

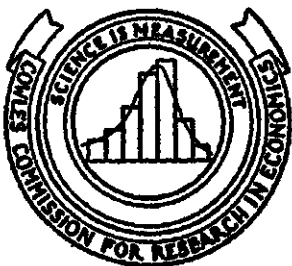
**COWLES COMMISSION FOR
RESEARCH IN ECONOMICS
Monograph No. 10**

STATISTICAL INFERENCE
IN DYNAMIC ECONOMIC MODELS

By
COWLES COMMISSION RESEARCH STAFF MEMBERS
AND GUESTS

Edited by
TJALLING C. KOOPMANS

With Introduction by
JACOB MARSCHAK



JOHN WILEY & SONS, INC., NEW YORK
CHAPMAN & HALL, LIMITED, LONDON
1950

COPYRIGHT, 1950
BY
COWLES COMMISSION FOR RESEARCH IN ECONOMICS

All Rights Reserved

*This book or any part thereof must not
be reproduced in any form without
the written permission of the publisher.*

PRINTED IN THE UNITED STATES OF AMERICA

CONTRIBUTORS TO THIS VOLUME

- RICHARD L. ANDERSON
North Carolina State College
- THEODORE W. ANDERSON, JR.
Cowles Commission and Columbia University
- TRYGVE HAAVELMO
Cowles Commission and University of Oslo
- HAROLD HOTELLING
University of North Carolina
- LEONID HURWICZ
Cowles Commission and University of Illinois
- TJALLING C. KOOPMANS
Cowles Commission, The University of Chicago
- ROY B. LEIPNIK
Cowles Commission (formerly), Institute for Advanced Study
- JACOB MARSCHAK
Cowles Commission, The University of Chicago
- HENRY B. MANN
Ohio State University
- HERMAN RUBIN
Cowles Commission (formerly), Stanford University
- ABRAHAM WALD
Columbia University

ACKNOWLEDGMENTS

The publication of this study, and the research of Cowles Commission staff members incorporated in it, has been made possible by generous grants received from the Rockefeller Foundation over an extended period.

Thanks are due to the guest authors who have made their contributions available for inclusion in this volume.

Mr. Dickson H. Leavens, Mr. B. A. de Vries, and the staff of the Cowles Commission, especially Mrs. Jane Novick, Editorial Secretary, Mr. Gerhard Stoltz, and Mrs. Lorraine Kirk, have put much care into editing and preparing the manuscript. Mr. Marvin Schuler has displayed skill, ingenuity, and perseverance in varityping the manuscript for reproduction by photo-offset.

TABLE OF CONTENTS

Chapter	Page
Principles of Notation	xiii
INTRODUCTION	
I. Statistical Inference in Economics: An Introduction By J. MARSCHAK	1
PART ONE: SIMULTANEOUS EQUATION SYSTEMS	
II. Measuring the Equation Systems of Dynamic Economics By T. C. KOOPMANS, H. RUBIN, and R. B. LEIPNIK	53
1. Description of the Systems Considered	54
2. The Identification of Economic Relations	69
3. Estimation of the Parameters	110
4. Computation of the Maximum-Likelihood Estimates	153
PROBLEMS OF IDENTIFICATION	
III. Note on the Identification of Economic Relations By A. WALD	238
IV. Generalization of the Concept of Identification By L. HURWICZ	245
V. Remarks on Frisch's Confluence Analysis and Its Use in Econometrics By T. HAAVELMO	258
PROBLEMS OF STRUCTURAL AND PREDICTIVE ESTIMATION	
VI. Prediction and Least Squares By L. HURWICZ	266
VII. The Equivalence of Maximum-Likelihood and Least-Squares Estimates of Regression Coefficients By T. C. KOOPMANS	301
VIII. Remarks on the Estimation of Unknown Parameters in Incomplete Systems of Equations By A. WALD	305
IX. Estimation of the Parameters of a Single Equation by the Limited- Information Maximum-Likelihood Method By T. W. ANDERSON, JR.	311
PROBLEMS OF COMPUTATION	
X. Some Computational Devices By H. HOTELLING	323

PART TWO: PROBLEMS SPECIFIC TO TIME SERIES

Chapter	TREND AND SEASONALITY	Page
XI.	Variable Parameters in Stochastic Processes: Trend and Seasonality By L. HURWICZ	329
XII.	Nonparametric Tests against Trend By H. B. MANN	345
XIII.	Tests of Significance in Time-Series Analysis By R. L. ANDERSON	352
ESTIMATION PROBLEMS		
XIV.	Consistency of Maximum-Likelihood Estimates in the Explosive Case By H. RUBIN	356
XV.	Least-Squares Bias in Time Series By L. HURWICZ	365
CONTINUOUS STOCHASTIC PROCESSES		
XVI.	Models Involving a Continuous Time Variable By T. C. KOOPMANS	384
PART THREE: SPECIFICATION OF HYPOTHESES		
XVII.	When Is an Equation System Complete for Statistical Purposes? By T. C. KOOPMANS	393
XVIII.	Systems with Nonadditive Disturbances By L. HURWICZ	410
XIX.	Note on Random Coefficients By H. RUBIN	419
	References	423
	Index	429

PRINCIPLES OF NOTATION

MATHEMATICAL SYMBOLS

- $\alpha, \beta, \eta, \chi, \dots$ Greek characters denote known or unknown constants (parameters).
- a, b, γ, z, \dots Latin characters denote quantities subject to a probability distribution, or exogenous variables assumed given.
- $\alpha, a; \pi, p; \dots$ Unknown parameters and their estimates are denoted as much as possible by corresponding (Greek and Latin) characters.
- $A, A; B, B; \dots$ Greek and Latin capitals can be distinguished by the fact that all Greek characters are vertical, all Latin characters italicized.
- $g = 1, \dots, G; \dots$ Latin characters (lower case) are also used as subscripts for numbering of variables, equations or observation periods. In such cases the range of the subscript is often from 1 to a maximum denoted by the corresponding Latin capital letter.
- $f, F, \varphi, \Phi, \dots$ Greek and Latin characters are used without distinction to denote distribution functions.
- a, b, A, B, \dots Vectors are denoted by lower case type, matrices by capitals. Subscripts preceding a vector or matrix denote the deletion of elements, rows, columns, or sets thereof.
- $a_{12}, \beta_{gh}, x_k, \dots$ Scalar elements of matrices or vectors are denoted in the normal manner by subscripts following the lower case character corresponding to the character denoting the matrix or vector.
- $\alpha_g, y_{II}, M_{yz}, \dots$ In some cases, explained in the text, subscripts denote sets of elements, rows or columns, to be included in a submatrix or subvector of the matrix or vector denoted by the symbol to which these subscripts are attached.
- $I A, III a, [1] \Phi, \dots$ Preceding subscripts similarly denote elements, rows or columns excluded in forming subvectors

or submatrices.

$\mathcal{E}, \mathcal{P}, \mathcal{Q}, \mathcal{R}, \dots$ Script letters denote operators.

$\mathbb{G}, \mathbb{S}, \dots$ German letters denote sets.

REFERENCES

Items in the list of references at the end of this volume are referred to by author's name and where necessary by year of presentation or publication in square brackets: [Haavelmo, 1943].

Articles contained in this volume are referred to by Roman numerals in square brackets: [III].

Sections of articles are referred to by italicized Arabic numerals: [I-2.5].

Formulae are referred to by vertical Arabic numerals in parentheses: (26).

**STATISTICAL INFERENCE IN ECONOMICS:
AN INTRODUCTION**

I. STATISTICAL INFERENCE IN ECONOMICS: AN INTRODUCTION

BY J. MARSCHAK

	Page
0. Purpose of the Volume	1
0.1. The Problem	1
0.2. The Discussion	4
1. Nonstochastic Models	6
1.1. The Model	6
1.2. Structural Changes and Policies	8
1.2.1. Structure and reduced form	8
1.2.2. Use of observations	9
1.2.3. Structural changes and policies	11
1.3. Identification	14
1.4. Use of Experiments	17
2. Stochastic Models	18
2.1. Random Disturbances	18
2.2. Shock Model and Structure	22
2.3. Identification	30
2.4. Dynamic Models	32
2.5. Estimation	34
2.6. The Choice of Model	44
3. Plan of the Volume	45

0. PURPOSE OF THE VOLUME

0.1. *The Problem*

Quantitative economic study has a threefold basis: it is necessary to formulate economic hypotheses, to collect appropriate data, and to confront hypotheses with data. The latter task, statistical inference in economics, was discussed at a Cowles Commission conference held at the University of Chicago from January 27 to February 1, 1945. Staff members of the Cowles Commission prepared,

and circulated in advance, some of the papers; others were delivered by the Commission's guests.¹

In its *Annual Report* for 1944 and for subsequent years, the Commission stressed the importance of adapting statistical methods to the peculiarities of the data and the objectives of economic research. The economist's objectives are similar to those of an engineer but his data are like those of a meteorologist. The economist is often required to estimate the effects of a given (intended or expected) change in the "economic structure," i.e., in the very mechanism that produces his data. None of these changes can he produce beforehand, as in a laboratory experiment; and since some of the changes envisaged have never happened before, the economist often has to estimate the results of changes that he has never observed.

The economist can do this if his past observations suffice to estimate the relevant structural constants prevailing before the change. Having estimated the past structure the economist can estimate the effects of varying it. He can thus help to choose those variations of structure that would produce - from a given point of view - the most desirable results. That is, he can advise on policies (of a government or a firm).

Thus, practical considerations bring about the economist's concern with economic structure. Hypotheses about economic structure are also known as economic theories. They try to state relations that describe the behavior and environment of men and determine the values taken at any time by economic variables such as prices, output, and consumption of various goods and services, and the prices and amounts of various assets. As there are several variables the economic structure must involve several simultaneous relations to determine them. In this, economic theory is analogous to theories used in experimental science.

Also, economic variables as well as those of experimental science are, in principle, random (stochastic) variables: that is, their properties are described by probability distributions. In particular, the stochastic character of the observed data can often be ascribed to their dependence on stochastic nonobservable variables: such nonobservable variables are random "errors" in the observation of single variables or random "shocks" suffered by the relations connecting

¹The guests included R. L. Anderson, T. Haavelmo, H. Hotelling, W. G. Madow, H. B. Mann, G. Tintner, and A. Wald. Staff members who participated in the conference were L. Hurwicz, L. R. Klein, T. C. Koopmans, R. Leipnik, J. Marschak, and H. Rubin. T. W. Anderson and T. Haavelmo joined the staff at a later date.

them. However, an experimenter could replace the natural conditions by laboratory conditions. To study one of the several relations, the experimenter observes the random values taken by one variable when the other observables that determine it are made reasonably free of the influence of errors and shocks. The economist cannot thus control variables and isolate relations. His data are produced by the existing economic structure, as described by a system of simultaneous relations between these random variables: the observables themselves, the errors, and the shocks. To use such data for the estimation of the system - "structural estimation" - is a new statistical problem.

This new statistical problem is thus forced upon the economist by the occurrence or consideration of structural changes (including policy changes on which his advice is sought), and by his inability to make experiments of either of two kinds: experiments producing in advance the considered change in structure (e.g., wind-tunnel experiments on airplanes), and experiments in which some of the random variables of nature are given fixed values (e.g., experiments to test fundamental laws of physics). See [J.Marschak, 1947B].

The role of simultaneous equations is familiar to economic theorists. But it has often been forgotten by economic statisticians who tried to estimate a single stochastic relation as if no other such relations had taken part in determining the observed values of the variables. On the other hand, economic theorists are apt to forget that the observed economic variables are, in general, stochastic. To be susceptible of empirical tests an economic hypothesis must be formulated as a statistical one, i.e., be specified in terms of probability distributions.

The statistical problem of the economist is complicated by the fact that many an economic relationship connects current and past values of the same or other variables involved. The economic structure determines, accordingly, not a set of constant values, one for each variable, but a set of probable paths, one for each variable, provided certain initial values are given. This dynamic character of economic structure creates, in the absence of experiments, further statistical difficulties: many economic data have the form of time series in which successive items are not independent. Statistical inference from time series of this kind involves further new problems.

Thus, economic data are generated by systems of relations that are, in general, stochastic, dynamic, and simultaneous. Occurring jointly, these three properties give rise to unsolved problems of statistical inference from the observed data to the relations.

Yet these very relations constitute economic theory and knowledge of them is needed for economic practice.

0.2. *The Discussion*

All these difficulties, under names like "pitfalls in demand and supply analysis" (Ragnar Frisch), "lack of independence in economic time series," etc., have caused uneasiness for a long time. Of the many attempts to grapple with the problem, Ragnar Frisch's contributions [1929, 1931, 1933, 1934, 1938] were probably the most stimulating ones. However, he did not take full account of the random disturbances (shocks) in the economic relations, nor of the simultaneous character of these relations. Moreover, Frisch's hypotheses on random disturbances (errors) in variables were not specified in probability terms. The latter, but not the former, defect was corrected in the early work of Koopmans [1937] and Wald [1940]. A new milestone was reached in 1943 when two articles were published in *Econometrica* by Haavelmo [1943], and by Mann and Wald [1943]. Haavelmo formulated the economist's simultaneous-equations model as a statistical hypothesis by assuming a random disturbance (shock) in each equation, in addition to random errors in each observable variable, and by specifying the distribution of these (unobservable) random quantities [Haavelmo, 1944, esp. Chapter III]. Mann and Wald outlined a solution of the estimation problem arising from the new formulation, though only for the case of large samples, and omitting the observation errors. For a set of short time series of interrelated variables, the contemporary (and incomplete) work on time series in a single variable has to be utilized as a start; and important suggestions can also be expected from the study of continuous random processes that is being developed currently in the service of other sciences. As to combining shocks and errors in one equation system, recent investigations of T. W. Anderson and L. Hurwicz [1947] were stimulated by discussions with G. Tintner.

For its quantitative studies of economic behavior, the Cowles Commission had to expect much from the criticism of statisticians who had contributed to the estimation of simultaneous equations and to the theory of time series. Such was the object of the conference. Earlier drafts of articles II, V, VI, XI, XII, XIV, XV, XVI, XVII, XVIII, XIX, were prepared for the conference. Articles III, VIII, X, XIII, are discussions that were contributed in the conference and written up shortly afterwards.

Articles or parts of articles added or substantially expanded later are I, II-2.3, II-3.3, II-4, IV, VI, VII, IX, XI-10.2, XIV, XV. The revision was helped by discussions with persons not present at the conference: especially on the problem of computations [II-4] with A. Adrian Albert and John von Neumann. An alternative method of structural estimation, that of "limited information," suggested by M. A. Girshick and worked out by T. W. Anderson and H. Rubin [1949], is briefly presented in this volume by T. W. Anderson [IX]. Other conference contributions, by L. R. Klein [1946 A] and by W. G. Madow [1945], have been published elsewhere. Madow's subject was explored further by R. B. Leipnik [1947].

The manuscript of the present volume was completed early in 1947, but publication has been delayed by typographical and other printing difficulties.

In the next two sections of the present introductory paper, the author has drawn freely on the results attained in the papers that follow and on suggestions made in the daily work and discussion within the Cowles Commission. His debt to Leonid Hurwicz and Tjalling C. Koopmans is particularly heavy.

We have tried to achieve conformity in terminology and, to some extent, in notation. However, since the several contributions differ in purpose as well as in emphasis, rigorous uniformity is neither possible nor desirable.

Most of the contributions to this volume presuppose on the part of the reader a general knowledge of mathematical principles of statistics; to explain these principles to a more general reader would take more space than is available. The present introduction, in summarizing the purpose and the main results of the studies collected in this volume, is addressed to the mathematically-minded economist rather than to the statistician. Hence - the use made of nonstochastic models (section 1) and the attention paid, in the stochastic case, to the properties of populations (section 2.5.1) as distinct from samples. This summary, too, has to be terse. For a less compressed treatment and further economic illustrations and applications the reader is referred to the following publications of Cowles Commission staff members: [Girshick and Haavelmo, 1947], [Haavelmo, 1943, 1944, 1947 A, 1947 B], [Hurwicz, 1947], [Klein, 1946 B, 1947, 1950], [Koopmans, 1945, 1949], [Marschak and Andrews, 1944], [Marschak, 1947 A, 1947 B]. The group continues to work on statistical inference in economics, both in general and with respect to specific economic models.¹ It is hoped that the present volume

¹See the *Annual Report* of the Cowles Commission, in particular the *Five-Year Report* for 1942-46.

will stimulate further cooperation of mathematical statisticians and economists in solving the many problems that have been indicated but not solved in this volume. Plans are under way for a parallel monograph (No. 12) of a more expository character, in which emphasis is placed on a discussion of the main ideas and techniques developed in this volume with the help of simple illustrative models, rather than on formal mathematical proof.

1. NONSTOCHASTIC MODELS

Economic relations are, in general, stochastic. They involve variables whose properties are described with the aid of probability distributions. Moreover, the estimates of parameters of these relations, obtained by statistical methods from a limited number of observations, are also random variables.

However, important distinct properties of empirical economics can be brought out even if, for simplicity, we assume the data to be measured exactly and to satisfy exactly the relations of theory. The equations or inequalities serving to determine the parameters of such relations from observations are free from random variables. The problem of estimation degenerates into that of determination. This simplifies the study of certain "prestatistical" problems facing the economic statistician - in particular that of identification (section 1.3), and also helps to see why, as indicated in section 0.1, these problems originate in the need for policy decisions in the absence of experiments.

In the present section we shall deal with this special, or degenerate, case to meet in particular the habits of readers with economic rather than statistical background.

1.1. *The Model*

1.1.0. Denote the observable variables (or observables) by a vector $x \equiv (x_1, \dots, x_y)$. The first, second, ..., T th observations on x , succeeding each other in time, or arranged in space or in any other way, form a matrix:

$$X^o \equiv [x_n(t)], \quad n = 1, \dots, N, \quad t = 1, \dots, T;$$

or

$$X^o \equiv \begin{bmatrix} x(1) \\ \dots \\ x(T) \end{bmatrix}.$$

We call *a priori information* all statements (either true or false) arrived at independently of any knowledge of X^0 . We call *model* \mathfrak{G} the a priori information on a system of mutually consistent and independent equations

$$(1.1) \quad \varphi_g(x, \alpha_{(g)}) = 0, \quad g = 1, \dots, G,$$

where $\alpha_{(g)}$ is a vector of P_g parameters. We shall denote the vector of all parameters of the system by $\alpha \equiv (\alpha_{(1)}, \dots, \alpha_{(G)})$, and write $\sum P_g = P$; we shall also write¹ $\bar{\varphi} \equiv (\varphi_1, \dots, \varphi_G)$.

1.1.1. We shall assume throughout that \mathfrak{G} defines a) the form of $\bar{\varphi}$, and b) the "a priori restrictions," i.e., equations or inequalities in parameters α .

1.1.2. To provide this information we must make full use of our independent knowledge of existing production conditions (technology, legal statutes, etc.) and of plausible, if not necessarily rational, individual behavior. The equations of the model must refer to individual agents in specified markets (as consumers or manufacturers of certain goods, or as workers, bankers, landlords, etc.).

1.1.3. However, to reduce the model to a manageable size, variables referring to single individuals in finely subdivided markets must be grouped into aggregates. Suppose the value of some variable relevant to a practical decision (see below, section 1.2.3) is calculated on the basis of such an aggregative model. This value will contain an error inasmuch as it will deviate from the corresponding value calculated on the basis of a true, detailed model, with separate equations for each commodity and individual. Optimum aggregation should combine highest manageability (e.g., shortest computations) with smallest error. This aggregation problem (which includes that of index numbers) has not been solved or even formulated in detail, nor will it be studied in this volume.

1.1.4. The model \mathfrak{G} is called *self-contained* if $G = N$; *sectional* if $G < N$.

We call a model *complete* if it is either self-contained or has the following property: a subset of x containing $K = N - G$ elements - call this subset $z \equiv (z_1, \dots, z_K)$ - is determined by K

¹Subscripts indicate the elements of a vector which are scalars or functions: x_n, φ_g . In the present article, subvectors formed from the elements of a vector are indicated by subscripts in parentheses: $\alpha_{(g)}$. In addition, it will here prove convenient to use a barred letter for a vector whose elements are functions: $\bar{\varphi}$.

possibly unknown "subsidiary" equations

$$(1.2) \quad \varphi_{G+k}(z) = 0, \quad k = 1, \dots, K,$$

which are independent of the equation (1.1); the equations (1.2) do not contain any elements of x that are not in z . Denote the latter elements by vector $y \equiv (y_1, \dots, y_G)$. Then, assuming differentiability,

$$(1.3) \quad \frac{\partial \varphi_{G+k}}{\partial y_g} = 0, \quad g = 1, \dots, G, \quad k = 1, \dots, K.$$

A model that is not complete is called incomplete (or partial).

The observables z are called *exogenous* and the observables y *endogenous*, with respect to the model \mathcal{G} .

Equation (1.1) can be rewritten as

$$(1.4) \quad \varphi_g(y, z; \alpha_{(g)}) = 0, \quad g = 1, \dots, G.$$

1.2. Structural Changes and Policies

1.2.1. Structure and reduced form.

1.2.1.1. We call *structure* S all properties of the equations (1.4), including the properties not known a priori. Any model \mathcal{G} is a class of structures. Each structure is defined by the functional forms of the equations and the values of the parameters occurring in them. We can write¹

$$(1.5) \quad S \equiv (\bar{\varphi}, \alpha).$$

When the equations (1.4) are thus fully specified we call them *structural equations*. We call α the structural parameters. Concepts defined in section 1.1.4 with reference to the model ("complete," "sectional" model; variables exogenous to a model; etc.) will also be applied to a structure without causing ambiguity.

1.2.1.2. Given the structure S , equations (1.4) can be solved for y in terms of z , involving new parameters which we shall denote

¹In what follows quantities depending on S will be introduced. They are, properly speaking, "functionals" with respect to the argument $\bar{\varphi}$ and "functions" with respect to the argument α , although we shall denote and refer to them as "functions" of S .

by vector π :

$$(1.6) \quad y = \bar{\eta}(z; \pi),$$

say, where $\bar{\eta}$ is a vector whose elements are functions. Equation system (1.6) is called *reduced form*.¹ $\bar{\eta}$ depends on $\bar{\varphi}$; π is obtained by applying a transformation to the parameters α in (1.4). This transformation itself depends on the functions $\bar{\varphi}$. If we call this transformation $\bar{\pi}_\varphi$, we can write, accordingly,

$$(1.7) \quad \pi = \bar{\pi}_\varphi(\alpha) = \bar{\pi}(S),$$

say; furthermore, we can also write $\bar{\eta}_\varphi$ instead of $\bar{\eta}$ to emphasize the dependence of the functions $\bar{\eta}$ on the functions $\bar{\varphi}$. Thus

$$(1.8) \quad y = \bar{\eta}_\varphi[z; \bar{\pi}_\varphi(\alpha)].$$

1.2.1.3. If the structural functions φ are linear, the functions $\bar{\eta}$ and $\bar{\pi}_\varphi$ will also be linear, and the set of parameters $\pi = \bar{\pi}_\varphi(\alpha)$ of the reduced form corresponding to a given structure S will be unique. If the functions φ are nonlinear, several sets of parameters α may be compatible with the structure S . We shall neglect here this complication although it does occur in economic theory, as in the case of "multiple equilibrium" [Marshall, app. H].

1.2.1.4. Apart from very special cases, the reduced forms (each one characterized by a function set $\bar{\eta}$ and a parameter set π) compatible with a given structure S will be finite in number, or at least denumerable; in a linear model the reduced form is unique.

On the other hand the number of structures compatible with a given reduced form $(\bar{\eta}, \pi)$ may or may not be (nondenumerably) infinite. Even if the model is linear, the structure may or may not be uniquely determined by the reduced form. (See also below, section 1.3.)

1.2.2. Use of observations.

1.2.2.1. Suppose the period (or, say, the geographical area) of observations consists of elements 1, ..., T , and is so chosen that the structure maintains throughout it the value

$$S^0 \equiv (\bar{\varphi}^0, \alpha^0).$$

Then the following equations are satisfied by the observations X^0 :

¹See [IX] where $\bar{\varphi}$ and consequently $\bar{\eta}$ are linear functions.

$$(1.9) \quad \varphi_g^o[x(t); \alpha_{(g)}^o] = 0, \quad g = 1, \dots, G, \quad t = 1, \dots, T.$$

The structure S^o is called *observational structure*. We can similarly consider the equations

$$(1.10) \quad y(t) = \bar{\eta}^o[z(t), \pi^o], \quad t = 1, \dots, T,$$

where

$$(1.11) \quad \bar{\eta}^o \equiv \bar{\eta}_{\varphi^o}^o$$

and

$$(1.12) \quad \pi^o \equiv \bar{\pi}_{\varphi^o}(\alpha^o) \equiv \bar{\pi}(S^o);$$

the "observational reduced form" is given by $\bar{\eta}^o; \pi^o$. It follows from definitions that (1.10) is satisfied by observations X^o .

1.2.2.2. The form of the functions $\bar{\eta}_{\varphi^o}$ and $\bar{\pi}_{\varphi^o}$ is determined by the functions $\bar{\varphi}^o$ of the model. If the latter are linear so are the former. The model also provides restrictions on α^o , and these restrictions can be transformed, by (1.12), into equations or inequalities in π^o . If the model is linear, a unique set of parameters π^o will correspond to the set of structural parameters α^o . Furthermore, if T , the number of observations, is sufficiently large, the equations (1.9) together with the restrictions, determine the parameters π^o of the observational reduced form, given the observations X^o . We shall denote this operation by $\bar{\rho}_1$, so that

$$(1.13) \quad \bar{\rho}_1 X^o = \pi^o.$$

1.2.2.3. As stated in section 1.2.1.4, a given reduced form may or may not determine the structure uniquely. Therefore, although, as just stated, the observational reduced form can be determined from the model, given the observations X^o (if their number T is sufficiently large), it may be impossible to determine a unique observational structure from the model and the observations *however large their number T* . It may sometimes be possible, in other words, to replace a given structure by any one of an infinite number of other structures without contradicting the observations. There exists thus a problem of *identifying a structure* (treated below in section 1.3) but no problem of identifying a reduced form.

1.2.3. Structural changes and policies.

1.2.3.1. The observational structure S^0 may be different from some structure S valid for a different period (or geographical area, etc.). Suppose we know the "structural changes" \mathfrak{J} , a transformation that carries S^0 into S : $S = \mathfrak{J}S^0$. If, in addition, we know S^0 from observations, we can obtain $S \equiv (\varphi, \alpha)$ and hence also the functions $\bar{\eta}_\varphi$ and $\bar{\pi}_\varphi$ and therefore the parameters π of the reduced form. It is then possible to evaluate¹ y for a given z .

1.2.3.2. We shall distinguish between two kinds of structural changes: the *controllable* ones, \mathfrak{J}_c ; and the *uncontrollable* ones, \mathfrak{J}_u . The former ones are also called "structural policy;" the introduction, or abolition, of price control is an example. We can distinguish similarly between two sets of exogenous variables: the controllable ones, z_c ; and the uncontrollable ones, z_u . The fixing of z_c is called "nonstructural" (or "routine") policy: for example the annual revision of tax rates [Marschak, 1947A].

1.2.3.3. The policy-maker (on behalf of a government, of an individual firm, etc.) tries to maximize the "gain," or "welfare," ω , a certain function of the observables which, in principle, must be supposed to be known to him:

$$(1.14) \quad \begin{aligned} \omega &= \omega(y, z) = \omega\{\bar{\eta}_\varphi[z; \bar{\pi}_\varphi(\alpha)], z\} \\ &= \omega_\varphi[z; \bar{\pi}_\varphi(\alpha)], \end{aligned}$$

say. Thus the gain (welfare) function ω and the functions φ of the model combine to determine the function ω_φ of exogenous variables which is to be maximized, given the structure $(\bar{\varphi}, \alpha)$.

1.2.3.4. Suppose structural changes are neither intended nor expected for the period (or area) to which policy is to be applied, compared with the period (or area) of observation. (The policy consists, in this case, in fixing the value of z_c only: it is "nonstructural.") In this case, \mathfrak{J}_c as well as \mathfrak{J}_u is the identical transformation and we have $\varphi = \varphi^0$, $\alpha = \alpha^0$, and

$$(1.15) \quad \pi = \pi_{\varphi^0} \quad (\text{a set of constants}).$$

Thus the gain

$$(1.16) \quad \omega = \omega_{\varphi^0}(z; \pi) = \omega_{\varphi^0}(z_c, z_u; \pi)$$

¹But see section 1.2.1.3 for a qualification.

can be affected only by variations in z_c . Best policy is the value \hat{z}_c of z_c that gives ω its maximum value $\hat{\omega}$:

$$(1.17) \quad \hat{\omega} = \max_{z_c} \omega_{\varphi^0}(z_c, z_u; \pi^0) = \omega_{\varphi^0}(\hat{z}_c, z_u; \pi^0).$$

By comparing the values of ω for varying z_c at fixed z_u , the best policy \hat{z}_c can be determined for any given value of the uncontrollable variables z_u , provided the parameters π^0 of the (observational) reduced form are known. But π^0 can indeed be found from the observations X^0 (section 1.2.2.2). The operation (1.13),

$$\mathfrak{P}_1 X^0 = \pi^0,$$

depends on the model only and is called "predictive determination when structure is unchanged."¹ Operation \mathfrak{P}_1 provides, then, the parameters π^0 , to be used for the choice of policy under unchanged structure.

1.2.3.5. As a rule, however, some structural changes will be intended or expected, or both. That is, neither \mathfrak{V}_c nor \mathfrak{V}_u will, in general, be the identical transformation. We have $S = \mathfrak{V}S^0$; where $\mathfrak{V} = \mathfrak{V}_c \mathfrak{V}_u$ (neglecting the question of the order in which the transformations are applied). Further, by (1.14) and (1.15), we have

$$(1.18) \quad \omega = \omega_{\varphi} [z; \pi_{\varphi}(\alpha)] = \omega^*(z; S),$$

say, where the form of the function ω^* depends on the form of the gain function ω only. Further,

$$(1.19) \quad \begin{aligned} \omega^*(z, S) &= \omega^*(z_c, z_u; \mathfrak{V}S^0) \\ &= \omega^{**}(z_c, z_u; \mathfrak{V}_c, \mathfrak{V}_u; \varphi^0, \alpha^0). \end{aligned}$$

Best policy is defined by values $\hat{z}_c, \hat{\mathfrak{V}}_c$ that jointly maximize the gain. Let the maximum gain be

$$(1.20) \quad \hat{\omega} = \omega^{**}(\hat{z}_c, z_u; \hat{\mathfrak{V}}_c, \mathfrak{V}_u; \varphi^0, \alpha^0).$$

¹See [VI].

By comparing the values of ω for varying z_c , \mathfrak{U}_c , given the values of the uncontrollable, exogenous variables and given the uncontrollable changes in structure, one can determine the best policies $(\hat{z}_c, \hat{\mathfrak{U}}_c)$, provided the observational structure S^o is known. The practical procedure is, in principle, as follows: from S^o and any given $\mathfrak{U} = \mathfrak{U}_c \mathfrak{U}_u$ derive the new structure $S = \mathfrak{U} S^o$ for the period of policy application; from S derive the parameters $\pi = \pi(S)$; and compute the variables y and the gain ω as in (1.14).

1.2.3.6. An operation $\mathfrak{S} X^o = S^o$ determining S^o from the observations X^o will be called *structural determination*; the question of its existence will be discussed in section 1.3. If structural determination is possible, it is possible to derive the parameters π by applying in succession the operations \mathfrak{S} and \mathfrak{U} :

$$(1.21) \quad \pi = \bar{\pi}(S) = \bar{\pi}(\mathfrak{U} \mathfrak{S} X^o) = \mathfrak{P}_y X^o,$$

say. The operation \mathfrak{P}_y thus defined is called *predictive determination when structure undergoes a given change \mathfrak{U}* .

1.2.3.7. Structural determination provides a master key for predictive determination, and for the calculation of alternative gains, for any of the various possible structural changes. The structural changes to be considered are seldom known long in advance. Therefore, although formally $\pi = \bar{\pi}(\mathfrak{U} \mathfrak{S} X^o)$ can be computed without a stop, it is preferable and often essential to pause at the step $\mathfrak{S} X^o = S^o$. The knowledge of observational structure means greater flexibility with regard to various alternative policies. This is one reason why people are interested in any kind of theory!

1.2.3.8. However, the considered transformations \mathfrak{U} (structural policies and uncontrollable structural changes) and the subset of variables relevant to the evaluation of the gain, may happen to be such as to make the knowledge of all parameters α^o of the structure S^o unnecessary. A partial knowledge of S^o – some elements of α^o , or perhaps some functions of them – may be all that is needed.

1.2.3.9. It was required in section 1.1.2 that the equations of the model describe plausible behavior of specified economic agents, thus making full use of our a priori knowledge of behavior (rational or otherwise). We now see a further practical reason for shunning relations that do not refer to specified

economic agents. We call such relations "anonymous." Consider changes in human behavior, institutions, technology. The gain (personal or social) due to any such intended or expected change cannot be evaluated unless behavior, institutions, and technology are explicitly stated; such statements must be provided by the form, and by the values of parameters, of the equations of the model; that is, by the structure.

1.2.3.10. As an example, consider the following - admittedly oversimplified - model:

(1) demand for a commodity depends on its price and on national income;

(2) supply depends on price;

(3) demand equals supply;

and it is assumed that

(4) we can treat national income as exogenous.

Suppose we want to evaluate the effect of replacing free demand of the public by a fixed demand determined by the government; that is, relation (1) is replaced by

(1') demand = a constant.

If we know the form and parameters of (2) we can evaluate, with the help of (1') and (3), the price the suppliers will ask and the government will have to pay. Suppose, however, that instead of (1), (2), (3), we had at our disposal the following relation obtained by elimination of demand from (1), using (3), (2):

(5) price depends on income.

This "anonymous" relation (which, in this case, is a reduced form) can be computed from observations (section 1.2.2.2) but cannot help to evaluate the effect of structural policy, i.e., the effect of replacing (1) by (1'). For another example, see [Marschak, 1947B].

1.3. Identification

1.3.1. It was remarked in section 1.2.2.3 that a given set of observations, however numerous, or a given reduced form, will not, in general, determine a unique structure. For a rigorous introduction to this *problem of identification* in the case of stochastic models we refer to [I-2] and [IV]. For the purposes of the present paper, it is convenient to approach the problem by using the concept of reduced form and studying it first in the nonstochastic case.

Consider equations (1.8) and the structure $S^o \equiv (\varphi^o, \alpha^o)$.

Denote by \mathfrak{S}_*^o the class of all structures $S_*^o \equiv (\varphi_*^o, \alpha_*^o)$ that are compatible with the model \mathfrak{S} and yield the same reduced form as the structure S^o . The latter condition means that the following equations are satisfied identically in z :

$$(1.22) \quad y = \bar{\eta}_{\varphi^o}[z; \bar{\pi}_{\varphi^o}(\alpha^o)] = \bar{\eta}_{\varphi_*^o}[z; \bar{\pi}_{\varphi_*^o}(\alpha_*^o)].$$

S^o is said to be *uniquely identifiable* by the model \mathfrak{S} if S^o is the only element of \mathfrak{S}_*^o , i.e., if every $S_*^o = S^o$. Furthermore, S^o is said to be *incompletely identifiable* by the model \mathfrak{S} if the class \mathfrak{S}_*^o contains a nondenumerably infinite number of elements. Only in the former case is it possible to determine α^o uniquely from the parameters of the reduced form, $\pi^o = \bar{\pi}_{\varphi^o}(\alpha^o)$, provided the functions φ^o and the a priori restrictions on α^o are given.

1.3.2. The concept of identifiability can be easily extended to any subset of α^o , say $\alpha_{(i)}^o$ (*partial identification*): for example, all or some of the parameters of some of the equations (1.4) may be uniquely determinable from the parameters of the reduced form.

1.3.3. If a subset of α^o is nonidentifiable it is impossible to determine it from X^o *however large the number of observations* T . If α^o is completely and uniquely identifiable, it can be obtained from X^o using equations (1.9) jointly with the a priori restrictions, provided T is sufficiently large: this is structural determination, denoted in section 1.2.3.6 by \mathfrak{S} . A similar statement applies to any subset $\alpha_{(i)}^o$.

1.3.4. As an example, let $x = (y_1, y_2, z_1)$ be the coordinates of a point and let $\varphi_g(x) = 0$, $g = 1, 2$, be a pair of equations of two distinct planes. The observations on x will yield a set of collinear points which will determine a straight line – corresponding to (1.6) – permitting the prediction of y_1 or y_2 for a given z_1 . But it will be impossible to reconstruct any particular pair of planes. A (partial) identification is, however, possible if, for example, it is known a priori that one of the planes is vertical; in which case this plane (but not the other one) is identifiable.

1.3.5. This corresponds to the economic example of section 1.2.3.10 with $y_1 =$ demand (=supply), $y_2 =$ price, $z_1 =$ income. The parameters of (2) but not those of (1) are identifiable. If

(2) were modified into

(2') supply depends on price and wage rate,

with wage rate assumed exogenous [Koopmans, 1945 p.451] all parameters of the structure would become identifiable.

1.3.6. The conditions of identifiability in a linear stochastic model are studied in [II-2]; their application to a nonstochastic model is easily derived. The most important criteria are supplied by the presence or absence of variables in each of the equations of the system (1.4). In particular, the occurrence of different exogenous variables in different equations contributes to identifiability.

1.3.7. Let the number of endogenous variables, $G = 1$; assume φ_1 is linear, and normalize the parameters α by choosing the coefficient of y_1 to be equal to 1:

$$(1.23) \quad y_1 - \sum \alpha_k z_k = 0.$$

The *linear uniequational structure* (1.23) is always identifiable, and $\alpha = \pi$, provided there exist no linear relations between the z 's. The proposition can be easily extended to nonlinear uniequational structures, apart from trivial modifications.

1.3.8. *Predictive determination when structure is unchanged*, i.e., the operation $\mathcal{P}_1 X^0$ in (1.13) is possible, regardless of whether or not the observational structure S^0 is identifiable.

1.3.9. *Predictive determination when structure undergoes a known change* \mathcal{U} , i.e., the operation $\pi = \mathcal{P}_Y X^0$ in (1.21) is possible when all parameters α^0 are identifiable, and impossible when none of them are identifiable.

1.3.10. However, the structural change \mathcal{U} and the gain function ω may be such as to require the knowledge of some but not all elements of α^0 (section 1.2.3.8); in this case partial identifiability (section 1.3.2) is all that is needed. Also, \mathcal{U} and ω may be such as to require the knowledge, not of the parameters α^0 themselves, but of some functions of them; in which case unique identifiability of every single parameter separately is not necessary for the choice of best policy.

1.3.11. An incomplete model has more endogenous variables than it has equations (section 1.1.4). For example a uniequational model involving two or more endogenous variables is incomplete. The parameters of such a model form a subset of the parameters of some complete model. The parameters of the incomplete model may or may

not be determinable from observations since the structure may or may not be partially identifiable (section 1.3.2), by the complete model, with respect to the particular subset of its parameters.

1.4. Use of Experiments

1.4.0. We shall now distinguish between: *original* structure S^{00} , *future* structure $S = \mathfrak{J}S^{00}$, and *observational* structure $S^0 = \mathbb{Q}S^{00}$. In the preceding section observations were supposed to be made on the original structure, so that \mathbb{Q} was the identical transformation. This is indeed the situation in nonexperimental science. Experiments, on the other hand, consist in applying certain transformations that change the original structure and in getting observations from the new structure $S^0 = \mathbb{Q}S^{00}$ that is thus obtained. The transformations \mathbb{Q} are chosen in such a way as to permit predictive determination under future structure S without determining either S or the (possibly nonidentifiable) original structure S^{00} . According to the nature of the transformation \mathbb{Q} , there are various types of experiments. Two particular types of experiments deserve our attention [Marschak, 1947 B].

1.4.1. *Experiments of type I: imitation of future structure* (example: wind tunnels for testing airplanes). Here $\mathbb{Q} = \mathfrak{J}$, hence $S^0 = S$ and $\pi^0 = \pi$. The structural change occurs between S^{00} and S . To evaluate y for a given z under the future structure S , and thus to evaluate the gain ω , it suffices to determine π^0 by the operation $\mathfrak{P}_1 X^0 = \pi$ (section 1.2.2.2). This does not require identifiability of either S^{00} or S .

1.4.2. *Experiments of type II: creation of uniequational complete structure* (example: controlled experiments in the physical laboratory). Let the original structure S^{00} correspond to a model

$$(1.24) \quad \varphi_g^{00}(y, z; \alpha_{(g)}^{00}) = 0, \quad g = 1, \dots, G.$$

Let \mathbb{Q}_1 be the operation of replacing all equations (1.24) but the first by the equations

$$(1.25) \quad y_g = z_{K+g}, \quad g = 2, \dots, G,$$

where the quantities $(z_{K+1}, \dots, z_{K+G}) \equiv$ vector z_* are fixed by the experimenter at various values $z_*(1), \dots, z_*(T)$, these

values being mutually independent. Denote the resulting observations by matrix X_1^o ; they can be regarded as produced by an un-equational complete structure involving a vector of parameters $\alpha_{(1)}^o$, as follows:

$$\mathbb{Q}_1 S^{oo} = S_1^o = (\varphi_1^o, \alpha_{(1)}^o);$$

that is,

$$(1.26) \quad \varphi_1^o(y_1, z, z_*; \alpha_{(1)}^o) = 0,$$

where $\varphi_1^o = \varphi_1^{oo}$, and $\alpha_{(1)}^o = \alpha_{(1)}^{oo}$. The vector $\alpha_{(1)}^o$ of structural parameters is related by trivial transformations (e.g., dividing by a constant if φ_1^o is linear) to the vector of parameters $\pi_{(1)}^o$ ($= \alpha_{(1)}^{oo}$) of the reduced form

$$(1.27) \quad y_1 = \eta_{\varphi o}(z, z_*; \pi_{(1)}^o).$$

Hence the subset $\alpha_{(1)}^{oo}$ ($= \alpha_{(1)}^o$) can be determined (section 1.3.7) even though the set α^{oo} may be nonidentifiable.

If it is desired to determine the whole set α^{oo} , experiments $\mathbb{Q}_1, \mathbb{Q}_2, \dots, \mathbb{Q}_G$ can be applied in succession.

2. STOCHASTIC MODELS

2.1. *Random Disturbances*

2.1.1. The stochastic character of economic relations will be recognized even by an out-and-out determinist. His world is, in principle, ruled by a set of very many equations in very many variables (both economic and other). But to make his theory verifiable by observation he will have to shorten the set considerably. The numerous causes that determine the error incurred in measuring a variable are not listed separately; instead, their joint effect is represented by the probability distribution of the error, a random variable. Also, the numerous causes that determine, say, the velocity of a gas particle are conveniently represented by the probability distribution of this velocity, a random variable. The economist acts similarly. He allows for

random errors of observation; and he represents the vagaries of, say, changing fashion by random "disturbances" or "shocks" that obey certain probability distributions; he thus cuts short the complicated causal explanation of why tastes fluctuate in the way they do.

The nondeterminist, on the other hand, will find it unnecessary to give any justification for the presence of random elements in economic models except by appealing directly to the "erratic," "unpredictable" character of certain types of events including human behavior; though he, too, will have to assume those events to be bound by certain probability distributions: even if actions are unpredictable, certain actions remain more probable than others.

2.1.2. Denote by $w \equiv (w_1, \dots, w_j)$ the vector of nonobservable random disturbances affecting economic observations, and by $x \equiv (x_1, \dots, x_n)$ the vector of observable variables. We shall call a stochastic model \mathfrak{S} the a priori information on a system of equations

$$(2.1) \quad \varphi_g(x, w; \alpha_{(g)}) = 0, \quad g = 1, \dots, G,$$

and on the joint distribution density

$$(2.2) \quad f(w; \varepsilon),$$

where ε as well as $\alpha_{(g)}$ denote parameter vectors. As in section 1.1.0 we shall write $\alpha \equiv (\alpha_{(1)}, \dots, \alpha_{(G)})$. We shall assume, in particular, that a priori information exists a) on the forms of the functions $\bar{\varphi} \equiv (\varphi_1, \dots, \varphi_G)$ and f ; and b) on some equations or inequalities in the parameters of these functions. We shall call *structure* all properties of the equations (2.1) and of the distribution (2.2) including the properties not known a priori.

If $w = 0$, we have the nonstochastic case treated in section 1 of this article.

2.1.3. If it is possible to substitute for w from (2.1), the distribution density function f of the disturbances can be transformed into the distribution density function of the observable vector x :

$$(2.3) \quad g_x(x),$$

say. Given the form of the functions φ and f and given their parameters α , ε , the distribution function g_x and its parameters are determined. Certain parts of the present volume [II-2 and IV] deal with the converse problem: given the distribution (2.3) of the observables, determine the properties of the equations (2.1) and the properties of the distribution (2.2) of the disturbances. Well-known methods permit estimation of the distribution of the observables (2.3) from the observations. But to use the knowledge of that distribution to determine the properties of (2.1) and (2.2) raises a new problem in inference, that of identification. The relation between the stochastic model and the distribution of the observables will show certain analogies with the relation, already studied, between the nonstochastic model and the reduced form.

2.1.4. The random disturbances w may include as a subset a vector of additive *disturbances in variables* (additive errors of observations, or briefly, errors) $v = (v_1, \dots, v_n)$. Thus the model (2.1), (2.2) can be rewritten, slightly changing notation,

$$(2.4) \quad \begin{aligned} \varphi_g(x - v, w') &= 0 \\ f(v, w'), & \quad g = 1, \dots, G, \end{aligned}$$

where the vector w' is complementary to v in w , i.e., $w \equiv (v, w')$.

2.1.5. There are also random *disturbances in relations*, especially of the additive kind, which we call for brevity *shocks*, $(u_1, \dots, u_G) \equiv u$, a subset of w' . If we denote by vector w'' the complement of u in w' so that $w \equiv (u, v, w'')$, we can rewrite the model (again slightly changing the meaning of φ):

$$(2.5) \quad \begin{aligned} \varphi_g(x - v, w'') &= u_g, \\ f(u, v, w''), & \quad g = 1, \dots, G. \end{aligned}$$

2.1.6. If w'' is empty the model may be called "simple shock-and-error model" [T. W. Anderson and Hurwicz, 1947]. In the present volume the contributions of the Cowles Commission staff are confined to the more special case where w'' and v are empty: the "simple shock model." Other studies, [Frisch and Mudgett, 1931], [Koopmans, 1937], [Wald, 1940], [Tintner, 1946], [Reiersøl, 1945], and [Geary, 1942], may be said to deal with another special case: the "simple error model," in which both w'' and u are empty.

2.1.7. Although even the simpler stochastic models present considerable difficulties, economists will probably be right in calling the statisticians' attention to more complicated models, i.e., those in which w'' is not empty, and this for three reasons:

(1) there may be nonadditive errors of observation;

(2) there may be nonadditive shocks, since an economic relation can be disturbed in a variety of ways (for example, a linear relation, say a demand curve, can fluctuate owing to random changes in its constant term as well as in any of the coefficients, see [XVIII], [XIX]);

(3) there occur "prospective"¹ variables, such as prospective profits, prices, etc., that affect people's behavior and must enter the equations of the model.

As a rule the economist cannot observe these (except by difficult questioning of a sample of people). But he may have hypotheses describing the determination of each of these variables in the minds of people: such "forecast equations" [Hurwicz, 1946, p.130 ff.] would relate the "prospective" variables to certain "actual" variables (such as the current or lagged national income, etc.) or to the observations on the "actual" variables. These additional structural equations will be themselves subject to psychological fluctuations, expressed as random disturbances, additive or non-additive. By use of the forecast equations, the prospective variables can be eliminated. The model thus derived can be used in predicting effects of structural changes occurring in equations other than the forecast equations. The derived model will involve functions of the random disturbances that occurred in the forecast equations. Even if they occurred additively in the forecast equations, these disturbances will enter the model, generally, in a nonadditive fashion, i.e., as elements of w'' in (2.5).

However, the remaining parts of this article and the bulk of the monograph itself will deal with simple shock models only, though some of the statements might be generalized to apply to the general model (2.5).

2.2. Shock Model and Structure

2.2.1. The introduction of shocks, i.e., additive random

¹The word "anticipated" variables is sometimes suggested. But this would seem to exclude expectations that do not come true. The term "expected" might be confused with the statistical term, expected (= mean) values.

disturbances in relations (excluding other random disturbances), can be considered as weakening the nonstochastic hypothesis (1.1): the observable variables $x \equiv (x_1, \dots, x_N)$ are related by a system of equations

$$(2.6) \quad \varphi_g(x, \alpha_{(g)}) = u_g, \quad g = 1, \dots, G,$$

the distribution function (probability density) of the random vector $u \equiv (u_1, \dots, u_G)$ being denoted by

$$(2.7) \quad f(u; \varepsilon)$$

where $\alpha_{(g)}$ and ε are parameter vectors. Thus each equation (other than definitional equations that can be properly eliminated from the system) is subject to a disturbance.

A priori information on the functions and parameters in (2.6), (2.7) constitutes a shock model, \mathcal{G} . A shock structure, S , consists of all properties of those functions and parameters:

$$S \equiv (\varphi, f, \alpha, \varepsilon),$$

where vector $\alpha \equiv (\alpha_{(1)}, \dots, \alpha_{(G)})$. If, as will be assumed, the model supplies the forms of the functions (in addition to some a priori restrictions on the parameters), it is permissible also to write, more briefly, $S \equiv (\alpha, \varepsilon)$, provided φ, f do not change. The structural properties to be estimated from the observations will then be the numerical values of α, ε ; or, more generally, some relations, not known a priori, between the parameters.

2.2.2. We can apply the earlier definition (section 1.1.4) of a self-contained model ($G=N$) and a sectional model ($G<N$). But the definition of a complete model and of exogenous vs. endogenous variables must be supplemented. The "subsidiary" system of equations will, in general, involve random shocks of its own, so that (1.2) becomes

$$(2.8) \quad \varphi_{G+k}(z) = u_k^{(k)}, \quad k = 1, \dots, K = N - G.$$

Rewriting the earlier condition (1.3) as

$$(2.9) \quad -\frac{\partial \varphi_{G+k}}{\partial y_g} = 0, \quad g = 1, \dots, G, \quad k = 1, \dots, K,$$

we see that this condition is not sufficient to make the variables z determinable outside of the system (2.6). To make them thus determinable, and the considered model complete, we must add the condition that the shocks u of the considered model and the shocks $u^{(z)} \equiv (u_1^{(z)}, \dots, u_K^{(z)})$ of the subsidiary model be distributed independently:

$$(2.10) \quad f_x(u, u^{(z)}) = f(u) \cdot f_z(u^{(z)}),$$

where f_x, f, f_z are density functions. If both conditions (2.9) and (2.10) are fulfilled, the model

$$(2.11) \quad \begin{aligned} \varphi_g(y, z; \alpha_{(g)}) &= u_g, \\ f(u; \epsilon), & \quad g = 1, \dots, G, \end{aligned}$$

is said to be complete, with y as the endogenous and z the exogenous set of observables.

2.2.3. The distribution (2.3), $g_x(x)$, of the observable vector can always be represented as a product of a conditional and a marginal distribution

$$(2.12) \quad g_x(x', x'') \equiv g(x' | x'') \cdot g_{x''}(x''),$$

where x', x'' is any pair of mutually complementary subsets of x . Suppose now that the model is complete; put $x' = y$, $x'' = z$, and denote by vector λ the parameters of the conditional distribution $g(y|z)$:

$$(2.13) \quad g_x(y, z) \equiv g(y | z; \lambda) \cdot g_z(z),$$

say. We observe that, by section 2.2.2, z is determined by the subsidiary model and is independent of u , the random disturbances of the considered model. The parameters of the ("subsidiary") marginal distribution $g_z(z)$ are not related to the parameters λ of the other factor in (2.13), the conditional distribution of the endogenous variables

$$(2.14) \quad g(y | z; \lambda).$$

This distribution is determined by the structure S , as is seen by solving the structural equations in (2.11) for y . The *reduced form*,

$$(2.15) \quad y = \bar{\eta}_\varphi[z, u; \bar{\pi}_\varphi(\alpha)],$$

differs from that of the nonstochastic case (1.8) inasmuch as it involves the random vector u . Given z , the value taken by $u \equiv (u_1, \dots, u_G)$ determines the value taken by $y \equiv (y_1, \dots, y_G)$.

The distribution of y depends on the elements of the structure $S \equiv (f, \varphi, \alpha, \varepsilon)$; this can be expressed thus:

$$(2.16) \quad g(y | z; \lambda) = g_{f, \varphi}(y | z; \lambda),$$

$$(2.17) \quad \lambda = \bar{\lambda}(\alpha, \varepsilon) = \bar{\lambda}_{f, \varphi}(\alpha, \varepsilon) = \bar{\lambda}(S).$$

Thus the structure S determines the form and the parameters of the conditional distribution $g(y | z; \lambda)$.

On the other hand, S does not determine the distribution $g_z(z)$, the other factor in the product in (2.13). For the purpose of estimating S from observations, we have to consider not the joint distribution $g_x(y, z)$ but merely the conditional distribution $g(y|z)$: that is, we can disregard the possibly random character of z , and treat z as if it were fixed in repeated samples of y .

2.2.4. The reduced form (2.15) is a system of stochastic equations. In the nonstochastic case, the prediction of y from z involved knowledge of π , the parameters of the reduced form. In the stochastic case, the prediction of y consists in estimating the parameters (mean, variance, etc.) of the conditional distribution of y , given z . That is, λ , the parameters of the conditional distribution $g(y|z; \lambda)$ are to be estimated. Now consider (2.17), and let

$$(2.18) \quad \lambda^0 = \bar{\lambda}(S^0) = \bar{\lambda}(\alpha^0, \varepsilon^0)$$

be the parameters of the conditional distribution during the observation period, that is, of the distribution

$$(2.19) \quad g^0[y(t) | z(t); \lambda^0], \quad t = 1, \dots, T.$$

It is possible to compute from observations an estimate l^0 of the vector λ^0 by an operation

$$(2.20) \quad \mathcal{P}_1 X^0 = l^0;$$

operation \mathcal{P}_1 is called *predictive estimation when structure is unchanged*. The notation is analogous to but not identical with that of (1.13) of the nonstochastic case. There we dealt with predictive determination of the reduced form, here with predictive estimation of the distribution of observables.

2.2.5. If changes of structure are intended or expected, the operation \mathcal{P}_1 will, in general, not suffice for the prediction of y from z .

Suppose it is possible to compute estimates (a^0, e^0) of the parameter vectors $(\alpha^0, \varepsilon^0)$ of the observational structure S^0 . Call this operation \mathcal{S} :

$$(2.21) \quad \text{est } S^0 \equiv (a^0, e^0) = \mathcal{S}X^0.$$

\mathcal{S} thus denotes *structural estimation* (analogous to the structural determination of section 1.2.3.6). Now let the observational structure undergo a change \mathcal{U} : i.e., $\mathcal{U}S^0 = S$. The estimates l of the parameters λ of the new distribution of observables are obtained by the operation

$$(2.22) \quad \begin{aligned} l &= \bar{\lambda}(a, e) = \bar{l}[\mathcal{U}(a^0, e^0)] \\ &= \bar{\lambda}(\mathcal{U} \mathcal{S}X^0) = \mathcal{P}_{\mathcal{U}} X^0, \end{aligned}$$

say; the operation $\mathcal{P}_{\mathcal{U}}$ is called *predictive estimation when structure undergoes a given change \mathcal{U}* .

Whether optimal properties of the estimates a^0, e^0 are preserved under subsequent transformations $\mathcal{U}, \bar{\lambda}$, depends on the estimating methods; in particular, the function ψ of a maximum-likelihood estimate of a parameter θ is the maximum-likelihood estimate of the function $\psi(\theta)$.

2.2.6. One can interpret the operations (2.20) and (2.22) more generally. l^0 and l may be used to denote, not the point estimates of the parameters λ^0 and λ , but the estimates of the parameters of the joint distribution of those estimates or the estimates of their confidence regions. However, this interpretation need not be applied in the present article.

2.2.7. Structural estimation is needed, in particular, when

structural changes are intended, i.e., when *structural policies* are being undertaken. This was shown in section 1.2.3.5 with regard to structural determination when the model was nonstochastic. Section 2.5.1.3 illustrates the stochastic case.

We have to modify the earlier treatment since y and consequently the gain (welfare) ω as previously defined are random variables. For example, the real income (of a firm or of a nation, as the case may be) is the quotient of money income over price level, i.e., a function of endogenous variables y . Real income therefore depends on the value of exogenous variables z and on the distribution $g(y | z; \lambda)$, and is a random variable. We must, accordingly, redefine the gain (welfare) ω , i.e., the quantity which the policy-maker tries to maximize. The policy-maker will prefer a high mean value of real income to a low mean value; he may possibly at the same time prefer a small variance of income to a large variance. In short, he will prefer certain probability distributions of income to others. He will maximize a quantity ω that is a functional of the probability distribution function of income.¹ But the probability distribution of income depends on the value of the exogenous variables z and on the distribution $g(y | z; \lambda)$. Since the form of g is fixed a priori (depending on the function forms φ, f of the model), the gain ω is a known function of the distribution parameters λ and of z only. By (2.17) the parameters λ depend on the structure S in a way uniquely determined by the model. The structure S is, in turn, (analogously to section 1.2) the result of controlled or uncontrolled changes applied to the observational structure. Thus

$$(2.23) \quad \lambda = \bar{\lambda}(S) = \bar{\lambda}(\mathfrak{U}_c \mathfrak{U}_u S^o).$$

Furthermore, the exogenous variables can be either controlled (z_c) or uncontrolled (z_u) by the policy-maker. His aim is thus to get

$$(2.24) \quad \hat{\omega} = \omega(\hat{z}_c, z_u, \hat{\mathfrak{U}}_c \mathfrak{U}_u S^o),$$

where ω is a known gain function, and $\hat{z}_c, \hat{\mathfrak{U}}_c$ are, respectively, the optimal values of controllable exogenous variables and of

¹Compare [Hurwicz, 1946, p.132], [Tintner, 1942]. This statement can be shown to be implied in the statement that the mean value of utility of real income is maximized.

controllable structural changes. That is, $\hat{\omega}$ is a maximum of ω with respect to policies (z_c, \mathfrak{J}_c) .

To estimate the best policy $(\hat{z}_c, \hat{\mathfrak{J}}_c)$, for a given set (z_u, \mathfrak{J}_u) , of uncontrolled exogenous variables and uncontrolled structural changes, one has to have an estimate of the observational structure S^0 , i.e., to have the estimates (a^0, e^0) of the observational structural parameters $(\alpha^0, \varepsilon^0)$.

In practice, the maximization of ω with respect to policies consists in estimating for each considered policy (z_c, \mathfrak{J}_c) , the parameters λ of the conditional distribution of y given z , and the consequent probability distribution of (say) real income. The policy which gives the probability distribution that is preferable to all others is then chosen.

2.2.8. *Relation between regression, reduced form, and structural equations.* Sometimes the series of alternative policies considered will be such as to make it unnecessary to know the joint conditional G -variate distribution (2.14) of all endogenous variables after the change in structure. It may suffice to know the distribution of some particular subset, say $y_{(a)}$ of y , that is, the distribution

$$(2.25) \quad g_{(a)}(y_{(a)}, z; \lambda_{(a)}).$$

Such conditional distribution is easily derived from the general conditional distribution (2.14)¹. The set $y_{(a)}$ is called the *predicand*; the *predictor* is, in this case, the whole exogenous set z . Practical importance may be attached in particular to the simplest case, viz., the univariate conditional distributions

$$(2.26) \quad g_1(y_1 | z; \lambda_1), \dots, g_G(y_G | z; \lambda_G).$$

In particular, one may be interested in the first moments of such distributions, e.g., in the expectation

$$(2.27) \quad \mathfrak{E}(y_1 | z) = R(z; \chi_{(1)}).$$

R is the regression function of y_1 on z . Its parameters, denoted here by the vector $\chi_{(1)}$, are called regression coefficients if R happens to be linear, so that, writing $\chi_{(1)} \equiv (\chi_{10}, \chi_{11}, \dots, \chi_{1K})$,

¹See [VI].

$$(2.28) \quad \mathcal{E}(y_1 | z) = \sum_1^K \chi_{1k} z_k + \chi_{10}.$$

A sufficient condition for the function R in (2.27) to be linear is that the functions φ , f of the model be linear and normal, respectively. We can then conveniently split the structural coefficients α into two subsets $\beta \equiv \{\beta_{gh}\}$ and $\gamma \equiv \{\gamma_{gk}\}$, ($g, h = 1, \dots, G$; $k = 0, \dots, K$), denoting respectively the coefficients of endogenous and of exogenous variables; further we can choose the constant terms γ_{g0} so as to make $\mathcal{E}u = 0$. The model (2.11) becomes:

$$(2.29) \quad \sum_{h=1}^G \beta_{gh} y_h + \sum_{k=1}^K \gamma_{gk} z_k + \gamma_{g0} = u_g, \quad g = 1, \dots, G,$$

$$f(u),$$

where $f(u)$ is normal. In this case the reduced form (2.15) will also be linear; for example the first equation of the reduced form will be

$$(2.30) \quad y_1 = \sum_1^K \pi_{1k} z_k + \pi_{10} + \text{a linear function of } u.$$

Taking the expectation of y_1 given z , we find that the regression coefficients $\chi_{(1)}$ and the coefficients $\pi_{(1)}$ of the reduced form coincide: see (2.28).

Furthermore, suppose our complete model consists of one equation only ($G=1$); if we choose the units of the z 's so as to make $\beta_{11} = 1$, we can write the unique structural equation as

$$(2.31) \quad y_1 - \sum_{k=0}^K \gamma_{1k} z_k - \gamma_{10} = u_1.$$

In this particular case the coefficients π_1 of the reduced form will coincide not only with the regression coefficients $\chi_{(1)}$, but also with the structural coefficients $\gamma_{(1)}$. [This applies, of course, also to the case where a complete model involving $G > 1$ variables can be split into G uniequational complete models such as (2.31); in each of which, one endogenous variable depends on

exogenous variables only.] In general, however, there will be no such coincidence. An example will be given in section 2.5.1.2.

2.2.9. *Predicting exogenous variables.* If the conditional distribution (2.19) estimated from past observations is known to remain valid for $t > T$, this distribution, or any of its practically relevant conditional distributions (2.25) can be used, respectively, as the future distribution of y , or of a subset $y_{(d)}$, for given future values of z .

However, it may also be desirable to make predictions about exogenous variables, in particular the uncontrolled ones, z_u . It is seen from section 2.2.7 that the choice of best policy presupposes the knowledge of uncontrolled factors: z_u as well as \mathcal{J}_u . It is possible to estimate future values of the variables z_u if they are related to observable variables that do not enter the model considered. For example, suppose prediction is done on behalf of a firm to help it in the choice of its policies. The firm's sales is an endogenous variable in the model describing the demand and production conditions for the firm's product. This model contains national income as an exogenous variable which the firm cannot control. The national income itself is an endogenous variable of another model; this may also include exogenous variables such as foreign crops (which affect the demand for this country's exports of manufactured goods, and hence affect this country's income). Thus the firm will be interested in predicting national income, using distributions such as (2.26), where y_1 would denote national income, and z will include foreign crops z_1 . However, another variable endogenous to the national economy model - e.g., the imports y_2 - can be usefully included in the predictor set, if, in the future, information on national income becomes available later than information on imports. The prediction of y_1 for a given z and a specified value of y_2 is more accurate than the prediction of y_1 from z with the value of y_2 unspecified. Thus, it may be useful to derive from the observational distribution $g^o(y | z)$ a distribution

$$(2.32) \quad g_{(d,r)}^o(y_{(d)} | y_{(r)}, z),$$

say; this distribution is more general than (2.25) in that its predictor set (y_r, z) includes endogenous variables, $y_{(r)}$. In particular, we may derive a regression equation to estimate

$$(2.33) \quad \mathcal{E}(y_1 \mid y_2, z).$$

We may also want to estimate a variable exogenous not only to the relevant model (e.g., the model of the firm) but also to the model (e.g., of the national economy) from the variables of which it has to be estimated. For example, foreign crops z_1 may be of direct interest to the firm and also may influence national imports y_2 ; but observations on factors determining foreign crops may be unavailable. If data on y_2 become known earlier than those on z_1 , z_1 can be estimated from y_2 . This case differs from the preceding one inasmuch as z_1 may be a nonstochastic variable. However, it may be possible to use the past distribution $g^0(y \mid z; \lambda^0)$ of the observables to derive a confidence region for z_1 (looked upon as a parameter): values of z_1 such as would give rise to the known values of y_2 only with a small probability (lower than a preassigned significance level) will lie outside the confidence limits for z_1 .

2.2.9.1. These remarks are of some importance in view of many attempts to use endogenous variables as predictors; e.g., to predict national income from contemporaneous imports, or retail sales, etc. It may be possible to determine such a relation from past observations. But its usefulness is limited to the cases just mentioned. To choose between various national policies, it cannot be useful to be able to predict national income from imports, since the latter are themselves affected by any policy chosen. For purposes of private policies, on the other hand, such prediction can be useful. But such prediction is possible only inasmuch as information on imports is available earlier than that on national income; and only if, in addition, the relevant aspects of the national economic structure can be assumed unchanged between the period of observation and the time when the private policy is going to be applied. The same considerations are valid for the case when, because of a lag in the available data, variables endogenous to the national economy are used to estimate exogenous ones.

2.3. Identification

2.3.1. In section 1.3 it was shown for nonstochastic models

that structure can be determined to the extent that the a priori conditions (i.e., the model) suffice to make the structure uniquely identifiable. In this case, only one structure will correspond to a given reduced form. Or, in terms of parameters: if only one set α^o corresponds to a given set π^o , then, since π^o can always be determined from (sufficiently numerous) observations, the set α^o can also be so determined.

In the stochastic case, a structure is said to be uniquely identifiable by the model if only one structure compatible with the model corresponds to a given conditional distribution of the observations. Or, in terms of parameters: if only one set α^o, ε^o , corresponds to a given set λ^o , then, since λ^o can be estimated from the observations, the structural parameters α^o, ε^o can also be so estimated.

2.3.2. The structural transformations $\mathfrak{D}_c, \mathfrak{D}_u$ and the gain functional ω may be such as to make identifiability of all parameters α^o, ε^o unnecessary for the choice of best policy. Section 1.3.10 on partial identifiability applies accordingly.

2.3.3. If the structure, or a part of it, is not identifiable its estimation is not possible, however numerous the observations on the variables treated as observables in the model. However, observations on other variables may provide additional information (which is a priori with respect to the structure considered; see section 1.1.0) such as to make the structure, or its relevant part, identifiable. The failure of a model to identify the structure is not a ground for rejecting the model; rather, it calls for additional information, to be provided by a new type of observations. Suppose, for example, that the relation describing the investment behavior of the aggregate of American firms is not identifiable within a model that involves, in addition to this relation, relations describing the behavior of consumers, lenders, etc. This would make it necessary to add a new type of information, based, e.g., on records of single firms, or on interviews with businessmen.

2.3.4. The a priori information provided by the model may involve the parameters α of the structural equations as well as the parameters ε of the shock distribution. If no a priori information on the parameters ε of the shock distribution exists, and if a nonstochastic structure

$$(2.34) \quad \varphi_g(y, z; \alpha_{(g)}) = 0, \quad g = 1, \dots, G,$$

is not completely identifiable, then the parameters α of the stochastic structure

$$(2.35) \quad \begin{aligned} \varphi_g(y, z; \alpha_{(g)}) &= u_g, \\ f(u; \epsilon), & \quad g = 1, \dots, G, \end{aligned}$$

will also not be completely identifiable; but information on the distribution $f(u; \epsilon)$, if available, may make (2.35) completely identifiable. Thus the example in section 1.3.5 may or may not hold for the stochastic case if information on the shock distribution is available: such as the knowledge that the shocks on the demand side are not correlated with those on the supply side; or knowledge of the ratio between the variances, etc.: [Frisch, 1933], [Mann and Wald, esp. p.219], [Marschak and Andrews, §§ 18-22], and sections 2.5.1.1, 2.5.1.4 below.

2.3.5. In particular, a uniequational complete model is completely identifiable, apart from trivial transformations (section 1.3.7). Hence the remarks in section 1.4.2, on the role of experiments (of type II) in making the observational structure identifiable, will apply. If, in addition (section 2.2.8); the uniequational complete model generated by an experiment is linear with shocks distributed normally, the structural coefficients will not only be identifiable but will coincide with the regression coefficients. The model can then also be regarded as one in which one (the "dependent") variable is subject to measurement errors while all others are free from such errors: a familiar case in the history of the application of statistics to experiments. It arises in nonexperimental science only if the mechanism producing the observations can be adequately represented by a model involving only one nonlagged endogenous variable.

2.4. Dynamic Models

2.4.1. A model is called dynamic if it has at least one of the following two properties: 1) at least one observable variable occurs in the structural equations with values taken at various points of time (this includes the case of time derivatives, differences, and integrals over time); 2) at least one equation contains functions of time (trend, seasonal fluctuations, etc.). If the first property is present the model is called *multitemporal*; if both properties are absent, it is called *unitemporal*.

2.4.2. We obtain a *discrete multitemporal* shock model if the equations in (2.11) are replaced by

$$(2.36) \quad \varphi_g[y'(t), \dots, y'(t-\tau_{y'}), z'(t), \dots, z'(t-\tau_{z'})]; \alpha_{(g)}] = u_g(t),$$

and the distribution of shocks in (2.11) by

$$(2.37) \quad f[u(t), \dots, u(t-\tau_{y'})]; \varepsilon].$$

The time interval between two successive observations is chosen as a time unit; time lags smaller than 1 are not admitted by the model; $\tau_{y'}$ and $\tau_{z'}$ denote the largest time lag with which the corresponding variables occur; and t takes all integral values through the time interval during which the model is supposed to be valid.

2.4.3. A case of great practical importance arises when successive shocks are mutually independent; that is,

$$(2.38) \quad \begin{aligned} f[u(t), u(t-1), \dots, u(t-\tau_{y'})] \\ = f[u(t)] \cdot f[u(t-1)] \cdot \dots \cdot f[u(t-\tau_{y'})], \end{aligned}$$

say. If this condition is fulfilled then, as shown in [I-1] and [XVII], lagged variables can be treated as additional variables in judging the identifiability of a structure. Lagged endogenous variables can then be treated as if they were fixed in repeated samples (i.e., like exogenous variables). The following model can then be regarded as a complete one:

$$(2.39) \quad \varphi_g[y(t), z(t); \alpha_{(g)}] = u_g(t), \quad g = 1, \dots, G,$$

$$(2.40) \quad f[u(t); \varepsilon],$$

where the notation in (2.36) has been changed as follows:

$$y'(t) \equiv y(t) \equiv \text{"jointly dependent variables,"}$$

$$(2.41) \quad [y'(t-1), \dots, y'(t-\tau_{y'}), z'(t), \dots, z'(t-\tau_{z'})] \equiv z(t) \\ \equiv \text{"predetermined variables,"}$$

following Koopmans' terminology [XVII].

2.4.4. It is plausible to assume, as in (2.38), the statistical independence of successive shocks if the time interval between two successive observations is not too short. It would be more realistic to study models in which at least some of the observables $y(t)$, $z(t)$, and also some of the shocks $u(t)$ are functions of a continuous time variable. Some properties of such continuous stochastic models are outlined in [XVI]. However, the remainder of this volume deals with discrete models and most of the time assumes condition (2.38) to be valid.

2.4.5. The reduced form defined in (2.15) applies to the multitemporal model with variables $y(t)$ and $z(t)$ defined by (2.41). The variables on the right-hand side of the reduced form (the predictor set) occur with lags no higher than those in the structural equations. However, other kinds of predictor sets can be considered. In the *separated form* each nonlagged variable is expressed as depending only on its own lagged values and on the exogenous variables with or without lags; this is possible because, by "shifting time back," one can obtain enough structural equations to eliminate the lagged values of other endogenous variables. One can go further and eliminate the lagged values of all endogenous variables, leaving only the exogenous variables (with and without lags) as the predictor set: the *resolved form*.

As shown by Tinbergen, the separated form (his "final equation") can be used to build up, year by year, the path of a single endogenous variable beginning with given initial values and giving effect every year to changes in exogenous variables; that is, the solution of each equation of the reduced form (looked upon as a difference equation) expresses the predictand variable as a function of time ("cyclical fluctuations") and of exogenous variables.

Corresponding to the reduced, the separated, and the resolved forms in which a multitemporal model can be written, there are three kinds of distributions with which predictive estimation under unchanged structure can be concerned, and which can be estimated (in principle) regardless of whether the structure is or is not identifiable. However, as will be shown below (section 2.5.3), only in the case of reduced form, but not in the case of separated and of resolved form, is the estimation amenable to known methods.

2.5. Estimation

2.5.0. Before we summarize the results obtained in this

volume with regard to the estimation of relevant parameters from finite samples, we shall restate and illustrate certain population properties. We shall then discuss the estimation of regression equations, and of complete and incomplete structures, and add a few remarks on the choice of models.

2.5.1. *Population properties.* Important differences between, on the one hand, the structural parameters α , and, on the other, the parameters π of the reduced form and the regression parameters χ , will be briefly restated. (It is understood that the parameters to be treated here are generated by the observational structure S^o ; but the superscript o is omitted from the symbols where no ambiguity arises.) To fix the ideas we shall assume the model to be linear with normally distributed shocks and shall give examples when certain elements of χ or π are or are not equal to some elements of α .

While χ and π can always be estimated from observations, the same is not true of the structural parameters α unless they are identifiable. Furthermore (section 2.2.8) the coefficients $\pi_{(1)}$, of that equation of the reduced form which determines the endogenous variable y_1 , are equal to the coefficients $\chi_{(1)}$ of the corresponding regression equation. If there exists a structural equation relating y_1 to the exogenous subset z and containing no endogenous variables, its coefficients $\alpha_{(1)}$ will also be equal to $\chi_{(1)}$; provided this structural equation, if stated together with the distribution $-f_1(u_1)$, say - of its shock-variable, constitutes a complete model. The latter condition implies (section 2.2.2) that u_1 is distributed independently of the shock variables u_2, \dots, u_G entering the remaining structural equations. If the structural equation considered contains lagged endogenous variables, and if successive shocks are independent, the statements just made can be extended so as to include in z not only the exogenous but also the lagged endogenous variables (section 2.4.3).

For example, structural and regression coefficients will coincide in each of the following complete models (each is assumed to be linear, with successive shocks independent and normally distributed; the choice of the "dependent" variable in the regression will be indicated in each case by the phrase describing the model). Of these four models, the first three are uniequational, while the last one can be partitioned into two uniequational complete models!

¹A more complicated model partitionable into four uniequational models is given in [Bentzel and Wold, p.104].

(1) Current yield per acre depends on current temperature and current rainfall.

(2) Current yield per acre depends on current temperature and rainfall, and also on the amount of fertilizers fixed by the experimenter.

(3) Current yield per acre depends on current and past temperature and rainfall, and on the decision farmers made in the previous year regarding the amount of fertilizers.

(4) Market price of a nonstorable good depends on supply (demand equation); supply depends on previous year's price (supply equation); the behavior of buyers and that of sellers undergo mutually independent shocks.

2.5.1.1. In the following example¹; on the other hand, each structural equation contains more than one endogenous variable, and cannot therefore (taken together with the distribution of its shock variable) be regarded as a complete model in itself; but the two equations together constitute a complete model. We shall show that the structural coefficients will not be equal to any of the regression coefficients. Let y_1 (national product, identical with national income, or the supply of all goods) and y_2 (demand for all goods) obey the following equations for any time t :

$$y_2(t) - \beta_1 y_1(t) - \beta_0 = u_1(t) \quad (\text{behavior of buyers}),$$

$$y_2(t) - y_1(t) = u_2(t) \quad (\text{behavior of producers}).$$

The model does not contain predetermined variables. Its random shocks are $u_1(t)$ ("shift of demand") and $u_2(t)$ ("failure to adjust production to demand"). Each pair of values $u_1(t), u_2(t)$ determines

$$y_1(t) = [u_1(t) - u_2(t) + \beta_0] / (1 - \beta_1),$$

$$y_2(t) = [u_1(t) - \beta_1 u_2(t) + \beta_0] / (1 - \beta_1).$$

Suppose successive shocks are independent and the joint distribution $f[u_1(t), u_2(t)]$ of the shocks is normal and independent of time.

¹For more fully developed examples from the same branch of economics see [Haavelmo, 1947 A].

Let its moments be

$$(2.42) \quad \mathcal{E} u_1(t) u_2(t) = 0; \quad \mathcal{E} u_1^2(t) = \sigma_{11}; \quad \mathcal{E} u_2^2(t) = \sigma_{22}.$$

Then the distribution $g(y_1, y_2)$ of the observables is normal and has moments

$$(2.43) \quad \begin{aligned} \mathcal{E} y_1 &\equiv \mathcal{E} y_2 = \beta_0 / (1 - \beta_1); & \sigma_{y_1 y_2} &= (\sigma_{11} + \beta_1 \sigma_{22}) / (1 - \beta_1)^2; \\ \sigma_{y_1 y_1} &= (\sigma_{11} + \sigma_{22}) / (1 - \beta_1)^2; & \sigma_{y_2 y_2} &= (\sigma_{11} + \beta_1^2 \sigma_{22}) / (1 - \beta_1)^2. \end{aligned}$$

These are four mutually independent equations (not counting the identity $\mathcal{E} y_1 \equiv \mathcal{E} y_2$) to determine the four unknown structural parameters $S = (\beta_1, \beta_0, \sigma_{11}, \sigma_{22})$ from the parameters of the distribution $g(y_1, y_2)$ of the observables. The structural parameters are completely identifiable.

2.5.1.2. Consider now the regression of y_2 on y_1 :

$$\mathcal{E}(y_2 | y_1) = \chi_1 y_1 + \chi_0,$$

say, and

$$(2.44) \quad \chi_1 = \sigma_{y_1 y_2} / \sigma_{y_1 y_1} = (\sigma_{11} + \beta_1 \sigma_{22}) / (\sigma_{11} + \sigma_{22}),$$

which approaches β_1 or 1 as $\sigma_{11} / \sigma_{22} \rightarrow 0$ or $\rightarrow \infty$, respectively. Thus β_1 , the "marginal propensity to spend," is distinct from the regression coefficient χ_1 of spending (y_2) on income (y_1), except in a limiting case when the behavior of buyers is not subject to random shocks (while the behavior of producers is).

2.5.1.3. This example also illustrates the different practical purposes of estimating β_1 , β_0 or χ_1 , χ_0 . Suppose a firm expects a given change in the economic structure; for example, a rise in the "marginal propensity to spend" whereby β_1 will be replaced by $k\beta_1$; suppose the other three structural parameters are known to stay unchanged. If the firm wants to use old observations on y_1 , y_2 to predict the distribution of these variables under the new circumstances, it will have first to obtain β_1 , β_0 , σ_{11} , σ_{22} from observations, and then insert the new set of structural parameters

$(k\beta_1, \beta_0, \sigma_{11}, \sigma_{22})$ into (2.43). Furthermore, if the firm wants to estimate national spending from national income (because the latter is published earlier than the former; section 2.2.9.1), it will have to use, not the old regression coefficient (2.44), but a new regression coefficient, $(\sigma_{11} + k\beta_1\sigma_{22})/(\sigma_{11} + \sigma_{22})$. This cannot be obtained from χ_1 , but can be obtained from the old structural parameters $(\beta_1, \beta_0, \sigma_{11}, \sigma_{22})$ since k is known. Only if the structure is known to remain unchanged can χ_1 be used.

2.5.1.4. If the structure were not identifiable, there would be no way to predict effects of known structural changes. This would be the case, for example, if the a priori assumption - in (2.42) - of noncorrelated shocks could not be admitted. We should then have five structural parameters to determine; the four equations (2.43) would therefore not suffice. On the other hand, the structure would be identifiable, even with the noncorrelation assumption dropped, if a different predetermined variable had been introduced into each of the two equations of the model.

2.5.2. *Estimation of regression coefficients.* When structure is known to remain unchanged, estimates of regression coefficients χ , together with other parameters of the distribution of observables, help to estimate the future distribution of predictand variables for given values of others, and hence to estimate future positions of maximum gain. When, on the other hand, structure is known to undergo a given change, the estimation of future values of variables requires knowledge of the structural parameters (α, ϵ) . We have seen (section 2.2.8) that in the special case of certain uniequational complete models (and of models that can be decomposed into such uniequational complete models), $\alpha = \chi$. In this case, if $q = q(X^0)$ is a function of observations X^0 that is an unbiased estimate of the parameters χ , then $\mathcal{E}q = \chi = \alpha$; that is, q can serve as an unbiased estimate of structural parameters. In general, however, there is no equality between α and χ . Therefore a function of observations that is an unbiased estimate of χ cannot be an unbiased estimate of α , even for infinitely large samples. If such a function is used as an estimate of α for predictive estimation under a structure subjected to given changes, the future endogenous variables, and hence the future gains, will be estimated with a bias; consequently, other than optimal policies will be chosen.

This does not make the estimation of regression coefficients useless: uniequational complete models of the appropriate type may

exist; and structures do not always change. Hence our interest in reconsidering the properties of regression coefficients.

In a general multitemporal model, the variables are connected by a set of equations

$$\varphi_g[y_1(t), \dots, y_1(t-\tau_1), y_2(t), \dots, y_G(t-\tau_G), z] = u_g,$$

$$g = 1, \dots, G,$$

where z denotes exogenous variables (lagged as well as nonlagged). The regression equation for y_1 - denoting by $(y_{(r)}, z)$ the whole predictor set - is

$$\mathcal{E}[y_1(t) | y_{(r)}, z] = R(y_{(r)}, z),$$

which has not been studied in general. The simplest case is that of the linear regression equation with $y_{(r)}$ empty:

$$(2.45) \quad \mathcal{E}[y_1(t) | z] = \sum_1^K \chi_{1k} z_k + \chi_{10};$$

it is satisfied when $y_1(t)$ and z are connected by an unequational unitemporal model, cf. (2.31):

$$(2.46) \quad y_1(t) - \sum_1^K \gamma_{1k} z_k - \gamma_{10} = u_1(t),$$

provided successive shocks $u_1(t), u_1(t+1), \dots$, are normally and independently distributed. Certain optimal properties of the least-squares estimate of χ ($=\gamma$) are well established for this case. The case has to be generalized in two important directions (possibly, but not necessarily, preserving the assumptions of linearity and normality of the structure): 1) the complete model may be made multiequational; 2) the complete model may be made multitemporal.

In the case of multiequational (but unitemporal) model, the properties of the least-squares estimate and the maximum-likelihood estimate of the regression coefficients, and in particular the conditions of equivalence of these two estimates, have been studied in [VI] and [VII] by Hurwicz and Koopmans for large as well as small

samples.

For the case of a uniequational multitemporal model with only one lag,

$$(2.47) \quad y_1(t) = \alpha y_1(t-1) + u_1(t),$$

where $u_1(t)$ is distributed normally with a constant variance and with zero covariances $E u_1(t) u_1(t')$, $t \neq t'$, Hurwicz [XV] shows that the least-squares estimate of the regression (and structural) coefficient α has a bias. For $T = 20$ (a length of time series common to economic studies based on annual data of the interwar period), the bias approaches 9 per cent as $\alpha \rightarrow 0$. The bias seems to disappear as the sample increases or as $\alpha \rightarrow 1$.

2.5.3. This bias in the regression estimates from short time series will in general exist in every multitemporal model. If the model is multiequational, there arises, in addition, the question of the *choice of predictors*. Each predictand variable can be expressed as a function of the predictor variables and of the shock variables as in (2.15) and (2.30). Only contemporary shock variables will be involved provided the predictor variables occur with the same lags as those of the structural equations. Such is the case with the reduced form (see section 2.4.5). The predictand will then be related to the predictor by an equation such as (2.30), where the z 's would stand for exogenous as well as lagged endogenous variables; and where the successive random terms will be independent. In this case, least-squares large-sample estimates of the coefficients of the variables will have the usual optimal properties. But it is different in the case of the "resolved" and the "separated" forms (section 2.4.5). In these cases, structural shocks relating to various points of time will be contained in the same equation, as the result of eliminating certain endogenous variables after replacing t by $t-1$, $t-2$, ..., or "shifting the time back." Therefore successive random terms of a "resolved" or a "separated" form are not independent; and the least-squares estimates of corresponding regression equations will, in general, be biased.¹

2.5.4. As stated in section 2.3.2, expected or intended structural changes may be of such a nature that the estimation of all the parameters of a complete model is not necessary for purposes of

¹ Pointed out by Haavelmo. See [Klein, 1946 B, p.303 ff.].

prediction and policy: hence the practical importance of *partial structural estimation*, i.e., the estimation of a selected set of parameters. Such estimation is possible if the set of parameters in question is identifiable, even though some or all of the remaining parameters may not be identifiable. In particular, the estimation of the parameters of a few selected structural equations, or even of a single structural equation (e.g., demand equation) has an obvious practical interest.

On the other hand, the expected or intended changes of structure may not be known long beforehand. It is therefore often desirable to have estimated the structure completely.

2.5.5. The method of *complete structural estimation* most fully discussed in the volume is that of maximum likelihood. The joint probability density of all the observed values of the variables is regarded as a function (the likelihood function) of the structural parameters. Those values of the parameters for which this function attains its highest value are called maximum-likelihood estimates. In important classes of cases these estimates are "consistent" (they converge with probability 1 to the true values in the limit for infinitely large samples) and "efficient" (they have, in large samples, variances that never exceed those of any other normally distributed estimate). In [Mann and Wald] are discussed the maximum-likelihood estimates obtained when all available a priori information is used and when the model is as follows: a completely identifiable complete system of linear difference equations with no exogenous variables; shocks mutually independent; the model generates a stationary process, i.e., the observable variables would converge in time to constant values if shocks were absent. These authors proved the consistency (but not the efficiency) of the estimates under these conditions; they also showed that, in the absence of a priori restrictions, the structure is not identifiable. In [II-3] Koopmans and Rubin estimate the parameters of a complete linear model that is more general owing to the introduction of exogenous variables. Their discussion also covers the estimation of the parameters of some identifiable structural equations if other equations of the complete system are not identifiable. They give a proof of the consistency of the estimates which takes account of these two generalizing assumptions. In [XIV] Rubin extends the proof to a simple case of a nonstationary ("explosive") process.

There is an alternative maximum-likelihood method for obtaining consistent estimates of structural parameters: each equation of the structure is estimated separately via the reduced form as

described in section 2.5.6. For each equation, this procedure has to use at least as much a priori information as is necessary to make that equation identifiable; but it leaves unused a (possibly large) part of the remaining a priori information. This method will therefore lead, in general, to less efficient estimates (i.e., larger sampling variances of the estimates) than the maximum-likelihood method using all information described above, although both kinds of estimates have the consistency property. For brevity, we call the two estimation methods "information-preserving maximum-likelihood estimation" and "limited-information maximum-likelihood estimation," respectively. The latter method is also known as the "method of reduced forms."

The practical usefulness of the consistency property of estimates diminishes if the sample becomes small. Small sample properties of structural estimates have not been studied, except [XV] for the unequational multitemporal model, equation (2.47) of the present article. The bias found in the estimate in this case suggests that in general both proposed methods of estimating multitemporal structures are biased if applied to short time series.

2.5.6. *Incomplete (partial) structural estimation.* A complete model (and structure) was defined above, in sections 1.1.4, 1.2.1.1, and 2.2.2. Obviously the limited-information maximum-likelihood method just described can be used to estimate complete as well as incomplete structures. At the time when this introduction is being written this method appears to be best developed. But other suggestions for estimating incomplete structures will also be discussed below (section 2.5.6.1) after adding a few more remarks on the reduced forms method.

The coefficients π of a linear reduced form (2.30) are regression coefficients of a nonlagged endogenous variables on exogenous and lagged endogenous variables. Their least-squares estimates have the consistency property. At the same time, they are functions of the structural parameters. This suggests that consistent estimates of structural coefficients can be obtained by applying appropriate transformations to the least-squares estimates of the coefficients of the reduced form. Similarly, parameters (e.g., the variance) of the estimated distribution of the random terms of the reduced form can be transformed into the parameters of the distribution of the shocks in the structural equations. The suggestion has been familiar for some time¹ but has been applied rigorously for the first

¹See, for example, [Mann and Wald, p.219], [Haavelmo, 1944, pp. 103 - 104].

time by T. W. Anderson, M. A. Girshick, and H. Rubin; their joint work is summarized in [IX]. Since an equation of the reduced form has as many unknown coefficients as there are predetermined variables in the system, it will, in general, have more parameters than the structural equations (or equation) to be estimated. Hence equations connecting the unknown estimates of structural parameters with estimates of parameters of the reduced form may be more numerous than the number of these unknowns. To avoid this overdeterminacy, part of the available information has to be dropped. In particular, the method of reduced forms, which has been applied so far only to one equation at a time - in [Girshick and Haavelmo], [Haavelmo, 1947 A], [Klein, 1947, 1950] - does not use the observations on jointly dependent variables outside of this equation, though it does use the observations on predetermined variables of the system. The only a priori information this method uses are the linear restrictions on the parameters of the equation in question (including the prescription as to which variables enter this equation).

2.5.6.1. Other procedures of incomplete structural estimation were suggested by Koopmans [1945] and by Wald [VIII]. To estimate the parameters of F ($< G$) structural equations, Koopmans proposed the following approximation method. The $G-F$ "complementary" equations of the model are "sketched in," e.g., by using admittedly biased single-equation least-squares estimates or some a priori guesses. Then proceed with the estimation of the parameters of the F equations to be estimated. The estimates of the remaining parameters of the complete model will then be improved compared with what they would be if the complementary part of the model were entirely neglected.

Wald's suggestion, [VIII], is different: even if we do not know anything about the complementary part of the model, our a priori knowledge about the F equations that interest the investigator may be sufficient to exclude hosts of originally admissible hypotheses about these equations, and thus to construct confidence regions for parameter estimates. For example, if we know that successive shocks are independent, then a set of values of the structural parameters must be rejected whenever the estimate of shocks computed from observations (the "residuals") fails to have approximately the characteristics of a random series. The elaboration of the method will consist of showing how to use a priori knowledge concerning an incomplete model to construct shortest confidence regions. Essentially, the difference between Koopmans' and Wald's suggestions on the estimation of incomplete models consists in

attaching different weight to our knowledge of the complementary part of the entire model.

The method of reduced forms makes use of the maximum-likelihood principle, since least-squares estimates of regression coefficients are maximum-likelihood estimates (under normally distributed shocks); however, it leaves unused certain types of available a priori knowledge. In Wald's approach to incomplete systems, still fewer known restrictions are used.

For a given number $F (\leq G)$ of equations that are to be estimated, the choice of method will depend on two considerations: the mathematical and computational simplicity on the one hand, and the degree of use of available information on the other. For $F = 1$, the reduced forms method is less laborious than the approach via the estimation of the complete system (Wald's method has not yet been studied in this respect). But the complete estimation has one advantage over the other methods: it utilizes in full the a priori information on the model as well as the observations on all variables of the system.

2.6. *The Choice of Model*

There are many competing sets of a priori restrictions that can be imposed upon the structural parameters without contradicting what we know of human behavior and environment; and there is a wide variety of functional forms to specify the relations and distributions involved. It is often asserted that the choice is much wider in economics than in other empirical sciences. The usual testing considers only one hypothesis (and its negation) at a time. This is an inadequate procedure when a number of hypotheses classifiable according to a large number of attributes are in competition. In this volume no attempt is made to approach this problem.¹ In fact, the present volume is little concerned with the testing of hypotheses. Yet, the following remarks implied in the basic ideas of this volume seem appropriate.

A completely identifiable structure is said to be "just identifiable" by the model if the omission of one of the a priori restrictions of the model makes the structure incompletely identifiable; a completely identifiable structure that is not just identifiable is called overidentifiable. If several alternative overidentifying models are acceptable on a priori grounds, each

¹See [Wald, 1942, pp. 8 - 9], [Brookner]. This is the problem of "multiple (as distinguished from dual) decisions."

of them (and its negation) can be tested against data by the existing methods; this is, in fact, attempted in [IX-6]. However, new methods are needed to test the whole set of such models simultaneously.

Regarding the great variety of functions equally appropriate, on a priori grounds, to describe structural economic relations, one may expect some help from the statisticians' recent attempts at nonparametric estimation of distribution functions [Wald and Wolfowitz]. Certain weak a priori restrictions on the structural relations, such as the sign of certain partial derivatives, the independence of successive shocks, etc., the economist can assert with better conscience than the restrictions upon, say, the degree of polynomials chosen to describe the structural relations. If confidence limits for joint probability densities of the variables could be estimated on the basis of such weak restrictions, predictive estimation under properly defined structural changes might become possible without introducing stronger but less justifiable hypotheses.

3. THE PLAN OF THE VOLUME

3.1. This volume is concerned with empirical inference in economics. The bulk of its contents is determined by the *stochastic* character of economic models. Their character as *systems of simultaneous equations* requires certain modifications of the usual inference method. The usual method is appropriate when no changes in structure are expected or intended, so that no structural estimation is needed; or when structural estimation is based on experiments which isolate single equations that constitute the complete model by making all variables but one predetermined; or when experiment is used to reproduce the structural change considered (section 1.4 above). The modifications of the inference method that are necessary when none of these conditions is satisfied are treated in Part One: "Simultaneous Equation Systems." The *dynamic* character of economic models calls for modifications in the technique of inference, especially drastic in the treatment of small samples; these modifications are discussed in Part Two: "Problems Specific to Time Series." Part Three, on "Specification of Hypotheses," discusses the construction of models of the type analyzed in the earlier parts of this article, and other more general ones.

3.1.1. In particular, Part One deals with three subjects separately: identification, estimation, and computation. Koopmans, Rubin, and Leipnik treat these problems extensively in [II], and the discussion that follows is organized according to the three headings.

3.1.2. Koopmans and Rubin establish criteria for identification of linear dynamic systems. Their extensive treatment of this especially important case is followed by comments of three authors: Hurwicz, in [IV], provides the logical basis of the problem in its general formulation; Wald, in [III], contributes certain general criteria for identifiability; and Haavelmo, in [V], confronts Ragnar Frisch's technique of bunch maps (devised to establish the presence or absence of multicollinearity in the data of a model based on errors without specified distribution) with the theory of models of simultaneous equations involving shocks with specified distributions.

3.1.3. The logical foundations of structural estimation, compared with predictive estimation under changed or unchanged structure, is given by Hurwicz, in the introductory sections of [VI]. The more technical parts of this paper study the properties of least-squares estimates of regression coefficients when the complete model is multiequational. For one of Hurwicz's results - the equivalence of least-squares and maximum-likelihood estimates of regression coefficients - Koopmans gives an alternative proof in [VII].

The procedure of *complete* structural estimation by the maximum-likelihood method is developed in this volume more fully than the methods of *incomplete (partial)* structural estimation. To the former, the extensive article [II-3] by Koopmans, Rubin, and Leipnik is devoted. Methods of incomplete structural estimation, summarized above in section 2.5.6, are discussed only briefly: by Wald, [VIII], who has sketched a nonparametric approach, and by T. W. Anderson, [IX], who presents the results obtained in a forthcoming larger article [T. W. Anderson and Rubin, 1949] on the method of reduced forms.

3.1.4. A difficult task, for the first time explored in this volume, is to develop appropriate techniques for the computation of maximum-likelihood estimates that utilize all a priori information contained in the multiequational model. When the model consists of a single equation, the maximization of the likelihood function with respect to structural parameters yields a system of

"normal equations" (familiar from multiple-regression theory) that is linear in the maximum-likelihood estimates: the usual single-equation least-squares method of estimating parameters is a special case of the maximum-likelihood method. When, however, the model has two or more equations, the "normal equations" have to be replaced by more general equations, from which to compute maximum-likelihood estimates of the structural parameters. (If the structure is not uniquely identifiable, these equations will be dependent and the estimates indeterminate in some degree.) The maximum-likelihood equations are nonlinear in the estimates. The direct solution of these equations is therefore troublesome. Koopmans, Rubin, and Leipnik [I-4] have provided iterative methods to solve these equations with any desired degree of approximation. On this problem, much help was derived from the advice of von Neumann and from Hotelling [X]. More recently, the methods involved were discussed and explored from a more general point of view by Chernoff [1949].

Since the likelihood function will have, in general, several maxima, it would be important to know their number and to ensure, by a proper choice of initial values, that successive iterations approach that set of values of parameters which corresponds to the highest of the likelihood maxima. These problems are as yet unsolved.

3.2. The fact that economic models are, in general, dynamic and have the nature of "stochastic processes" of a time variable gives rise to special problems treated in Part Two. Here the existence of simultaneous equations is provisionally relegated to the background and the study is occasionally confined to a uni-equational model.

3.2.1. The first division of Part Two ("Trend and Seasonals") may be said to specify the general (possibly nonlinear) model by using certain information on the possible role of the exogenous variable, time. A specified economic model of this kind is studied in Hurwicz's article [XI] on variable parameters in stochastic processes; structural estimation is the author's objective. In the two other articles of the section, given by guests of the Cowles Commission, Mann and R. L. Anderson, less reliance is placed on knowledge from economic theory; correspondingly, no attempt is made to estimate economic structure. Rather than estimating structural parameters, these two articles discuss how to test the hypothesis that no trend, or that no seasonal fluctuation, is present. Mann [XII] has devised a new nonparametric test

against trend (applicable in situations like quality control of successive products of a machine); R. L. Anderson [XII] gives account of the analysis of variance as used in testing for the presence of seasonality in a unitemporal (lagless) system.

3.2.2. The second division of Part Two ("Estimation Problems") considers the consequences of relaxing certain assumptions that in Part One were shown to be sufficient to make maximum-likelihood estimates unbiased: in particular, the assumption that the samples are large, and the assumption that the model describes a stationary process.

Problems of short time series were discussed at the conference by Madow and Hurwicz. The former's paper, [Madow, 1945], contained an exact formula through which the distribution of the autocorrelation coefficient of a discrete series when its population value differs from zero is connected with the corresponding distribution when the population value equals zero; the latter was derived in [Koopmans, 1942], [Dixon], and elsewhere. This formula applies to a circular series, i.e., a series in which the first and last elements are identical: a condition often satisfied in space series (as in the sampling of families around a block). The removal of this last condition is thus the only remaining step for the full solution of the problem. Hurwicz [XV] shows that, in short one-variable series, the least-squares estimate of the autoregression coefficient has a very considerable bias. This is a case when the regression coefficient equals the structural coefficient - hence the presence of bias in estimated structural coefficients is also proved. Moreover, in one of the cases treated (fixed initial value of the variable) the least-squares estimate equals the maximum-likelihood estimate. Hence, maximum-likelihood as well as least-squares estimates of both structural and regression coefficients are generally biased. This is an important warning since the time series that we have to use to estimate dynamic economic structures are short.

The assumption of stationary processes is discussed by Rubin [XIV]. He shows that, in the case of one variable, this assumption is not necessary for the consistency of maximum-likelihood estimates. This is important since economic "explosions" (speculative panics, etc., often described by economists as "cumulative processes") are known to occur.

3.2.3. The third division of Part Two, unlike the preceding two divisions, derives its importance not so much from specifying or relaxing the model of Part One as from correcting it. Koopmans'

article [XVI], which constitutes this division, suggests treating variables as functions of a continuous time variable instead of breaking up the continuous economic process into arbitrary periods of finite length. In addition to being more realistic the introduction of a continuous time variable helps to solve the following important paradox. In Part One of this volume, the answer to the question of whether a variable occurred with or without a time lag could, under certain circumstances, decide whether or not a system was identifiable and hence susceptible to estimation of all its parameters. The presence or absence of a time lag could also, under certain other circumstances, decide whether or not the estimation of some parameters by the single-equation least-squares method was subject to bias. Thus a lag of one day would appear to make an enormous difference for the statistician. Koopmans shows that this is not the case: as the time unit is being diminished, the assumption of independent successive disturbances becomes less and less permissible; and it is on the latter assumption that the above statements of Part One were based. A similar result is obtained for the discrete case by Hurwicz [XI]. If proper mathematical treatment of stochastic models can be developed, such models promise to be a more accurate and more flexible tool for inference in economics than the discrete models used heretofore.

3.3. Quite naturally, certain "maintained hypotheses" (as distinct from the "considered hypotheses" that are being tested on the basis of data) had to be specified throughout the volume. For example, linear relations were assumed, the effects of the presence of predetermined variables discussed, continuous and discrete models compared. Part Three is called "Specification of Hypotheses," with a view to more particular needs of economics.

It is introduced by Koopmans' article [XVII], in which he gives formal criteria for the successive delimitation of the exogenous and the predetermined variables. These criteria are not based on the distinction between "economic" and "noneconomic" variables; they refer, instead, to the algebraic role played by a given variable in the model. For example: unless significant lags exist between the economic facts and the political facts that are affected by them and that, in turn, affect economic facts, economic causation of politics must be treated as one of the relationships of the complete system. This shows the importance of the questions raised in Parts One and Two: how to deal with incomplete systems, and how to determine whether or not lags are significant.

The "nonadditive random disturbances" treated by Hurwicz in [XVIII] are studied to take into account the economists' view that, for example, the erratic variations in taste - from year to year or from person to person - do not necessarily affect the demand equation's constant term only; e.g., a change in slope of a linear demand curve may accompany its shift. This complicates the mathematics of statistical inference quite considerably. A beginning has been made in Rubin's brief article [XIX] that concludes the volume.

3.4. A number of major problems remain unsolved. It suffices to list some of those already mentioned in various places of this introduction: observation errors (and the problem of nonobservable variables in general); multiple decisions; problems of aggregation; continuous stochastic processes. Even where the statistical investigation has proceeded beyond the definition of the problem, the results obtained so far are subject to severe limitations. Structural estimation from small samples is in its beginnings; nonlinear relations, often called for in economic theory, have hardly been approached, etc. All these problems will require much further work on the part of statisticians.

PART ONE

SIMULTANEOUS EQUATION SYSTEMS

II. MEASURING THE EQUATION SYSTEMS OF DYNAMIC ECONOMICS

BY T. C. KOOPMANS, H. RUBIN, AND R. B. LEIPNIK¹

	Page
1. Description of the Systems Considered	
1.1. The Economic and Statistical Basis of a System of Equations	54
1.2. Specifying the Joint Distribution of all Variables . .	55
1.3. Exogenous and Endogenous Variables	56
1.4. The Disturbances in the Equations	56
1.5. Economic Interpretation of the Disturbances in the Equations	57
1.6. Errors of Measurement or Disturbances in the Variables	57
1.7. Nonsingularity of Σ	58
1.8. Jointly Dependent Variables and Predetermined Variables	59
1.9. Nonsingularity of B_0	60
1.10. Timing of the Variables	62
1.11. The Problem of Identification	62
1.12. The A Priori Restrictions	64
1.13. A Priori Restrictions on the Distribution of Disturbances	66
1.14. Inequalities as A Priori Restrictions	67
1.15. Rules of Normalization	68
1.16. Summary of Subsequent Sections	68
2. The Identification of Economic Relations	
2.1. The Concept of Identification	69
2.2. Identification of One Structural Equation under Linear Restrictions	78
2.3. Treatment of Unidentifiable Structural Equations by Linear Dummy Restrictions	85
2.4. Identification of a Set of Structural Equations under Linear and Bilinear Restrictions	94
2.5. Incompleteness of the Present Discussion of Identification Problems	106

¹Rubin contributed to sections 1 - 3 and to an early draft of the methods developed in section 4. Leipnik contributed to section 4 and directed most of the computations reported there. Further computations were directed by B. A. de Vries.

	Page
3. Estimation of the Parameters	
3.1. Properties of the Unrestricted Likelihood Function . .	110
3.2. Properties of the Restricted Likelihood Function . . .	120
3.3. Large-Sample Properties of the Maximum-Likelihood Estimates	133
4. Computation of the Maximum-Likelihood Estimates	
4.1. Introductory Remarks	153
4.2. A Complete Set of Unrestricted Parameters	160
4.3. The Case of Uncorrelated Disturbances	166
4.3.1. The nature of the problem	166
4.3.2. The methods $\hat{\rho}_1$, $\hat{\rho}_h$, and $\hat{\rho}_h^n$	168
4.3.3. Asymptotic convergence properties of $\hat{\rho}_1$, $\hat{\rho}_h$, and $\hat{\rho}_h^n$	172
4.3.4. Arrangement of computations for $\hat{\rho}_1$, $\hat{\rho}_h$, and $\hat{\rho}_h^n$	190
4.3.5. The Newton method	203
4.4. The Case of Unrestricted Correlations between the Disturbances	211
4.5. Concluding Remarks	230

1. DESCRIPTION OF THE SYSTEMS CONSIDERED

1.1. *The economic and statistical basis of a system of equations.* The analysis and explanation of economic fluctuations has been greatly advanced by the study of systems of equations connecting economic variables. The construction of such a system is a task in which economic theory and statistical method combine. Broadly speaking, considerations both of economic theory and of statistical availability determine the choice of the variables. Economic theory predominates in the definition of the "behavior equations" describing a certain type of economic decisions taken by a certain category of economic agents, and in the specification of the variables that may possibly enter each behavior equation (i.e., of the conditions that may affect that decision by that group of agents). "Institutional equations" describe behavior patterns set by law or rule. Technical knowledge enters into the definition, and selection of variables, of the "technical equations" expressing the physical relation between input and output

in production. A fourth group of equations, usually referred to as "identities" (like the savings-investment identity), which occupy a place in economic literature out of proportion to their theoretical triviality, should be classified as deriving directly from the definitions of the variables through the principles of economic accounting. Theoretical preconceptions, statistical evidence, and sometimes mere assumption or approximation, are intermingled in the determination of the form of each equation, as regards linearity and as regards the occurrence and length of time lags. All these things being determined, it is almost entirely left to statistical methods to estimate the numerical values of the coefficients in the equations, and to assess the possible degree of error in those estimates, subject to the assumptions made.

Several equation systems of this kind have been constructed by Tinbergen [1939] and others for different countries and periods. We shall in this article assume a general knowledge of the nature of those systems, and of the uses to which they are put.

Tinbergen gives ample consideration to the economic assumptions on which these systems are based. Only recently has attention been directed systematically to the specific problems of statistical method involved in estimating the coefficients of any equation that forms part of such a system of equations. Haavelmo [1943, 1944] has pointed out that the methods developed for the measurement of a single relationship under conditions of experimental control over - or at least independent determination of - all variables except the one "dependent" variable, are inadequate if we are faced with a system of simultaneous equations between the variables. He has indicated the general principles of a statistical method appropriate to the latter situation. Mann and Wald [1943] have applied these principles to give a statistical treatment of large samples of a number of variables which satisfy an equal number of linear difference equations.

1.2. Specifying the joint distribution of all variables. The main principle advanced by Haavelmo is that the measurement of a system of equations should be based on a specification of the joint probability distribution of all values of all variables involved. This principle has been generally accepted in other applications of statistical method. Probably, economic statisticians have largely been unaware of the fact that their methods did not satisfy this requirement. Actually, the probability distributions that were employed always referred to one equation taken in isolation, and distributions specified with regard to different equations were usual-

ly incompatible.

1.3. *Exogenous and endogenous variables.* This article is concerned with linear systems of difference equations of the following general form:

$$(1.1) \quad \sum_{i=1}^G \sum_{\tau=0}^{\tau^{\square}} \beta_{gi\tau} y_i(t-\tau) + \sum_{k=1}^K \sum_{\tau=0}^{\tau^{\square}} \gamma_{gk\tau} z_k(t-\tau) = u_g(t),$$

$$g = 1, 2, \dots, G; \quad t = 1, 2, \dots, T.$$

This form is slightly more general than that studied by Mann and Wald in that we consider G equations containing both G endogenous variables $y_i(t)$ and K exogenous variables $z_k(t)$. The latter are defined as variables that influence the endogenous variables but are not themselves influenced by the endogenous variables. It will be clear that at this stage the distinction between exogenous and endogenous variables is a theoretical, a priori distinction on which statistical evidence may or may not be obtained at a later stage. Because of the general interdependence of economic variables, exogenous variables are most likely to be found among noneconomic phenomena like temperature, rainfall, earthquake intensities, etc. Both endogenous and exogenous variables are assumed to be observable.

In the equations (1.1) the exogenous variables are treated as if they are given functions of time, the values of which remain the same in repeated samples. Another contribution to this volume, [XVII], is devoted to the justification of this procedure, which we shall here assume to be correct.

1.4. *The disturbances in the equations.* The distribution of the endogenous variables is then defined by means of the not directly observable disturbances $u_g(t)$. The latter are terms in the equations specified only to the extent that they are assumed to be subject to a joint probability distribution. Because of the presence of these terms, the system (1.1) is called a system of *stochastic equations*.

It will be assumed here that the $u_g(t)$ have a joint probability distribution of the form

$$(1.2) \quad \prod_{t=1}^T f(u_1(t), \dots, u_G(t)) du_1(t) \dots du_G(t).$$

The assumption (1.2) implies independence of disturbances in successive time intervals. It also implies that the disturbances are independent of the values of the exogenous variables. The conditions to be imposed on the distribution function f are not the same in the various sections of this article. In all cases, we shall assume that first-order moments exist and are equal to zero,

$$(1.3) \quad \mathcal{E} u_g(t) = 0,$$

and that second-order moments (variances and covariances),

$$(1.4) \quad \mathcal{E} u_g(t) u_h(t) = \sigma_{gh},$$

also exist. In certain sections of this article we shall go further and assume that $f(u_1, \dots, u_G)$ represents the joint normal distribution function of G variables.

The assumed distribution of the $u_g(t)$ defines the joint distribution of all values $y_i(t)$ of the endogenous variables for which $t = 1, \dots, T$, provided we specify in addition that any values $y_i(t)$ for which $t \leq 0$ and which occur in (1.1) are regarded as given numbers that remain the same in repeated samples.

1.5. Economic interpretation of the disturbances in the equations. In each behavior equation, the disturbance is interpreted as representing the joint effect, on the behavior described by that equation, of all variables of minor individual importance that have not been explicitly introduced into the system of equations. For instance, random variation in consumers' tastes will lead to a certain amount of shifting in the curve of consumers' demand. Similarly, in the technical relations between input and output, a certain amount of random shifting in the relationships is due to a large number of minor causes of variation not explicitly studied. The term "random" is used here in the sense of the assumptions (1.2), (1.3), and (1.4), made regarding the disturbances in the equations.

1.6. Errors of measurement or disturbances in the variables. It is important to note that in the interpretation of disturbances just given, each disturbance is associated with an equation of the system, and not with a variable. This excludes the interpretation of the "disturbances in the equations" as errors of measurement. If errors of measurement occur to a marked degree, separate provision must be made for them in the probability dis-

tribution of observed variables by introducing additional "disturbances in the variables." In order to concentrate on the effect of disturbances in the equations, we shall assume in this study that all variables are measured without error. Systems in which "disturbances in the equations" (also called "shocks") and "disturbances in the variables" (also called "errors") occur side by side have been studied in [T. W. Anderson and Hurwicz].

1.7. *Nonsingularity of Σ .* For some purposes the mathematical treatment of systems like (1.1) is simplified if we can restrict ourselves to cases in which there is no *functional* relation (as distinct from *stochastic* dependence) between the G disturbances $u_g(t)$ for any t . This requires in particular that the matrix $\Sigma = [\sigma_{gh}]$ defined by (1.4) be nonsingular.

Now each "identity" that is present among the equations (1.1) makes all elements in the corresponding row and column of the matrix Σ vanish, because by their nature identities are not subject to disturbances. However, the variables entering into a given identity, and the coefficients with which they enter, are always known a priori (often the coefficients are +1 or -1). It is therefore possible, whenever the assumption of nonsingularity of Σ is desirable for mathematical reasons, to remove the identities from the system by elimination of as many variables as there are identities to be removed. For instance, the identity "volume of production equals real income" can be removed by replacing the variable "volume of production," wherever it occurs, by the variable "real income." A less trivial example: if the profit margin is conceived to be a determining factor of investment activity, the identity defining the profit margin can be removed by replacing the variable "profit margin" in the equation explaining investment activity by the linear combination "product price less the sum of factor prices per unit of product." The latter example shows that the elimination of variables defined by identities may introduce a priori proportionalities or other restrictions among the coefficients occurring in the remaining equations of the system. We shall revert to these a priori restrictions below.

In the case in which the identities have thus been disposed of, it is reasonable to assume that no functional relation exists between the disturbances in different behavior equations and technical equations. This can be seen if we ask ourselves what, for instance, would be implied in an exact proportionality of the disturbances in two given equations. This would mean, not only that

precisely the same minor causes would be operative in the random shifts of each of these relationships, but also that the relative strengths with which these causes operate in each equation are the same - an obvious impossibility. A similar, slightly more complicated argument applies to disqualify the assumption of a functional relation involving disturbances in more than two equations.

1.8. *Jointly dependent variables and predetermined variables.* Besides the distinction between endogenous and exogenous variables, it is desirable to introduce another classification of variables, which is based partly on the former distinction, and partly on the timing of each variable. That is, for the purposes of the classification now to be introduced, it will be necessary to regard for instance $y_i(t)$ and $y_i(t-1)$, and generally all variables measured with different time lags, as different variables.

The equations (1.1) for a given value of t are intended to describe the process of the formation of the endogenous variables $y_i(t)$ at time t , under the influence of earlier values $y_i(t-\tau)$, $\tau \geq 1$, of the endogenous variables, of the exogenous variables $z_k(t-\tau)$, $\tau \geq 0$, and of the disturbances $u_g(t)$. Generalizing a terminology of single-equation least-squares regression theory, the values $y_i(t)$, without time lag, of the endogenous variables may be called *jointly dependent variables at time t* . To bring out more clearly that the equations (1.1) describe the formation of the jointly dependent variables, these equations may be written in the form

$$(1.5) \quad \sum_{i=1}^G \beta_{gi0} y_i(t) = \\ - \sum_{i=1}^G \sum_{\tau=1}^{\tau^{\square}} \beta_{gi\tau} y_i(t-\tau) - \sum_{k=1}^K \sum_{\tau=0}^{\tau^{\square}} \gamma_{gk\tau} z_k(t-\tau) + u_g(t).$$

The right-hand member contains, besides the disturbances, a group of variables that we shall call *predetermined variables at time t* . The lagged values $y_i(t-\tau)$, $\tau \geq 1$, of the endogenous variables are predetermined in a temporal sense, in that their values $y_i(t-\tau)$ for a given value of t are determined by variables and disturbances relating to time intervals preceding t . In particular, they are unaffected by the disturbances $u_g(t)$ of the time interval t . The

exogenous variables $z_k(t)$ without time lag are predetermined in the logical sense that they are influenced only by causes outside the economic system studied, and are independent of all other variables and disturbances (measured at time t or earlier) included in the system of equations. The lagged exogenous variables $z_k(t-\tau)$, $\tau \geq 1$, are predetermined in both senses.

1.9. *Nonsingularity of B_0 .* Although the need for the foregoing classification will become clear only when we come to estimation problems, it is introduced here because of a related basic assumption that is of general importance. If the right-hand member in (1.5) represents a set of causal factors in the determination of the dependent variables $y_i(t)$ in the left-hand members, without any causal action in the opposite direction, then it is necessary to specify that the matrix

$$(1.6) \quad B_0 = \begin{bmatrix} \beta_{110} & \cdots & \beta_{1G0} \\ \cdot & \cdots & \cdot \\ \beta_{G10} & \cdots & \beta_{GG0} \end{bmatrix}$$

be nonsingular. For if B_0 were singular, there would be a set of numbers λ_g , $g = 1, \dots, G$, not all equal to zero such that

$$(1.7) \quad \sum_{g=1}^G \lambda_g \beta_{gi0} = 0.$$

Writing $w_g(t)$ for the right-hand member in (1.5), the validity of (1.7) would then entail a linear restriction,

$$(1.8) \quad \sum_{g=1}^G \lambda_g w_g(t) = 0,$$

on the expressions represented by $w_g(t)$. Such a restriction, however, is contrary to the assumed direction of causation, and is in particular incompatible with the fact that there is no linear functional dependence between the stochastic variables $u_g(t)$.

The foregoing argument needs amendment in case (1.5) contains identities, because in that case the corresponding quantities $u_g(t)$ vanish. Suppose that the equations (1.5) for $g = 1, \dots, G_0$ are identities, and that for $g = G_0 + 1, \dots, G$ the equations involve disturbances of positive variance. Then a restriction (1.8) is compatible with the assumed direction of causation if and only if

$$(1.9) \quad \lambda_{G_0+1} = \dots = \lambda_G = 0.$$

The only form of nonsingularity permissible in B_0 in the present case is therefore, according to (1.7), linear dependence among the first G_0 rows, each of which contains the coefficients of an identity. Such linear dependence is, of course, precluded by the simple reason that any linear dependence should be removed from the (fully known) identities before they are admitted to the system of equations.

There is, of course, no a priori reason why a number of behavior equations could not happen to be such that B_0 is singular. But there is good empirical evidence that this case can be ruled out, at least in dynamic equation systems. For if B_0 were singular (or even if its determinant value $\det B_0$ were very small compared with its term largest in absolute value), small disturbances in any direction in the space of the disturbance vector $u = (u_1, \dots, u_G)$ incompatible with the linear restriction (1.7), would lead to infinite (or very large) *simultaneous* changes in the variables $y_i(t)$. Such phenomena have not been observed. Although small causes occasionally have great effects in economic developments, in such cases time is required for the effects to materialize.

One could not have equal confidence in a statement that the matrix \bar{B} with elements $\bar{\beta}_{gi} = \sum_{\tau=1}^{\tau_0} \beta_{gi\tau}$, describing a corresponding static system

obtained through the neglect of all time lags, is far from being singular. It is easily seen, however, that if \bar{B} is singular for such a static system, then the deterministic dynamic system obtained from (1.5) by omitting the disturbances and giving arbitrary constant values to the exogenous variables does not in general have a solution asymptotically approaching a set of finite equilibrium values.

1.10. *Timing of the variables.* The nonsingularity of B_0 in particular excludes the possibility that one of the variables $y_i(t)$ would not occur at all in the equations (1.1) except with a positive time lag. A variable of that kind properly belongs among the exogenous variables because a past quantity cannot be influenced by present developments. An objection to this reasoning might be that the timing with respect to which the variables are defined may be varied at will by a transformation,

$$(1.10) \quad \begin{aligned} y_i^\oplus(t) &= y_i(t + \theta_i), & i &= 1, \dots, G, \\ z_k^\oplus(t) &= z_k(t + \theta'_k), & k &= 1, \dots, K. \end{aligned}$$

However, such transformations cannot be regarded as permissible for our purposes - except in the trivial case where all quantities θ_i and θ'_k are equal. The time variable is more than an index used to distinguish successive values of one and the same variable. It indicates historical time - the medium in which causation and interaction between economic and other variables takes place. Therefore the matrix B_0 must be such that the endogenous variables, which by their definition are in continuous and instantaneous interaction with each other, are all represented in the system (1.1) by simultaneous values with zero time lag ($\tau=0$). Further light is thrown on this important point in another contribution already referred to [XVII].

1.11. *The problem of identification.* The statistical measurement of a system of equations like (1.1) involves two logically distinct and successive problems, which have here been called the problem of the *identification* of each equation and the problem of the *estimation of the parameters* of each equation. Section 2 is devoted to the former of these problems, which arises especially with regard to data governed by more than one equation at the same time. It originates from the fact that, if a system like (1.1) is viewed *only* as a mathematical specification of the joint probability distribution of the observable variables, it can be written in many different ways. Any linearly independent system of G linear combinations of the equations (1.1) with a correspondingly transformed distribution of the disturbance terms will be a mathematically equivalent way of defining the probability distribution of the variables.

Let a "way of writing" the system be called a *representation*

of the distribution of the variables. Two representations are called *observationally equivalent* if they define the same probability distribution of the variables. Haavelmo [1944, p. 91] uses the expression "indistinguishable on the basis of the observations" to describe two equivalent representations, because even if the probability distribution of the observations were fully known - the best that can be expected from statistical methods - there would still be no way to distinguish observationally equivalent representations. The distribution of the variables can be looked upon as determining the set of all observationally equivalent representations of it, and is completely defined by any of these representations. Mathematically speaking, it is immaterial which representation is employed, except that it will be desirable to choose a simple one. Economically, however, different representations of the same system are not at all equivalent.

The study of a system of equations like (1.1) derives its sense from the postulate - already implicit in earlier parts of this section - that there exists one and only one representation in which each equation corresponds to a specified law of behavior (attributed to a specified group of economic agents), to a specified technical law of production, or to a specified identity. Let us call these particular equations the *structural equations*, because they are the elements of which the dynamic economic structure of society is composed. The representation composed of the structural equations may be called the *representation according to economic structure*, or briefly the *structural representation*. Any discussion of the effects of changes in economic structure, whether brought about by gradual trends or by purposive policies, is best put in terms of changes in the structural equations. For those are the elements that can, at least in theory, be changed one by one, independently. For this reason, it is important to have the system (1.1) in a form in which the greatest possible number of its equations can be identified and recognized as structural equations.

Suppose for a moment that the structural representation be known to investigator A, who as a mathematical exercise derives another representation from it by taking linear combinations. In that process, the economic identity of the structural equations is lost, and when A hands the derived representation over to B without disclosing its source or method of computation, B is faced with the problem of identifying among all linear combinations of the equations of the representations given to him, the structural equations that alone reflect specified laws of economic behavior, of the

technique of production, or of economic accounting.

1.12. *The a priori restrictions.* The position of our investigator B corresponds exactly to the position of the econometrician who sets out to measure a system of economic relations. Statistical observation will in favorable circumstances permit him to estimate, with a precision again subject to estimation, the characteristics of the probability distribution of the variables. Under no circumstances whatever will passive statistical observation permit him to distinguish between different mathematically equivalent ways of writing down that distribution. Because he has no experimental control over economic variables, the simultaneous validity of all the structural equations prevents him from isolating and individually observing any one of them on a statistical basis alone. The only way in which he can hope to identify and measure individual structural equations implied in that system is with the help of a priori specifications of the form of each structural equation.

The most important instrument of identification is a specification as to *which variables may enter into which structural equations with which possible time lags*. Assuming now that the system (1.1) is the structural representation, this can be expressed mathematically by putting equal to zero all coefficients of terms that do not enter into the respective equations,

$$(1.11) \quad \begin{aligned} \beta_{g_r i_r \tau_r} &= 0, & r &= 1, 2, \dots, R_\beta^{(1)}, \\ \gamma_{g_r k_r \tau_r} &= 0, & r &= R_\beta^{(1)} + 1, \dots, R_\beta^{(1)} + R_\gamma^{(1)}, \quad R_\beta^{(1)} + R_\gamma^{(1)} \equiv R_\alpha^{(1)}. \end{aligned}$$

Sometimes it is useful to state these restrictions on the coefficients in a slightly more general form which includes (1.11) as a special case,

$$(1.12) \quad \sum_{i=1}^G \sum_{\tau=0}^{\tau^\square} \chi_{ri\tau}^{(g_r)} \beta_{g_r i\tau} + \sum_{k=1}^K \sum_{\tau=0}^{\tau^\square} \psi_{rk\tau}^{(g_r)} \gamma_{g_r k\tau} = 0,$$

$$r = 1, 2, \dots, R_\alpha^{(1)}.$$

The special case (1.11) will be distinguished from (1.12) as the

case of single-parameter restrictions. The form (1.12) of restrictions, in which the quantities $\chi_{ri\tau}^{(g_r)}$ and $\psi_{rk\tau}^{(g_r)}$ are a priori known constants, permits inclusion of cases where the ratio of two coefficients in the same equation, or another linear relation between the coefficients of an equation, is given a priori. Examples of this type of restriction have already been given.

It will be noted that each condition (1.12) connects only coefficients that occur in the same structural equation. There is a further type of restrictions involving coefficients occurring in different equations. This can again be illustrated with the example of the profit margin. Assume that the profit margin enters as such into at least two behavior equations, whereas none of its constituents enters explicitly. Suppose further that the definition of the profit margin, with the help of which we wish to eliminate that variable, contains an unknown parameter. (This may happen, for instance, if the conversion factor by which the price of any given factor of production is related to the unit of product is not known.) The type of restriction arising from such a situation is one in which two coefficients, of the variables y_{i_r} and $y_{i'_r}$, respectively, are required to have the same ratio in two different structural equations (numbered g_r and g'_r):

$$(1.13) \quad \begin{bmatrix} \beta_{g_r i_r \tau_r} & \beta_{g_r i'_r \tau_r} \\ \beta_{g'_r i_r \tau_r} & \beta_{g'_r i'_r \tau_r} \end{bmatrix} = 0.$$

Similar restrictions may arise from the approximation of a distributed time lag by a linear combination of terms with discrete lags. If a variable y_{i_r} is supposed to occur with the same lag distribution in two equations numbered g_r and g'_r , this leads to a restriction of the type

$$(1.14) \quad \begin{bmatrix} \beta_{g_r i_r \tau_r} & \beta_{g_r i_r \tau'_r} \\ \beta_{g'_r i_r \tau_r} & \beta_{g'_r i_r \tau'_r} \end{bmatrix} = 0.$$

While the restrictions (1.12) are linear in the unknown coefficients

$\beta_{gi\tau}$, $\gamma_{gk\tau}$, restrictions like (1.13) and (1.14) are bilinear, and lead to greater mathematical complications in what follows. We shall assume that there are $R_{\alpha}^{(2)}$ bilinear restrictions of the types (1.13) and (1.14), possibly involving coefficients $\beta_{gi\tau}$, $\gamma_{gk\tau}$ of both endogenous and exogenous variables.

The restrictions (1.11) or (1.12), (1.13), and (1.14) - and such similar restrictions as we may wish to add later - will be called the *a priori restrictions*. In section 2 we investigate necessary and sufficient conditions under which the *a priori* restrictions suffice to identify a given equation (1.1) as a specified structural equation. On this basis we shall distinguish, and include in subsequent sections, the case, not covered by Mann and Wald [1943], in which one or more but not all of the structural equations can be identified within the system.

It will be seen that, even in the case where the *a priori* restrictions are insufficient in number and variety to permit identification of all structural equations, there may be among the *a priori* restrictions one or more that can be omitted without thereby removing further equations from the list of identifiable ones. "*A priori*" restrictions of this kind are in principle subject to statistical testing (on the basis of the remaining *a priori* restrictions). For this reason, statistical evidence was quoted, in the opening paragraph of this article, as one of the bases for a determination of the form of the structural equations. If restrictions supported to a degree by statistical evidence are nevertheless imposed *a priori*, this will in general reduce the sampling variances of the estimates of some or all parameters subject to estimation. The use of *a priori* restrictions should therefore be resorted to whenever the theoretical grounds are strong enough. To make possible a formal mathematical treatment according to established procedures of statistical inference, we shall in this article regard the "*a priori* restrictions" strictly as given *a priori* and imposed without reference even to the possibility of statistical test. It goes without saying that we should eliminate from consideration sets of *a priori* restrictions that are mutually incompatible or mutually dependent.

1.13. *A priori* restrictions on the distribution of disturbances. It may well happen that one or more specified structural equations cannot be identified on the basis of such *a priori* restrictions of the forms (1.12), (1.13), and (1.14) as are considered theoretically justified. We shall therefore study further a

priori restrictions on the matrix Σ of the variances and covariances of the disturbances, which require that the covariance between the disturbances in two specified equations shall vanish. The elements of Σ that are thereby required to vanish may or may not follow a regular pattern. One particular pattern is of interest both because of its special mathematical consequences and because the assumptions involved may present a fair approximation to reality. In this pattern (which is here formulated for systems from which all identities have been removed) it is supposed that the G equations can be classified into N groups of G_1, G_2, \dots, G_N equations respectively, with $G_1 + G_2 + \dots + G_N = G$, such that the disturbances $u_g(t)$ of equations in different groups are independent. In that case, the matrix Σ can be partitioned into the form

$$(1.15) \quad \Sigma = \begin{bmatrix} \Sigma_1 & 0 & \dots & 0 \\ 0 & \Sigma_2 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \Sigma_N \end{bmatrix},$$

where each of the matrices $\Sigma_1, \dots, \Sigma_N$ is positive definite. Each element σ_{gh} , $g < h$, of Σ that is thereby prescribed to be zero gives rise to an a priori restriction, which turns out (see section 2) to be equivalent to a bilinear restriction in the elements of the corresponding rows of the coefficient matrix $[B \ \Gamma]$. For this reason we shall denote by R_σ the total number of such restrictions, and write $R_\alpha^{(1)} + R_\alpha^{(2)} + R_\sigma = R$. The particular choice of vanishing elements indicated by the partitioning in (1.15) is simpler than other choices because it is invariant under inversion of Σ .

1.14. Inequalities as a priori restrictions. A further class of a priori restrictions that can often be based on economic considerations is inequalities. Frequently, the sign of coefficients $\beta_{gi\tau}$ or $\gamma_{gk\tau}$ is known beforehand. Sometimes it may be possible to prescribe the sign of, or set another limit to, the correlation of the disturbances in the structural equations. In the present article we do not study the question of how to give effect to restrictions of this kind.

1.15. *Rules of normalization.* The equations (1.1) as well as the restrictions (1.12), (1.13), and (1.14) are homogeneous in the coefficients $\beta_{gi\tau}$ and $\gamma_{gk\tau}$ and the parameters σ_{gh} for each value of g , i.e., they are unaffected by a change in scale

$$(1.16) \quad \beta_{gi\tau}^{\oplus} = v_g \beta_{gi\tau}, \quad \gamma_{gk\tau}^{\oplus} = v_g \gamma_{gk\tau}, \quad \sigma_{gh}^{\oplus} = v_g \sigma_{gh} v_h,$$

of each equation (1.1). It will be useful sometimes to fix the scale factors v_g by imposing a *normalization rule* on each equation. The precise form of the normalization rule is obviously a matter of choice, and different normalization rules are most convenient in different problems. We shall consider the following two of many possible alternative sets of G normalizing restrictions:

$$(1.17) \quad \begin{cases} (1.17a) & \beta_{gi_g 0} = 1, & g = 1, 2, \dots, G, \\ (1.17b) & \sigma_{gg} = 1, & g = 1, 2, \dots, G. \end{cases}$$

In the case of the first rule (1.17a) there should of course be no conflict with (1.11) or (1.12). The second rule (1.17b) still leaves open the choice of the sign of one of the nonvanishing coefficients "β" or "γ" in each equation.

Because of the trivial nature of the question of normalization, we shall sometimes omit specification of a normalization rule. It is therefore useful to introduce the convention that the g th equation can be called completely identified by the a priori restrictions even if its scale has not been fixed, provided the ratios between all its coefficients and the quantities $\sigma_{gg}^k, \sigma_{gh}, h \neq g$, are determinate. In case normalization rules are specified, they will be comprised in the term "a priori restrictions."

1.16. *Summary of subsequent sections.* In section 2, we discuss conditions for the identifiability of a given structural equation under a priori restrictions of the type (1.12), (1.13), (1.14), or (1.15). Necessary and sufficient conditions for identifiability under the restrictions (1.12) are derived (section 2.2). In section 2.4, the problem of extending these conditions to cases where restrictions of the types (1.13), (1.14), (1.15) are added is discussed but not solved. Means are indicated in section 2.3 to make possible the estimation of certain identifiable structural equations, even if certain other equations remain

unidentifiable. General observations indicating the incomplete state of our knowledge with regard to identification problems conclude this section.

Sections 3.1 and 3.2 deal with those properties of the likelihood function, before and after imposition of a priori restrictions, which are relevant to the maximum-likelihood method of estimation. It is found that whenever restrictions are imposed on structural equations that are not indispensable for the identification of those equations, the likelihood function is prevented from reaching its unrestricted absolute maximum, except in a set of samples of probability zero. Under such restrictions, maximum-likelihood estimates using all a priori information can only be obtained by computational procedures essentially more complicated than the least-squares method applied to the "reduced form" without regard to restrictions. Section 3.3 discusses and generalizes results regarding the limiting distribution of the maximum-likelihood estimates reached by earlier writers.

In section 4, iterative computation methods for the maximum-likelihood estimates are developed and discussed for the two cases in which the covariance matrix Σ of the disturbances is diagonal (section 4.3), and unrestricted (section 4.4). Unsolved problems connected with these methods are indicated.

Sections preceded by the symbol * can be passed over in a first reading without seriously affecting the understanding of the remaining parts of the article.

2. THE IDENTIFICATION OF ECONOMIC RELATIONS

2.1. *The Concept of Identification*

2.1.1. *Earlier discussions of the identification problem.* The first systematic discussion of the problem of identification was given by Frisch in an unpublished memorandum [1938]. Frisch's terminology is rather different from that employed here, and the concepts are slightly different in that the disturbances and their distribution are not explicitly introduced in his formulae. Nevertheless, the underlying ideas are to a large extent the same, and the present authors desire to acknowledge their indebtedness, and to emphasize the support found in Frisch's memorandum for the discussion of the problem of identification in this article.

Frisch indicates that what is here called the identification

problem arises from the passive nature of economic observations. There is no possibility of independently varying the several factors entering a given behavior equation. The only observations available are those which by assumption satisfy all structural equations simultaneously.

The same point is emphasized by Haavelmo, who has continued and extended Frisch's work in a very general discussion [Haavelmo, 1944, pp. 91-98] of one central problem in identification: the formulation of conditions under which *all* structural relations of the system can be identified. Haavelmo also does not use the term identification, but describes the above mentioned problem as the "problem of confluent relations" or, alternatively, as the "problem of arbitrary parameters," and classifies it under the heading "estimation." As regards this classification, it appears to the present authors that the identification problem is concerned with the unambiguous definition of the parameters that are to be estimated - a logical problem that precedes estimation. It is therefore not a problem in statistical inference, but a prior problem arising in the specification and interpretation of the probability distribution of the variables. As such it deserves separate classification.

Haavelmo's discussion of the "problem of confluent relations" is more general than the present discussion of identification problems in that he does not in any way restrict the functional form of the equations concerned. The conditions to be given below for the identifiability of *all* structural equations in a linear system could therefore be obtained as a specialization of Haavelmo's results, although we shall derive them directly. The present discussion, while restricted to linear systems, goes further in that we also discuss conditions under which any one particular structural equation can be identified.

2.1.2. *Notation.* It will be convenient to use a matrix notation for the equation system (1.5) in which the distinction between jointly dependent and predetermined variables introduced in section 1.8 is given explicit expression, whereas that between endogenous and exogenous variables is concealed. The variables and the disturbances will be represented by row vectors, and the coefficients will be regarded as the elements of matrices, as follows:

$$(2.1) \quad y(t) \equiv [y_1(t) \quad \cdots \quad y_G(t)],$$

...

$$\begin{aligned}
 z(t) &\equiv [y_1(t-1) \quad \cdots \quad y_G(t-\tau^{\square}) \quad z_1(t) \quad \cdots \quad z_K(t-\tau^{\square})], \\
 u(t) &\equiv [u_1(t) \quad \cdots \quad u_G(t)], \\
 (2.1) \quad B &\equiv \begin{bmatrix} \beta_{110} & \cdots & \beta_{1G0} \\ \cdot & \cdots & \cdot \\ \beta_{G10} & \cdots & \beta_{GG0} \end{bmatrix}, \\
 \Gamma &\equiv \begin{bmatrix} \beta_{111} & \cdots & \beta_{1G\tau^{\square}} & \gamma_{110} & \cdots & \gamma_{1K\tau^{\square}} \\ \cdot & \cdots & \cdot & \cdot & \cdots & \cdot \\ \beta_{G11} & \cdots & \beta_{GG\tau^{\square}} & \gamma_{G10} & \cdots & \gamma_{GK\tau^{\square}} \end{bmatrix}.
 \end{aligned}$$

Here the vector $y(t)$ of $G \equiv K_y$ elements comprises the variables jointly dependent at time t , and the vector $z(t)$ of K_z elements, say, comprises all variables predetermined at time t . The matrix B has previously been denoted by B_0 . In this notation, the equations (1.1) can be written as follows:

$$(2.2) \quad By'(t) + \Gamma z'(t) = u'(t),$$

where y' denotes the column vector which is the transpose¹ of the row vector y . The probability density function $f\{u_1(t), \dots, u_G(t)\}$ of the disturbances will be denoted by $f\{u(t)\}$. Occasionally, the argument t of y , z , u will be omitted.

For some purposes even the distinction between dependent and predetermined variables is irrelevant. Then we shall denote the equation system by

¹For reasons to be stated in the footnote on p. 81, we have reversed the more usual notation in which the transposition sign denotes a row vector, its absence a column vector.

$$(2.3) \quad \sum_{k=1}^{K_x} \alpha_{gk} x_k = u_g, \quad g = 1, \dots, G, \quad K_x \equiv K_y + K_z \equiv G + K_z,$$

or $Ax' = u'$, where $A = [B \quad \Gamma]$, $x = [y \quad z]$.

2.1.3. *Observationally equivalent structures.* The concept of identifiability is to be defined with reference to the joint distribution function F_T of all observations, as determined by (1.2), (1.1), and the requirement that all values $z_k(t)$ and those values of $y_i(t)$ for which $t \leq 0$ are given constants. In order to write F_T explicitly as a distribution function in terms of the variables $y_i(t)$, $i = 1, \dots, G$; $t = 1, \dots, T$, it is necessary to regard (1.1) as a transformation expressing the $u_g(t)$ in terms of the $y_i(t)$. It is well known (see for instance [Wilks, p. 28]) that the absolute value of the volume element transforms according to

$$(2.4) \quad \left| du_1(1) du_2(1) \cdots du_G(T) \right| = J_T \left| dy_1(1) dy_2(1) \cdots dy_G(T) \right|,$$

$$J_T \equiv \left| \frac{\partial \{u_1(1), \dots, u_G(1); u_1(2), \dots, u_G(2); \dots; u_1(T), \dots, u_G(T)\}}{\partial \{y_1(1), \dots, y_G(1); y_1(2), \dots, y_G(2); \dots; y_1(T), \dots, y_G(T)\}} \right|.$$

Here J_T represents the absolute value of the Jacobian determinant of the transformation (1.1). In evaluating the elements $\partial u_g(t) / \partial y_i(t')$ of the determinant in (2.4) we find that

$$(2.5) \quad \frac{\partial u_g(t)}{\partial y_i(t)} = \beta_{gi0}, \quad \frac{\partial u_g(t)}{\partial y_i(t')} = 0,$$

if $t' > t$. Viewed as a matrix, the Jacobian in (2.4) can therefore be partitioned as follows:

$$(2.6) \quad \begin{bmatrix} B & & & \\ 0 & B & & \\ \cdot & \cdot & \dots & \\ 0 & 0 & \dots & B \end{bmatrix},$$

if, as before, $B = [\beta_{g't_0}]$. It follows that the determinant value of this matrix is independent of the elements $\partial u_g(t)/\partial y_i(t')$ with $t' < t$ [which have been left blank in (2.6)], and equals

$$(2.7) \quad J_T(B) = |\det B|^T,$$

where the symbol "det" followed by a square matrix denotes the corresponding determinant. By assumption, the $\beta_{g't_0}$ are such that $\det B$ differs from zero (see section 1.9). The distribution function therefore equals

$$(2.8) \quad F_T = |\det B|^T \cdot \prod_{t=1}^T f\{B y'(t) + \Gamma z'(t)\}.$$

In all integrations over the whole or part of the sample space, this function must, of course, be multiplied with the volume element $dy_1(1) \cdots dy_G(T)$.

The nature of the identification problem has already been explained in section 1.1. We shall now formalize it by introducing the following concepts:

DEFINITION 2.1.3.1. *A structure S consists of a set of values of the coefficient matrices B and Γ (of which B is nonsingular), and a distribution function $f(u)$ of the vector u of disturbances.*

DEFINITION 2.1.3.2. *Two structures $S = (B, \Gamma, f)$ and $S^\oplus = (B^\oplus, \Gamma^\oplus, f^\oplus)$ are called (observationally) equivalent if they imply the same probability distribution of the observations, i.e., if, for all values of T , $y_i(t)$, $z_k(t)$,*

$$(2.9) \quad \det^T(B) \prod_{t=1}^T f\{B y'(t) + \Gamma z'(t)\} \equiv \det^T(B^\oplus) \prod_{t=1}^T f^\oplus\{B^\oplus y'(t) + \Gamma^\oplus z'(t)\}.$$

This equivalence is denoted by $S \sim S^\oplus$. It follows from these definitions that equivalence of structures is transitive: if $S \sim S^\oplus$ and $S^\oplus \sim S^{\oplus\oplus}$ then $S \sim S^{\oplus\oplus}$.

We shall derive necessary and sufficient conditions for the equivalence of two structures. For that purpose, we restate certain restrictions which both structures are required to satisfy, and which are implied in the assumptions made in sections 1.4 and 1.9 respectively:

ASSUMPTION 2.1.3.3.

$$\mathcal{E}\{u(t) \mid z(t)\} = \mathcal{E}u(t) = 0.$$

This assumption is a consequence of the independence of $u(t)$ from the exogenous variables as well as from previous values $u(t')$, $t' < t$, of the disturbance vector, which together with exogenous variables determine the predetermined (lagged) endogenous variables now included in $z(t)$.

ASSUMPTION 2.1.3.4.

$$\det B \neq 0.$$

Solving (2.2) for y' through premultiplication with B^{-1} we obtain

$$(2.10) \quad y'(t) = -B^{-1} \Gamma z'(t) + B^{-1} u'(t).$$

For any equivalent structure S^\oplus we likewise have

$$(2.11) \quad y'(t) = -B^{\oplus-1} \Gamma^\oplus z'(t) + B^{\oplus-1} u^{\oplus'}(t).$$

Since the identity (2.9) in Definition 2.1.3.2 should hold for all values of T , it implies the identity

$$(2.12) \quad |\det B| \cdot f\{B y'(t) + \Gamma z'(t)\} \equiv |\det B^\oplus| \cdot f^\oplus\{B^\oplus y'(t) + \Gamma^\oplus z'(t)\}$$

of the conditional probability distributions of $y(t)$ for given $z(t)$ according to the two structures. It follows that, upon taking conditional expectations in (2.10) and (2.11) for given values of $z(t)$, and using Assumption 2.1.3.3, the same function of the elements of $z(t)$ must result from both structures:

$$(2.13) \quad \mathcal{E}\{y'(t) \mid z(t)\} = -B^{-1} \Gamma z'(t) = -B^{\oplus-1} \Gamma^\oplus z'(t).$$

Consequently

$$(2.14) \quad B^{-1} \Gamma = B^{\oplus-1} \Gamma^\oplus.$$

The square matrix of order G

$$(2.15) \quad \Upsilon = B^\oplus B^{-1}$$

is nonsingular by Assumption 2.1.3.4, and satisfies, on account of (2.15) and (2.14),

$$(2.16) \quad B^\oplus = \Upsilon B, \quad \Gamma^\oplus = \Upsilon \Gamma, \quad u'^\oplus = \Upsilon u',$$

the last equality in (2.16) being obtained from the first two by (2.2) and its counterpart for the equivalent structure S^\oplus .

Conversely, let S be a given structure, and let S^\oplus now be a structure derived from S by the transformation (2.16) where Υ is any nonsingular matrix of order G . It is easily seen that S^\oplus is then equivalent to S . For (2.16) now implies successively the nonsingularity of B^\oplus , (2.14), and

$$(2.17) \quad B^{\oplus-1} u'^\oplus(t) = B^{\oplus-1} \Upsilon u'(t) = B^{-1} u'(t).$$

Hence (2.10) and (2.11) define the same conditional distribution (2.12) of the dependent variables, for any set of predetermined variables and for any distribution function $f(u)$ of u . The two structures S and S^\oplus thus define the same distribution function (2.9) of the observed variables and are accordingly equivalent.

It will be noticed that in (2.16) the coefficient matrices B

and Γ occur in the same manner. It is therefore appropriate to express the result just obtained in the notation introduced by (2.1):

THEOREM 2.1.3.5. *A necessary and sufficient condition for the equivalence of two structures $S = (A, f(u))$ and $S^\oplus = (A^\oplus, f^\oplus(u^\oplus))$ satisfying Assumptions 2.1.3.3 and 2.1.3.4 is that they are connected by a linear transformation*

$$(2.18\alpha) \quad A^\oplus = \Upsilon A,$$

$$(2.18u) \quad u'^\oplus = \Upsilon u',$$

with nonsingular matrix Υ .

2.1.4. *Two interpretations of the implied transformation of the parameters Σ .* We note that the transformation (2.18 α) implies a transformation for the covariance matrix Σ of the disturbances whenever that matrix exists. This transformation, together with (2.18 α), can be written in matrix form as follows:

$$(2.18) \quad \begin{cases} (2.18\alpha) & A^\oplus = \Upsilon A, \\ (2.18\sigma) & \Sigma^\oplus = \Upsilon \Sigma \Upsilon'. \end{cases}$$

It should be stressed that the present discussion of the identification problem is based on the assumption that all available knowledge regarding the distribution function $f(u)$ of the disturbances is expressed by Assumption 2.1.3.3. If additional information on the functional form of f were available, the possibility exists that such information could be used for identification purposes. However, it is easily seen that there is an alternative case, in which the conclusions as regards identifiability of structural equations are precisely the same as under the present assumptions. This is the case in which the very general Assumption 2.1.3.3 is replaced by the special assumption that $f(u)$ represents a nonsingular joint normal distribution of the disturbances u_1, \dots, u_g . In this case, the space of structures S becomes the space of the parameters α_{gh} , σ_{gh} , and (2.18 σ) supplants (2.18 u) in the definition of a linear transformation in the "parameter space." Whenever the transformation (2.18) is quoted in what follows, either of the two interpretations just given is applicable.

It so happens that any additional restrictions on $f(u)$ which we shall consider in what follows are restrictions on the parameters σ_{gh} . These restrictions have the same identifying effects whether $f(u)$ is previously restricted only by Assumption 2.1.3.3 and the assumption that Σ exists, or whether $f(u)$ is previously restricted to the form of the normal distribution. For this reason, it will be convenient from now on to discuss the identification problem in terms of points (A, Σ) in the parameter space rather than in terms of structures $S = (A, f(u))$. We therefore supplement Definition 2.1.3.3 by

DEFINITION 2.1.3.6. *Two points (A, Σ) and $(A^\oplus, \Sigma^\oplus)$ in the parameter space are called (observationally) equivalent if they are connected with equivalent structures.*

Theorem 2.1.3.5 can now also be interpreted as stating conditions for the equivalence of two points in the space of the parameters A and Σ .

2.1.5. *Equivalent points in the restricted parameter space.* The identification problem in this article consists in the study of the extent to which there exist nontrivial transformations (2.18) which preserve the a priori restrictions. The following definitions are helpful in developing the concept of identification.

DEFINITION 2.1.5.1. *By the restricted parameter space we understand the set of those points (A, Σ) in the parameter space that satisfy the a priori restrictions.*

It will further be useful to rule out as irrelevant certain trivial transformations that do not affect the economic identity of the equations whose identification is studied.

DEFINITION 2.1.5.2. *A transformation (2.18) is called trivial with respect to the g_0 th structural equation if it involves only a change of scale*

$$(2.19) \quad \begin{aligned} \alpha_{g_0}^\oplus k &= v_{g_0 g_0} \alpha_{g_0} k, & \sigma_{g_0 g_0}^\oplus &= v_{g_0 g_0}^2 \sigma_{g_0 g_0}, \\ \sigma_{g_0}^\oplus g &= v_{g_0 g_0} \sum_h \sigma_{gh} v_{gh}, & & g \neq g_0, \end{aligned}$$

in the parameters of that equation.

The concept of identifiability of a structural equation is now defined as follows:

DEFINITION 2.1.5.3. *The g_0 th structural equation in (2.3) is said to be identifiable by a set of a priori restrictions, in the point (A, Σ) of the restricted parameter space, if each point $(A^{\oplus}, \Sigma^{\oplus})$ in the restricted parameter space, equivalent to (A, Σ) , is obtainable from (A, Σ) by a transformation (2.18) which is trivial with respect to the g_0 th equation.*

DEFINITION 2.1.5.4. *The system (2.3) of structural equations is said to be identifiable by a set of a priori restrictions, in the point (A, Σ) of the restricted parameter space, if each of its equations is thereby identifiable.*

If the latter definition is applied with reference to a set of a priori restrictions that includes an unambiguous normalization rule for each structural equation (2.3), the definition of identifiability of the system (2.3) is equivalent to requiring that the set of points in the restricted parameter space, equivalent to (A, Σ) , consist only of the point (A, Σ) .

2.2. Identification of One Structural Equation under Linear Restrictions

2.2.1. Necessary and sufficient conditions for identifiability of a given structural equation under linear single-parameter restrictions. Let us first consider the case of the identification of the g_0 th equation by linear a priori restrictions of the single-parameter form (1.11), which require certain specified α_{gk} to vanish. It is useful to rearrange the conditions (1.11) in the order of the structural equations to which they apply:

$$(2.20) \quad \alpha_{gk_r} = 0, \\ r = \bar{R}_{g-1} + 1, \dots, \bar{R}_g, \quad \bar{R}_g - \bar{R}_{g-1} = R_g, \quad g = 1, \dots, G,$$

with $\bar{R}_0 \equiv 0$, $\bar{R}_G \equiv R_1 + R_2 + \dots + R_G \equiv R_{\alpha}^{(1)}$.

THEOREM 2.2.1. *A necessary and sufficient condition for the*

identifiability, under Assumptions 2.1.3.3 and 2.1.3.4 of the g_0 th structural equation in (2.3) by the a priori restrictions (2.20), is that the matrix

$$(2.21) \quad A^{(g_0)} \equiv [\alpha_{g_0 k_r}],$$

$$g = 1, \dots, G, \quad r = \bar{R}_{g_0-1} + 1, \dots, \bar{R}_{g_0},$$

obtained from the complete matrix A of the coefficients α_{gk} by selecting those columns $k = k_r$ for which $\alpha_{g_0 k}$ is required to vanish, is of rank¹ $G - 1$.

Obviously, $A^{(g_0)}$ cannot have a rank higher than $G - 1$ because its g_0 th row consists of zeros only. Stated in other words, the condition in Theorem 2.2.1 requires that from the g_0 th row of the matrix $A = [\alpha_{gk}]$ we can select in at least one way $G - 1$ elements that the a priori restrictions require to be zero, such that the determinant obtained by combining the columns of those elements with all other rows differs from zero. If this theorem is true, the following corollary ensues.

COROLLARY TO THEOREM 2.2.1. *A necessary condition for the identifiability of the g_0 th structural equation by the a priori restrictions (2.20) is that the number R_{g_0} of these restrictions involving coefficients of the g_0 th equation be at least equal to the number G of structural equations less one.*

In order to prove the theorem let us first assume that $\Upsilon = [v_{gh}]$ defines a transformation of the type (2.18), which preserves those restrictions (2.20) for which $r = \bar{R}_{g_0-1} + 1, \dots, \bar{R}_{g_0}$. Then

$$(2.22) \quad \alpha_{g_0 k_r}^{\oplus} = \sum_{h=1}^G v_{g_0 h} \alpha_{h k_r} = 0, \quad r = \bar{R}_{g_0-1} + 1, \dots, \bar{R}_{g_0}.$$

Because of (2.20) the term with $h = g_0$ can be omitted from the

¹A matrix X is said to be of rank ρ , if at least one of the determinants of order ρ , and none of the determinants of order $\rho + 1$, that can be formed from the elements of X , is different from zero. Obviously, ρ cannot exceed the number of rows or columns in X .

summation. As a consequence of the condition specified in Theorem 2.2.1, the system of homogeneous equations (2.22), in which the $v_{g_0 h}$ with $h \neq g_0$ are now regarded as $G-1$ unknowns, has the rank $G-1$ (for the omission of a row of zeros from a matrix does not affect its rank). We can therefore in at least one way select from (2.22) $G-1$ equations, with $r = r_1^{(g_0)}, \dots, r_{G-1}^{(g_0)}$, say, in which the determinant $\Delta(g_0; r_1^{(g_0)}, \dots, r_{G-1}^{(g_0)})$ of the coefficients of the unknowns differs from zero. It follows that

$$(2.23) \quad v_{g_0 h} = 0$$

for $h \neq g_0$, and the transformation Υ can only be of the type (2.19) admitted in Definition 2.1.5.2. This proves that the condition stated in Theorem 2.2.1 is sufficient for identifiability of the g_0 th structural equation.

That this condition is also necessary is seen if we now assume that the g_0 th equation is identifiable and that therefore the only nonsingular transformations Υ satisfying (2.22) if (2.20) holds are those that satisfy (2.23). Suppose that at the same time $A^{(g_0)}$ has a rank lower than $G-1$. Then (2.22) would possess at least two linearly independent solutions $v_{g_0 h}^{(1)}$ and $v_{g_0 h}^{(2)}$, say, of which the first can be taken to satisfy (2.23). The more general solution $v_{g_0 h} = \lambda_1 v_{g_0 h}^{(1)} + \lambda_2 v_{g_0 h}^{(2)}$ then satisfies (2.23) only if $\lambda_2 = 0$. Since λ_1 and the other rows ($g \neq g_0$) of Υ can always be selected so that Υ coincides with the identical transformation (with unit matrix) when $\lambda_2 = 0$, there are values $\lambda_2 \neq 0$ of λ_2 (e.g., in a neighborhood of $\lambda_2 = 0$) for which Υ is nonsingular, but as stated does not satisfy (2.23). This contradicts the assumption made at the beginning of this paragraph. Therefore a rank $G-1$ of $A^{(g_0)}$ is also necessary.

2.2.2. *Identifiability conditions under more general linear restrictions.* Theorem 2.2.1 can easily be generalized to the case where the linear a priori restrictions take the form (1.12). It will be useful now to write these restrictions in matrix form:

$$(2.24) \quad \alpha(g)\Phi'_g = 0, \quad g = 1, \dots, G.$$

Here $\alpha(g)$ is a row vector¹ containing the elements of the g th row of A . The matrix Φ'_g is the transpose of a matrix Φ_g of rank R_g , containing in R_g rows and K_x columns the coefficients "χ" and "ψ" of those restrictions (1.12) that refer to the g th structural equation in (2.3). (It is permissible to take the rank of Φ_g equal to the number of its rows, since otherwise the restrictions expressed by (2.24) would not be independent.) Again

$$(2.25) \quad R_1 + \dots + R_G = R_\alpha^{(1)}.$$

We consider only such transformations (2.18) that preserve the restrictions (2.24). Hence, if v_g is the g th row of $Y \equiv Y_{yy}$,

$$(2.26) \quad \alpha^\oplus(g)\Phi'_g = v_g A \Phi'_g, \quad g = 1, \dots, G.$$

By a repetition of the previous reasoning, it is seen that the g_0 th condition (2.26) will then and only then require the elements of v_{g_0} (except $v_{g_0 g_0}$) to vanish, if $A \Phi'_{g_0}$ (of which the g_0 th row consists of zeros only) has the rank $G - 1$. We therefore have

THEOREM 2.2.2. *A necessary and sufficient condition for the identifiability, under the Assumptions 2.1.3.3 and 2.1.3.4, of the g_0 th structural equation by the a priori restrictions (2.24) is that $A \Phi'_{g_0}$ has the rank $G - 1$.*

Since the rank of Φ_g is assumed equal to the number R_g of its rows (the number of independent restrictions imposed on the g th structural equation), and since the rank of a matrix cannot increase through premultiplication with another matrix, we have

COROLLARY TO THEOREM 2.2.2. *A necessary condition for the*

¹Vectors are here considered as one-row matrices rather than the more commonly used one-column matrices in order to be able to treat rows of A as vectors, with A in the form corresponding to the natural way (1.1) of writing the structural equations.

identifiability, under the assumptions of Theorem 2.2.2, of the g_0 th structural equation by the a priori restrictions (2.24) is that the number of independent restrictions expressed by Φ_{g_0} (i.e., the number of rows of Φ_{g_0}) be at least $G - 1$.

2.2.3. *Identifiability almost everywhere in the parameter space.* It is of interest to note that the identifiability of the g_0 th structural equation depends only on the matrix Φ_{g_0} expressing the restrictions on that particular equation, and on the coefficient matrix A . In practice, however, the elements of A are unknown before estimation, and are not known exactly after estimation. Uncertainty with regard to the rank of $A\Phi'_{g_0}$ may therefore remain even after estimation. A question that can be answered before estimation, however, is whether, in cases where $\rho(\Phi_{g_0}) \geq G - 1$, the a priori restrictions (2.24) with $g \neq g_0$, i.e., those restricting structural equations other than the one whose identification is studied, do or do not reduce the rank of $A\Phi'_{g_0}$ below $G - 1$ identically, i.e., for all values of A permitted by (2.24). If the a priori restrictions are in the form (2.20), this question can be decided in a simple way by exhaustive study of a diagram of the elements of the matrix A , in which a zero is entered for every element required to be zero by (2.20), and a cross for every other element. A determinant Δ extracted from A is then not identically zero if at least one term of the determinant can be found that is the product of elements represented by crosses only.

The generalization of this technique to the case of restrictions in the form (2.24), although mathematically interesting, is somewhat complicated in operation. It appears preferable, for this particular purpose, to reduce to a minimum the number of restrictions not in the form (2.20). This can be achieved by retaining the identities as part of the equation system. This is a possible procedure because the assumption of nonsingularity of Σ has not been used in the present discussion of identification problems, and identities are therefore admissible to the equation system of which identification properties are studied. In this case, of course, no identification problem arises with respect to the identities themselves, since these are completely known already. If a few restrictions (2.24) remain that cannot be reduced to the form (2.20), it is advisable first to carry out the analysis indicated while ignoring

those particular restrictions, investigating thereafter whether or not the conclusions are changed by the presence of these restrictions.

Assuming that most situations arising in practice can be dealt with in such manner, we shall not attempt to formulate a general theorem covering all cases under which the rank of $A \Phi'_{g_0}$ may be *identically* below that of Φ_{g_0} . In section 2.2.4, however, we shall discuss some interesting special cases in which this occurs. Meanwhile, it is already possible to make the following generalization: All criteria for identifiability formulated above are in terms of ranks of matrices. Since a determinant is a linear function of any element, and of the elements in any row, a matrix that under linear restrictions of the type (2.20) or (2.24) attains a required rank in one point of the parameter space, attains the required rank in all points except for a set of measure zero. Thus a structural equation that is not identically unidentifiable under such restrictions, is identifiable almost everywhere in the parameter space.

*2.2.4. *Cases where the rank of $A \Phi'_{g_0}$ is identically below that of Φ'_{g_0} .* Let us now consider some special cases in which the rank of $A \Phi'_{g_0}$ is identically less than that of Φ'_{g_0} on account of the restrictions imposed on the other structural equations. Since the rank of Φ_{g_0} is at the same time the number R_{g_0} of rows in Φ_{g_0} , every nonvanishing vector λ_{g_0} containing R_{g_0} elements satisfies

$$(2.27) \quad \Phi'_{g_0} \lambda'_{g_0} \neq 0.$$

Let $\rho(X)$ represent the rank of any matrix X . Then, if

$$(2.28) \quad \rho(A \Phi'_{g_0}) < \rho(\Phi'_{g_0}) = R_{g_0}$$

for all values of A satisfying (2.24), there exists for every such value of A a nonvanishing vector λ_{g_0} such that

$$(2.29) \quad A \Phi'_{g_0} \lambda'_{g_0} = 0, \text{ or } A \xi' = 0, \text{ where } \xi' \equiv \Phi'_{g_0} \lambda'_{g_0} \neq 0.$$

Now there are two possible cases. It may occur that a *constant* vector λ exists for which (2.29) is satisfied by all values of A permitted by (2.24). In this case the restriction

$$(2.30) \quad \alpha(g) \xi' = 0,$$

ξ constant, must hold for each row $\alpha(g)$, $g = 1, \dots, G$, of A , as a consequence of (2.24), i.e., there exist vectors λ_g such that

$$(2.31) \quad \Phi'_g \lambda'_g = \xi', \quad g = 1, \dots, G.$$

This means that the restrictions on the individual structural equations imply at least one restriction (2.30), which is common to all of them. The coefficients ξ of this common restriction do not depend on A , but can be determined or selected on the basis of the Φ_g , $g = 1, \dots, G$, alone. It follows that the number K_x of variables x_k can be reduced by one without changing the nature of the problem studied. This is obvious in the special case that ξ has only one nonvanishing element, say ξ_1 , because then the a priori restrictions imply that the variable x_1 does not actually occur in any one of the structural equations. The same conclusion can be drawn as follows if ξ is any other constant vector. Let Ξ be a nonsingular square matrix of order K_x containing ξ as its first row. Then the linear transformation

$$(2.32) \quad A \Xi \equiv A^\oplus, \quad \Xi^{-1} x' \equiv x'^\oplus,$$

of variables x and coefficients A leads to a situation where the a priori restrictions imply $\alpha_{1g}^\oplus = 0$, $g = 1, \dots, G$, that is, where the variable x_1^\oplus does not actually occur.

Alternatively, it may occur that (2.29) can be satisfied only by a vector ξ which itself depends on A . A simple example is that of $G = 3$ equations where the a priori restrictions require a submatrix consisting of the k th and l th columns of A to be as follows:

$$(2.33) \quad \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ \alpha_{3k} & \alpha_{3l} \end{bmatrix},$$

with α_{3k} , α_{3l} , and the remaining elements of $\alpha(1)$ and $\alpha(2)$ unrestricted. Here we must take

$$(2.34) \quad \xi_k = -\alpha_{3l}, \quad \xi_l = \alpha_{3k}, \quad \xi_m = 0 \quad \text{if } m \neq k, l,$$

to obtain a restriction (2.30) common to all structural equations. However, the first $G - 1 = 2$ structural equations have two independent constant restrictions in common¹, and are not identifiable.

*2.3. Treatment of Unidentifiable Structural Equations by Linear Dummy Restrictions

*2.3.1. *Dummy restrictions to produce formal identifiability.*
In case one or more structural equations are not identifiable, this fact should not be allowed to interfere with the estimation of such other equations as are identifiable. We may even wish to go further and estimate such linear functions of the parameters of the unidentifiable equations as are not affected by the lack of identification. From the point of view both of estimation and computation, therefore, it is an important problem to write the parameters (A, Σ) as functions of two sets of parameters, θ_1 and θ_2 say, of which the first uniquely defines a set of points in the restricted (A, Σ) space, while variation of the second set of parameters θ_2 only causes the point (A, Σ) to vary within a set of equivalent points in that space. Using Wald's extension of the identifiability concept to individual parameters [III], we may say that the parameters (A, Σ) are written as functions of a set of identifiable parameters θ_1 and a set of unidentifiable parameters θ_2 .

A device whereby the separation of θ_1 and θ_2 can be carried out is the addition to the original a priori restrictions of dummy restrictions such that the so augmented set of restrictions ensures identifiability. The dummy restrictions are then made to contain as many parameters θ_2 as are required to define a point within a set of equivalent points. It will be clear that this device is applicable only within a region of the restricted parameter space, in which the matrices $A \Phi'_g$, $g = 1, \dots, G$, have constant ranks. For the sake of simplicity, we shall only discuss the case of a region in which the rank of each $A \Phi'_g$ equals that of the corresponding

¹It is again true here that the variables x_k and x_l occur only in the linear combination $\alpha_{3k}x_k + \alpha_{3l}x_l$, but this does not permit us to reduce the number of variables because α_{3k} and α_{3l} are unknown.

restriction matrix Φ'_g . That is, we shall disregard the case in which the rank of an $A\Phi'_g$ is identically depressed on account of restrictions on the structural equations other than the g th.

*2.3.2. *A lemma.* In preparation for a theorem stating what can be achieved by dummy restrictions, we shall first prove

LEMMA 2.3.2. *If Φ and Ψ are two matrices with an equal number of rows, such that¹*

$$(2.35) \quad \rho(\Phi \ \Psi) < \rho(\Phi) + \rho(\Psi),$$

then there exist two nonvanishing vectors λ and μ such that

$$(2.36) \quad \Phi \lambda' + \Psi \mu' = 0, \quad \Phi \lambda' \neq 0.$$

Let $c(\Phi)$ denote the number of columns of Φ , and assume first as a special case that

$$(2.37) \quad \rho(\Phi) = c(\Phi), \quad \rho(\Psi) = c(\Psi).$$

(It is only with respect to this special case that the lemma is used in the present section; the further case (2.39) is added for later use, see sections 3.2.5 and 4.3.4.6.) Then, from (2.35) and (2.37),

$$(2.38) \quad \rho(\Phi \ \Psi) < c(\Phi \ \Psi).$$

Hence there exists a nonvanishing vector $[\lambda \ \mu]$ satisfying the equality in (2.36). However, this cannot be a vector such that λ vanishes, because then the equality in (2.36) would imply $\Psi \mu' = 0$ with $\mu' \neq 0$, which is precluded by the second condition in (2.37). Similarly, the first condition in (2.37) precludes the vanishing of $\Phi \lambda'$ now that $\lambda \neq 0$.

Now assume, more generally, that

$$(2.39) \quad \rho(\Phi) \leq \dot{c}(\Phi), \quad \rho(\Psi) \leq c(\Psi).$$

Then it is possible, wherever an inequality sign in (2.39) applies,

¹Square brackets $[\]$ denoting matrices are omitted when a matrix appears as argument of the functions $\rho(\)$, $r(\)$, $c(\)$.

to delete one or more columns from Φ or Ψ or both, so as to obtain matrices $\bar{\Phi}$ and $\bar{\Psi}$ respectively such that

$$(2.40) \quad \rho(\Phi) = \rho(\bar{\Phi}) = c(\bar{\Phi}), \quad \rho(\Psi) = \rho(\bar{\Psi}) = c(\bar{\Psi}),$$

and, from (2.35) and (2.40),

$$(2.41) \quad \rho(\bar{\Phi} \quad \bar{\Psi}) \leq \rho(\Phi \quad \Psi) < \rho(\bar{\Phi}) + \rho(\bar{\Psi}).$$

The conditions (2.40) and (2.41) are equivalent to (2.35) and (2.37) with Φ and Ψ replaced by $\bar{\Phi}$ and $\bar{\Psi}$. Hence there exist nonvanishing vectors $\bar{\lambda}$ and $\bar{\mu}$ such that

$$(2.42) \quad \bar{\Phi} \bar{\lambda}' + \bar{\Psi} \bar{\mu}' = 0, \quad \bar{\Phi} \bar{\lambda}' \neq 0.$$

By adding zero elements in the proper places, these can be completed to vectors λ and μ satisfying (2.36).

**2.3.3. The number and type of dummy restrictions required.*

We shall now study the imposition of dummy restrictions on the g th structural equation in the neighborhood of such a point A_0 in the space of the parameters A in which

$$(2.43) \quad \rho(A_0 \cdot \Phi'_g) = \rho(\Phi'_g).$$

As before, the restriction matrix Φ_g is so chosen that

$$(2.44) \quad \rho(\Phi'_g) = c(\Phi'_g) = R_g$$

equals the number of independent restrictions imposed on the g th structural equation. Since the g th row of $A_0 \Phi'_g$ consists of zeros only, we may as well operate with a matrix ${}_g A_0$ obtained from A_0 by deleting the g th row, and write instead of (2.43)

$$(2.45) \quad \rho({}_g A_0 \cdot \Phi'_g) = \rho(\Phi'_g).$$

For reasons of notational symmetry, we shall in the remainder of this section 2.3 occasionally use the symbol K_y , introduced in

section 2.1.2 as synonymous with G (the number of structural equations and of dependent variables). Suppose now that

$$(2.46) \quad R_g < G - 1 \equiv K_y - 1,$$

that is, that the g th structural equation is not identified by the original restrictions (2.24). We shall continue to refer to the restrictions (2.24) as the a priori restrictions, and to regard them as the basis for the concept of equivalence according to Definition 2.1.3.6. We now wish to achieve identification of the g th equation by the addition to (2.24) of dummy restrictions which we denote

$$(2.47) \quad \alpha(g) \cdot \bar{\Phi}_g^{(0)} = 0, \quad c(\bar{\Phi}_g^{(0)}) = \bar{R}_g.$$

We shall of course require that the dummy restrictions are independent of each other and of the a priori restrictions, i.e.,

$$(2.48) \quad \rho(\Phi'_g \quad \bar{\Phi}_g^{(0)}) = R_g + \bar{R}_g.$$

THEOREM 2.3.3. *If in a point (A_0, Σ_0) of the parameter space the following conditions are satisfied*

- (a) *the a priori restrictions (2.24),*
- (b) *the rank condition (2.45),*
- (c) *the insufficiency (2.46) of the number R_g of independent (2.44) a priori restrictions on the g th structural equation to identify that equation,*

then there exists a neighborhood N of (A_0, Σ_0) in the a priori restricted parameter space, and a matrix $\bar{\Phi}_g^{(0)}$ defining

$$(2.49) \quad \bar{R}_g = G - 1 - R_g$$

dummy restrictions (2.47) on the g th structural equation, which are independent, mutually and from the a priori restrictions, such that

- (i) *to each point (A, Σ) in N can be found at*

least one equivalent point satisfying the dummy restrictions (2.47),

- (ii) the a priori and dummy restrictions taken in combination identify the g th structural equation.

The statement (i) in this theorem assures that no probability distribution of the observations permitted by the a priori restrictions is deprived of representation by imposing dummy restrictions.

Introducing the notation

$$(2.50) \quad \bar{\Phi}'(g)^{(0)} \equiv [\Phi'_g \quad \bar{\Phi}'_g{}^{(0)}],$$

we shall first show that there exists a matrix $\bar{\Phi}'_g{}^{(0)}$ with the number of columns \bar{R}_g as given by (2.49), such that

$$(2.51) \quad \rho({}_g A_0 \cdot \bar{\Phi}'(g)^{(0)}) = G - 1 = K_y - 1 = \rho(\bar{\Phi}'(g)^{(0)}).$$

The first step is to choose an arbitrary orthogonal complement of ${}_g A_0$, that is, a matrix $\bar{\bar{\Phi}}_g^{(0)}$ such that both

$$(2.52) \quad \rho(\bar{\bar{\Phi}}_g^{(0)}) = c(\bar{\bar{\Phi}}_g^{(0)}) = K_x - K_y + 1 = K_z + 1,$$

say, and

$$(2.53) \quad {}_g A_0 \cdot \bar{\bar{\Phi}}_g^{(0)} = 0.$$

However this choice is made, we must have

$$(2.54) \quad \rho(\Phi'_g \quad \bar{\bar{\Phi}}_g^{(0)}) = c(\Phi'_g \quad \bar{\bar{\Phi}}_g^{(0)}) = R_g + K_z + 1.$$

For otherwise, according to Lemma 2.3.2, nonvanishing vectors λ_g and $\bar{\lambda}_g$ exist such that

$$(2.55) \quad \Phi'_g \cdot \lambda'_g + \bar{\bar{\Phi}}_g^{(0)} \cdot \bar{\lambda}'_g = 0,$$

and, on account of (2.53),

$$(2.56) \quad {}_g A_0 \cdot \Phi'_g \cdot \lambda'_g = 0,$$

which is incompatible with (2.44) and (2.45).

Because of (2.54), we can, as the next step, choose a matrix $\bar{\Phi}'^{(0)}$ with a number of rows \bar{R}_g as given by (2.49), such that

$$(2.57) \quad \bar{\Phi}'(g) \equiv [\bar{\Phi}'_g \quad \bar{\Phi}'_g^{(0)} \quad \bar{\Phi}'_g^{(0)}] \equiv [\bar{\Phi}'(g)^{(0)} \quad \bar{\Phi}'_g^{(0)}]$$

is a nonsingular square matrix of order $K_x \equiv K_y + K_z$. The nonsingularity of $\bar{\Phi}'(g)^{(0)}$ ensures that (2.48) is satisfied. It follows further that

$$(2.58) \quad \rho({}_g A_0 \cdot \bar{\Phi}'(g)^{(0)}) = \rho({}_g A_0),$$

because postmultiplication with a nonsingular square matrix does not affect rank. On the other hand, from (2.53) and (2.57),

$$(2.59) \quad \rho({}_g A_0 \cdot \bar{\Phi}'(g)^{(0)}) = \rho({}_g A_0 \cdot \bar{\Phi}'(g)^{(0)}).$$

Finally, since the structural equations are independent [see also (1.6)],

$$(2.60) \quad \rho(A_0) = K_y = 1 + \rho({}_g A_0).$$

Combining (2.58), (2.59), and (2.60), we have completed the proof of the first equality in (2.51). The second equality in (2.51) follows directly from (2.48), (2.49), and the nonsingularity of $\bar{\Phi}'(g)^{(0)}$

Since a determinant is a continuous function of its elements, it follows from (2.51) that

$$(2.61) \quad \rho({}_g A \cdot \bar{\Phi}'(g)^{(0)}) = K_y - 1 = \rho(\bar{\Phi}'(g)^{(0)})$$

in a neighborhood ${}_g N_\alpha$ of the point ${}_g A_0$ in the space of the parameters ${}_g A$ subject to the a priori restrictions (2.24).

Let N be a neighborhood of the point (A_0, Σ_0) in the restricted parameter space such that for all points (A_0, Σ_0) in N , the coor-

dinates ${}_g A$ define a point ${}_g A$ in ${}_g N_\alpha$. We shall now show that N and $\Phi_g^{(0)}$ have the properties required in the theorem.

To satisfy condition (i) we associate with any point (A, Σ) of N , a point

$$(2.62) \quad \bar{A} = \Upsilon A, \quad \bar{\Sigma} = \Upsilon \Sigma \Upsilon',$$

by the following choice of the transformation Υ ,

$$(2.63) \quad \Upsilon = \begin{bmatrix} 1 & \dots & 0 & 0 & 0 & \dots & 0 \\ \cdot & \dots & \cdot & \cdot & \cdot & \dots & \cdot \\ 0 & \dots & 1 & 0 & 0 & \dots & 0 \\ v_{g1} & \dots & v_{g,g-1} & 1 & v_{g,g+1} & \dots & v_{gG} \\ 0 & \dots & 0 & 0 & 1 & \dots & 0 \\ \cdot & \dots & \cdot & \cdot & \cdot & \dots & \cdot \\ 0 & \dots & 0 & 0 & 0 & \dots & 1 \end{bmatrix},$$

in which the vector

$$(2.64) \quad {}_g v(g) = [v_{g1} \quad \dots \quad v_{g,g-1} \quad v_{g,g+1} \quad \dots \quad v_{gG}]$$

of as yet unspecified elements in the g th row $v(g)$ of Υ is determined by

$$(2.65) \quad {}_g v(g) \cdot {}_g A \cdot \bar{\Phi}'(g)^{(0)} = -\alpha(g) \cdot \bar{\Phi}'(g)^{(0)}.$$

As a consequence of (2.61), there is always one and only one solution ${}_g v(g)$ to this condition. Rewriting (2.65) as

$$(2.66) \quad \alpha(g) \cdot \bar{\Phi}'(g)^{(0)} = v(g) \cdot A \cdot \bar{\Phi}'(g)^{(0)} = 0,$$

we see, in connection with (2.50), that the point (2.62) satisfies both the a priori conditions (2.24) and the dummy restrictions (2.47).

Finally the uniqueness of the solution ${}_g v(g)$ of (2.65) implies that any other point $(\bar{A}^\oplus, \bar{\Sigma}^\oplus)$ that is obtainable from (A, Σ) by a linear transformation (2.62) and satisfies the combined (a priori and dummy) restrictions is obtainable from $(\bar{A}, \bar{\Sigma})$ by a transformation that is trivial with respect to the g th structural equation (see Definition 2.1.5.2). For, the combined restrictions on the g th equation are expressed by (2.66), which permit free choice only of the diagonal element v_{gg} in $v(g)$. This shows that condition (ii) of the theorem is also satisfied.

*2.3.4. *The degree of indeterminacy of an a priori unidentifiable structural equation.* It is of interest now to consider the reverse problem that arises after estimation has been carried out subject to dummy restrictions added to obtain formal identifiability: From a parameter point $(\bar{A}, \bar{\Sigma})$ that satisfies both the a priori and the dummy restrictions, reconstruct the set of all equivalent points (A, Σ) in the restricted parameter space; i.e., all points obtained through linear transformations inverse to (2.62) and satisfying the a priori restrictions but not necessarily the dummy restrictions. Since the identification problem under the restrictions (2.24) can be studied for each equation separately, it is sufficient to study those transformations for which A differs from \bar{A} only as regards the g th row $\alpha(g)$ or

$$(2.67) \quad {}_g A = {}_g \bar{A}.$$

This means that we can confine ourselves to transformations

$$(2.68) \quad A = \Upsilon^{-1} \bar{A}, \quad \Sigma = \Upsilon^{-1} \bar{\Sigma} \Upsilon'^{-1},$$

with a matrix Υ^{-1} inverse to a matrix Υ of type (2.63). It is easily seen that Υ^{-1} itself then is of the same type, with nondiagonal elements in the g th row equal to

$$(2.69) \quad v^{gh} = -v_{gh}, \quad h \neq g.$$

For the g th row of A this gives us in particular

$$(2.70) \quad \alpha(g) = \bar{\alpha}(g) - {}_g v(g) \cdot {}_g \bar{A},$$

and, through postmultiplication with $\Phi_g'^{(0)}$ and $\bar{\Phi}_g'^{(0)}$ respectively, using (2.24), (2.50), and (2.66),

$$(2.71) \quad \begin{cases} (2.71r) & 0 = - {}_g v(g) \cdot {}_g \bar{A} \cdot \Phi_g'^{(0)}, \\ (2.71\bar{r}) & \alpha(g) \cdot \bar{\Phi}_g'^{(0)} = - {}_g v(g) \cdot {}_g A \cdot \bar{\Phi}_g'^{(0)}. \end{cases}$$

Of these conditions, (2.71r) expresses the restrictions on Υ^{-1} , arising through (2.69) from the fact that $\alpha(g)$ must satisfy the a priori restrictions. The left-hand member in (2.71 \bar{r}) arises from the fact that $\alpha(g)$ is not bound by the dummy restrictions. It is easily seen that precisely those linear combinations of the elements of $\alpha(g)$ which are the elements of the vector

$$(2.72) \quad \bar{\theta}_g \equiv \alpha(g) \bar{\Phi}_g'^{(0)}$$

can be chosen arbitrarily before, according to (2.61), (2.67), and (2.69), Υ^{-1} and therewith $\alpha(g)$ is fully determined. Therewith we have proved:

THEOREM 2.3.4. *If under the conditions of Theorem 2.3.3 the parameter point $(\bar{A}, \bar{\Sigma})$ satisfies both the a priori restrictions (2.24) and the dummy restrictions (2.47) identifying the g th structural equation, the set of points (A, Σ) , equivalent to $(\bar{A}, \bar{\Sigma})$ but satisfying only the a priori restrictions, is obtained from the latter point by a transformation (2.68) with a matrix Υ^{-1} of which the g th row is through (2.69) determined from (2.71), with arbitrary choice of the vector $\bar{\theta}_g$ of dummy parameters (2.72).*

The theorem does not stipulate that the remaining rows of Υ^{-1} correspond to the form (2.63), since no assumptions were made as to the identifiability of the remaining structural equations. It will be clear, however, that the transformed point (2.68) must satisfy the a priori restrictions (2.24) on all structural equations.

2.4. Identification of a Set of Structural Equations under Linear and Bilinear Restrictions

2.4.1. *Problems arising from additional types of restrictions.*
We shall now discuss identification problems that arise if two fur-

ther types of a priori restrictions are added, each of which binds the parameters of one structural equation to those of another. The first type is given by (1.13) or (1.14), which is rewritten in slightly different notation in (2.73 α). The second type (2.73 σ) expresses absence of correlation (or, under the normality assumption, independence) between the disturbances in two equations.

$$(2.73) \left\{ \begin{array}{l} (2.73\alpha) \quad \begin{bmatrix} \alpha_{g_r k_r} & \alpha_{g_r l_r} \\ \alpha_{h_r k_r} & \alpha_{h_r l_r} \end{bmatrix} = 0, \quad r = 1, \dots, R_\alpha^{(2)}, \\ (2.73\sigma) \quad \sigma_{g_r h_r} = 0, \quad \begin{array}{l} g_r < h_1, \\ r = R_\alpha^{(2)} + 1, \dots, R_\alpha^{(2)} + R_\sigma. \end{array} \end{array} \right.$$

For the time being, we shall assume no particular pattern for the occurrence of zeros among the elements of Σ . Later we shall say a few words concerning the special pattern (1.15).

If a given structural equation can already be identified on the basis of the restrictions (2.24) alone, the imposition of additional restrictions (2.73) of course does not detract from the identifiability of that equation. The only identification problem of interest in connection with the restrictions (2.73) is therefore under what circumstances an equation not identifiable on the basis of (2.24) alone can be identified if the restrictions (2.73) are added.

Each of the restrictions (2.73) connects two structural equations, numbered g_r and h_r , which we shall call the two equations *referred to* in that restriction. [Similarly, we shall speak of the one equation referred to by any one of the restrictions (2.24).] The restrictions (2.73) link up the identification problem of individual equations, and statements regarding identifiability will therefore in general relate either to the whole set of structural equations (2.3) or to subsets thereof which only in special cases may consist of one single equation.

2.4.2. *The additional restrictions are bilinear.* If the parameters A , Σ satisfy the a priori restrictions (2.24) and (2.73), to which we shall add the normalization rules (1.17a), the requirement that the transformed parameters (2.18) shall satisfy the same restrictions leads to the following expression of the

$R_\alpha^{(1)} + R_\alpha^{(2)} + R_\sigma + G$ a priori restrictions in terms of the rows v_g , $g = 1, \dots, G$, of the transformation matrix Y :

$$(2.74l) \quad \left\{ \begin{array}{l} (2.74lh) \quad v(g) \cdot A \cdot \Phi'_g = 0, \quad R \text{ restrictions,} \\ (2.74ln) \quad v(g) \cdot A \cdot \iota(i_g) = 1, \quad \text{one restriction,} \end{array} \right. \\ \left. \begin{array}{l} (R_\alpha^{(1)} + R_\sigma \\ \text{restrictions}) \end{array} \right\} \quad g = 1, \dots, G.$$

$$(2.74) \quad \left\{ \begin{array}{l} (2.74b) \quad \left[\begin{array}{ll} v(g_r) \cdot A \cdot \iota'(k_r) & v(g_r) \cdot A \cdot \iota'(l_r) \\ v(h_r) \cdot A \cdot \iota'(k_r) & v(h_r) \cdot A \cdot \iota'(l_r) \end{array} \right] = 0, \\ (R_\alpha^{(2)} + R_\sigma \\ \text{restrictions}) \quad \left. \begin{array}{l} (2.74b\alpha) \\ (2.74b\sigma) \end{array} \right\} \quad \begin{array}{l} r = 1, \dots, R_\alpha^{(2)}, \\ v(g_r) \cdot \Sigma \cdot v'(h_r) = 0, \\ r = R_\alpha^{(2)} + 1, \dots, R_\alpha^{(2)} + R_\sigma. \end{array} \end{array} \right.$$

Here $\iota(k)$ is the k th row of the unit matrix of order K_x ; i.e., a vector of which the k th element is 1 and all other elements are 0. The normalization rules (2.74l) are nonhomogeneous in the elements of Y .

All other restrictions in (2.74) will be referred to as the homogeneous restrictions. The two types of restrictions under (2.74b), while being quadratic in the elements of Y , are linear in the elements of any row of Y . For this reason we shall refer to (2.74b) as the *bilinear* restrictions.

The occurrence of bilinear restrictions greatly complicates the identification problem. The following discussion should be regarded as a first exploration of the field, and does not lead to firm criteria such as were derived for linear restrictions only.

2.4.3. *The solution $Y = I$ is always present.* It is important to note that, because the parameters A and Σ are assumed to satisfy the a priori restrictions (1.17a), (2.24), and (2.73), the system of restrictions (2.74) and any of its subsystems always

permit of one particular solution Υ , viz., the identical transformation

$$(2.75) \quad \Upsilon = I,$$

where I represents the unit matrix of order K . For this reason there can never be too many compatible restrictions for identification. There can only be either too few or a sufficient number.

2.4.4. *Unique, multiple, and complete identification.* We shall discuss the identifiability of a given subset S of the structural equations (2.3), containing the H equations for which $g = g_1, \dots, g_H$, on the basis of a given subset R of the restrictions (2.74). This discussion is concerned with matrices of the type

$$(2.76) \quad \Upsilon^S \equiv \begin{bmatrix} u_{g_1 1} & \cdots & u_{g_1 G} \\ \cdot & \cdots & \cdot \\ u_{g_H 1} & \cdots & u_{g_H G} \end{bmatrix}$$

combining the rows, corresponding to the equations of S , of a solution Υ of the subset R of the restrictions (2.74).

DEFINITION 2.4.4.1. *A subset S of the structural equations will be said to be uniquely identifiable by a subset R' of the restrictions (2.74) that includes all normalization rules (2.74ln) relevant to S , if for all solutions Υ of R' we have*

$$(2.77) \quad \Upsilon = I_{[G_S, G]},$$

where $I_{[G_S, G]}$ represents a unit matrix of order G_S adjoined to a zero matrix of G columns. We shall speak of multiple identification if to all solutions Υ of R' there corresponds a finite number of different matrices Υ^S that exceeds one, and of incomplete identification if the number of different matrices Υ^S is infinite. Complete identification means either unique or multiple identification.

Incomplete identifiability is, of course, synonymous with un-

identifiability. In order to obtain conformity with Definition 2.1.5.3 we add:

DEFINITION 2.4.4.2. *A subset S of the structural equations will be said to be uniquely, multiply, or incompletely identifiable with respect to a subset R of the homogeneous restrictions in (2.74) if, after addition to R of the normalization rules relevant to S , it is so identifiable in the sense of Definition 2.4.4.1.*

The possibility of complete but multiple identification arises, of course, from the presence of quadratic restrictions in (2.74). A simple example is that of three equations in three variables, in which the a priori restrictions assume the form

$$(2.78) \begin{cases} (2.78a) & \alpha_{11} = \alpha_{22} = \alpha_{33} = 0, \\ (2.78b) & \sigma_{12} = \sigma_{23} = \sigma_{31} = 0, \quad \sigma_{11} = \sigma_{22} = \sigma_{33} = 1. \end{cases}$$

If for convenience we impose the normalization condition in (2.78b) only on the original matrix Σ but not on Σ^\oplus , we find that the remaining conditions (2.78) for the transformed matrices A^\oplus and Σ^\oplus permit, not only of any transformation with a diagonal matrix Υ corresponding to a change of scales, but also of the transformation of which the elements are obtained from

$$(2.79) \begin{aligned} v_{11} &= \alpha_{12} \alpha_{13} (\alpha_{21}^2 + \alpha_{31}^2) + \alpha_{21} \alpha_{31} \alpha_{23} \alpha_{32}, \\ v_{12} &= \alpha_{31} (\alpha_{12} \alpha_{23} \alpha_{31} - \alpha_{13} \alpha_{32} \alpha_{21}), \\ v_{13} &= -\alpha_{21} (\alpha_{12} \alpha_{23} \alpha_{31} - \alpha_{13} \alpha_{32} \alpha_{21}), \end{aligned}$$

by cyclical permutation (followed again by any change of scales).

In the case of complete but multiple identification, the number of solutions can sometimes be reduced, or even unique identification can be achieved, through additional a priori restrictions in the form of inequalities (see section 1). The use of such restrictions with respect to the elements of Σ in a case of incomplete identification has been demonstrated by Marschak and Andrews [1944].

In the remainder of this section we shall concentrate on the question of completeness or incompleteness of identifiability, irrespective of the number of solutions in the case of complete identification. We shall first make some remarks on the counting of

restrictions as an indication of identifiability. Thereafter, we shall discuss in which respects procedures based on counting alone may be insufficient to establish either identifiability or lack of identifiability.

*2.4.5. *Criteria of identifiability based on the counting of restrictions.* Apart from H normalization factors (one for each row), the matrix (2.76) contains $H(G - 1)$ unknowns. If R contains at least $H(G - 1)$ homogeneous restrictions, however, these may still be unevenly divided between the different rows of Υ^S , leaving some rows undetermined for lack of an adequate number of restrictions. In formulating principles for counting restrictions on individual rows, it is necessary to remember that each of the bilinear restrictions (2.74b) refers to two rows of Υ^S , and obviously should not be counted as a new restriction with regard to the identification of each of those two rows. These considerations lead to the following definition.

DEFINITION 2.4.5. *The subset R of the restrictions (2.74) will be said to be adequate in number and variety with respect to (the identification of) the subset S of the structural equations (2.3) if it is possible to assign each bilinear restriction (2.74) occurring in R to one of the two equations (2.3) to which it refers in such a way, that the number of homogeneous linear equations (2.74lh) in R referring to, plus the number of bilinear conditions (2.74b) in R assigned to, each equation of S is at least $G - 1$.*

*2.4.6. *The completed subset of structural equations.* The concept introduced by Definition 2.4.5 can be applied in particular to the identification of the set of all equations (2.3) by the set of all restrictions (2.74). If some of the equations (2.3) cannot be identified for lack of an adequate number and variety of a priori restrictions, it becomes necessary to develop criteria that are of assistance in finding the largest subset S that can be identified.

DEFINITION 2.4.6.1. *The subset R_S of a priori restrictions (2.74) associated with a given subset S of the structural equations consists of all homogeneous linear conditions (2.74lh) referring to an equation of S and all bilinear conditions (2.74b) referring to two equations of S .*

DEFINITION 2.4.6.2. *A subset S_0 of the structural equations will be called a completed subset if a) the subset R_{S_0} of restric-*

tions (2.74) associated with S_0 is adequate in number and variety with respect to S_0 and b) there exists no larger subset S' that contains all the equations of S_0 and one or more others besides, and with respect to which the associated subset $R_{S'}$ of the restrictions (2.74) is adequate in number and variety.

THEOREM 2.4.6. *There is at most one completed subset S_0 of the structural equations (2.3).*

Suppose there are two different subsets S_1 and S_2 satisfying Definition 2.4.6.2. Obviously S_1 cannot be a subset of S_2 or vice versa, because then either S_1 or S_2 would not be a completed subset. On the other hand, if S_1 and S_2 have no equations in common, the combination $S_1 + S_2$ of both sets would possess an associated subset $R_{S_1+S_2} = R_{S_1} + R_{S_2}$ of restrictions (2.74) which is adequate in number and variety with respect to $S_1 + S_2$, contrary to the assumption made about S_1 . In the third possible case, in which S_1 and S_2 have a set S_1S_2 in common, which differs from both S_1 and S_2 , it can likewise be inferred that, contrary to the assumption regarding S_1 , the set $R_{S_1+S_2}$ of restrictions (2.74) associated with the combination $S_1 + S_2$ is of the requisite number and variety with respect to $S_1 + S_2$. This is seen by assigning the bilinear restrictions (2.74b) in $R_{S_1+S_2}$ (which contains the combination $R_{S_1} + R_{S_2}$ of R_{S_1} and R_{S_2}) in the following way. The bilinear restrictions (A) (see Fig. 2.4.6) in R_{S_1} are assigned in the same way as they were assigned to meet the requirements of Definition 2.4.5 with respect to S_1 . The equations of S_1 are thereby provided with an adequate number and variety of a priori restrictions. Then the bilinear restrictions (B) in R_{S_2} but not in R_{S_1} are assigned as they were to meet the requirements of Definition 2.4.5 with respect to S_2 . This adequately provides the equations in $S_3 = S_2 - S_1S_2$, that is, the equations in S_2 but not in S_1 , because the restrictions in R_{S_1} now excluded from consideration do not refer to these

equations. It follows that, irrespective of how any member of the third group of bilinear restrictions in $R_{S_1 + S_2}$, namely those neither in R_{S_1} nor in R_{S_2} , are assigned, the requirements of Definition 2.4.5 are always met with respect to $S_1 + S_2$. The assumption of two different completed subsets S_1 and S_2 has therewith been disproved.

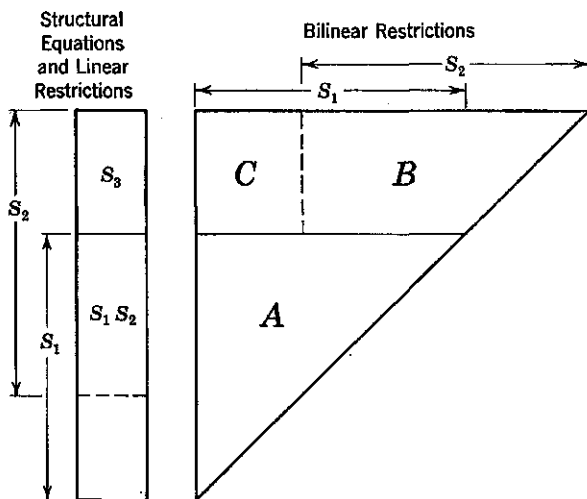


Figure 2.4.6

*2.4.7. *Construction of the completed subset.* It is of interest to give an example in which the completed subset can easily be constructed. This is the case in which the bilinear restrictions require Σ to be diagonal, while the structural equations can be ordered in such a way that there are $g-1$ homogeneous linear restrictions on the g th equation. Then the G th structural equation is subject to an adequate number of linear restrictions alone, and the $G-1$ bilinear restrictions

$$(2.80) \quad \sigma_{Gg} = 0, \quad g = 1, \dots, G-1,$$

referring to it can be assigned one by one to each of the first $G-1$ equations. This provides the $(G-1)$ th structural equation

with an adequate number of restrictions, and the $G-2$ remaining bilinear restrictions

$$(2.81) \quad \sigma_{G-1, g} = 0, \quad g = 1, \dots, G-2,$$

referring to it can now be assigned to the first $G-2$ equations. A repetition of this process shows that the completed subset S_0 contains all structural equations.¹

This example suggests a practical procedure for constructing the completed subset S_0 . First the set S_1 of those structural equations, for which the linear restrictions (2.74b) alone are adequate in number and variety, is included in S_0 . Then bilinear restrictions connecting the equations of S_1 with the remaining ones are assigned to equations outside S_1 , and the equations thereby provided with an adequate number of restrictions are included in S_0 . This process is repeated until it has become impossible to include in S_0 further equations one by one. Thereafter, it may still be possible to include small sets of three or more, counting also bilinear restrictions connecting the equations being included as a set.

*2.4.8. *Lack of sufficiency of criteria based on counting.* It has already been indicated that counting of restrictions alone is inconclusive in establishing identifiability. The condition that a structural equation belongs to the completed subset S_0 is neither necessary nor sufficient for its identifiability. Unfortunately, the formulation of necessary and sufficient conditions generalizing those established for linear restrictions is a task beset with considerable difficulties owing to the presence of nonlinear restrictions. Nevertheless, cases in which equations inside S are unidentifiable, or equations outside S_0 are identifiable, are "exceptional" in one sense or another, and we shall presently discuss the nature of the "exceptions."

¹An additional point of interest in this example is that the rows of Y can be obtained successively (from the bottom row up), each as the solution of a linear equation system. Therefore, if, under the a priori restrictions stated, the set of structural equations is completely identifiable at all, it is uniquely identifiable, in spite of the fact that among the restrictions imposed, $G(G-1)/2$ are bilinear.

A structural equation belonging to S_0 may fail to be identifiable (as was also found in the case of linear restrictions discussed earlier in this section) through a functional dependence of the relevant restrictions (2.74) on Y . Such functional dependence may come about because the parameters A, Σ happen to fall within a set of measure zero in the parameter space. Or it may even come about everywhere in the parameter space, as a result of the special nature of the a priori restrictions. Examples of the latter possibility were discussed in section 2.2.4. Another simple example is the case where the a priori restrictions are invariant for the interchange of two of the structural equations. These two equations are then inevitably unidentifiable, because the only case in which an interchange of two equations is innocuous, i.e., the case of complete equality of corresponding coefficients " α " and of corresponding covariances " σ " connecting the two equations with other equations of the system, is precluded by the assumed nonsingularity of B .

*2.4.9. *Lack of necessity of criteria based on counting.* A new element in the situation, which did not arise under linear restrictions only, is the fact that to belong to the completed subset is not even necessary for identifiability of a given structural equation. This is due to two restrictions on the parameter space which follow from the nature of the problem studied. The first of these is the restriction to real values of the parameters A, Σ . This restriction had no effect under linear a priori restrictions, since linear systems of equations with real coefficients only permit of real solutions. Quadratic or even bilinear systems possess no such property. Therefore, the possibility exists in the present case, that the a priori restrictions (including rules of normalization) confine Y to a point set in the complex space of all its elements, of which all *real* points are such that the g_0 th row of Y equals the corresponding row of the unit matrix - even though the g_0 th structural equation be outside of the completed subset S_0 . The condition that this shall occur is in the nature of a tangency condition.¹ We have not attempted either to prove (by the construction of an example) the possibility of such an occurrence under bilinear restrictions, or to prove its impossibility. One would expect such a tangency to occur only on a point set of measure zero in the parameter space, but the possibility cannot be excluded

¹An ellipsoid and a plane in three-dimensional space may have only one real point in common, although they represent only two inhomogeneous equations in three unknowns.

without proof that it could occur everywhere in the parameter space as a result of a special choice of a priori restrictions.

*2.4.10. *Restrictions requiring a certain partitioning of A and Σ .* The second relevant restriction on the parameter space is the nonsingularity of B, and therefore of Y. That this restriction may affect the identification problem appears from a constructed example which was kindly brought to our attention by A. Wald. The following formulation contains Wald's example as a special case.

Let the a priori restrictions be such that, if the structural equations are in a certain way exhaustively subdivided into two sets, S_I and S_{II} , and if at the same time the dependent variables are arranged in a certain order, the matrices B and Σ partition as follows:

$$(2.82) \quad B = \begin{bmatrix} B_{I I} & B_{I II} \\ 0 & B_{II II} \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{I I} & 0 \\ 0 & \Sigma_{II II} \end{bmatrix}.$$

In addition to the restrictions implied in (2.82), there may be further linear and bilinear restrictions involving the elements of $\Sigma_{I I}$, $\Sigma_{II II}$, $B_{I I}$, $B_{I II}$, $B_{II II}$ and the remaining elements of A.

In order that the transformation (2.16) shall preserve the partitioning (2.82), we must have

$$(2.83) \quad \begin{cases} (2.83\beta) & B_{II I}^{\oplus} = Y_{II I} B_{I I} = 0, \\ (2.83\sigma) & \Sigma_{II I}^{\oplus} = Y_{II I} \Sigma_{I I} Y_{I I} + Y_{II II} \Sigma_{II II} Y_{I II} = 0. \end{cases}$$

Now the nonsingularity of B requires that $B_{I I}$ be nonsingular. It follows from (2.83 β) that

$$(2.84) \quad Y_{II I} = 0.$$

Consequently, the nonsingularity of Y requires that $Y_{II II}$ be nonsingular. From this, (2.83 σ), (2.74), and the nonsingularity of $\Sigma_{II II}$, it also follows that

$$(2.85) \quad Y_{I II} = 0,$$

so that Υ partitions as follows

$$(2.86) \quad \Upsilon = \begin{bmatrix} \Upsilon_{\text{I I}} & 0 \\ 0 & \Upsilon_{\text{II II}} \end{bmatrix}.$$

This result means that, if no *further* bilinear restrictions connecting a structural equation of S_{I} with one of S_{II} are introduced, the identification problems of the two sets of structural equations have been separated. For the only transformations permitted by (2.86) are those within each set of equations.

Let there be G_{I} structural equations in S_{I} , G_{II} in S_{II} , with $G_{\text{I}} + G_{\text{II}} = G$. Then a counting criterion for the separate identifiability of the equations of S_{I} , or a subset thereof, on the basis of restrictions *additional* to (2.82) which refer exclusively to equations of S_{I} , can be formulated as follows: A completed subset $S_{\text{I}}^{(0)}$ of S_{I} is again defined by Definition 2.4.6.2, with this modification, that only restrictions additional to (2.82) are counted; and that their number and variety is deemed adequate if the requirements of Definition 2.4.5 are met with G_{I} substituted for G .

Now it is possible for $S_{\text{I}}^{(0)}$ to be nonempty, or even to contain all equations of S_{I} , even though the unmodified application of Definition 2.4.6.2 to the total set $S_{\text{I+II}}$ of G structural equations leads to an empty completed subset S_0 of $S_{\text{I+II}}$. This is seen most clearly in Wald's example, which takes $G_{\text{I}} = 1$, $G_{\text{II}} \geq 2$, and assumes no restrictions besides (2.82). Then $G_{\text{I}} - 1 = 0$, and $S_{\text{I}}^{(0)}$ consists of the one equation in S_{I} , whereas S_0 is empty. Similar examples with higher values of G_{I} can easily be constructed.

The modified counting criterion just indicated for the identifiability of structural equations in S_{I} , and a similar criterion with reference to S_{II} , can be subsumed under a more general criterion applicable to the full set S of structural equations. In this general criterion both the restrictions (2.82) and any additional restrictions are counted under the unmodified Definitions 2.4.5 and 2.4.6.2, but the $G_{\text{I}} G_{\text{II}}$ restrictions on Σ in (2.82) are to be counted as equivalent to the *linear* restrictions (2.85) on $\Upsilon_{\text{I II}}$, i.e., on the rows of Υ corresponding to the equations of

S_I , with G_{II} restrictions falling on each of the G_I equations of S_I . This procedure is justified because any attempt to construe the bilinear conditions (2.83 σ) arising from (2.82) as restrictions on the equations of S_{II} would require $Y_{I II}$ to have a positive rank, which was shown to be incompatible with the nonsingularity of Y . This general criterion can also be used if the additional restrictions contain further bilinear restrictions connecting an equation of S_I with one of S_{II} .

The foregoing discussion also applies in the more general case in which, instead of the partitioning (2.82) of B , we have (after some rearrangement of the variables x_k into two groups 1 and 2) a similar partitioning

$$(2.87) \quad A = \begin{bmatrix} A_{I1} & A_{I2} \\ 0 & A_{II2} \end{bmatrix}$$

of the rectangular matrix A , provided A_{I1} is square and nonsingular. The case (2.82) in which A_{I1} is a submatrix of B , however, is especially important, and will be studied further in section 3.2.7 in connection with estimation problems.

2.4.11. *Other special cases.* In special cases where the number of structural equations is moderate or the number of bilinear restrictions small, or both conditions hold, a more conclusive discussion of the identification problem than was given here for the general case, may be more easily possible. An example is the study of the measurement of production functions [Marschak and Andrews, 1944] already referred to.

Another example may be briefly indicated without attempting rigorous statement. This is the case in which no structural equation is subject to fewer than $G-2$ homogeneous linear restrictions, while for the set S of those equations that do not possess at least the adequate number $G-1$ of such restrictions, the deficiency is just made up by the bilinear restrictions between them. Further analysis then needs to be concerned only with the equations of S , and with the matrix Y^S containing the corresponding rows of Y . The elements of any row of Y^S can now be expressed as a linear function of two of them. Writing θ_g for the ratio of those elements in the row v_g , we easily see that the bilinear re-

lations take the form

$$(2.88) \quad \kappa_{gh} \theta_g \theta_h + \lambda_{gh} \theta_g + \mu_{gh} \theta_h + \nu_{gh} = 0.$$

Since by assumption the number and variety of linear and bilinear restrictions on the equations of S is just adequate, each variable θ_h enters into two restrictions (2.88) together with θ_g and θ_i respectively, say. The elimination of θ_h from these two restrictions leads to the same type of restriction between θ_g and θ_i . Continuation of this process of elimination must finally lead to a restriction

$$(2.89) \quad \kappa_g \theta_g^2 + \lambda_g \theta_g + \nu_g = 0$$

connecting θ_g with itself. If all bilinear restrictions have been eliminated in this process, there is unique, multiple, or incomplete identification, according as (2.89) has one, two, or infinitely many solutions. If there are one or more other sets of bilinear restrictions connecting other subsets of the rows of Υ^S , these must be investigated in the same manner, until all bilinear restrictions have been accounted for.

2.5. *Incompleteness of the Present Discussion of Identification Problems*

2.5.1. *Dummy restrictions if some a priori restrictions are bilinear.* It has already been pointed out that the present approach has not led to necessary and sufficient conditions for identifiability in the general case of linear and bilinear restrictions, and is not likely to do so without considerable further study. For that reason, no study has been made of the degree of indeterminacy of a priori unidentifiable structural equations with the help of dummy restrictions.

2.5.2. *Wald's criterion for identifiability.* In [III], by a different approach, Wald obtains a criterion in terms of ranks of matrices, which is both necessary and sufficient, applies to each parameter separately rather than to all parameters of a structural equation as a group, and permits a much more general class of a priori restrictions. Against these advantages must be set the fact

that the matrices whose rank must be examined generally have a much higher order, equal to or exceeding the square of the number K_x of variables in the structural equations¹ (2.3).

2.5.3. *Other indications of possibly incomplete identification.* Both Wald's criterion, and the criteria developed in this section, are inevitably stated in terms of the unknown values of the parameters. It has already been indicated that there are exceptional point sets (of measure zero) in the parameter space, in which the matrices involved in the criteria suffer a decrease in rank, and in which therefore the parameters are subject to a greater degree of indeterminacy than was already recognized by the general analysis of the identification problem. The practical question then arises, whether a parameter point within the exceptional set could have produced the actual sample of observations with any degree of likelihood. Fortunately, there are further indications of such an occurrence: the maximum of the likelihood function must then be very "flat" with respect to one or more particular permissible directions in the parameter space - permissible as regards the a priori restrictions. In the most extreme case the maximum is completely "flat," i.e., it is reached along a curve, surface, etc., rather than in a point. Such situations reveal themselves 1) through slow convergence of the iterative computation procedure (in the extreme case through more rapid convergence to a solution which depends on the initial values used, section 4.3.3.11), and 2) through very high (in the extreme case, infinite) values of the estimated sampling variances of parameters erroneously believed determinate. In this way deficiencies in the analysis of identification problems will come to light in later stages of the investigation. It will be clear, however, that the computational stage can be handled much more efficiently, if all indeterminacies in the parameter space have already been recognized through the study of identification problems.

2.5.4. *Identification should be based on a minimum of firmly established assumptions.* It is therefore important to make the prior analysis of identification problems as complete and general as possible. In particular, one should avoid as much as possible employing assumptions that might not be satisfied by the data, and which are at the same time essential to the conclusions reached

¹In the equations (2.3), it will be remembered, values of the same economic variable measured with different time lags are to be considered as different variables.

regarding the identifiability of structural equations. For this reason, the present discussion of identification problems has been made independent of the assumption of normality of the distribution of disturbances. This is in contrast with those parts in subsequent sections, especially the evaluation of sampling variances of maximum-likelihood estimates, in which the normality assumption was already known to be relatively harmless, even if the data do not strictly correspond to it.

For the same reason, the present authors are inclined to rely more firmly and more extensively on restrictions involving the coefficients A that have a good basis in economic considerations, than on restrictions on the covariance matrix Σ of disturbances, at least until the nature of the disturbances has been more fully analyzed by theory and observation.

2.5.5. *Linearity of the structural equations.* The question should be raised whether the assumption that the structural equations are linear does not conceal from view possible further cases of indeterminacy in the measurement of economic relations that need not be strictly linear. To answer this question it is necessary to formulate what is the alternative to linearity. If the alternative is the addition of higher-degree terms in the variables to obtain polynomial expansions, the answer is that the present analysis can be extended to cover such cases as follows. First the present analysis is applied to the equation system obtained by omitting all nonlinear terms to find for any point in the restricted parameter space a set T_L of transformations preserving the a priori restrictions on the linear terms. As long as the vanishing of all nonlinear terms is not excluded a priori, the set T preserving, everywhere in the relevant part of the parameter space, any a priori restrictions on the nonlinear terms in addition to those on the linear terms, can only be a subset of T_L . Because of the algebraic independence of terms of different degrees under linear transformations, T_L can thus be narrowed down to T by successively applying the restrictions, if any, on the terms of each of the higher degrees. This reasoning indicates that the admission of nonlinear terms does not lead to new indeterminacies unsuspected in the linear case.

Another possible situation is that in which the a priori restrictions prescribe linearity for some structural equations, and for some other equations, a type of relationship that excludes linearity (e.g., hyperbolic or exponential). While the general mathe-

mathematical treatment of identification problems in such cases might be more difficult, we conjecture that again the set T would in some sense be narrower than in the corresponding case where all equations are linear. For this reason, we believe, "mixed" prescriptions of this type, as regards the form of the equations, are more likely to conceal than to reveal cases of indeterminacy of economic parameters, except where indubitable a priori evidence exists as regards the validity of such prescriptions.

2.5.6. *Transformations in the parameter space involving shifts in time.* It should be pointed out that among the assumptions on which the present discussion of identification problems is based, there is still at least one of the undesirable type against which we have just put in a word of warning. That is, there is one assumption that may not be too well fulfilled by the data, whereas its removal may open up new possibilities of indeterminacy. This is the assumption of independence between disturbances in successive time units.

The proof of Theorem 2.1.2 is based on that assumption. One example is sufficient to show that this basic theorem does not hold without the independence assumption. Suppose that we admit serial correlation between disturbances relating to successive time points, but do not think it justified to impose any particular mathematical form on the autocorrelation function¹. Consider the system of two equations

$$(2.90) \quad \begin{array}{rcl} \alpha_{110} x_1(t) & + \alpha_{111} x_1(t-1) + \alpha_{121} x_2(t-1) & = u_1(t), \\ \alpha_{210} x_1(t) + \alpha_{220} x_2(t) & & = u_2(t), \end{array}$$

in which the open spaces indicate the coefficients prescribed to be zero. Each of these equations is identifiable under the assumptions of Theorem 2.1.2. But under the present assumptions, the transformation

$$(2.91) \quad \begin{array}{l} u_1^{\oplus}(t) = u_1(t) + \lambda u_2(t-1), \\ u_2^{\oplus}(t) = u_2(t), \end{array}$$

¹Hurwicz demonstrates in [XI-10.2] that certain specific assumptions regarding the form of the autocorrelation functions restore identifiability.

preserves the form of the equations (2.90), and the assumed form of the distribution of the disturbances, which now permits correlation between $u_1^\oplus(t)$ and $u_2^\oplus(t-1)$. The transformation (2.91) affects the coefficients of the first equation according to

$$(2.92) \quad \alpha_{110}^\oplus = \alpha_{110}, \quad \alpha_{1g1}^\oplus = \alpha_{1g1} + \lambda \alpha_{2g0}, \quad g = 1, 2.$$

The first equation (2.90) has thus ceased to be identifiable.

The transformation (2.91) permits one of the structural equations to be shifted in its timing before it is linearly combined with another equation. If such transformations are permissible, the study of identification problems is greatly complicated even if the a priori restrictions are linear. It is argued in [XVI] that these problems can perhaps be studied more fruitfully if at the same time the time variable is made continuous rather than discrete.

3. ESTIMATION OF THE PARAMETERS

3.1. Properties of the Unrestricted Likelihood Function

3.1.1. Maximum-likelihood estimation using all a priori restrictions. We now turn to the problem of estimating the parameters $\beta_{gt\tau}$, $\gamma_{gk\tau}$, σ_{gh} of the distribution function (2.8). It is assumed that the study of identification problems has shown whether or not the various structural equations on which this distribution is built are uniquely identified by the a priori restrictions. It is further assumed that this analysis has revealed the extent and nature of the indeterminacy in the parameters of those equations that are not uniquely identified.

We shall now make a more restrictive assumption on the nature of the distribution function $f(u)$ of the disturbances, at least for the purpose of constructing estimates of the parameters:

ASSUMPTION 3.1.1. *The disturbances u_g have a joint normal distribution function with nonsingular covariance matrix $\Sigma \equiv [\sigma_{gh}] \equiv [\sigma^{gh}]^{-1}$,*

$$(3.1) \quad f(u) = (2\pi)^{-\frac{1}{2}G} \det^{-\frac{1}{2}} \Sigma \exp -\frac{1}{2} \sum_{g,h=1}^G u_g \sigma^{gh} u_h .$$

We shall use as estimates those functions of the observations which for this choice of $f(u)$ constitute *maximum-likelihood estimates* of the parameters. If (3.1) is inserted in (2.8), the probability density (2.8) in any particular sample point, i.e., for any particular set of observations $y(t), z(t), t = 1, \dots, T$, is a function of the parameters B, Γ, Σ , known as the likelihood function. The maximum-likelihood estimates here considered are those values

$$(3.2) \quad B, C, S$$

of the parameters for which, subject to all the a priori restrictions, the likelihood function reaches its highest maximum. Following Mann and Wald, the properties of these estimates can then be studied both under the same normality assumption for the distribution of the disturbances, and under some less restrictive assumption.

3.1.2. Maximum-likelihood estimates under partial disregard of a priori information. T. W. Anderson and Rubin have indicated¹ other estimates based on a suggestion of M. A. Girshick. These estimates are obtained by mathematically simpler, and in most cases less laborious, computational methods. These estimates are maximum-likelihood estimates under disregard of a suitably chosen part of the a priori information available. The simplification of computational problems is obtained at a cost of increased sampling variances of the estimates (reduced efficiency of the method of estimation). Further comparison with this elegant method, called the "reduced-form method" will be made in section 3.2.1. In the remainder of this article, the term "maximum-likelihood estimates" will be used for such estimates obtained with the aid of all a priori information available. Where a distinctive expression is needed, the term "information-preserving maximum-likelihood method" will be used for the method of estimation here applied.

3.1.3. Classification of the variables. For most of the present section, the relevant distinction is that between "jointly dependent" and "predetermined" variables, made in the introduction. The importance of this distinction is based on the fact that the coefficients of the jointly dependent variables enter the Jacobian (2.2) of the transformation (2.1), whereas those of the predetermined variables do not. In the equations defining the maximum-likelihood estimates, and in the formulae for their estimated as-

¹[1949] and unpublished manuscript. See also [IX].

ymptotic sampling variances and covariances, the position of the jointly dependent variables has similarities with that of the one dependent variable in the single-equation least-squares method.

On the contrary, in these equations and formulae the predetermined variables occur without any distinction as to whether they are exogenous variables, or lagged values of endogenous variables. The latter distinction is relevant in the present context only in one instance: in the proof of consistency of the maximum-likelihood estimates. That the distinction is irrelevant elsewhere, is, of course, connected with the fact that the present study is confined to large-sample approximations.

3.1.4. *Notation.* We shall therefore continue to use the notation introduced in section 2.1.2. We restate the partitioning of coefficients and variables

$$(3.3) \quad A \equiv [B \quad \Gamma], \quad x \equiv [y \quad z],$$

and the equation system (1.1) in this notation,

$$(3.4) \quad Ax'(t) = By'(t) + \Gamma z'(t) = u'(t),$$

where $x'(t)$ is the transpose of x , and

$$(3.5) \quad u(t) = [u_1(t) \quad \cdots \quad u_G(t)], \quad G \equiv K_y.$$

We shall more fully use the symmetric notation $K_y \equiv G$, K_z , and $K_x \equiv K_y + K_z$ for the number of jointly dependent, of predetermined, and of all variables respectively.

If (3.1) is substituted in the likelihood function (2.8) and logarithms are taken, we obtain in the new notation¹

$$(3.6) \quad \begin{aligned} \frac{1}{T} \log F \equiv L \equiv L(A, \Sigma) = & -\frac{1}{2} K_y \log 2\pi + \log \det B \\ & -\frac{1}{2} \log \det \Sigma - \frac{1}{2} \text{tr} \Sigma^{-1} A M_{xx} A'. \end{aligned}$$

Here M_{xx} is the observed symmetric "moment" matrix

¹tr X, the trace of a square matrix X, denotes the sum of all diagonal elements of X.

$$(3.7) \quad M_{xx} \equiv \frac{1}{T} \left[\sum_{t=1}^T x_k(t) x_l(t) \right] = \frac{1}{T} \sum_{t=1}^T x'(t) x(t),$$

$k, l = 1, \dots, K_x,$

of all variables $x_k(t)$, and partitions according to

$$(3.8) \quad M_{xx} = \begin{bmatrix} M_{yy} & M_{yz} \\ M_{zy} & M_{zz} \end{bmatrix}.$$

3.1.5. Positive definiteness of M_{xx} . In what follows we shall assume that the sample obtained is one of those, occurring with probability one, for which M_{xx} is nonsingular and therefore positive definite. (A symmetric matrix M_{xx} is called positive definite if $a M_{xx} a' > 0$ for every nonvanishing vector (one-row matrix) a . A nonsingular moment matrix is positive definite because $a M_{xx} a' = \frac{1}{T} \sum_t a x'(t) x(t) a'$ is a sum of squares of the vector products $a x'(t)$, $t = 1, \dots, T$, while $a M_{xx} a' = 0$ for some nonvanishing a would entail the singularity of M_{xx} .) It follows that M_{yy} and M_{zz} have ranks equal to their respective orders K_y and K_z .

3.1.6. The reduced form of the structural equations and its parameters. We shall first study the maximum properties of L without imposing any a priori conditions on the parameters A, Σ . From Theorem 2.1.3.5, we know that any maximum properties of the likelihood function should be preserved by a transformation of the type (2.18), which we rewrite in the new notation

$$(3.9) \quad A^\oplus = \Upsilon A, \quad \Sigma^\oplus = \Upsilon \Sigma \Upsilon'.$$

A unique representation of each set of mutually equivalent points (see Definition 2.1.3.6) in the unrestricted parameter space is obtained by writing $\Upsilon = B^{-1}$ in (3.9), which makes

$$(3.10) \quad B^{\oplus} = I,$$

if I denotes the unit matrix of order K_y . The fact that each equivalent point set contains one and only one point satisfying (3.10) corresponds to the fact that each system (3.4) can be written in one and only one way in what has been called the *reduced form* [Mann and Wald, p. 201, equation (85)]:

$$(3.11) \quad y_i(t) = \sum_{k=K_y+1}^{K_x} \pi_{ik} z_{k-K_y}(t) + v_i(t), \quad i = 1, \dots, K_y,$$

or

$$y'(t) - \Pi z'(t) = v'(t),$$

in which each equation contains only one of the independent variables, $y_k(t)$, $k = 1, \dots, K_y$, with unity as coefficient. Here we have written $-\Pi$ for the value of A^{\oplus} in (3.9) corresponding to (3.10), and we shall write Ω for the corresponding value of Σ^{\oplus} ,

$$(3.12) \quad \omega_{kl} = \mathcal{E} v_k(t) \cdot v_l(t), \quad k, l = 1, \dots, K_y,$$

or

$$\Omega = \mathcal{E} v'_y(t) \cdot v_y(t).$$

It follows that Ω , and therefore also Ω^{-1} , are symmetric and positive definite, since B is nonsingular. In this notation, (3.9) becomes

$$(3.13) \quad [I \quad -\Pi] = B^{-1} A, \quad \Omega^{-1} = B' \Sigma^{-1} B.$$

Since

$$(3.14) \quad \log \det B - \frac{1}{2} \log \det \Sigma = \frac{1}{2} \log \det (B' \Sigma^{-1} B),$$

the function L in (3.6) can now be written as

$$(3.15) \quad L = L(-\Pi, \Omega) = -\frac{1}{2} G \log 2\pi + \frac{1}{2} \log \det \Omega^{-1} \\ - \frac{1}{2} \text{tr} \{ \Omega^{-1} (M_{yy} - M_{yz} \Pi' - \Pi M_{zy} + \Pi M_{zz} \Pi') \}.$$

Because of the uniqueness of the reduced form, there is no transformation in the (Π, Ω) space, other than the identical transformation, which preserves the form of (3.15).

3.1.7. *Rules for the differentiation of functions of matrices.* Before proceeding to maximize (3.15) it will be useful to state a few rules regarding the differentiation of a matrix $X(\xi) = [x_{mn}(\xi)]$ with respect to a scalar parameter ξ . If X is square or rectangular, and if Y is a constant matrix with the same number of rows and columns respectively, we have, because of the linearity of the trace operation,

$$(3.16) \quad \frac{d}{d\xi} \operatorname{tr}(XY') = \operatorname{tr}\left(\frac{dX}{d\xi} Y'\right).$$

If X is square, let X^{mn} denote the cofactor of x_{mn} , such that the typical element of the inverse X^{-1} of X is $x^{mn} = X^{nm} / \det X$. Then we have

$$(3.17) \quad \frac{d}{d\xi} \log \det X = \frac{\frac{d}{d\xi} \det X}{\det X} = \frac{\sum_{m,n} X^{mn} \frac{dx_{mn}}{d\xi}}{\det X} = \operatorname{tr} X^{-1} \frac{dX}{d\xi}.$$

An expression for $\frac{dX^{-1}}{d\xi}$ is derived as follows:

$$(3.18) \quad XX^{-1} = I, \quad \frac{dX}{d\xi} X^{-1} + X \frac{dX^{-1}}{d\xi} = 0, \quad \frac{dX^{-1}}{d\xi} = -X^{-1} \frac{dX}{d\xi} X^{-1}.$$

Therefore, if Y is also square, of the same order, and constant,

$$(3.19) \quad \frac{d}{d\xi} \operatorname{tr}(X^{-1} Y') = -\operatorname{tr}\left(X^{-1} \frac{dX}{d\xi} X^{-1} Y'\right).$$

3.1.8. *Maximizing the likelihood function with respect to Π .* The study of the maximum properties of $L(-\Pi, \Omega)$ is facilitated by

maximizing L in two successive steps as follows: First we consider Ω as a given matrix, and determine the value P of Π at which the quadratic form $L(-\Pi, \Omega)$ in the elements of Π has a maximum (for variations of Π only).

Writing

$$(3.20) \quad \Pi = P + \varepsilon \underset{\rightarrow}{P}$$

in (3.15), we have

$$(3.21) \quad \begin{aligned} \frac{dL(-\Pi, \Omega)}{d\varepsilon} &= -\frac{1}{2} \text{tr} \{ \Omega^{-1} (M_{yz} \underset{\rightarrow}{P}' + \underset{\rightarrow}{P} M_{zy} - \Pi M_{zz} \underset{\rightarrow}{P}' - \underset{\rightarrow}{P} M_{zz} \Pi') \} \\ &= -\text{tr} \{ \Omega^{-1} (M_{yz} - \Pi M_{zz}) \underset{\rightarrow}{P}' \} \end{aligned}$$

in connection with the symmetry¹ of M_{xx} and Ω^{-1} . Furthermore

$$(3.22) \quad \frac{d^2 L}{d\varepsilon^2} = -\text{tr} (\Omega^{-1} \underset{\rightarrow}{P} M_{zz} \underset{\rightarrow}{P}')$$

identically in ε , and in particular for $\varepsilon = 0$.

A necessary condition for a maximum of $L(-\Pi, \Omega)$ in $\Pi = P$ is that $\left(\frac{dL}{d\varepsilon} \right)_{\varepsilon=0}$ shall vanish for all possible values of $\underset{\rightarrow}{P}$. It is seen from (3.21) that this requires

$$(3.23) \quad M_{yz} - P M_{zz} = 0, \quad \text{or} \quad P = M_{yz} M_{zz}^{-1},$$

using the nonsingularity of M_{zz} . There is therefore only one extremum of L , which is reached in a point $\Pi_{yz} = \underset{\rightarrow}{P}_{yz}$ which proves to be independent of Ω . Moreover, this extremum is actually a maximum (and since we are dealing with a quadratic function, an absolute maximum) because (3.22) is a negative definite quadratic form in the elements of $\underset{\rightarrow}{P}$. For the matrices Ω^{-1} and M_{zz} , being positive definite, can be decomposed (see [9], p. 246) according to

$$(3.24) \quad \Omega^{-1} = \Psi' \Psi, \quad M_{zz} = R R',$$

¹Use has been made of the properties $\text{tr} X Y' = \text{tr} Y' X = \text{tr} Y X'$ which follows directly from the definition of the trace.

where Ψ and R are real. Thereby (3.22) turns into the negative sum of squares¹

$$(3.25) \quad \frac{d^2 L}{d\varepsilon^2} = - \operatorname{tr}\{(\Psi \vec{P} R)(\Psi \vec{P} R)'\}.$$

The reader will have noticed that the elements

$$(3.26) \quad \hat{\pi}_{ik} = \sum_{l=K_y+1}^{K_x} m_{il} m_{(zz)}^{lk}, \quad \text{with} \quad [m_{(zz)}^{lk}] = M_{zz}^{-1},$$

$$l, k = K_y + 1, \dots, K_x,$$

of the i th row of the matrix P represent the coefficients of the elementary regression of the dependent variable $y_i(t)$, $1 \leq i \leq K_y$, on the predetermined variables $z_k(t)$, $k = K_y + 1, \dots, K_x$, i.e., the coefficients estimated by the single-equation least-squares method.

3.1.9. *Maximizing the likelihood function with respect to Ω .* The second step is to insert (3.23) in the expression (3.15) for L , which on account of the symmetry of M_{zz} becomes

$$(3.27) \quad L(\Omega) \equiv -\frac{1}{2} K_y \log 2\pi + \frac{1}{2} \log \det \Omega^{-1} - \frac{1}{2} \operatorname{tr}(\Omega^{-1} \cdot {}^z M_{yy}),$$

where

$$(3.28) \quad {}^z M_{yy} \equiv M_{yy} - M_{yz} \cdot M_{zz}^{-1} \cdot M_{zy}$$

is again positive definite, because it is the moment matrix of the

"residuals" $v_i^\oplus(t) \equiv y_i(t) - \sum_{l=1}^{K_x} \hat{\pi}_{i, K_y+l} z_l(t)$, $i = 1, \dots, K_y$,

from the elementary regressions of each of the dependent variables separately, on all predetermined variables. We shall now maximize (3.27) with respect to the variations of Ω or of Ω^{-1} . If we write

$$(3.29) \quad \Omega^{-1} = V + \eta \vec{V}, \quad V \equiv W^{-1}, \quad \vec{V} = \vec{V}',$$

¹Use is made of the properties $(XY)' = YX'$ and $\operatorname{tr}(XY) = \sum_{n, n} (x_{nn})^2$.

the assumed symmetry of \vec{V} and V and the positive definiteness of V will ensure positive definiteness of Ω^{-1} in an η -neighborhood of $\eta=0$ for any particular value of \vec{V} . For, if ω_1^{-1} and ω_1^{-1} represent the absolute minima of $v \Omega^{-1} v'$ and $v V v' \equiv v W^{-1} v'$ respectively, under the restriction $v v' = 1$, and $|\eta_1|$ the similarly restricted absolute maximum of $|v \vec{V} v'|$, we have $\omega_1^{-1} > 0$, $\omega_1^{-1} \geq \omega_1^{-1} - |\eta| |\eta_1| > 0$ for sufficiently small values of $|\eta|$. Using (3.16) and (3.17), we have from (3.27)

$$(3.30) \quad \frac{dL(\Omega)}{d\eta} = \frac{1}{2} \operatorname{tr}(\Omega \vec{V}) - \frac{1}{2} \operatorname{tr}(\vec{V} \cdot {}^z M_{yy}) = \frac{1}{2} \operatorname{tr}\{(\Omega - {}^z M_{yy}) \vec{V}\},$$

and, using (3.19),

$$(3.31) \quad \left(\frac{d^2 L}{d\eta^2} \right)_{\eta=0} = -\frac{1}{2} \operatorname{tr}(W \vec{V} W \vec{V}).$$

From (3.30) we see that $L(\Omega)$ has one stationary value, which is reached if Ω equals

$$(3.32) \quad W \equiv {}^z M_{yy}.$$

This stationary value is a maximum because the quadratic form (3.31) can be shown to be negative definite in the elements of \vec{V} , by an argument similar to that used in the case of (3.22), and using the symmetry of \vec{V} .

There are various ways of proving that (3.32) indicates the absolute maximum of $L(\Omega)$ in the space of symmetric and positive definite matrices Ω^{-1} . Perhaps the most elementary proof is as follows: If $\Omega^{-1}(1)$ is a matrix in that space different from W^{-1} , the matrix

$$(3.33) \quad \Omega^{-1}(\theta) \equiv \theta \Omega^{-1}(1) + (1 - \theta)W^{-1}, \quad 0 \leq \theta \leq 1,$$

is easily shown to belong to the same space. The function

$$(3.34) \quad \bar{L}(\theta) \equiv L\{\Omega^{-1}(\theta)\}$$

possesses continuous first and second derivatives with respect to θ for $0 \leq \theta \leq 1$. These derivatives satisfy the two conditions

$$(3.35) \quad \left(\frac{d\bar{L}(\theta)}{d\theta} \right)_{\theta=0} = 0, \quad \frac{d^2\bar{L}(\theta)}{d\theta^2} < 0 \quad \text{for} \quad 0 \leq \theta \leq 1.$$

The first condition is satisfied because a stationary value $L(\Omega)$ is reached for $\Omega^1 = \Omega^1(0)$. The second condition is satisfied because the negative definiteness of (3.31) is not dependent on W satisfying the maximum conditions (3.32). It follows from (3.35) by use of Taylor's theorem, that

$$(3.36) \quad L\{\Omega(1)\} = \bar{L}(1) = \bar{L}(0) + \frac{1}{2} \left(\frac{d^2\bar{L}(\theta)}{d\theta^2} \right)_{\theta=\theta'} < \bar{L}(0) = L(W^{-1}),$$

where θ' represents some number between 0 and 1.

3.1.10. The absolute maximum of the likelihood function. By inverting the transformation (3.13) we can summarize the maximum properties of the likelihood function in the following

THEOREM 3.1.10. *In the absence of any a priori restrictions the logarithmic likelihood function (3.6) has one and only one maximum value*

$$(3.37) \quad L_{\max} = -\frac{1}{2} K_y (1 + \log 2\pi) - \frac{1}{2} \log \det(M_{yy} - M_{yz} M_{zz}^{-1} M_{zy}),$$

which is an absolute maximum. This maximum is reached in each point

$$(3.38) \quad \begin{aligned} B &= \text{any nonsingular square matrix of order } K_y, \\ \Gamma &= -B P, \\ \Sigma &= B W B, \end{aligned}$$

of the set of points equivalent to the point

$$(3.39) \quad B = I, \quad \Gamma = -P = -M_{yz} M_{zz}^{-1}, \quad \Sigma = W = M_{yy} - M_{yz} M_{zz}^{-1} M_{zy}.$$

This theorem establishes the uniqueness of the maximum of the likelihood function in the unrestricted parameter space, in the sense that there is one and only one set of mutually equivalent points on which the maximum is reached.

We add an expression for the likelihood function that is derived from (3.15) with the help of (3.39)

(3.40)

$$L = -\frac{1}{2}K \log 2\pi + \frac{1}{2} \log \det \Omega^{-1} - \frac{1}{2} \text{tr} \left(\Omega^{-1} \{ W + (P - \Pi) M_{zz} (P - \Pi)' \} \right),$$

and that brings out clearly the significance of the statistics P and W established by Theorem 3.1.10.

3.2. Properties of the Restricted Likelihood Function

3.2.1. *The case in which the restricted likelihood function can attain its absolute maximum.* In the case where a priori restrictions are introduced, a somewhat weaker theorem can be formulated as long as the a priori restrictions do not prevent the likelihood function from attaining its absolute maximum (3.37).

THEOREM 3.2.1. *Under a priori restrictions of any kind that permit the likelihood function to attain its absolute maximum in some point (A, Σ) , this maximum is attained only in all points $(A^{\oplus}, \Sigma^{\oplus})$ of the restricted parameter space that are equivalent to the point (A, Σ) ,*

This theorem follows immediately from Theorem 3.1.10 and Definition 2.1.5.1. It should be noted that Theorem 3.2.1 does not preclude the existence of one or more relative maxima where the likelihood function attains a value lower than (3.37).

The question of whether or not the a priori restrictions permit the likelihood function to attain its absolute maximum is important for two reasons. In the first place this question is connected with the relations between the reduced-form method¹ based essentially on single-equation least-squares procedures, and the maximum-likelihood method preserving all a priori information, as applied in this article. According to Theorem 3.1.10, as long as the absolute maximum of the likelihood function can be reached, the information-preserving maximum-likelihood method of estimation is mathematically equivalent to the single-equation least-squares

¹See section 3.1.2 and also [IX].

method applied to each equation of the reduced form (3.11). For the respective rows of P in (3.39) are identical with the estimates obtained for the coefficients of the corresponding equations (3.11) by the latter method. After P and W have been determined from (3.39), it is then possible to determine the transformation (3.38) so as to satisfy the a priori restrictions.

The second reason is connected with the computation of maximum-likelihood estimates, and is a consequence of the first reason. In case L attains the value L_{\max} , the procedure just described always leads to the absolute maximum of the likelihood function. In case L cannot attain L_{\max} , the maximum-likelihood equations are essentially nonlinear, and the only practicable methods of computation available are iterative methods. So far we do not know with certainty under what conditions each of these methods converges to the absolute maximum. One may possibly be led to a relative maximum, depending on the initial values chosen at the start of the iterative procedure, and the particular method of iteration used. As far as our present results reach, therefore, the case where L cannot attain the value L_{\max} is subject to an uncertainty which is absent when L_{\max} can be attained.

3.2.2. *Attainability of the absolute maximum under linear and bilinear restrictions.* For these reasons, to which another will be added in section 4.3.3.4, it is important to know under which conditions L_{\max} can be attained, that is, under which conditions (3.38) is compatible with the a priori restrictions. If the latter consist of the linear restrictions (2.24) combined with the bilinear restrictions (2.73), this question must be answered from an equation system obtained by inserting (3.38) in these restrictions:

$$(3.41lh) \quad \beta(g)[-I \ P] \Phi'_g = 0, \quad g = 1, \dots, K_y,$$

$$(3.41ln) \quad \beta(g)[-I \ P] \iota'(i_g) = 1,$$

$$(3.41b\alpha) \quad \left[\begin{array}{cc} \beta(g_r)[-I \ P] \iota'(k_r) & \beta(g_r)[-I \ P] \iota'(l_r) \\ \beta(h_r)[-I \ P] \iota'(k_r) & \beta(h_r)[-I \ P] \iota'(l_r) \end{array} \right] = 0, \quad r = 1, \dots, R_\alpha^{(2)},$$

$$(3.41b\sigma) \quad \beta(g_r) W \beta'(h_r) = 0, \quad r = R_\alpha^{(2)} + 1, \dots, R_\alpha^{(2)} + R_\sigma,$$

where $\iota(k)$ is again the k th row of the unit matrix of order K_x . We shall, for convenience, refer to these equations as follows:

$$(3.41) \quad \left\{ \begin{array}{l} (3.41l) \left\{ \begin{array}{l} (3.41lh) \\ (3.41ln) \end{array} \right. \\ (3.41b) \left\{ \begin{array}{l} (3.41b\alpha) \\ (3.41b\sigma) \end{array} \right. \end{array} \right.$$

The matrix $[-I \quad P]$ in (3.41) is put together in analogy to

$$(3.42) \quad [-I \quad \Pi] = -B^{-1}[B \quad \Gamma] = -B^{-1}A,$$

as defined in (3.13).

The equations (3.41) are similar in form to the equations (2.74), and the present problem is therefore closely related to the identification problem. Nevertheless, there are two important differences in the two problems, one in the assumptions, and one in the question to be answered. In the identification problem, it is known by assumption that the equations (2.74) have at least one real solution $\Upsilon = I$, and the question to be answered is under what conditions there is only one, or a finite number of real solutions. In the present problem it is not known whether there is at all a real solution B to the equations (3.41), and the question is under what conditions there is at least one solution. To obtain an answer to the present question, the counting of equations and variables is even less conclusive than in the identification problem. For even in a case in which, on the basis of counting, the number of real or complex solutions B is believed to be finite almost everywhere in the space of P and W , we cannot without further analysis say that the part of the sample space in which all solutions are complex is of measure zero.

We have so far not succeeded in finding general conditions for the existence of at least one real solution. However, the iterative computation procedures for solution of the maximum-likelihood

equations to be described in section 4 lead to such solutions if they exist, provided the computation is started with suitable initial values - although again we do not know precisely which initial values are suitable.

3.2.3. *Attainability of the absolute maximum under linear restrictions only.* It is not difficult to state exact conditions for the attainability of the absolute maximum of L in the case where only linear a priori restrictions of the type (2.24) are introduced. This leads to the conditions (3.41lh) which we wish to be satisfied by at least one real solution B .

THEOREM 3.2.3.1. *A necessary condition for the attainability of the absolute maximum of the likelihood function under the homogeneous linear a priori restrictions (2.24) is that a) none of the matrices $P\Phi'_g$, $g = 1, \dots, K_y$, has a rank exceeding $K_y - 1$. A necessary and sufficient condition is that, in addition to a), b) the consequently nonempty set of solutions B of (3.41lh) contains at least one nonsingular solution ($\det B \neq 0$).*

THEOREM 3.2.3.2. *A necessary condition for the attainability, almost everywhere in the sample space, of the absolute maximum of the likelihood function under the homogeneous linear a priori restrictions (2.24) is that none of the matrices Φ_g , $g = 1, \dots, K_y$, has a rank exceeding $K_y - 1$.*

Theorem 3.2.3.1 follows directly from the conditions for the existence of a solution of a homogeneous system of linear equations. Theorem 3.2.3.2 follows because the condition that the rank of $P\Phi'_g$ shall fall below K_y if the rank of Φ'_g is at least K_y entails a restriction on P satisfied only on a point set of measure zero in the sample space.

3.2.4. *Connections between attainability of the absolute maximum of the likelihood function and identifiability of structural equations.* It is of interest to compare Theorem 3.2.3.2 with Theorem 2.2.2 and its corollary. The latter states that identifiability of the g th structural equation requires the rank of Φ_g to be at least $K_y - 1$. Theorem 3.2.3.2 states that attainability of the absolute maximum requires that rank to be at most $K_y - 1$ (for all values of G). Thus, if independent linear restrictions on the coefficients of the g th equation are added one by one (beginning in a situation where L_{\max} is attainable), the point at which

in general complete identification of the g th structural equation is attained almost everywhere in the parameter space, is at the same time the point beyond which no further restrictions referring to the equation can be added without preventing L_{\max} from being attainable almost everywhere in the sample space.

A similar situation is found under quite general a priori restrictions, which we shall denote

$$(3.43) \quad \varphi_r(A, \Sigma) = 0, \quad r = 1, \dots, R.$$

We shall assume these restrictions to be independent¹, compatible, and to imply normalization of all structural equations. We shall further assume that the functions φ_r possess continuous derivatives with respect to the parameters A, Σ .

DEFINITION 3.2.4.1. *By the restricted reduced parameter space we understand the space of the parameters*

$$(3.44) \quad \Pi = -B^{-1} \Gamma, \quad \Omega = B^{-1} \Sigma B'^{-1}, \quad \det \Omega \neq 0,$$

subject to such restrictions,

$$(3.45) \quad \psi_r(\Pi, \Omega) = 0,$$

if any, as are a consequence of (3.43).

The usefulness of this definition is based on the fact, recognized above, that the parameters Π, Ω of the reduced form uniquely specify the distribution of the observed variables. In other words, there is a one-to-one correspondence between the points of the reduced parameter space and the sets of mutually equivalent points in the restricted parameter space (see Definitions 2.1.4 and 2.1.5.1).

DEFINITION 3.2.4.2. *The space of the statistics M_{xx} will be called the moment space. The space of the statistics P and W defined by (3.39) and (3.42) will be called the reduced moment space. We exclude from these spaces any singular values of M_{xx} or W , which will be referred to as arising from "singular" samples.*

¹The concept of independence for the purposes of this discussion is sharply defined by Definition 3.2.5 below.

On the basis of the foregoing definitions, the condition (3.38) for attainment of the absolute maximum of the likelihood function can be rewritten as

$$(3.46) \quad \Pi = P, \quad \Omega = W,$$

for this is obtained if the expressions (3.38) for A and Σ are substituted in (3.44). We thus have:

THEOREM 3.2.4.1. *A necessary and sufficient condition for the attainability of the absolute maximum of the likelihood function for a given nonsingular sample of observations, is that, to the point P, W in the reduced moment space, there corresponds, by (3.46), a point Π, Ω that belongs to the restricted reduced parameter space, i.e., satisfies the restrictions (3.45).*

This theorem shows that the attainability of the absolute maximum of L under given restrictions depends only on the statistics P, W , not on the statistics M_{zz} on which, as shown in (3.40), the likelihood function also depends. Moreover, the attainability of L_{\max} does not even depend on P and W if the set of restrictions (3.45) is empty.

We can now state an important theorem, which indicates the connection between identifiability of structural equations and attainability of the absolute maximum of the likelihood function, alluded to at the beginning of the present section 3.2.4.

THEOREM 3.2.4.2. *Let N_m be a region of positive measure in the reduced moment space, in each point of which the likelihood function can attain its absolute maximum under the restrictions (3.43). Let N_η be the corresponding region, according to (3.46), in the reduced parameter space. Assume that in every point of N_η the structural equations belonging to a nonempty set S are completely identifiable. Then, the addition to (3.43) of one further restriction, which is independent in N_η (in the sense of Definition 3.2.5 below) of the original restrictions (3.43), and which refers to equations of S only, will prevent the likelihood function from attaining its absolute maximum almost everywhere in N_η .*

*3.2.5. *Proof of Theorem 3.2.4.2.* Denote by $\theta \equiv [\theta_S \quad \theta_{-S}]$ a vector (one-row matrix) containing, under the notation θ_p , $p = 1, \dots, P$, all elements of A and Σ . Let θ_S contain, under the notation

θ_p , $p = 1, \dots, P_S$, all elements of those rows of A corresponding to structural equations belonging to the set S and those elements of Σ referring to two equations of S . Let θ_{-S} with elements θ_p , $p = P_S + 1, \dots, P$, contain all remaining parameters. Let the vector $\eta = \eta(\theta)$ contain all P^\ominus elements of Π and Ω , and h those of P and \bar{W} .

The assumptions with regard to N_m and N_η respectively, stated in the theorem, imply, owing to Theorem 3.2.4.1, that in every point η of N_η , the equation system (3.43), (3.44) admits at least one solution θ , and further that among those solutions, there is only a finite number of different values of θ_S . It follows from theorems regarding implicit functions, that (3.43) and (3.44) define θ_S as an implicit (multivalued) function of η , of which the derivatives are found from

$$(3.47) \quad \begin{aligned} \Phi_S \cdot \delta\theta'_S + \Phi_{-S} \cdot \delta\theta'_{-S} &= 0, \\ H_S \cdot \delta\theta'_S + H_{-S} \cdot \delta\theta'_{-S} &= \delta\eta', \end{aligned}$$

by elimination of $\delta\theta_{-S}$. Here

$$(3.48) \quad \Phi_S \equiv \begin{bmatrix} \frac{\partial\varphi_1}{\partial\theta_1} & \dots & \frac{\partial\varphi_1}{\partial\theta_{P_S}} \\ \cdot & \dots & \cdot \\ \frac{\partial\varphi_R}{\partial\theta_1} & \dots & \frac{\partial\varphi_R}{\partial\theta_{P_S}} \end{bmatrix}, \quad \Phi_{-S} \equiv \begin{bmatrix} \frac{\partial\varphi_1}{\partial\theta_{P_S+1}} & \dots & \frac{\partial\varphi_1}{\partial\theta_P} \\ \cdot & \dots & \cdot \\ \frac{\partial\varphi_R}{\partial\theta_{P_S+1}} & \dots & \frac{\partial\varphi_R}{\partial\theta_P} \end{bmatrix},$$

and H_S , H_{-S} are defined similarly with respect to the elements of $\eta = \eta(\theta)$.

The fact that (3.47) possesses at least one solution $[\delta\theta_S \quad \delta\theta_{-S}]$ for any $\delta\eta$ requires that

$$(3.49) \quad \rho \begin{pmatrix} \Phi_S & \Phi_{-S} \\ H_S & H_{-S} \end{pmatrix} = \rho(\Phi_S \quad \Phi_{-S}) + \rho(H_S \quad H_{-S}).$$

For the only alternative to (3.49) is that the right-hand member exceeds the left-hand member, in which case there would, according to Lemma 3.2.2, exist two nonvanishing vectors φ , η such that

$$(3.50) \quad [\bar{\varphi} \quad \bar{\eta}] \begin{bmatrix} \Phi_S & \Phi_{-S} \\ H_S & H_{-S} \end{bmatrix} = 0.$$

But then we could conclude from the existence of a solution of (3.47) that

$$(3.51) \quad \bar{\varphi} 0' + \bar{\eta} \delta\eta' = \bar{\eta} \delta\eta' = 0,$$

and values of $\delta\eta$ not satisfying (3.51) would not permit a solution of (3.47), contrary to the assumption made.

On the other hand, it is known that, after elimination of $\delta\theta_{-S}$, (3.47) has only a finite number of solutions $\delta\theta_S$, which in view of the linearity of the system (3.47) can only be one. The uniqueness of this solution is essentially a property of the homogeneous system of equations obtained from (3.47) by writing $\delta\eta = 0$. From the uniqueness of $\delta\theta_S$ it follows that

$$(3.52) \quad \rho \begin{pmatrix} \Phi_S & \Phi_{-S} \\ H_S & H_{-S} \end{pmatrix} = \rho \begin{pmatrix} \Phi_S \\ H_S \end{pmatrix} + \rho \begin{pmatrix} \Phi_{-S} \\ H_{-S} \end{pmatrix}.$$

For otherwise, according to Lemma 2.3.2, two nonvanishing vectors $\bar{\theta}_S$, $\bar{\theta}_{-S}$ would exist such that

$$(3.53) \quad \begin{bmatrix} \Phi_S & \Phi_{-S} \\ H_S & H_{-S} \end{bmatrix} \begin{bmatrix} \bar{\theta}'_S \\ \bar{\theta}'_{-S} \end{bmatrix} = 0,$$

and a scalar multiple of $\bar{\theta} \equiv [\bar{\theta}_S \quad \bar{\theta}_{-S}]$ could be added to the solution $\delta\theta = [\delta\theta_S \quad \delta\theta_{-S}]$ of (3.47) to produce other solutions

for the same value of $\delta\eta$, which differ in regard to $\delta\theta_S$. In addition, we have

$$(3.54) \quad \rho \begin{pmatrix} \Phi_S \\ H_S \end{pmatrix} = P_S,$$

which is the highest rank a matrix of P_S columns can attain. For, if the left-hand member in (3.54) were less than P_S , then a nonvanishing vector $\bar{\theta} = [\bar{\theta}_S \quad 0_{-S}]$ could be found, a scalar multiple of which could be added to the solution $\delta\theta = [\delta\theta_S \quad \delta\theta_{-S}]$ of (3.47) to produce other solutions for the same value of $\delta\eta$, which differ in regard to $\delta\theta_S$.

Now suppose that an additional a priori restriction

$$(3.55) \quad \varphi(\theta_S) = 0$$

is imposed such that N_η contains at least one point η_0 in which (3.43), (3.44), and (3.55) have a solution θ . (If no such point exists the theorem is already true.) We shall investigate what are the conditions to be satisfied by the derivatives of φ and φ_r in that point in order that the absolute maximum of the likelihood function is attainable everywhere in a neighborhood of the corresponding point $h_0 = \eta_0$ of the reduced moment space. For that to be so, it is necessary that if the row

$$(3.56) \quad [\varphi_S \quad 0_{-S}] \equiv \left[\frac{\partial\varphi}{\partial\theta_1} \quad \cdots \quad \frac{\partial\varphi}{\partial\theta_p} \quad 0_{P_S+1} \quad \cdots \quad 0_{P_S} \right]$$

is added to the matrix $[\Phi_S \quad \Phi_{-S}]$, the system (3.47) so enlarged or

$$(3.57) \quad \begin{aligned} \varphi_S \cdot \delta\theta'_S &= 0, \\ \Phi_S \cdot \delta\theta'_S + \Phi_{-S} \cdot \delta\theta'_{-S} &= 0, \\ H_S \cdot \delta\theta'_S + H_{-S} \cdot \delta\theta'_{-S} &= \delta\eta', \end{aligned}$$

satisfies the properties established for the original system (3.47). This would have the following consequences: From (3.54), applied both to the original and to the enlarged system, it follows that

$$(3.58) \quad \rho \begin{pmatrix} \Phi_S \\ H_S \end{pmatrix} = P = \rho \begin{pmatrix} \varphi_S \\ \Phi_S \\ H_S \end{pmatrix}.$$

From (3.52) applied to both members of (3.58) follows

$$(3.59) \quad \rho \begin{pmatrix} \Phi_S & \Phi_{-S} \\ H_S & H_{-S} \end{pmatrix} = \rho \begin{pmatrix} \varphi_S & 0 \\ \Phi_S & \Phi_{-S} \\ H_S & H_{-S} \end{pmatrix},$$

since the addition of a row of zeros to the last matrix in (3.52) does not change its rank. From (3.49) applied to both members of (3.59)

$$(3.60) \quad \rho \begin{pmatrix} \Phi_S & \Phi_{-S} \end{pmatrix} = \rho \begin{pmatrix} \varphi_S & 0 \\ \Phi_S & \Phi_{-S} \end{pmatrix}.$$

This, then, is a necessary condition for the existence of a solution of the enlarged system (3.57) for every value of $\delta\eta$. It will now be shown that if (3.60) is not satisfied, the only values of $\delta\eta$ permitting such a solution are those subject to a linear restriction of the type (3.51). If (3.60) is not true, we must have

$$(3.61) \quad \rho \begin{pmatrix} \Phi_S & \Phi_{-S} \end{pmatrix} = \rho \begin{pmatrix} \varphi_S & 0 \\ \Phi_S & \Phi_{-S} \end{pmatrix} - 1.$$

Suppose then that $\delta\eta$ is such that a solution of the enlarged

system (3.57) exists. The validity of (3.58) and (3.59) is not affected by the fact that $\delta\eta$ now represents one particular value, instead of all possible values. For the proof of the equalities for the enlarged system, equivalent to (3.52) and (3.54) respectively, depends only on the uniqueness of the solution with regard to $\delta\theta_S$, and this was already recognized as being a property of the homogeneous system obtained from (3.57) by taking $\delta\eta = 0$.

From (3.59), (3.61), and (3.49), we conclude

$$(3.62) \quad \rho \begin{pmatrix} \varphi_S & 0 \\ \Phi_S & \Phi_{-S} \\ H_S & H_{-S} \end{pmatrix} = \rho \begin{pmatrix} \varphi_S & 0 \\ \Phi_S & \Phi_{-S} \end{pmatrix} + \rho \begin{pmatrix} H_S & H_{-S} \end{pmatrix} - 1,$$

from which, as before, we can derive the existence of a linear restriction on $\delta\eta$ of the type (3.51).

It is well known that, if (3.60) is satisfied everywhere in N_η , then $\varphi(\theta_S) \equiv \varphi([\theta_S \ 0_{-S}])$ is a function of the remaining functions $\varphi_r(\theta)$, and (3.55) is either dependent on or incompatible with (3.43). The present theorem could probably be proved on the assumption that (3.60) holds in N_η only on a set of measure zero. We shall make a somewhat different assumption, which is better adapted to this particular proof, and is sufficient for our purposes:

DEFINITION 3.2.5. *The restriction (3.55) is called independent in N_η of the restrictions (3.43), if (3.60) is not satisfied in any point in N_η in which (3.43), (3.44), and (3.45) permit a solution θ (or Λ, Σ).*

If this is the case, the set of points η_0 in N_η in which (3.43), (3.44), and (3.55) permit a solution θ can only be of measure zero, because from any one such point, any neighboring points can now only be reached by variations $\delta\eta$ subject to a linear restriction.

A special case in which the additional restriction satisfies Definition 3.2.5 is, of course, that in which (3.60) is not satisfied in any point in N_η .

***3.2.6. Tabular summary of possible cases.** We shall now apply Theorem 3.2.4.2 to the case of linear and bilinear a priori restrictions (2.24) and (2.27). It may be useful to set out the various

TABLE 3.2.6

CONNECTION BETWEEN IDENTIFIABILITY OF STRUCTURAL EQUATIONS AND ATTAINABILITY OF THE ABSOLUTE MAXIMUM OF THE LIKELIHOOD FUNCTION

Note: This classification excludes point sets of measure zero in the parameter space (col. 3) and in the sample space (col. 4) and is subject to other exceptions discussed in sections 2.4.8 - 10. It is assumed that the a priori restrictions are compatible and that they are mutually independent in the sense of Definition 3.2.5.

Possible Cases		Statements relating to these cases	
(1)	(2)	(3)	(4)
The completed subset ¹ S_0 of the structural equations	The a priori restrictions in the associated subset ² R_{S_0} of (3.41) are with respect to S_0	The following structural equations are completely identifiable:	Can the likelihood function attain its absolute maximum?
(A) is empty.		none.	Only if among the solutions B of (3.41) there is a real solution. ⁴
(B) is not empty but does not contain all structural equations.	(1) just adequate in number and variety. ³	only those of S_0 .	Only if among the solutions B of (3.41) there is a real solution. ⁴
	(2) more than just adequate in number and variety. ³	only those of S_0 .	No.
(C) contains all structural equations.	(1) just adequate in number and variety. ³	all structural equations.	Only if among the solutions B of (3.41) there is a real solution. ⁴
	(2) more than just adequate in number and variety. ³	all structural equations.	No.

¹See Definition 2.4.6.2. ²See Definition 2.4.6.1. ³See Definition 3.2.6.

⁴If the a priori restrictions consist of linear restrictions only, this clause can be replaced by an unqualified "yes." It may be stated without

cases as to identifiability and attainability of the absolute maximum of L in a tabular form based on the counting of restrictions, even though the validity of this criterion is subject to exceptions already noted. In connection with Thm. 3.2.4.2, it is desirable to supplement Definition 2.1.5.5 by

DEFINITION 3.2.6. *A subset R of the a priori restrictions (2.28) will be said to be just adequate in number and variety with respect to (the identification of) a subset S of the structural equations if it is adequate in the sense of Definition 2.1.5.5 but loses that property if any of the restrictions in R are omitted.*

3.2.7. A factorization of the likelihood function. A further remark may be made about the case where the a priori restrictions imply a simultaneous partitioning (2.82) or

$$(3.63) \quad B = \begin{bmatrix} B_{I I} & B_{I II} \\ 0 & B_{II II} \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{I I} & 0 \\ 0 & \Sigma_{II II} \end{bmatrix},$$

of the matrices B and Σ . It is easily seen that this entails a factorization of the likelihood function, expressed by the following splitting-up of its logarithm (3.6) into two terms

$$(3.64) \quad L(A, \Sigma) = L_I(A_I, \Sigma_{II}) + L_{II}(A_{II}, \Sigma_{II II}),$$

where

$$(3.65) \quad \begin{aligned} L_I &= -\frac{1}{2}K_I \log 2\pi + \log \det B_{I I} - \frac{1}{2} \log \det \Sigma_{I I} \\ &\quad - \frac{1}{2} \operatorname{tr}(\Sigma_{I I}^{-1} \cdot A_I \cdot M_{xx} \cdot A_I'), \\ L_{II} &= -\frac{1}{2}K_{II} \log 2\pi + \log \det B_{II II} - \frac{1}{2} \log \det \Sigma_{II II} \\ &\quad - \frac{1}{2} \operatorname{tr}(\Sigma_{II II}^{-1} \cdot A_{II} \cdot M_{xx} \cdot A_{II}'). \end{aligned}$$

In [XVII] this factorization property of the likelihood function

proof that the "yes" applies even if the number of linear restrictions referring to each structural equation is $K_y - 1$.

is used to justify the concept of exogenous variables. Here it is sufficient to remark that, if no further a priori restrictions connect the structural equations with coefficients A_I with those having coefficients A_{II} , the estimation problems of these two subsystems of the system of structural equations have been effectively separated. For the two terms in (3.64) then depend on entirely independent sets of parameters, and the function (3.64) can only reach its maximum if each of the two terms reaches its own maximum.

It should be noted that the expressions (3.65) for L_I and L_{II} are of precisely the same form as the original logarithmic likelihood function (3.6). The variables indicated by the subscript II occur as dependent variables in L_{II} while the variables corresponding to the subscript I do not occur in L_{II} (the moments $M_{y_1 x}$ are multiplied into vanishing coefficients). The variables corresponding to the subscript II occur as predetermined variables in L_I , in which the variables corresponding to the subscript I represent the dependent variables.

3.3. Large-Sample Properties of the Maximum-Likelihood Estimates

3.3.1. Assumptions. In this section 3.3 we shall discuss the large-sample properties of the maximum-likelihood estimates of the parameters of the system (1.1) of structural equations. Following Mann and Wald, [1943, p. 192], we shall assume that the equation system is stable. Reverting to the notation of section 1, we express this by the following two assumptions:

ASSUMPTION 3.3.1.1. All roots ρ of the equation

$$(3.66) \quad \det \left[\sum_{\tau=0}^{\tau^{\square}} B(\tau) \rho^{\tau^{\square}-\tau} \right] = \det \left[\sum_{\tau=0}^{\tau^{\square}} \beta_g i_{\tau} \rho^{\tau^{\square}-\tau} \right] = 0,$$

$g, i = 1, \dots, K_y,$

satisfy

$$(3.67) \quad |\rho| < 1.$$

ASSUMPTION 3.3.1.2. If

$$(3.68) \quad m_{z_k z_l}(\tau, \tau', T) \equiv \frac{1}{T} \sum_{t=1}^T z_k(t - \tau) z_l(t - \tau'),$$

$k, l = 1, \dots, K,$

is a moment of two exogenous variables $z_k(t)$ and $z_l(t)$, there exists a finite limit

$$(3.69) \quad \lim_{T \rightarrow \infty} m_{z_k z_l}(\tau, \tau', T) = \infty_{z_k z_l}(\tau, \tau')$$

for every k, l, τ , and τ' .

Regarding the distribution of the disturbances we shall make two alternative assumptions:

ASSUMPTION 3.3.1.3. *The distribution function $f(u_1, \dots, u_{K_y})$ of the disturbances possesses finite $(4 + \epsilon)$ th-order moments for some $\epsilon > 0$. Its first-order moments vanish and its second-order moments form a nonsingular matrix Σ .*

Alternatively, we shall specify a particular distribution admitted under Assumption 3.3.1.3.

ASSUMPTION 3.3.1.4. *The disturbances u_1, \dots, u_{K_y} have a joint normal distribution (3.1) with mean zero and nonsingular second-order moment matrix Σ .*

Values $y(t)$, $t \leq 0$, of endogenous variables, with a timing preceding the period $1 \leq t \leq T$ during which the dependent variables are observed, are treated (together with the values of the exogenous variables) as given constants which remain the same in repeated samples.

3.3.2. *Quasi-maximum-likelihood estimates.* Under Assumption 3.3.1.4, the distribution function of the observations $x_1(1), \dots, x_{K_y}(T)$ is

$$(3.70) \quad F(M_{xx}, A, \Sigma) = (2\pi)^{-\frac{1}{2} \cdot K_y \cdot T} \cdot \det^T B \cdot \det^{-\frac{1}{2} \cdot T} \Sigma \cdot \exp\left\{-\frac{1}{2} \text{tr}(\Sigma^{-1} A M_{xx} A')\right\}.$$

As a function of the parameters A, Σ , we have called (3.70) the likelihood function, and defined maximum-likelihood estimates as the values of the parameters that, subject to the a priori restrictions, maximize this function. Under the wider Assumption 3.3.1.3, (3.70) has no necessary connection with the distribution of the observations. Nevertheless, we can use the function (3.70) to define estimates of the parameters by the same maximizing procedure. In these circumstances, we shall call (3.70) the quasi-likelihood function, and call the maximizing values of its parameters quasi-maximum-likelihood estimates. We shall also discuss some large-sample properties of these estimates.

3.3.3. *Results of Mann and Wald.* A very thorough analysis of large-sample properties of the quasi-maximum-likelihood estimates has been given by Mann and Wald [1943]. The system considered by these authors satisfies Assumption 3.3.1.1 and a slightly more restrictive version of Assumption 3.3.1.3. Their system does not contain exogenous variables $z_k(t)$ (except a constant term in each equation). Finally, they assume that each equation is completely identified. Our main concern in the present section 3.3 is to indicate that Mann and Wald's results can be extended to the case where exogenous variables satisfying Assumption 3.3.1.2 are present, and to the case where some but not all of the structural equations are identifiable. We shall first discuss the large-sample properties of the moment matrix M_{xx} , and thereafter those of the quasi-maximum-likelihood estimates.

3.3.4. *Asymptotic distribution of the moments.* Extended to include systems with exogenous variables, Mann and Wald's results regarding the moments can be stated as follows:

THEOREM 3.3.4. *Under Assumptions 3.3.1.1, 3.3.1.2, and 3.3.1.3, the expected value*

$$(3.71) \quad \mathcal{E}M_{xx} \equiv M_{xx}$$

of the moment matrix M_{xx} possesses the properties a) that

$$(3.72) \quad \lim_{T \rightarrow \infty} M_{xx} \equiv M_{cc}$$

exists and is finite and b) that those elements of

$$(3.73) \quad M_{xx} - M_{xx}$$

which are subject to sampling variation have a joint asymptotically normal distribution with a variance-covariance matrix of order T^{-1} .

The matrix M_{xx} comprises the square and cross moments of the endogenous variables y_i (with and without time lags), the cross moments between the endogenous variables y_i and the exogenous variables z_k (with and without time lags), and the square and cross moments (3.68) of the exogenous variables. Since the latter variables are treated as given functions of time (see [XVII]) not subject to a probability distribution, the elements (3.68) of M_{xx} are likewise given functions of T , equal to the corresponding elements of M_{xx} .

We shall not indicate in detail the incorporation of exogenous variables in Mann and Wald's proof, since a somewhat different proof including exogenous variables will be published by one of the present authors [Rubin, 1948].

3.3.5. *A property of the logarithmic quasi-likelihood function.* In the remainder of this section 3.3, we shall notationally combine in one vector θ all elements of Λ and Σ , and we shall write M , \mathbf{M} instead of M_{xx} , M_{xx} . Occasionally, in particular in the present section 3.3.5, we shall distinguish notationally between the true values θ of these parameters, and the argument $\bar{\theta}$ of the quasi-likelihood function (3.70). The logarithmic quasi-likelihood function (divided by T),

$$(3.74) \quad L(M, \bar{\theta}) = \frac{1}{T} \log F(M, \bar{\theta})$$

as written out in (3.6), is linear in the moment matrix M . Its expected value therefore equals

$$(3.75) \quad \mathcal{E}L(M, \bar{\theta}) = L(\mathbf{M}, \bar{\theta}),$$

a function we shall refer to as the expected logarithmic quasi-likelihood function. The expected moment matrix \mathbf{M} occurring in (3.75) depends, for any given T , only on the fixed values of the exogenous variables $z_k(t)$, on the requisite number of initial values $x_k(t)$, $t \leq 0$, of all variables, and on the true values θ of

the parameters. We express the last-mentioned dependence by

$$(3.76) \quad M = M(\theta).$$

The function (3.75) possesses the following important property:

THEOREM 3.3.5. *Under Assumption 3.3.1.3, the expected logarithmic likelihood function*

$$(3.77) \quad L\{M(\theta), \bar{\theta}\}$$

reaches its (unrestricted) absolute maximum with respect to the parameters $\bar{\theta}$ in the point

$$(3.78) \quad \bar{\theta} = \theta.$$

We shall first prove this theorem under the normality Assumption 3.3.1.4. In that case the function $F(M, \theta)$ serves both as distribution function of the observations, and as a function defining the maximum-likelihood estimates. Therefore, if dx stands for $dx_1(1) \cdots dx_K(T)$, and $\int dx$ for integration over the whole sample space, we have $\int F(M, \theta) dx = 1$, and

$$(3.79) \quad \begin{aligned} 0 &= \frac{\partial}{\partial \theta} \int F(M, \theta) dx = \int \frac{\partial F}{\partial \theta} dx = \int F \frac{\partial \log F}{\partial \theta} dx \\ &= \int F(M, \theta) \left(\frac{\partial \log F(M, \bar{\theta})}{\partial \bar{\theta}} \right)_{\bar{\theta}=\theta} dx \\ &= \left[\frac{\partial}{\partial \bar{\theta}} \int F(M, \theta) \log F(M, \bar{\theta}) dx \right]_{\bar{\theta}=\theta} = \left[\frac{\partial}{\partial \bar{\theta}} \mathcal{E} \log F(M, \bar{\theta}) \right]_{\bar{\theta}=\theta}, \\ &= T \left[\frac{\partial}{\partial \bar{\theta}} \mathcal{E} L(M, \bar{\theta}) \right]_{\bar{\theta}=\theta} = T \left[\frac{\partial}{\partial \bar{\theta}} L\{M(\theta), \bar{\theta}\} \right]_{\bar{\theta}=\theta}, \end{aligned}$$

using (3.75) in the last equality. The differentiations with respect to θ and $\bar{\theta}$ are performed without regard to the a priori restrictions. On the other hand, we know from Theorem 3.3.10, that the function (3.77) of an unrestricted $\bar{\theta}$ is stationary only in points where its absolute maximum is reached. It follows from

(3.79) that that maximum is reached for $\bar{\theta} = \theta$.

The moments M entering in the definition (3.71) of the function (3.76) can be expressed, for any value of T , as quadratic or linear functions of the disturbances $u_g(t)$. This is seen most readily by repeated substitution of the right-hand member of the reduced form (3.11) for the $y_i(t)$ in the definition (3.7) of the moments, taking $t = T, T-1, \dots, 1$, successively. It follows that the function $M(\theta)$ remains the same under the more general Assumption 3.3.1.3 regarding the distribution of the disturbances. Consequently (3.79), and therewith Theorem 3.3.5, are also valid under Assumption 3.3.1.3.

3.3.6. Consistency of quasi-maximum-likelihood estimates of identifiable parameters. It will be clear that any statement regarding consistency¹ of quasi-maximum-likelihood estimates can relate only to the estimation of parameters that are uniquely identifiable in a neighborhood of the true parameter point θ . Since maximum-likelihood estimation is invariant for functional transformation in the parameter space, we can achieve greater generality and flexibility by formulating our statements in terms of identifiable functions of the parameters θ , defined as follows:

Let the a priori restrictions be denoted, as in (3.43), by

$$(3.80) \quad \varphi(\theta) \equiv \left[\varphi_1(\theta) \quad \cdots \quad \varphi_R(\theta) \right] = 0.$$

The restrictions (3.80) define the restricted parameter space, within which as before (section 2.1.5) we distinguish sets of mutually equivalent points θ , or briefly "equivalent point sets."

DEFINITION 3.3.6.1. A parameter $\zeta = \zeta(\theta)$ is called uniquely identifiable in a region \mathcal{N} of the restricted parameter space if it is constant, within \mathcal{N} , on any set of mutually equivalent points.

Let $\eta(\theta) \equiv \left[\eta_1(\theta) \quad \cdots \quad \eta_{P^\oplus}(\theta) \right]$ represent, as before, the P^\oplus parameters Π, Ω of the reduced form (3.11) of the structural equations.

¹An estimate t_q of a parameter θ_q , derived from a sample of size T , is called consistent if, for any $\varepsilon > 0$, $\lim_{T \rightarrow \infty} P(|t_q - \theta_q| > \varepsilon) = 0$, where $P(E)$ denotes the probability of an event E . This relationship of t_q and θ_q is also denoted by $\text{plim}_{T \rightarrow \infty} t_q = \theta_q$.

DEFINITION 3.3.6.2. *The a priori restrictions (3.80) will be called regular in the parameter point θ if in a neighborhood N_θ of that point in the unrestricted parameter space the following three conditions are satisfied:*

- (i) *the functions $\varphi_r(\theta)$ possess continuous third derivatives,*
 - (ii)
- $$(3.81) \quad \rho \left(\frac{d\varphi}{d\theta'} \right) = R,$$
- (iii)
- $$(3.82) \quad \rho \left(\frac{d\eta}{d\theta'}, \quad \frac{d\varphi}{d\theta'} \right) = P^\oplus + R^{\oplus\oplus},$$

say, is constant.

On the basis of these definitions, we shall prove:

THEOREM 3.3.6. *Let $\zeta(\theta) \equiv [\zeta_1(\theta) \quad \cdots \quad \zeta_Q(\theta)]$ be a set of Q parameters that*

- (i) *are uniquely identifiable in a neighborhood ${}^T N_\theta$ of the true parameter point θ in the parameter space as restricted by means of a priori restrictions (3.80) regular in that point,*
- (ii) *possess continuous third derivatives in a neighborhood N_θ of θ in the unrestricted parameter space containing ${}^T N_\theta$, and*
- (iii) *in N_θ satisfy*

$$(3.83) \quad \rho \left(\frac{d\zeta}{d\theta'}, \quad \frac{d\varphi}{d\theta'} \right) = Q + R.$$

Then the quasi-maximum-likelihood estimates $\hat{\zeta}$ of ζ are consistent, and have an asymptotically normal distribution with variance-covariance matrix of order T^{-1} .

*3.3.7. *Three lemmas.* In order to prove this theorem, we shall first establish the following three lemmas.

LEMMA 3.3.7. *If the restrictions (3.80) are regular in the point θ , they imply exactly*

$$(3.84) \quad R^{\oplus} \equiv R - R^{\oplus\oplus}$$

independent restrictions

$$(3.85) \quad \psi(\eta) \equiv \left[\psi_1(\eta) \quad \cdots \quad \psi_{R^{\oplus}}(\eta) \right] = 0, \quad \rho \left(\frac{d\psi}{d\eta'} \right) = R^{\oplus},$$

on the parameters $\eta(\theta)$ of the reduced form, in a neighborhood N_{θ} of the point θ .

Since the parameters of the reduced form are independent,

$$(3.86) \quad \rho \left(\frac{d\eta}{d\theta'} \right) = P^{\oplus},$$

in any region of the parameter space (excluding, of course, the points with $\det B = 0$). It follows from (3.81), (3.82), (3.84), and (3.86) that

$$(3.87) \quad 0 \leq R^{\oplus\oplus} \leq R, \quad \text{so} \quad 0 \leq R^{\oplus} \leq R.$$

If $R^{\oplus} > 0$ and hence $R^{\oplus\oplus} < R$, it follows from (3.81), (3.82), and (3.86) that there exists a vector function of R^{\oplus} elements $\psi(\eta, \varphi) \equiv \left[\psi_1(\eta, \varphi) \quad \cdots \quad \psi_{R^{\oplus}}(\eta, \varphi) \right]$ such that in N_{θ}

$$(3.88) \quad \psi\{\eta(\theta), \varphi(\theta)\} = 0, \quad \rho \left(\begin{array}{c} \frac{\partial \psi}{\partial \eta'} \\ \frac{\partial \psi}{\partial \varphi'} \end{array} \right) = R^{\oplus}.$$

Moreover, these functions must be such that, in the point set $N_{\eta, \varphi}$ on which N_{θ} is mapped through the functions $\eta(\theta)$ and $\varphi(\theta)$,

$$(3.89) \quad \rho \left(\frac{\partial \psi}{\partial \eta'} \right) = \rho \left(\frac{\partial \psi}{\partial \varphi'} \right) = R^{\oplus},$$

For, if for instance $\rho(\partial\psi / \partial\eta')$ were less than R^\oplus , there would exist a vector function $\kappa(\eta, \varphi)$ containing R^\oplus elements such that on $N_{\eta, \varphi}$,

$$(3.90) \quad \frac{\partial\psi}{\partial\eta'} \kappa' = 0, \quad \frac{\partial\psi}{\partial\varphi'} \kappa' \equiv \bar{\kappa}' \neq 0.$$

In this case the equations

$$(3.91) \quad \begin{bmatrix} \frac{d\eta}{d\theta'} & \frac{d\varphi}{d\theta'} \end{bmatrix} \begin{bmatrix} \frac{\partial\psi}{\partial\eta'} \\ \frac{\partial\psi}{\partial\varphi'} \end{bmatrix} = 0$$

obtained from (3.88) by differentiation with respect to θ' would possess a linear combination

$$(3.92) \quad \begin{bmatrix} \frac{d\eta}{d\theta'} & \frac{d\varphi}{d\theta'} \end{bmatrix} \begin{bmatrix} \frac{\partial\psi}{\partial\eta'} \\ \frac{\partial\psi}{\partial\varphi'} \end{bmatrix} \kappa' = \frac{d\varphi}{d\theta'} \bar{\kappa}' = 0, \quad \bar{\kappa}' \neq 0,$$

in contradiction with the regularity condition (3.81) on the a priori restrictions. Writing now

$$(3.93) \quad \psi(\eta) \equiv \psi(\eta, 0)$$

(3.85) follows from (3.88) and (3.89).

We note for later use that in ${}^rN_\theta$, as a consequence of (3.91) and (3.93)

$$(3.94) \quad \rho \left(\begin{array}{cc} \frac{d\eta}{d\theta'} & \frac{d\psi(\eta)}{d\eta'} \\ \frac{d\varphi}{d\theta'} & \end{array} \right) = \rho \left(\begin{array}{cc} -\frac{d\varphi}{d\theta'} \frac{\partial\psi(\eta, \varphi)}{\partial\varphi'} & \frac{d\varphi}{d\theta'} \\ \frac{d\varphi}{d\theta'} & \end{array} \right) = \rho \left(\frac{d\varphi}{d\theta'} \right).$$

LEMMA 3.3.7.2. *If Z and Ξ are two matrices with equal numbers of rows and columns respectively, and Φ is a third matrix with an equal number of rows, such that*

$$(3.95) \quad \rho(Z - \Xi \quad \Phi) = \rho(\Phi),$$

then

$$(3.96) \quad \rho(Z \quad \Phi) = \rho(\Xi \quad \Phi).$$

Proof: It follows from (3.95) that there exists a matrix Π such that

$$(3.97) \quad Z - \Xi = \Phi \Pi.$$

Hence

$$(3.98) \quad \rho(Z \quad \Phi) = \rho(\Xi + \Phi \Pi \quad \Phi) = \rho(\Xi \quad \Phi).$$

LEMMA 3.3.7.3. *If Ξ, Ψ , and Φ are three matrices with an equal number of rows, such that*

$$(3.99) \quad \rho(\Xi \quad \Phi) = \rho(\Xi) + \rho(\Phi), \quad \rho(\Xi) = c(\Xi), \quad \rho(\Phi) = c(\Phi),$$

and

$$(3.100) \quad \rho(\Psi \quad \Phi) = \rho(\Phi),$$

then

$$(3.101) \quad \rho(\Xi \quad \Psi) = \rho(\Xi) + \rho(\Psi).$$

Proof: It follows from (3.100) that there exists a matrix P such that

$$(3.102) \quad \Psi = \Phi P.$$

Now, if the left-hand member in (3.101) were smaller than the right-hand member, there would according to Lemma 2.3.2 exist two vectors λ and μ such that

$$(3.103) \quad \Xi \lambda' + \Psi \mu' = \Xi \lambda' + \Phi P \mu' = 0, \quad \Xi \lambda' \neq 0.$$

Regarding $P \mu'$ as a new vector $\bar{\mu}'$, this is in contradiction with (3.99), since the second condition in (3.103) precludes the vanishing of $[\lambda \quad \bar{\mu}]$. It is easily seen that the last two conditions in (3.99) only facilitate the proof, and can be dispensed with if necessary.

*3.3.8. *First part of the proof of Theorem 3.3.6.* It was noted in section 3.1.6 that the parameter vector $\eta(\theta)$ of the reduced form is constant on each equivalent point set in the unrestricted parameter space, and that $\eta(\theta)$ assumes different values on any two different equivalent sets. Consequently, the same is true in the restricted parameter space. It follows from Definition 3.3.6.1 that $\zeta(\theta)$ is in ${}^rN_\theta$ a one-valued function $\xi(\eta)$ of $\eta(\theta)$:

$$(3.104) \quad \zeta(\theta) = \xi\{\eta(\theta)\} \quad \text{whenever} \quad \varphi(\theta) = 0.$$

Since $\zeta(\theta)$, $\eta(\theta)$, and $\varphi(\theta)$ have continuous third derivatives with respect to the elements of θ in an unrestricted neighborhood N_θ of θ , $\xi(\eta)$ must have continuous third derivatives with respect to the elements of η . Therefore, (3.104) implies that in a restricted neighborhood ${}^rN_\theta$ of θ , viz., in the set ${}^rN_\theta$ of those points in N_θ for which $\varphi(\theta) = 0$,

$$(3.105) \quad \rho \left(\left(\frac{d\zeta}{d\theta'} - \frac{d\eta}{d\theta'} \frac{d\xi}{d\eta'} \right) \quad \frac{d\varphi}{d\theta'} \right) = \rho \left(\frac{d\varphi}{d\theta'} \right) = R.$$

Thus the matrices $Z \equiv d\zeta/d\theta'$, $\Xi \equiv (d\eta/d\theta')(d\xi/d\eta')$ and $\Phi \equiv d\varphi/d\theta'$ satisfy the condition of Lemma 3.3.7.2, and, from (3.83) and (3.96), we have

$$(3.106) \quad \rho \left(\frac{d\eta}{d\theta'} \frac{d\xi}{d\eta'} \quad \frac{d\varphi}{d\theta'} \right) = Q + R.$$

Since this is the maximum possible rank for a matrix of $Q + R$ columns, we also have in ${}^rN_\theta$

$$(3.107) \quad \rho \left(\frac{d\eta}{d\theta'} \frac{d\xi}{d\eta'} \right) = Q.$$

According to (3.81), (3.101), (3.106), and (3.107), the matrices $\Xi \equiv (d\eta/d\theta')(d\xi/d\eta')$, $\Psi \equiv (d\eta/d\theta')(d\psi/d\eta')$, and $\Phi \equiv d\varphi/d\theta'$, satisfy the conditions of Lemma 3.3.7.3. It follows, using (3.89), (3.101), and (3.107), that in ${}^rN_\theta$

$$(3.108) \quad \rho \left(\frac{d\eta}{d\theta'} \left[\frac{d\xi}{d\eta'} \quad \frac{d\psi}{d\eta'} \right] \right) = \rho \left(\frac{d\eta}{d\theta'} \frac{d\xi}{d\eta'} \quad \frac{d\eta}{d\theta'} \frac{d\psi}{d\eta'} \right) = Q + R^{\oplus}.$$

Since the rank of a matrix product does not exceed the rank of either of the two factors, we must have

$$(3.109) \quad \rho \left(\frac{d\xi}{d\eta'} \quad \frac{d\psi}{d\eta'} \right) \geq Q + R^{\oplus}$$

in the point set ${}^{\tau}N_{\eta}$ on which ${}^{\tau}N_{\theta}$ is mapped by the function $\eta(\theta)$. However, $Q + P^{\oplus}$ is also the number of columns of the matrix in (3.109). Hence, in ${}^{\tau}N_{\eta}$

$$(3.110) \quad \rho \left(\frac{d\xi}{d\eta'} \quad \frac{d\psi}{d\eta'} \right) = c \left(\frac{d\xi}{d\eta'} \quad \frac{d\psi}{d\eta'} \right) = Q + R^{\oplus},$$

and, because of the continuity of the functions involved, (3.110) holds also in a neighborhood N_{η} , of the point set ${}^{\tau}N_{\eta}$, in the space of the unrestricted parameters η - provided ψ is regarded as that function $\psi(\eta)$ of η only, defined by (3.93).

It follows from (3.110) that

$$(3.111) \quad S^{\oplus} \equiv P^{\oplus} - Q - R^{\oplus} \geq 0.$$

The equality sign holds only if ζ represents a complete set of identifiable parameters. However, whenever $S^{\oplus} > 0$ we can because of (3.110) choose a vector function $\chi(\eta) \equiv [\chi_1(\eta) \quad \cdots \quad \chi_{S^{\oplus}}(\eta)]$ in N_{η} having S^{\oplus} elements, with continuous third derivatives, such that in N_{η}

$$(3.112) \quad \rho \left(\frac{d\xi}{d\eta'} \quad \frac{d\psi}{d\eta'} \quad \frac{d\chi}{d\eta'} \right) = Q + R^{\oplus} + S^{\oplus} = P^{\oplus}.$$

The matrix in (3.112) has thus been made square and nonsingular, and there exists in N_{η} an inverse function $\eta(\xi, \psi, \chi)$, i.e., a one-valued function with continuous third derivatives such that identically in N_{η} ,

$$(3.113) \quad \eta\{\xi(\eta), \psi(\eta), \chi(\eta)\} = \eta, \quad \xi(\eta) = \zeta.$$

3.3.9. *Second part of the proof of Theorem 3.3.6.* We shall summarize the results reached in the previous section 3.3.8. At the same time, we shall revert to the notation used in section 3.3.5, whereby the argument $\bar{\theta}$ of the likelihood function, and functions $\bar{\eta} = \bar{\eta}(\bar{\theta})$, $\bar{\zeta} = \bar{\zeta}(\bar{\theta})$, etc., of $\bar{\theta}$, are distinguished from the true parameter point θ and the corresponding functional values $\eta = \eta(\theta)$, $\zeta = \zeta(\theta)$, etc., by placing bars on the former quantities.

It has been found that the P^\oplus parameters $\bar{\eta}$ of the reduced form of the structural equation can be expressed in a neighborhood N_η of η in the unrestricted space of the parameters $\bar{\eta}$, as one-valued and uniquely invertible functions $\bar{\eta} = \bar{\eta}(\bar{\zeta}, \bar{\psi}, \bar{\chi})$ possessing continuous third derivatives, of

(i) the Q identifiable parameters $\bar{\zeta} = \bar{\zeta}(\bar{\theta})$,

(ii) the R^\oplus functions $\bar{\psi}(\bar{\eta})$ expressing the restrictions (3.85) on $\bar{\eta}$ arising from the a priori restrictions (3.80) on $\bar{\theta}$.

(iii) S^\oplus auxiliary parameters $\bar{\chi} = \bar{\chi}(\bar{\theta})$, with $S^\oplus \equiv P^\oplus - Q - R^\oplus \geq 0$.

We go on to describe maximum-likelihood estimation of the parameters ζ under the restrictions (3.80) in terms of the functions that have been introduced. We start from the likelihood function (3.74) in the reduced form (3.15), now to be denoted

$$(3.114) \quad L = L^\oplus(M, \bar{\eta}),$$

a function possessing continuous derivatives of all orders. In this function we substitute

$$(3.115) \quad \bar{\eta} = \bar{\eta}(\bar{\zeta}, 0, \bar{\chi})$$

– thus automatically satisfying the restrictions (3.80) – and maximize with respect to $\bar{\chi}$ for any constant $\bar{\zeta}$. Let the maximizing value of $\bar{\chi}$ be denoted

$$(3.116) \quad \hat{\chi} = \hat{\chi}(M, \bar{\zeta}).$$

This function is one-valued in a neighborhood of $M = M$, $\bar{\zeta} = \zeta$, because of (3.112), and of the negative definiteness of $\partial^2 L(M, \bar{\eta}) / \partial \bar{\eta}' \partial \bar{\eta}$ in the point $\bar{\eta} = \eta$, to be shown below in (3.121). Furthermore, it possesses continuous second-order derivatives because

of the continuity of the third-order derivatives of $\bar{\eta}(\bar{\zeta}, 0, \bar{\chi})$. We insert the value (3.116) in (3.115),

$$(3.117) \quad \bar{\eta} = \bar{\eta}(\bar{\zeta}, 0, \hat{\chi}(M, \zeta)),$$

and write

$$(3.118) \quad L = L^{\oplus}[M, \bar{\eta}(\bar{\zeta}, 0, \hat{\chi}(M, \bar{\zeta}))] \equiv L^{\oplus\oplus}(M, \bar{\zeta})$$

for the function so obtained. It follows from the invariance of any maximizing process for continuous functional transformation of the parameters that the value $\hat{\zeta}$ of $\bar{\zeta}$ maximizing $L^{\oplus\oplus}(M, \bar{\zeta})$ represents the maximum-likelihood estimate of the parameter vector ζ . Explicitly,

$$(3.119) \quad \hat{\zeta} = \bar{\zeta}(\hat{\theta}),$$

if $\hat{\theta}$ is that value of $\bar{\theta}$ maximizing the likelihood function in its original form (3.6), subject to the restrictions (3.80).

It was shown in the proof of Theorem 3.1.10 that the matrix

$$(3.120) \quad \frac{\partial^2 L^{\oplus}(M, \bar{\eta})}{\partial \bar{\eta}' \partial \bar{\eta}}$$

is negative definite in any point $\bar{\eta} = \bar{\eta}(\bar{\theta})$ such that the original likelihood function $L(M, \bar{\theta})$ reaches its unrestricted absolute maximum in $\bar{\theta}$. Theorem 3.3.5 states that $L(M, \bar{\theta})$ reaches its unrestricted absolute maximum in the true parameter point θ . Since $\eta = \bar{\eta}(\theta)$, it follows that

$$(3.121) \quad \Lambda^{\oplus} \equiv \left(\frac{\partial^2 L^{\oplus}(M, \bar{\eta})}{\partial \bar{\eta}' \partial \bar{\eta}} \right)_{\bar{\eta}=\eta}$$

is negative definite. We shall now prove that in consequence

$$(3.122) \quad \Lambda^{\oplus\oplus} \equiv \left(\frac{\partial^2 L^{\oplus\oplus}(M, \bar{\zeta})}{\partial \bar{\zeta}' \partial \bar{\zeta}} \right)_{\bar{\zeta}=\zeta}$$

is also negative definite, where $\zeta = \bar{\zeta}(\theta)$ is the true value of the parameter vector ζ .

If we insert M for M in (3.117) and regard M as constant, $\bar{\eta}$ is expressed as a function

$$(3.123) \quad \bar{\eta} \equiv \bar{\eta}(\bar{\zeta}, 0, \hat{\chi}(M, \bar{\zeta})) \equiv \eta(\bar{\zeta}),$$

say, of $\bar{\zeta}$ alone, which possesses continuous first and second derivatives, the first being

$$(3.124) \quad \left(\frac{\partial \bar{\eta}}{\partial \bar{\zeta}'} + \frac{\partial \hat{\chi}}{\partial \bar{\zeta}'} \frac{\partial \bar{\eta}}{\partial \bar{\chi}'} \right)_{\bar{\psi}=0, \bar{\chi}=\bar{\chi}(M, \bar{\zeta})} \equiv \frac{d\bar{\eta}}{d\bar{\zeta}'},$$

say. In particular, owing to Theorem 3.5.5,

$$(3.125) \quad \bar{\eta}(\zeta) = \eta.$$

Differentiating (3.118) with respect to $\bar{\zeta}$, after substituting M for M , we have,

$$(3.126) \quad \frac{\partial L^{\oplus\oplus}(M, \bar{\zeta})}{\partial \bar{\zeta}'} = \frac{d\bar{\eta}}{d\bar{\zeta}'} \left(\frac{\partial L^{\oplus}(M, \bar{\eta})}{\partial \bar{\eta}'} \right)_{\bar{\eta}=\bar{\eta}(\bar{\zeta})}.$$

Because of Theorem 3.3.5, the quantities $\partial L^{\oplus} / \partial \bar{\eta}'$ in (3.126) vanish for $\bar{\zeta} = \zeta$. Therefore, and because continuous second derivatives of $\bar{\eta}(\bar{\zeta})$ exist, we have, using (3.125),

$$(3.127) \quad \Delta^{\oplus\oplus} = H \Delta^{\oplus} H',$$

where

$$(3.128) \quad H \equiv \left(\frac{d\bar{\eta}}{d\bar{\zeta}'} \right)_{\bar{\zeta}=\zeta}.$$

From (3.127) we conclude that the quadratic form

$$(3.129) \quad z \Delta^{\oplus\oplus} z' = z H \Delta^{\oplus} H' z' = y \Delta^{\oplus} y',$$

say, is equal to a negative definite quadratic form, in which, owing to (3.112), $y = z H$ vanishes only if z vanishes. Hence $\Delta^{\oplus\oplus}$ is negative definite, and therefore nonsingular.

It follows from Theorem 3.3.5 and from the definition of the

maximum-likelihood estimate $\bar{\zeta}$ of ζ , that the vector function

$$(3.130) \quad l^{\oplus\oplus}(\bar{M}, \bar{\zeta}) \equiv \frac{\partial}{\partial \bar{\zeta}} L^{\oplus\oplus}(\bar{M}, \bar{\zeta})$$

vanishes in the point $\bar{M} = M$, $\bar{\zeta} = \zeta$. Let us consider the Taylor expansion

$$(3.131) \quad 0 = \overset{*}{l^{\oplus\oplus}}(M, \hat{\zeta}) - l^{\oplus\oplus}(M, \zeta) = \left[\frac{\partial}{\partial \bar{\zeta}} \operatorname{tr} \left\{ (M - M) \frac{\partial L^{\oplus\oplus}(M, \bar{\zeta})}{\partial M'} \right\} \right]_{\bar{\zeta}=\zeta} + (\hat{\zeta} - \zeta) \Lambda^{\oplus\oplus} + \dots$$

from which we can solve for $\hat{\zeta} - \zeta$ by postmultiplication with $(\Lambda^{\oplus\oplus})^{-1}$,

$$(3.132) \quad \hat{\zeta} - \zeta = - \left[\frac{\partial}{\partial \bar{\zeta}} \operatorname{tr} \left\{ (M - M) \frac{\partial L^{\oplus\oplus}(M, \bar{\zeta})}{\partial M'} \right\} \right] (\Lambda^{\oplus\oplus})^{-1} + \dots$$

Reference to the form (3.6) of the likelihood function $L(M, \theta)$ shows that the coefficients of the elements $M - M$ in (3.132) do not all vanish. It follows that $\hat{\zeta} - \zeta$ is of the same order of magnitude $T^{-1/2}$ as $M - M$. The usual analysis of the quadratic term, omitted from (3.131), which is to be taken in a point intermediate between the two points $(M, \hat{\zeta})$ and (M, ζ) , will show that this term is of order T^{-1} , because of the continuity of the second derivative of the likelihood function (3.114) and of the functions $\bar{\eta}(\bar{\zeta}, 0, \bar{\chi})$ and $\bar{\chi}(\bar{M}, \bar{\zeta})$. Therefore, $\hat{\zeta} - \zeta$ is linear in $M - M$ up to terms of order T^{-1} . Theorem 3.3.6 follows from this observation combined with Theorem 3.3.4.

3.3.10. Asymptotic sampling variances and covariances of the maximum-likelihood estimates $\hat{\zeta}$. Mann and Wald's analysis shows that the expressions for the asymptotic sampling variances and covariances of the maximum-likelihood estimates $\hat{\theta}$ are greatly simplified by the normality Assumption 3.3.1.4 regarding the distribution of the disturbances. We shall here deal only with that case. The following derivation differs from that given by Mann and Wald [1943, pp. 213-214] only in that it applies to any set of parameters $\bar{\zeta}$ uniquely identifiable in a neighborhood of ζ , rather than to a complete set of identifiable parameters $\bar{\theta}$.

Under Assumption 3.3.1.4, the function F in (3.70), now to be denoted $F(M, \bar{\theta})$, also represents the probability density in the sample space. Writing

$$(3.133) \quad [\bar{\zeta} \quad \bar{\chi}] \equiv \bar{\omega},$$

we define, analogously to (3.114),

$$(3.134) \quad F(M, \bar{\theta}) = F^{\oplus}(M, \bar{\eta})$$

and, somewhat differently from (3.118),

$$(3.135) \quad F^{\oplus}(M, \bar{\eta}(\bar{\zeta}, 0, \bar{\chi})) \equiv \tilde{F}(M, \bar{\omega}),$$

with similar formulae in terms of $L = (1/T) \log F$. Finally, we define

$$(3.136) \quad \tilde{l}(M, \bar{\omega}) \equiv \frac{\partial}{\partial \bar{\omega}} \tilde{L}(M, \bar{\omega}).$$

Then, in the point $\bar{\omega} = \omega \equiv [\zeta \quad \chi]$,

$$\begin{aligned} \mathcal{E} \tilde{l}'(M, \omega) \tilde{l}(M, \omega) &= T^{-2} \int \tilde{F}(M, \omega) \frac{\partial \log \tilde{F}(M, \omega)}{\partial \omega'} \frac{\partial \log \tilde{F}(M, \omega)}{\partial \omega} dx \\ &= T^{-2} \int \frac{\partial \tilde{F}(M, \omega)}{\partial \omega'} \frac{\partial \log \tilde{F}(M, \omega)}{\partial \omega} dx \\ &= T^{-2} \int \left[\frac{\partial}{\partial \omega'} \left\{ F(M, \omega) \frac{\partial \log \tilde{F}(M, \omega)}{\partial \omega} \right\} \right. \\ &\quad \left. - \tilde{F}(M, \omega) \frac{\partial^2 \log \tilde{F}(M, \omega)}{\partial \omega' \partial \omega} \right] dx \\ (3.137) \quad &= T^{-2} \frac{\partial}{\partial \omega'} \int \tilde{F}(M, \omega) \frac{\partial \log \tilde{F}(M, \omega)}{\partial \omega} dx \\ &\quad - T^{-2} \int \tilde{F}(M, \omega) \frac{\partial^2 \log \tilde{F}(M, \omega)}{\partial \omega' \partial \omega} dx \\ &= 0 - T^{-1} \mathcal{E} \frac{\partial^2 \tilde{L}(M, \omega)}{\partial \omega' \partial \omega} \end{aligned}$$

$$= -T^{-1} \frac{\partial^2 L(M, \omega)}{\partial \omega' \partial \omega} \equiv -T^{-1} \tilde{\Lambda},$$

say.

On the other hand, we have a Taylor expansion

$$(3.138) \quad \begin{aligned} -\tilde{l}(M, \omega) &= \tilde{l}(M, \hat{\omega}) - \tilde{l}(M, \omega) \\ &= (\hat{\omega} - \omega) \frac{\partial^2 L(M, \omega)}{\partial \omega' \partial \omega} + \dots, \end{aligned}$$

where the omitted term is of order $T^{-1/2}$ relative to the term shown, because of the continuity of the third derivatives of the likelihood function (3.114) and of the function $\eta(\zeta, 0, \chi)$. The nonsingularity of the matrix

$$(3.139) \quad \frac{\partial^2 \tilde{L}(M, \omega)}{\partial \omega' \partial \omega} \equiv \tilde{L},$$

in a neighborhood of the point $M = M$ follows directly from that of the matrix Λ^\oplus defined by (3.121) and from the continuity of the derivatives involved. Therefore,

$$(3.140) \quad \hat{\omega} - \omega = -\tilde{l}(M, \omega) \tilde{L}^{-1} + \dots,$$

and from (3.137)

$$(3.141) \quad \mathcal{E}(\hat{\omega}' - \omega')(\hat{\omega} - \omega) = -T^{-1} \tilde{L}^{-1} \tilde{\Lambda} \tilde{L}^{-1} + \dots,$$

in which the omitted term is of order $T^{-1/2}$ relative to the term shown.

Owing to Theorem 3.3.4, property a), and the continuity of the relevant derivatives of the likelihood function,

$$(3.142) \quad \lim_{T \rightarrow \infty} \tilde{\Lambda} \equiv \tilde{\Lambda}_\infty$$

exists and is finite. Furthermore, because of property b) of the same theorem,

$$(3.143) \quad \text{plim}_{T \rightarrow \infty} \tilde{L} = \tilde{\Lambda}_{\infty}.$$

It follows from (3.141), (3.142), and (3.143), that

$$\text{plim}_{T \rightarrow \infty} T \mathcal{E}(\hat{\omega}' - \omega')(\hat{\omega} - \omega) = -\tilde{\Lambda}_{\infty}^{-1}.$$

This is the desired result in case χ is empty, i.e., in case $\bar{\zeta} = \bar{\omega}$ represents a complete set of unrestricted parameters. If $\bar{\zeta}$ is not complete, we shall use subscripts ζ and χ to indicate the partitioning of matrices illustrated by

$$(3.144) \quad \tilde{\Lambda} \equiv \begin{bmatrix} \tilde{\Lambda}_{\zeta\zeta} & \tilde{\Lambda}_{\zeta\chi} \\ \tilde{\Lambda}_{\chi\zeta} & \tilde{\Lambda}_{\chi\chi} \end{bmatrix}.$$

Our problem then is to evaluate

$$(3.145) \quad \text{plim}_{T \rightarrow \infty} T \mathcal{E}(\hat{\zeta}' - \zeta')(\hat{\zeta} - \zeta) = -(\tilde{\Lambda}_{\infty}^{-1})_{\zeta\zeta}$$

in terms that permit estimation on the basis of the quantities $\bar{\zeta}$ only. For this purpose we shall use the identity [Hotelling, 1943-1, p. 4]

$$(3.146) \quad (\tilde{\Lambda}^{-1})_{\zeta\zeta} = (\tilde{\Lambda}_{\zeta\zeta} - \tilde{\Lambda}_{\zeta\chi} \tilde{\Lambda}_{\chi\chi}^{-1} \tilde{\Lambda}_{\chi\zeta})^{-1}.$$

We recall the function $\hat{\chi}(M, \bar{\zeta})$ defined in (3.116), which we now need only for the argument $M = \bar{M}$. Besides the possession of a sufficient number of derivatives, the only property of this function used in the proof of Theorem 3.3.6 is that

$$(3.147) \quad \hat{\chi}(M, \zeta) = \chi.$$

At present, we must also use the property that, owing to the definition of $\hat{\chi}(M, \bar{\zeta})$,

$$(3.148) \quad \left(\frac{\partial \tilde{L}(M, \bar{\omega})}{\partial \bar{\chi}'} \right)_{\bar{\chi} = \hat{\chi}(M, \bar{\zeta})} = 0.$$

We differentiate (3.148) with respect to $\bar{\zeta}$,

$$(3.149) \quad \left(\frac{\partial^2 \tilde{L}(M, \omega)}{\partial \bar{\chi}' \partial \bar{\zeta}} + \frac{\partial^2 \tilde{L}(M, \bar{\omega})}{\partial \bar{\chi}' \partial \bar{\chi}} \frac{\partial \hat{\chi}(M, \bar{\zeta})}{\partial \bar{\zeta}} \right)_{\bar{\chi}=\hat{\chi}(M, \bar{\zeta})} = 0,$$

substitute $\bar{\zeta} = \zeta$ using (3.147), and solve as follows:

$$(3.150) \quad \frac{\partial \hat{\chi}'(M, \zeta)}{\partial \zeta} = - \tilde{\Lambda}_{\chi\chi}^{-1} \tilde{\Lambda}_{\chi\zeta},$$

using the nonsingularity of $\tilde{\Lambda}_{\chi\chi}$ which follows from the negative definiteness of Λ^{\oplus} .

On the other hand, if we write

$$(3.151) \quad \tilde{L}(M, \bar{\omega}) \equiv \tilde{L}(M, \bar{\zeta}, \bar{\chi}),$$

we have, from a comparison with (3.118),

$$(3.152) \quad L^{\oplus\oplus}(M, \zeta) = \tilde{L}\{M, \bar{\zeta}, \hat{\chi}(M, \bar{\zeta})\},$$

and hence, using (3.148),

$$(3.153) \quad \frac{\partial L^{\oplus\oplus}(M, \bar{\zeta})}{\partial \bar{\zeta}} = \left(\frac{\partial \tilde{L}(M, \bar{\zeta}, \bar{\chi})}{\partial \bar{\zeta}} \right)_{\bar{\chi}=\hat{\chi}(M, \bar{\zeta})}.$$

Differentiating once more with respect to $\bar{\zeta}$ and substituting $\bar{\zeta} = \zeta$, we obtain, using (3.150),

$$(3.154) \quad \begin{aligned} \Lambda^{\oplus\oplus} &= \frac{\partial^2 L^{\oplus\oplus}(M, \zeta)}{\partial \zeta' \partial \zeta} = \frac{\partial^2 \tilde{L}(M, \zeta, \chi)}{\partial \zeta' \partial \zeta} + \frac{\partial^2 \tilde{L}(M, \zeta, \chi)}{\partial \zeta' \partial \bar{\chi}} \frac{\partial \hat{\chi}'(M, \zeta)}{\partial \zeta} \\ &= \tilde{\Lambda}_{\zeta\zeta} - \tilde{\Lambda}_{\zeta\chi} \tilde{\Lambda}_{\chi\chi}^{-1} \tilde{\Lambda}_{\chi\zeta}. \end{aligned}$$

Comparison with (3.145) and (3.146) now yields

$$(3.155) \quad \text{plim}_{T \rightarrow \infty} T \mathcal{E}(\xi' - \zeta')(\hat{\xi} - \zeta) = -(\Lambda_{\infty}^{\oplus \oplus})^{-1}.$$

In practice, the matrix $\Lambda_{\infty}^{\oplus \oplus}$ must be estimated from a large sample. According to (3.143), a consistent estimate of $\Lambda^{\oplus \oplus}$ is also a consistent estimate of $\Lambda_{\infty}^{\oplus \oplus}$. According to Theorems 3.3.4 and 3.3.6 and the continuity of the relevant derivatives, the former quantity can again be estimated consistently by substituting M for M and maximum-likelihood estimates $\hat{\zeta}$ for ζ . This completes the proof of

THEOREM 3.3.10. *Under Assumption 3.3.1.3 (normally distributed disturbances), the product of the number of observations T and the matrix of sampling variances and covariances of the maximum-likelihood estimates $\hat{\zeta}$ of a set of parameters satisfying the conditions of Theorem 3.3.6 is consistently estimated by*

$$(3.156) \quad \text{est } T \mathcal{E}(\xi' - \zeta')(\hat{\xi} - \zeta) = - \left(\frac{\partial^2 L^{\oplus \oplus}(M, \bar{\zeta})}{\partial \bar{\zeta}' \partial \bar{\zeta}} \right)_{\bar{\zeta} = \hat{\zeta}}$$

as defined further by (3.118).

In section 4.4.13 this theorem will be used to determine sampling variances and covariances of the estimates of the parameters A in cases where the sampling variances and covariances of the estimates of Σ are not required.

4. COMPUTATION OF THE MAXIMUM-LIKELIHOOD ESTIMATES

4.1. Introductory Remarks

4.1.1. Nature of the computation problem. Apart from special cases, the equations to be satisfied by the maximum-likelihood estimates of the parameters A , Σ are essentially nonlinear and of a type that does not lend itself easily to direct solution. We shall therefore study iterative methods in which a sequence of successive approximations to the solution is obtained in such a way that the essential step in the determination of each approximation constitutes a linear problem.

The present discussion is exploratory. In section 4.5 we men-

tion several important problems that are left unsolved.

The authors wish to acknowledge very valuable help received from J. von Neumann with respect to the present problem, in the form of suggestions and advice only partially acknowledged by specific reference in what follows. Much support was also found in analogies with Hotelling's iterative method [Hotelling, 1943] for inverting a matrix.

4.1.2. *Notation.* We shall follow the rule of denoting functions of the observations by italic characters, using that notation also for the maximum-likelihood estimates, A , S , and for successive approximations, A_n , S_n , to these estimates. This notation will also be used for the initial values A_0 , S_0 , even though the latter need not (but frequently will) be functions of the observations. We shall continue to use A , Σ for the arguments of the likelihood function in general. Occasionally we shall use Θ to denote an arbitrary matrix of the same number of rows and columns as A , but which is not necessarily subject to the restrictions imposed on A .

4.1.3. *Positive definiteness of M_{xx} .* As before, we shall assume throughout that the moment matrix M_{xx} of the observed variables is positive definite. This assumption fails to be fulfilled only in cases occurring with probability zero, provided all linear identities are eliminated from the structural equation beforehand.

4.1.4. *A special case of a priori restrictions.* Before discussing the computation problem under the most general types of a priori restrictions that we have studied, it may be useful in a special and simple case to indicate a heuristic principle which has led to the computation methods discussed in what follows. In this case we assume that there is no correlation between the disturbances in different structural equations, and that normalization is imposed by taking

$$(4.1) \quad \Sigma = I.$$

We shall further assume that the only restrictions on the coefficients of the structural equations are single-parameter restrictions prescribing that certain coefficients are zero, the number of such restrictions on each equation being sufficient for its unique identification everywhere in a region N of the parameter space that contains the highest restricted maximum of the likelihood function as an internal point. The logarithm of the latter function, from

(3.6) and (4.1), is found to be, after division by T ,

$$(4.2) \quad \frac{1}{T} \log F = L(A) = \text{const} + \log \det B - \frac{1}{2} \text{tr}(A M_{xx} A').$$

4.1.5. *The iterative procedure now involves only the coefficients of the dependent variables.* It will be noted that the coefficients Γ of the predetermined variables occur only in the last (quadratic) term in (4.2). It is therefore useful first to maximize the likelihood function with respect to the nonprescribed elements of Γ only: the maximizing values $\hat{\Gamma}$ so obtained being functions of the elements of B . The last terms in (4.2) can be written as a sum of G terms of the type

$$(4.3) \quad -\frac{1}{2} \alpha(g) \cdot M_{xx} \cdot \alpha'(g) = -\frac{1}{2} \{ \beta(g) \cdot M_{yy} \cdot \beta'(g) + 2 \beta(g) \cdot M_{yz} \cdot \gamma'(g) + \gamma(g) \cdot M_{zz} \cdot \gamma'(g) \},$$

each term containing only coefficients of the corresponding structural equation indicated by g . Let the vector $\alpha^g \equiv [\beta^g \quad \gamma^g]$ be obtained from the g th row $\alpha(g) = [\beta(g) \quad \gamma(g)]$ of A by deleting all elements that are prescribed to be zero. Let $M^g \equiv M_{xx}^g \equiv M_{xx}'^g$ be obtained from M_{xx} by deleting the corresponding rows and columns. Then we wish to maximize

$$(4.4) \quad -\beta^g \cdot M_{yz}^g \cdot \gamma'^g - \frac{1}{2} \gamma^g \cdot M_{zz}^g \cdot \gamma'^g$$

by variation of γ^g only. It is easily seen that the maximizing values $\hat{\gamma}^g$ of γ^g are

$$(4.5) \quad \hat{\gamma}^g = -\beta^g \cdot M_{yz}^g \cdot (M_{zz}^g)^{-1}.$$

When these values are inserted in (4.3), (4.2) becomes

$$(4.6) \quad L^*(B) = \text{const} + \log \det B - \frac{1}{2} \sum_{g=1}^G \beta^g \cdot {}^z M_{yy}^g \cdot \beta'^g,$$

where

$$(4.7) \quad {}^z M_{yy}^g = M_{yy}^g - M_{yz}^g \cdot (M_{zz}^g)^{-1} \cdot M_{zy}^g.$$

The problem has now been reduced to finding the maximizing value B of B in (4.6). After this has been determined, the corresponding maximizing value C of Γ can be evaluated from (4.5). The computational advantage of this procedure is that the elements of C do not need to be recomputed with each iteration in the determination of B , but can be found directly from the result of the last iteration determining B .

4.1.6. *Revision of a single row.* Let B_0 represent a suitable initial value from which, through successive improvements, we attempt to reach the value B maximizing (4.6). The heuristic principle referred to above consists in revising only one row of B_0 at a time, as follows:

We write

$$(4.8) \quad B_0 = B_{0,0},$$

and determine another matrix $B_{0,1}$, which equals $B_{0,0}$ in all elements except those of the first row, the latter being determined so as to maximize (4.6). This leads to the first-order condition¹

$$(4.9) \quad d_{0,1} \cdot \text{cof } b_{0,1}^1 - b_{0,1}^1 \cdot {}^z M_{yy}^1 = 0,$$

where $\text{cof } b_{0,1}^g$ stands for a row vector containing as elements the cofactors in $B_{0,1}$ of the corresponding elements of $b_{0,1}^g$, and where the scalar quantity $d_{0,1}$ equals

$$(4.10) \quad d_{0,1} \equiv \det^{-1} B_{0,1} = (b_{0,1}^g \cdot \text{cof}' b_{0,1}^g)^{-1},$$

according to the Laplace expansion of $\det B_{0,1}$. Since the elements of $\text{cof } b_{0,1}^g$ are independent of the quantities $b_{0,1}^g$ now regarded as unknowns, we have a system (4.9) of linear equations in the unknowns $b_{0,1}^g$, $d_{0,1}$ and one quadratic equation (4.10).

Because of the positive definiteness of M_{xx} , and therefore that of ${}^z M_{yy}^1$, the unknowns $b_{0,1}^1$ can be solved uniquely from (4.9) in

¹Second-order conditions will be discussed in the general case below, see section 4.3.3.3.

terms of $d_{0,1}$ as

$$(4.11) \quad b_{0,1}^1 = d_{0,1} \cdot \text{cof } b_{0,1}^1 \cdot ({}^z M_{yy}^1)^{-1}.$$

The one remaining unknown $d_{0,1}$ is found, from (4.10) and (4.11), to satisfy

$$(4.12) \quad (d_{0,1})^2 = \{ \text{cof } b_{0,1}^1 \cdot ({}^z M_{yy}^1)^{-1} \cdot \text{cof } b_{0,1}^1 \}^{-1},$$

and can if desired be computed as the positive or negative square root of the right-hand member of (4.12). This indeterminacy of the sign of $d_{0,1}$ was to be expected, since the normalization rule (4.1) admits simultaneous changes in sign of all elements in any row of A.

4.1.7. *Successive versus simultaneous revision of rows of B_0 .* There are two important alternative ways in which the principle of revision of a row of B_0 , just described in terms of the first row, can be applied to all rows. In the first alternative, to be called successive revision of the rows of B_0 , the next step is to find a matrix $B_{0,2}$ which equals $B_{0,1}$ in all elements except those of the second row, the latter being determined again so as to maximize (4.6). In this way all rows are modified successively, the result of revision of the last row

$$(4.13) \quad B_{0,g} = B_1$$

being considered the result of the first complete iterative revision of the initial matrix B_0 .

The row-by-row revision just described, by which B_1 is obtained from B_0 , is economical for relatively small orders G of B , so that the computation of new cofactors after the revision of each row is not too laborious. Where economical, the row-by-row revision has a special flexibility in that one may depart from strict successive revision to give a higher frequency of revision to slowly converging subsets of the structural equations.

For larger systems, however, economy in terms of both quantity and standardization of work favors an alternative definition of B_1 , which requires simultaneous revision of all rows of B_0 . In this procedure, $B_{0,g}$ is defined, for every g , as being obtained from $B_{0,0}$ by the same process described above for $B_{0,1}$. Finally B_1 is such

that its g th row,

$$(4.14) \quad b_1(g) = b_{0,g}(g),$$

is equal to the g th row of $B_{0,g}$. Although one would expect slower convergence per iteration of this procedure, the saving through simultaneous computation of the cofactors of all elements cuts down the work per iteration to an extent increasing with G . All procedures studied in what follows require simultaneous revisions of all rows of B_0 . An additional reason for this choice is that generalization of the procedure to the case where no restrictions are imposed on the covariance matrix Σ of the disturbances is easier and more natural if A_1 is defined by simultaneous revision of all rows of A_0 .

4.1.8. Use of arbitrary scale factors in the approximations. A slight further saving arises if we realize that the only nonlinear operation in the procedure, viz., the computation of $d_{0,1}$ from (4.12), which serves only to determine common scale factors for the elements of each row of B_1 , need not be carried out for any iteration except the last one. [The term "scale factor" or "scale" is used here as distinct from "normalization" because it applies only to a row $\alpha(g)$ of A , not to the corresponding row $\sigma(g)$ or column $\sigma'(g)$ of Σ . While normalization is a matter of choice, the scale factor adjusting the absolute value of $\alpha(g)$ to that of $\sigma(g)$ or σ'_{gg} is of course determined by maximizing the likelihood function. See also equation (4.112) below.]

The elements of any row of B_n enter linearly and homogeneously in all relevant operations, their ratios being the relevant unknowns. For the determination of these ratios it is therefore permissible to choose any suitable value for $d_{0,g}$, for instance, unity. In what follows, it is found most suitable to substitute the known quantity $\det^{-1} B_{0,0}$ for the unknown quantity $\det^{-1} B_{0,g}$, even though the latter would give better scale factors in successive approximations. The effect of this substitution on the scale factors of the approximations B_n will be studied below. Its advantage is that each iteration is thereby completely reduced to a linear process, which can now be written

$$(4.15) \quad b_1^g \cdot z M_{yy}^g = v(g) \cdot B_0^{-1} \cdot \Phi^g'$$

Here Φ^g is a matrix of which each element is either 0 or 1, such that $\alpha(g)\Phi^g$ is a vector containing only those elements of α which are not prescribed to be zero, and $\iota(g)$ is the g th row of the identity matrix of appropriate order.

4.1.9. Saving through factorization of the likelihood function.

In what follows we consider quite general homogeneous linear restrictions on the elements of A , each restriction involving elements of one row $\alpha(g)$ only. Before specializing the restrictions on Σ , we again draw attention to the possibility that the restrictions on A and Σ taken together imply a factorization of the likelihood functions as a result of the partitioning (3.63) of these matrices. In this case it is permissible to treat the two corresponding subsystems of the structural equations separately, and a considerable saving in computation work results.

The two factors of the likelihood function indicated by (3.65) and connected with the two subsystems are of the same general form as the likelihood function for the total system of structural equations. For this reason no loss of generality is involved if we assume that the equation systems considered in the remainder of this section cannot be further reduced to subsystems in this manner.

4.1.10. Two cases regarding the a priori restrictions on Σ .

In two subsections we shall consider successively the case (4.3) where the disturbances are uncorrelated, and hence Σ is diagonal, and the case (4.4) where no restrictions at all are imposed on Σ . The latter case has simpler mathematical properties, although some of the formulae contain more terms and for that and other reasons the computations are more laborious. The former case, which was treated in the second place in the discussion of identification problems, is now taken up first. This is done mainly because the application of the heuristic principle indicated above is more straightforward in the case where Σ is restricted to be diagonal, while the experience so gained will be helpful in extending the methods to the case of an unrestricted Σ .

4.1.11. Dummy restrictions to insure identifiability. We shall assume throughout that each equation of the system is uniquely identifiable within a region N of the parameter space of which the highest restricted maximum of the likelihood function is an internal point. If this condition is not met initially, it must be met by adding dummy restrictions on the parameters as described in section 2.3. It was proved there that this can always be

done without further restricting the distribution function of the variables, in the case where Σ is unrestricted. No such proof was given for the case where Σ is required to be diagonal, because conditions of identifiability in that case were not fully analyzed. It may nevertheless be possible to apply the present computation methods in individual cases, either because identifiability can already be established without using the diagonality of Σ , or because a moderate size of the system permits the analysis of identification by *ad hoc* methods.

Before proceeding to the two specializations of the restrictions on Σ here considered, we shall in section 4.2 introduce a new formalism for the treatment of the restrictions on A which will facilitate the discussion of computation problems.

4.2. A Complete Set of Unrestricted Parameters

4.2.1. *The basic matrices Φ^g .* In previous sections, linear homogeneous a priori restrictions on the rows of A were used in the explicit form (2.24), which we rewrite¹

$$(4.16) \quad \alpha(g) \Phi_g' = 0, \quad \rho(\Phi_g) = r(\Phi_g) = R_g, \quad g = 1, \dots, G.$$

It will now be preferable to give implicit effect to the a priori restrictions on A by expressing all elements of A as linear functions of a basic set of unrestricted parameters. The most convenient choice of basic parameters is made separately for each row $\alpha(g)$ of A , by the use of an orthogonal complement Φ^g of Φ_g . The matrix Φ^g is defined, except for premultiplication by a nonsingular square matrix, through

$$(4.17) \quad \Phi^g \Phi_g' = 0, \quad \rho(\Phi^g) = r(\Phi^g) = Q_g = K_z - R_g.$$

We shall call the Φ_g , $g = 1, \dots, G$, the *restriction matrices*, and the Φ^g the *basic matrices*. It follows from (4.16) and (4.17) that

$$(4.18) \quad \Phi(g) \equiv \begin{bmatrix} \Phi^g \\ \Phi_g \end{bmatrix} \equiv [X(g) \quad \Psi(g)],$$

¹ $\rho(X)$ indicates the rank, $r(X)$ and $c(X)$ the number of rows and columns, respectively, of a matrix X .

where $X(g)$ has K_y and $\Psi(g)$ has K_z columns, is square and nonsingular and that the transformation

$$(4.19) \quad \theta(g) \equiv \theta(*g) \Phi(g) \equiv \theta^g \Phi^g + \theta_g \Phi_g$$

establishes a one-to-one correspondence between the space of the vector $\theta(g)$ and that of the vector

$$(4.20) \quad \theta(*g) \equiv [\theta^g \quad \theta_g]$$

of an equal number K_x of elements. If a vector $\alpha(g)$ satisfying the restriction (4.16) is substituted for $\theta(g)$ in (4.19), we find through postmultiplication with Φ'_g , using (4.17), that

$$(4.21) \quad \alpha_g \Phi_g \Phi'_g = 0$$

and hence, from the rank condition in (4.16), that $\alpha_g = 0$. Thus $\alpha(g)$ is expressible as

$$(4.22) \quad \alpha(g) = \alpha^g \Phi^g.$$

Conversely, any vector so expressible satisfies the restrictions (4.16). The components of the G vectors

$$(4.23) \quad \alpha^g, \quad g = 1, \dots, G,$$

represent an unrestricted set of parameters¹ except for such rules of normalization as it may be convenient to impose in certain cases.

The freedom of premultiplication by a nonsingular matrix in choosing each Φ^g should be used to make its form as simple as possible for computational purposes. Often the restrictions on A arise from the elimination of variables connected by identities to the variables retained in the system. Such elimination leads directly to the matrices Φ^g , without need for prior evaluation of the matrices Φ_g .

¹The notation α^g employed in section 4.1 represents a special case of the present notation. If matrices Φ^g were constructed for that special case, they would contain elements 1 and 0 only, with at most one 1 to each row or column. An example is contained in formula (4.128) below.

For certain purposes, especially in presenting theory, it is convenient to choose the rows of Φ^g orthogonal to each other, so that, after suitable normalization,

$$(4.24) \quad \Phi^g \Phi'^g = I.$$

From a computational point of view, such orthogonalization is often not necessary, and if carried out may increase the computational labor.

4.2.2. *Normalization.* If, as in subsection 4.3, we assume that Σ is diagonal, it is convenient for most purposes to normalize by equating the diagonal elements σ_{gg} of Σ to unity,

$$(4.25) \quad \Sigma = I.$$

For some purposes, however, it may be convenient not to restrict the σ_{gg} , but to normalize on one element of each α^g , by

$$(4.26) \quad \alpha^g \iota'(1) = 1, \quad g = 1, \dots, G,$$

say, if as before $\iota(1)$ indicates a vector of the appropriate order, of which the first element is 1 and all other elements are 0. These purposes include the calculation of sampling variances and covariances of the estimates a^g of the α^g . The normalization (4.26) will at any rate be applied in subsection 4.4 where Σ is unrestricted.

4.2.3. *The matrix A treated as a vector.* A set of parameters θ_{gk} , $g = 1, \dots, K_y$ ($\equiv G$), $k = 1, \dots, K_x$, can be considered either as a matrix

$$(4.27) \quad \Theta = \begin{bmatrix} \theta(1) \\ \vdots \\ \theta(K_y) \end{bmatrix}$$

k of K_y rows and K_x columns, or as a vector θ defined by

$$(4.28) \quad \theta \equiv \text{vec } \Theta \equiv [\theta(1) \ \theta(2) \ \cdots \ \theta(K_y)].$$

For certain operations, notably forming the determinant and the inverse of B, the matrix representation is convenient. Other operations, in particular those connected with the a priori restrictions, are simpler in the vector representation (4.28), because the restrictions are different for different rows of A. Finally, if we define M by¹

$$(4.29) \quad M \equiv \begin{bmatrix} M_{xx} & 0 & \dots & 0 \\ 0 & M_{xx} & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & M_{xx} \end{bmatrix} = I_{[G]} \otimes M_{xx},$$

and, if H, η are connected in a manner analogous to (4.28), the relations

$$(4.30) \quad \text{tr } \Theta H' = \Theta \eta', \quad \text{vec}(\Theta M_{xx}) = \Theta M, \quad \text{tr } \Theta M_{xx} \Theta' = \Theta M \Theta',$$

connect expressions of a simple type in both representations. The formal framework of the following analysis of computation problems will be an alternating use of the matrix and vector representations of the parameter space, taking advantage of the special properties of each.

4.2.4. *Projection of a matrix on the restricted parameter space.* We define the matrix

$$(4.31) \quad \Phi(*) \equiv \begin{bmatrix} \Phi^* \\ \Phi_* \end{bmatrix} \equiv \begin{bmatrix} \Phi^1 & 0 & \dots & 0 \\ 0 & \Phi^2 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \Phi^G \\ \Phi_1 & 0 & \dots & 0 \\ 0 & \Phi_2 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \Phi_G \end{bmatrix},$$

¹The symbol \otimes denotes the "direct product" or "Kronecker product" of two matrices; see [MacDuffee, p. 81].

which, owing to (4.16) and (4.17), is nonsingular and satisfies

$$(4.32) \quad \Phi^* \Phi'_* = 0.$$

Φ^* is called the *basic matrix*, Φ_* the *restriction matrix*. The transformation (4.19) for the individual vectors $\theta(g)$, $g = 1, \dots, G$, can now be summarized in

$$(4.33) \quad \theta = \theta(*) \Phi(*) = \theta^* \Phi^* + \theta_* \Phi_*,$$

where

$$(4.34) \quad \theta(*) = [\theta^* \theta_*], \quad \begin{cases} \theta^* = [\theta^1 \ \theta^2 \ \dots \ \theta^G], \\ \theta_* = [\theta_1 \ \theta_2 \ \dots \ \theta_G]. \end{cases}$$

In particular, if $\alpha = \text{vec } A$ arises from a matrix A satisfying the a priori restrictions (4.16), which in the new notation take the form

$$(4.35) \quad \alpha \Phi'_* = 0,$$

we must have

$$(4.36) \quad \alpha^* = [\alpha^1 \ \alpha^2 \ \dots \ \alpha^G], \quad \alpha_* = 0.$$

Thus, under the normalization rule (4.25), the elements of the vector α^* constitute a complete set of unrestricted parameters through which the original restricted parameters are expressed by

$$(4.37) \quad \alpha = \alpha^* \Phi^*.$$

Under the normalization (4.26), the complete set of unrestricted parameters consists of those elements of α^* not prescribed by (4.26), plus the other diagonal elements σ_{gg} , or all elements σ_{gh} , $h \geq g$, of the symmetric matrix Σ , according to the case considered.

Through (4.33) the arbitrary vector θ is expressed uniquely as the sum of two vectors, the first $\theta^* \Phi^*$ lying within the restricted parameter space, the second orthogonal to that space. This decomposition plays an essential role in what follows. Since it will also be applied to cases where θ is given as a matrix

product, it is convenient to introduce, in addition to the notations used in (4.33), the operator notations

$$(4.38) \quad \theta(*) \equiv \text{vec}^* \theta, \quad \theta^* \equiv \text{vec}^* \theta, \quad \theta_* \equiv \text{vec}_* \theta.$$

In these terms (4.33) runs

$$(4.39) \quad \theta = \text{vec} \theta = (\text{vec}^* \theta) \Phi^* + (\text{vec}_* \theta) \Phi_*,$$

from which we can solve for $\text{vec}^* \theta$ through postmultiplication by $\Phi'^*(\Phi^* \Phi'^*)^{-1}$ using (4.32), thus obtaining

$$(4.40) \quad \text{vec}^* \theta = (\text{vec} \theta) \Phi'^*(\Phi^* \Phi'^*)^{-1}.$$

Conversely, we define

$$(4.41) \quad \theta \equiv \text{mat} \theta \equiv \text{mat}^* \theta(*) \equiv \text{mat}^* \text{vec}^* \theta,$$

and

$$(4.42) \quad \text{mat}^* \theta^* \equiv \text{mat}^* [\theta^* \quad 0_*], \quad \text{mat}^* \theta_* \equiv \text{mat}^* [0^* \quad \theta_*].$$

The decomposition (4.39) can thus be written in matrix coordinates as

$$(4.43) \quad \theta = \mathcal{Q} \theta + \mathcal{R} \theta,$$

where

$$(4.44) \quad \mathcal{Q} \theta \equiv \text{mat}^* \text{vec}^* \theta, \quad \mathcal{R} \theta \equiv \text{mat}^* \text{vec}_* \theta.$$

The operation $\mathcal{Q} \theta$ can be regarded as a projection of the matrix θ on the restricted parameter space. In particular, the a priori restrictions (4.35) on the matrix A are equivalent to

$$(4.45) \quad A = \mathcal{Q} A, \quad \text{or} \quad \mathcal{R} A = 0.$$

We shall operate mostly in the vector space of θ^* , but occasionally return to the restricted matrix space of $\mathcal{Q} \theta$ through the transformation mat^* .

If the matrix H satisfies the restrictions $\eta_* = 0$, we have the

important property that

$$(4.46) \quad \text{tr } \Theta H' = \Theta \eta' = (\Theta^* \Phi^* + \Theta_* \Phi_*) \Phi'^* \eta'^* = \Theta^* \Phi^* \Phi'^* \eta'^*$$

because of (4.32) does not depend on Θ_* . In particular we have:

LEMMA 4.2.4. *A necessary and sufficient condition that $\text{tr } \Theta H' = 0$ for all values of H satisfying the a priori restrictions $\eta_* = \text{vec}_* H = 0$ is that $\Theta^* \equiv \text{vec}^* \Theta = 0$.*

The proof follows from (4.46) and the nonsingularity of $\Phi^* \Phi'^*$ due to the rank condition in (4.17).

4.3. The Case of Uncorrelated Disturbances

4.3.1. The nature of the problem.

4.3.1.1. *The maximum-likelihood equations.* In the present subsection, the matrix Σ of variances and covariances of the disturbances is assumed to be diagonal. Unless otherwise stated, normalization will be based throughout on (4.25) where Σ is equated to the unit matrix.

We shall now write down the first-order conditions for a maximum of the logarithmic likelihood function (4.2) which we rewrite:

$$(4.47) \quad L(A) = \text{const} + \log \det B - \frac{1}{2} \text{tr}(A M_{xx} A').$$

Let A_0 be a trial value of A satisfying the restrictions $\alpha_* = 0$, and write $(\delta A_0$ here denoting a *finite* change in A_0)

$$(4.48) \quad A = A_0 + \delta A_0, \quad \text{vec}_* \delta A_0 = 0,$$

which insures that A again satisfies the restrictions. The Taylor expansion of $L(A)$ in the neighborhood of $\delta A_0 = 0$ then contains the following constant and linear terms, derived with the use of (3.16) and (3.17),

$$(4.49) \quad \begin{aligned} L(A) &= L(A_0) + \text{tr}(B_0'^{-1} (\delta B_0)') - \text{tr}(A_0 M_{xx} (\delta A_0)') + \dots \\ &= L(A_0) + \text{tr} \{ (B_0'^{-1} I_{[k_y, k_x]} - A_0 M_{xx}) (\delta A_0)' \} + \dots \end{aligned}$$

(Here $I_{[K_y K_x]}$ is a submatrix, of K_y rows and K_x columns, of the unit matrix of order K_x , as follows: $I_{[K_y K_x]} = [I_{[K_y]} \quad 0]$, the first matrix in the right-hand member being the unit matrix of order K_y). According to Lemma 4.2.4, the necessary first-order condition for A_0 to coincide with a restricted maximum A of the likelihood function is therefore

$$(4.50) \quad \begin{cases} (4.50q) & \text{vec}^* (B'^{-1} I_{[K_y K_x]} - A M_{xx}) = 0, \\ (4.50r) & a_* \equiv \text{vec}_* A = 0. \end{cases}$$

These are the maximum-likelihood equations that are to be solved by an iterative process. If desired, the operators vec^* and vec_* can be replaced by \mathcal{Q} and \mathcal{R} , thus reverting to a matrix form.

It will be noted that the number of conditions equals the number of unknowns. If identification is incomplete, the equations become interdependent.

4.3.1.2. *Solutions without restrictions on A.* Further light is thrown on the mathematical nature of this problem if we first consider the case where no restrictions at all are imposed on the matrix A . Then a_* has no elements and the symbol vec^* can be omitted in (4.50q), so that

$$(4.51) \quad \begin{cases} (4.51y) & B'^{-1} - A M_{xy} = 0, \\ (4.51z) & A M_{xz} = 0. \end{cases}$$

Of these equations (4.51z) is solved by expressing A linearly in terms of an orthogonal complement of $M'_{xz} = M_{zx}$ as follows:

$$(4.52) \quad A = B [I_{[K_y]} \quad -M_{yz} M_{zz}^{-1}],$$

with B an arbitrary nonsingular matrix of order K_y ; Substituting this result in (4.51) we obtain as the condition on B

$$(4.53) \quad B'^{-1} - B(M_{yy} - M_{yz} M_{zz}^{-1} M_{zy}) = 0,$$

or

$$(4.54) \quad B' B = (M_{yy} - M_{yz} M_{zz}^{-1} M_{zy})^{-1},$$

which is solved, but for an arbitrary orthogonal matrix O , by

$$(4.55) \quad B = O(M_{yy} - M_{yz} M_{zz}^{-1} M_{zy})^{-1/2},$$

provided the inverse square root is taken to be symmetric. The right-hand member in (4.54) can also be denoted by $(M_{xx}^{-1})_{yy}$. The solutions (4.52) and (4.55) can also be obtained from the conditions (3.58) for the absolute maximum of the likelihood function by choosing a value B of B which will make Σ equal to the unit matrix as at present required.

It thus appears that in the absence of all restrictions on A , our problem is of the nature indicated by (4.54), leaving an arbitrary orthogonal matrix O in the solution. If identifying restrictions are now added gradually, more and more restrictions are imposed on O . The existence of a solution A of (4.51) within the restrictions is the necessary and sufficient condition for the likelihood function to be able to attain its absolute maximum. Since we assume here that the total set of restrictions (a priori and dummy combined) identifies each structural equation, there is just one special case in which a solution of (4.51) is still possible. This is the case in which the total set of restrictions is just adequate in number and variety, in accordance with Definition 3.6, to identify all structural equations. In that case a possible computation procedure would be to find one particular solution of (4.54) and then to determine O in such a way that A satisfies the restrictions. However, even in this case, computational economy may still favor the iterative methods developed below.

As soon as the restrictions are more than adequate in number and variety with respect to at least one structural equation, equation (4.54) cannot in general be satisfied any longer. We then have a more general problem where (4.54) must, in some sense, be satisfied as nearly as possible within the linear restrictions on B arising from those on A .

4.3.2. The methods \mathcal{P}_1 , \mathcal{P}_h , and \mathcal{P}_n .

4.3.2.1. *Choice of the linear path toward the next approximation.* Suppose now that in A_0 a restricted stationary value of the

likelihood function is not reached, but in point A in the neighborhood of A_0 a restricted maximum is reached. If we write as the next approximation

$$(4.56) \quad A_1 = A_0 + h \Delta A_0,$$

the matrix ΔA_0 indicates the direction on the linear path on which the next approximation is sought and the scalar h determines the distance traveled along that path. In sections 4.3.2 to 4.3.4 inclusive, we shall be concerned with methods based on the following choice of ΔA_0 :

$$(4.57) \quad \begin{cases} (4.57^*) & \text{vec}^* (\Delta A_0 \cdot M_{xx}) = \text{vec}^* (B_0'^{-1} I_{[K_y, K_x]} - A_0 M_{xx}), \\ (4.57_*) & \Delta a_{0*} = \text{vec}_* \Delta A_0 = 0. \end{cases}$$

We shall first show that a sufficiently small value of h will always lead to an increase in the likelihood function. If $h \Delta A_0$ is substituted for δA in (4.49), we have, on account of (4.49), (4.57), and (4.46),

$$(4.58) \quad \begin{aligned} L(A_1) - L(A_0) &= h \text{tr} \{ (B_0'^{-1} I_{[K_y, K_x]} - A_0 M_{xx}) (\Delta A_0)' \} + \dots \\ &= h \{ \text{vec}^* (B_0'^{-1} I_{[K_y, K_x]} - A_0 M_{xx}) \} \Phi^* \cdot \Phi'^* \cdot \Delta a_0^* + \dots \\ &= h \{ \text{vec}^* (\Delta A_0 \cdot M_{xx}) \} \Phi^* \cdot \Phi'^* \cdot \Delta a_0^* + \dots \\ &= h \text{tr} (\Delta A_0 \cdot M_{xx} \cdot \Delta A_0') + \dots \end{aligned}$$

Because of the positive definiteness of M_{xx} , the coefficient of h in this expansion is positive unless ΔA_0 vanishes identically, in which case a restricted stationary value of the likelihood function would already have been reached in A_0 . If h is sufficiently small but positive, therefore, the linear term in (4.58) will exceed in absolute value the sum of all subsequent terms, and $L(A_1) > L(A_0)$.

4.3.2.2. *Choices of the next approximation on the linear path selected. The process \mathcal{P}_1 .* We shall postpone until section 4.3.2.5 the proof that (4.57) always admits of one and only one solution ΔA_0 , and now discuss possible choices of h . We shall first show that the process obtained through the choice

$$(4.59) \quad h = 1,$$

now to be called the process \mathcal{P}_1 , is equivalent to the process demonstrated in a special case in section 4.1. For that choice of h , A_1 is defined by

$$(4.60) \begin{cases} (4.60^*) & \text{vec}^* (A_1 M_{xxx}) = \text{vec}^* [B_0'^{-1} \quad 0], \\ (4.60_*) & a_{1*} = \text{vec}_* A_1 = 0. \end{cases}$$

Furthermore, in the case referred to, all restrictions are of the single-parameter type which require certain elements of A to vanish. In this case, a suitable choice for each Φ^g in (4.22) is obtained by deleting from the unit matrix $I_{[K_x]}$ all rows corresponding to elements in $\alpha(g)$ that are required to vanish. We shall return (4.60) to matrix form by applying mat^* to all members:

$$(4.61) \begin{cases} (4.61\mathcal{L}) & \mathcal{L}(A_1 M_{xxx}) = \mathcal{L}[B_0'^{-1} \quad 0], \\ (4.61\mathcal{R}) & \mathcal{R}A_1 = 0. \end{cases}$$

because in the present case the operation \mathcal{L} consists simply in replacing by zero all elements of the matrix on which \mathcal{L} operates corresponding to elements of A that are required to vanish by (4.61 \mathcal{R}). In this case, therefore, if the partitioning of Θ corresponding to $A = [B \quad \Gamma]$ is denoted momentarily by $\Theta = [H \quad Z]$, the definition of \mathcal{L} can be extended to submatrices of Θ through

$$(4.62) \quad \mathcal{L}\Theta \equiv [\mathcal{L}H \quad \mathcal{L}Z],$$

and (4.61 \mathcal{L}) partitions into

$$(4.63\mathcal{L}) \begin{cases} (4.63\mathcal{L}_y) & \mathcal{L}(A_1 M_{xy}) = \mathcal{L}B_0'^{-1}, \\ (4.63\mathcal{L}_z) & \mathcal{L}(A_1 M_{xz}) = 0. \end{cases}$$

Of these conditions, (4.63 \mathcal{L}_z) is equivalent to (4.5), and (4.63 \mathcal{L}_y) to (4.15). This result establishes a presumption that \mathcal{P}_1 will have satisfactory convergence properties, except perhaps with respect to a common scale factor for each row of A_n .

4.3.2.3. *The process $\mathcal{P}_{1/2}$.* In a special borderline case the

choice $h = \frac{1}{2}$ leading to the process $\mathcal{P}_{1/2}$ has superior convergence properties. This is the case where there are no predetermined variables and no restrictions on the matrix B. (This, of course, implies dropping, for the present example, the previous assumption that each equation is completely identified.) In this case all matrices involved are square matrices of order G , and $A \equiv B$, $M_{xx} \equiv M_{yy}$. The process $\mathcal{P}_{1/2}$ now runs

$$(4.64) \quad A_1 = \frac{1}{2}(A_0'^{-1} \cdot M_{xx}^{-1} + A_0).$$

Simple calculations will show that this process possesses the property

$$(4.65) \quad A_1 M_{xx} A_1' - I = \frac{1}{4}(A_0 M_{xx} A_0' - I)(A_0 M_{xx} A_0')^{-1}(A_0 M_{xx} A_0' - I).$$

Since under the present assumptions the maximum-likelihood equations (4.51) to be solved iteratively are equivalent to

$$(4.66) \quad A M_{xx} A' = I,$$

(4.65) implies a high rate of convergence of $\mathcal{P}_{1/2}$, once convergence is obtained initially. The extent to which A_1 fails to satisfy (4.66), as measured by $A_1 M_{xx} A_1' - I$, is of second order compared with the corresponding quantity $A_0 M_{xx} A_0' - I$ in terms of A_0 . It follows that the number of decimal places which is correct in the n th iteration increases geometrically with n .

It will be clear that in cases where $\mathcal{P}_{1/2}$ possesses this very desirable property, a process \mathcal{P}_h based on any other constant value of h will produce a lower-order speed of convergence.

Of course, the solution of (4.66) is determined but for an orthogonal transformation, and it depends on the initial value A_0 which particular solution of (4.66) is approached by successive iterations. If the order G of the matrix A is reduced to one, the indeterminacy disappears, and (4.64) is specialized to a well-known iterative procedure,

$$(4.67) \quad a_1 = \frac{1}{2}\left(\frac{m^{-1}}{a_0} + a_0\right),$$

to obtain the square root of a scalar m^{-1} .

4.3.2.4... *The process* \mathcal{P}_{h_n} . Von Neumann has suggested determining h afresh with each iteration from the requirement that the sum of the linear and quadratic terms in the Taylor expansion of the likelihood function with respect to h shall have a vanishing first derivative. Extending the expansion (4.58), with the aid of (3.18), to

$$(4.68) \quad \begin{aligned} L(A_1) - L(A_0) = & h \operatorname{tr}\{(B_0'^{-1} I_{[K_y, K_x]} - A_0 M_{xx})(\Delta A_0)'\} \\ & + \frac{1}{2} h^2 \operatorname{tr}\{-B_0'^{-1} (\Delta B_0)' B_0'^{-1} (\Delta B_0)' \\ & - (\Delta A_0) M_{xx} (\Delta A_0)'\} + \dots, \end{aligned}$$

and using, as in (4.58), the definition of ΔA_0 given by (4.57), we find that the value,

$$(4.69) \quad h_0 = \frac{\operatorname{tr} \Delta A_0 \cdot M_{xx} \cdot (\Delta A_0)'}{\operatorname{tr} \left\{ (B_0'^{-1} (\Delta B_0)')^2 + \Delta A_0 \cdot M_{xx} \cdot (\Delta A_0)' \right\}},$$

of h satisfies the criterion mentioned. This procedure, which we denote by \mathcal{P}_{h_n} , may be expected to have an asymptotic speed of convergence superior to that obtained by any constant value of h . In particular, in cases like the preceding example where one particular constant value of h leads to a higher-order speed of convergence than all other constant values, the value h_n according to (4.69) may be expected to converge to that constant for $n \rightarrow \infty$.

4.3.3. Asymptotic convergence properties of

$$\mathcal{P}_1, \mathcal{P}_h, \mathcal{P}_{h_n}$$

4.3.3.1. *Orthogonalizing the basic matrix* Φ^* . The foregoing tentative discussion of various processes can be made more definite by studying their properties in the neighborhood of a restricted maximum A of the likelihood function. For this purpose, it is useful to assume that the basic matrix Φ^* defined by (4.31) is orthogonal according to (4.24). If the restriction matrix Φ_* is similarly orthogonalized, we have

$$(4.70) \quad \Phi(*) \Phi'(*) = I, \quad \text{so} \quad \Phi^{-1}(*) = \Phi'(*) .$$

The relationships (4.30) can then be extended in a simple way to the θ -space. We register for future use:

$$(4.71) \quad \begin{aligned} \text{tr}(\Theta H') &= \theta(*) \eta'(*), & \text{vec}*(\Theta M_{xx}) &= \theta(*) M(*), \\ \text{tr}(\Theta M_{xx} \Theta') &= \theta(*) M(*) \theta'(*), \end{aligned}$$

where

$$(4.72) \quad M(*) = \Phi(*) M \Phi'(*) = M'(*) = \begin{bmatrix} M^{**} & M^*_* \\ M^*_* & M_{**} \end{bmatrix} .$$

In particular we have, if H satisfies the restrictions $\eta_* = 0$, the important identities

$$(4.73) \quad \begin{aligned} \text{tr} \Theta H' &= \theta^* \eta'^* \\ \text{vec}*(H M_{xx}) &= \eta^* M^{**}, \\ \text{tr}(H M_{xx} H') &= \eta^* M^{**} \eta'^* . \end{aligned}$$

Incidentally, it follows from the second equality in (4.73) and the nonsingularity of M^{**} that (4.57) always admits of one and only one solution $\Delta \bar{A}_0$, a statement previously made without proof. This conclusion remains true even when identification is incomplete, since the nonsingularity of M^{**} is not affected thereby.

4.3.3.2. *Analysis of asymptotic convergence properties.* Let successive iterations be

$$(4.74) \quad A_n = A + \bar{A}_n, \quad \bar{a}_{n*} = \text{vec}_* \bar{A}_n = 0, \quad n = 0, 1, \dots .$$

We shall only consider linear and sometimes quadratic terms in Taylor expansions with respect to \bar{A}_0, \bar{A}_1 . In terms of \bar{A}_n the iterative process (4.56) runs

$$(4.75) \quad \bar{A}_1 = \bar{A}_0 + h(\Delta \bar{A}_0),$$

where $\Delta \bar{A}_0$ is identical with ΔA_0 as defined by (4.57). The Taylor expansion of (4.57*) now becomes, because of (4.50q),

$$(4.76) \quad \text{vec}^* (\Delta \bar{A}_0 M_{xx}) = \text{vec}^* (-B'^{-1} \bar{B}'_0 B'^{-1} I_{[K_y K_x]} - \bar{A}_0 M_{xx}) + \dots$$

We shall study this relation in the space of the unrestricted vectors $\bar{a}_0^* \equiv \text{vec}^* \bar{A}_0$ and $\Delta \bar{a}_0^* \equiv \text{vec}^* \Delta \bar{A}_0$. The term shown in the right-hand member of (4.76) is most easily understood in relation to the quadratic term in \bar{A}_0 in the Taylor expansion of the likelihood function $L(A_0)$ based on the point A . We now write for the latter term

$$(4.77) \quad \begin{aligned} \frac{1}{2} L_{(2)}(A_0) &\equiv \frac{1}{2} \text{tr}(-B'^{-1} \bar{B}'_0 B'^{-1} \bar{B}'_0 - \bar{A}_0 M_{xx} \bar{A}'_0) \\ &\equiv \frac{1}{2} \bar{a}_0^* L^* \bar{a}_0^*, \end{aligned}$$

thereby uniquely defining a symmetric matrix L^* . The explicit evaluation of L^* is immaterial at this stage, and will be demonstrated below in formulae (4.183) and (4.188) dealing with an analogous matrix L_0^* .

Differentiating the middle member of (4.77) with respect to \bar{a}_0^* , we have on the one hand, using the first relation (4.73), a row vector

$$(4.78) \quad \begin{aligned} \frac{1}{2} \frac{dL_{(2)}}{d\bar{a}_0^*} &= \frac{\frac{1}{2} \text{tr} \left(\frac{dL_{(2)}}{d\bar{A}_0} d\bar{A}'_0 \right)}{d\bar{a}_0^*} \\ &= \frac{\text{tr} \{ (-B'^{-1} \bar{B}'_0 B'^{-1} I_{[K_y K_x]} - \bar{A}_0 M_{xx}) d\bar{A}'_0 \}}{d\bar{a}_0^*} \\ &= \frac{\text{vec}^* (-B'^{-1} \bar{B}'_0 B'^{-1} I_{[K_y K_x]} - \bar{A}_0 M_{xx}) d\bar{a}_0^*}{d\bar{a}_0^*} \\ &= \text{vec}^* (-B'^{-1} \bar{B}'_0 B'^{-1} I_{[K_y K_x]} - \bar{A}_0 M_{xx}). \end{aligned}$$

Comparing this with the result of differentiating the last member of (4.77), we find that (4.76) is equivalent to

$$(4.79) \quad \Delta \bar{a}_0^* M^{**} = \bar{a}_0^* L^* + \dots,$$

in which the first member is obtained through the second relation (4.73). Through one further transformation

$$(4.80) \quad \begin{aligned} M^{**} &\equiv R^* R'^*, & L^\dagger &\equiv (R^*)^{-1} L^* (R^*)'^{-1}, \\ \Delta \bar{a}_0^\dagger &\equiv \Delta \bar{a}_0^* R^*, & \bar{a}_0^\dagger &\equiv \bar{a}_0^* R^*, \end{aligned}$$

this goes over into

$$(4.81) \quad \Delta \bar{a}_0^\dagger = \bar{a}_0^\dagger L^\dagger + \dots,$$

so that

$$(4.82) \quad \bar{a}_1^\dagger = \bar{a}_0^\dagger (I + hL^\dagger) + \dots$$

In the absence of higher-order terms, this iterative process has been studied for constant h by Hotelling [1933, 1936]. Its properties depend on the characteristic values k_q , $q = 1, \dots, Q$, of the matrix

$$(4.83) \quad K^\dagger \equiv I + hL^\dagger.$$

The choice of R^* in (4.80) is determined except for postmultiplication by an orthogonal matrix, and this freedom can be used to make L^\dagger , and hence K^\dagger , diagonal. The diagonal elements of these matrices then equal the characteristic values l_q and k_q , respectively, which are connected by

$$(4.84) \quad k_q = 1 + hl_q.$$

We are free to arrange the values (4.84) in descending order of algebraic magnitude¹ through suitable choice of R^* . The elementary vectors

$$(4.85) \quad v(q) = [0_1 \quad \dots \quad 0_{q-1} \quad 1_q \quad 0_{q+1} \quad \dots \quad 0_Q], \quad q = 1, \dots, Q,$$

¹Since we restrict ourselves to positive values of h , the descending order applies simultaneously to l_q and k_q .

form an orthogonal set of corresponding characteristic vectors, i.e.,

$$(4.86) \quad \begin{aligned} \iota(q)K^\dagger &= k_q \iota(q), & \iota(q)L^\dagger &= l_q \iota(q), \\ \iota(q_1) \iota'(q_2) &= \begin{cases} 0 & \text{if } q_1 \neq q_2, \\ 1 & \text{if } q_1 = q_2. \end{cases} \end{aligned}$$

If

$$(4.87) \quad \bar{a}_n^\dagger = \sum_{q=1}^Q \{\bar{a}_n^\dagger \iota'(q)\} \iota(q),$$

the iterative process (4.82) consists - to the first order of magnitude - of a multiplication

$$(4.88) \quad \bar{a}_1^\dagger \iota'(q) = k_q \{\bar{a}_0^\dagger \iota'(q)\}$$

of the q th component $\bar{a}_0^\dagger \iota'(q)$ underlying \bar{A}_0 by the factor k_q . Apart from the effect of higher-order terms (which is smaller, the nearer the initial value A_0 is to the solution A), convergence is assured if all characteristic values k_q are smaller than unity in absolute value.

4.3.3.3. *Identifiability of structural equations and nonsingularity of L^\dagger .* We have assumed that in A a maximum of the likelihood function is reached under completely identifying restrictions. In connection with the definition (4.77) of L^* , the fact that A is a maximum precludes the possibility that any characteristic value l_q of L^\dagger be positive. It does not necessarily preclude singularity of L^\dagger , since a maximum might be reached in a point where second derivatives vanish in a certain direction (e.g., the function $-x^4$ in the point $x = 0$). However, unless all third derivatives of the likelihood function vanish in A , complete identifiability of structural equations implies that L^* will be negative definite. If that is the case, all characteristic values l_q are negative, and a sufficiently small value of h will insure that $|k_q| < 1$ for all $q = 1, \dots, Q$, and will thus insure convergence from initial values A_0 sufficiently near to A . It is also clear from (4.84) that too small a value of h will make convergence quite

slow.

4.3.3.4. *Importance of the case in which the maximum of the likelihood function is not depressed by the restrictions.* In one important special case it is not difficult to make further statements about the characteristic values of L^\dagger , and therefore about those of the matrix (4.83). This is the case, discussed in section 3.2.2, in which the a priori restrictions do not depress the likelihood function (4.47) below its absolute maximum. It was shown above that this will be the case if and only if (4.51) permits a solution within the restrictions.

It was also noted that, since we are assuming complete identification of each structural equation, this case can occur identically in the sample space only if the restrictions are just adequate in number and variety for complete identification. This will rarely be so in systems of appreciable size, unless the investigator chooses to ignore the excess of a priori information over the minimum essential for identification. If excess information is available and is used, the case of a nondepressed likelihood function can still occur with probability zero in the sample space if a sample is drawn with "exceptional" values of M_{xx} that permit a restricted solution A of (4.51). Again this remark would be of little practical value, were it not that, in a large majority of sufficiently large samples, values of M_{xx} are obtained that are not far removed from a value M_{xx} that permits the likelihood function to reach its absolute maximum, provided the a priori information embodied in the restrictions is actually valid in the population. For in that case, Theorems 3.3.5 and 3.3.4 apply. In the first place, if the expectation,

$$(4.89) \quad M_{xx} = \mathcal{E}M_{xx},$$

of the moment matrix is inserted in the likelihood function, the absolute maximum can be reached by inserting the true values of the parameters A , Σ , values which obviously tally with any valid a priori information about these parameters. Secondly, we have

$$(4.90) \quad \text{plim}(M_{xx} - M_{xx}) = 0.$$

Thus any statement about the characteristic values of L^\dagger that is based on the assumption that (4.51) possesses a solution under valid restrictions is approximately true with high probability in sufficiently large samples.

4.3.3.5. *Characteristic values of L^* if the maximum of the likelihood function is not depressed.* It is well known that the characteristic values l_q of L^\dagger are the stationary values of the quadratic forms in \bar{a}^\dagger ,

$$(4.91) \quad L_{(2)} = \bar{a}^\dagger L^\dagger \bar{a}'^\dagger = \text{tr}\{-(B'^{-1} \bar{B}')^2 - \bar{A} M_{xx} \bar{A}'\},$$

under the restrictions on \bar{a}^\dagger ,

$$(4.92) \quad 1 = \bar{a}^\dagger \bar{a}'^\dagger = \text{tr} \bar{A} M_{xx} \bar{A}'.$$

Since a one-to-one correspondence has been established between the vectors \bar{a}^\dagger and the matrices \bar{A} satisfying the a priori restrictions, the values l_q are also the stationary values of the last member of (4.91) subject to the a priori restrictions on \bar{A} plus the restrictions (4.92).

Let us now make use of the assumption that A satisfies (4.51) to supplement it to a nonsingular matrix,

$$(4.93) \quad H \equiv \begin{bmatrix} A \\ D \end{bmatrix} \equiv \begin{bmatrix} B & C \\ 0 & F \end{bmatrix}, \quad \text{where} \quad F' F = M_{zz}^{-1},$$

which is such that

$$(4.94) \quad H M_{xx} H' = I, \quad \text{or} \quad M_{xx} = H^{-1} H'^{-1}.$$

If we now transform \bar{A} uniquely by

$$(4.95) \quad \bar{A} \equiv \tilde{A} H \equiv \tilde{B} A + \tilde{C} D, \quad \text{or} \quad \tilde{A} \equiv \bar{A} H^{-1} \equiv [\tilde{B} \quad \tilde{C}],$$

the linear a priori restrictions on \bar{A} entail similar restrictions on \tilde{A} . We shall refer to the latter restrictions as the \tilde{A} -restrictions, it being implied in the use of this expression that the a priori restrictions on A permit the likelihood function to reach its absolute maximum. Upon inserting (4.95) in (4.91) and (4.92) and using (4.93) and (4.94), we find that the l_q are the stationary values of the quadratic form

$$(4.96) \quad L_2 = \text{tr}(-\tilde{B}'^2 - \tilde{A} \tilde{A}') = - \sum_{g,h=1}^{K_y} \tilde{a}_{gh} \tilde{a}_{hg} - \sum_{g=1}^{K_y} \sum_{k=1}^{K_x} \tilde{a}_{gk}^2,$$

subject to the \tilde{A} -restrictions and the additional restriction

$$(4.97) \quad \text{tr} \tilde{A} \tilde{A}' = \sum_{g=1}^{K_y} \sum_{k=1}^{K_x} \tilde{a}_{gk}^2 = 1,$$

where \tilde{a}_{gk} denotes the element of \tilde{A} in the g th row and k th column.

Let us first revert to the case where no a priori restrictions are imposed on \tilde{A} and hence none on \tilde{A}' . Then the values $l_q, q = 1, \dots, P$, are those values of l that permit a solution \tilde{A} of

$$(4.98) \quad \frac{1}{2} \frac{d}{d\tilde{A}} (L_{(2)} - l \text{tr} \tilde{A} \tilde{A}') = -\tilde{B}' I_{[K_y K_x]} - (1 + l)\tilde{A} = 0.$$

The following table, whose entries are easily verified by substitution, contains all possible solutions, since the sum of the multiplicities of the characteristic values (determined as the corresponding number of linearly independent characteristic "vectors" A) equals $P = K_y K_x$.

	(a)	(b)	(c)	(d)
	Value of l	Value of $k=1+kl$	Multiplicity	Characteristic "vectors" \tilde{A} satisfy
(4.99)	0	1	$\frac{1}{2} K_y (K_y - 1)$	$\tilde{B} = -\tilde{B}', \tilde{C} = 0$
	-1	$1 - k$	$K_y K_x$	$\tilde{B} = 0, \tilde{C}$ arbitrary
	-2	$1 - 2k$	$\frac{1}{2} K_y (K_y + 1)$	$\tilde{B} = \tilde{B}', \tilde{C} = 0$

Since the values $l = 0$ and $l = -2$ are the extrema of the form (4.96) under the sole restriction (4.97), we have

$$(4.100) \quad -2 \leq L_{(2)} \leq 0 \quad \text{if} \quad \text{tr} \tilde{A} \tilde{A}' = 1.$$

Consequently, if a priori restrictions are now introduced that do not restrict the maximum of the likelihood function, the new values

l_q and k_q must satisfy

$$(4.101) \quad -2 \leq l_q \leq 0, \quad 1 - 2h \leq k_q \leq 1, \quad q = 1, \dots, Q.$$

Any particular characteristic value in (4.98) will remain a characteristic value under the a priori restrictions if the ensuing \tilde{A} -restrictions permit the corresponding condition in (4.99), column (d), to be satisfied. Its multiplicity then equals the number of independent "vectors" \tilde{A} satisfying that condition and the restrictions.¹

4.3.3.6. *Exclusion of the characteristic value $l = 0$ through complete identification.* Under the present restrictions (4.25) on Σ , but in the absence of any restrictions on A , the transformations

$$(4.102) \quad A^\oplus = \Upsilon A$$

preserving the form of the likelihood function are orthogonal:

$$(4.103) \quad \Upsilon \Upsilon' = I.$$

Let $A + \tilde{A}$ be obtained from A by such a transformation. Then

$$(4.104) \quad \tilde{C} = 0, \quad \tilde{B} = \Upsilon - I,$$

and, up to the first-order terms in \tilde{B} ,

$$(4.105) \quad \Upsilon \Upsilon' = (I + \tilde{B})(I + \tilde{B}') = I + \tilde{B} + \tilde{B}' + \dots = I,$$

so that in first approximation $\tilde{B} = -\tilde{B}'$. Thus the characteristic vectors (4.99) associated with $l = 0$ represent only directions of change from A corresponding to orthogonal transformations of A .

In the case of complete (unique or multiple) identification of all structural equations, all such transformations leading from a point A satisfying the restrictions to a point $A + \tilde{A}$ in a sufficiently small neighborhood of A are excluded by the restrictions on $A + \tilde{A}$. In that case, therefore, the characteristic value $l = 0$ is no longer present, and

¹Except that, with probability zero in the sample space, a "new" characteristic value $l = -1$ may be added, with a "vector" that is a linear combination of "vectors" corresponding to $l = -2$ and $l = 0$, respectively, in (4.99).

$$(4.106) \quad -2 \leq l_q < 0, \quad 1 - 2h \leq k_q < 1, \quad q = 1, \dots, Q.$$

4.3.3.7. *Complete identification through restrictions on A permits only trivial solutions with $l = -2$.* Let us now assume that identification is now complete on the basis of the restrictions on A alone, i.e., without recourse to (4.25). Then all points

$$(4.107) \quad A + \bar{A} = (I + \tilde{B})A, \quad (\tilde{C} = 0),$$

permitted by the restrictions on $A + \bar{A}$ are such that

$$(4.108) \quad \tilde{B} \text{ is diagonal.}$$

Since one can select only K_y linearly independent diagonal matrices of order K_y , the value $l = -2$ now has its multiplicity reduced to K_y ,

$$(4.109) \quad l_q = -2, \quad q = Q - K_y + 1, \dots, Q.$$

The accompanying characteristic "vectors" (4.108) correspond to the application of arbitrary scale factors $1 + c_{gg}$ to the rows $a(g)$ of A.

It follows that the remaining characteristic values now satisfy

$$(4.110) \quad -2 < l_q < 0, \quad 1 - 2h < k_q < 1, \quad q = 1, \dots, Q - K_y,$$

and that the corresponding characteristic "vectors" satisfy $\tilde{C} \neq 0$.

4.3.3.8. *Asymptotic properties of \mathcal{P}_1 .* The foregoing analysis suggests that, among processes \mathcal{P}_h employing a constant value of h , \mathcal{P}_1 is suitable if we have a large sample from a distribution satisfying known restrictions on the matrix A that are sufficient for its identification. For all relevant characteristic values are smaller in absolute value than unity, and convergence can be expected if the initial value A_0 is sufficiently close to the solution A. The only characteristic of A that does not participate in the convergence is a set of scale factors for the respective rows $a_n(g)$, $g = 1, \dots, G$. It is useful to define the scales of any approximation A_n by

$$(4.111) \quad m\{a_n(g)\} \equiv \{a_n(g) M_{xx} a'_n(g)\}^{1/2}$$

because the corresponding expressions for A satisfy

$$(4.112) \quad m\{a(g)\} \equiv \{a(g) M_{xx} a'(g)\}^{1/2} = 1,$$

as is easily proved by proportional variations of all elements of $a(g)$ in the likelihood function (4.47). Now if $m\{a_n(g)\}$ is sufficiently different from unity, successive approximations of \mathcal{P}_1 will, under the present assumptions, exhibit scales $m\{a_n(g)\}$ differing from unity by an amount asymptotically constant in absolute value but alternating in signs since the corresponding characteristic value satisfies $k = 1 + l = -1$. If desired for aesthetic or practical reasons, this oscillation of scales can be reduced by occasional modification of the scales so as to make $m\{a_n(g)\} = 1$, or by occasional application of $\mathcal{P}_{1/2}$ instead of \mathcal{P}_1 , whereby the characteristic values concerned become $1 + \frac{1}{2}l = 0$.

It has already been shown in section 4.1.8 that the alternating behavior of scales does not affect the speed of convergence of the ratios between the elements of any row $a_n(g)$. The asymptotic value of that speed is determined by the largest of the quantities $|k_q|$, $q = 1, \dots, Q - K_y$. If knowledge of these values were available beforehand, it would be possible to choose another constant value of h that minimizes the largest of the $|k_q|$, $q = 1, \dots, Q - K_y$. In the absence of such knowledge, the following considerations favor the use of \mathcal{P}_1 in the circumstances at present assumed:

(a) The cost of computation per iteration for \mathcal{P}_1 is below that for any other constant value of h , and considerably below that for the process \mathcal{P}_{h_n} with the variable value (4.69) of h_n .

(b) The interval to which l_q is confined according to (4.110) is such that values of h larger than unity may lead to divergent processes. Of the admissible values $0 < h \leq 1$, only the value $h = 1$ leads to an interval (4.110) for k_q of which 0 is the midpoint.

The second consideration is based on complete ignorance of the range of relevant values l_q inside the interval $(-2, 0)$, and may lose its weight if more experience about these values is accumu-

lated from economic data, or if the interplay of the restrictions with the conditions in the last column of (4.99) is analyzed theoretically more completely than is done here. In any case, the risk that some values l_q might be close to -2 presents an additional reason for interspersing the iterations of \mathcal{P}_1 with an occasional application of $\mathcal{P}_{1/2}$, which cuts down those components $\bar{a}^\dagger \cdot v'(q)$ of \bar{A} corresponding to characteristic values l_q nearest to -2 .

4.3.3.9. *Case where identification depends on the diagonality restriction on Σ .* If the restrictions on A alone are insufficient for identification of all structural equations, there exist nondiagonal matrices \tilde{B} such that $A + \bar{A}$ in (4.107) again satisfies the restrictions. The possibility exists that among those matrices \tilde{B} at least one symmetric matrix can be found. In that case, at least one of the characteristic values l_q , $q = 1, \dots, Q - K_y$, not associated with trivial scale factors, equals -2 , and if \mathcal{P}_1 is applied iteratively, the corresponding component $\bar{a}_0^\dagger \cdot v'(q)$ in the initial value is not reduced to zero in successive iterations. It may be possible, by criteria similar to those developed in section 2, to determine whether or not the restrictions on A permit at least one of the nondiagonal matrices \tilde{B} in (4.107) to be symmetric.¹ Alternatively, one may apply \mathcal{P}_h with a constant value of $h < 1$, say $h = 3/4$, or one may insert $\mathcal{P}_{1/2}$ at regular intervals. One initial application of $\mathcal{P}_{1/2}$ would indeed cut out, up to the first order of magnitude, those components $\bar{a}_0^\dagger \cdot v'(q)$ of \bar{A} that cannot be reduced by \mathcal{P}_1 , but the presence of higher-order terms in (4.88) may reintroduce those components to some extent, thus requiring that $\mathcal{P}_{1/2}$ now be inserted with regularity.

4.3.3.10. *Asymptotic properties of \mathcal{P}_{h_n} .* Alternatively, whether nontrivial characteristic values $l = -2$ are present or not, one may pay the higher cost per iteration of the process \mathcal{P}_{h_n} already described in order to obtain at each stage a value h_n of h that is already optimal (in a first-order sense) in relation to the relative sizes of the various components $\bar{a}_n^\dagger \cdot v'(q)$ present in \bar{A}_n . By transformation of ΔA_0 in (4.69) to the space of $\Delta \bar{a}_0^\dagger$ we obtain

¹Such an inquiry might also throw further light on questions left unanswered in section 2.4.

$$\begin{aligned}
 h_0 &= -\frac{\Delta \bar{a}_0^\dagger \cdot \Delta \bar{a}'_0^\dagger}{\Delta \bar{a}_0^\dagger \cdot L^\dagger \cdot \Delta \bar{a}'_0^\dagger} = -\frac{\bar{a}_0^\dagger (L^\dagger)^2 \bar{a}'_0^\dagger}{\bar{a}_0^\dagger (L^\dagger)^3 \bar{a}'_0^\dagger} + \dots \\
 (4.113) \quad &= -\frac{\sum_{q=1}^Q l_q^2 \{\bar{a}_0^\dagger v'(q)\}^2}{\sum_{q=1}^Q l_q^3 \{\bar{a}_0^\dagger v'(q)\}^2} + \dots
 \end{aligned}$$

It is easily seen that the lowest-order term shown in the last member of (4.113) represents that value of h_0 which minimizes the lowest-order term shown in the last member of

$$\begin{aligned}
 \bar{a}_1^\dagger \bar{a}'_1^\dagger &= \sum_{q=1}^Q \{\bar{a}_1^\dagger v'(q)\}^2 \\
 (4.114) \quad &= \sum_{q=1}^Q \{1 + h_0 l_q\}^2 \{\bar{a}_0^\dagger v'(q)\}^2 + \dots,
 \end{aligned}$$

in keeping with the principle from which the process \mathcal{P}_{h_n} is derived.

Another interesting property of \mathcal{P}_{h_n} may be mentioned. If we postmultiply (4.81) by $(L^\dagger)^{-1}$ and insert the resulting expression for \bar{a}_n^\dagger in (4.82)¹, we obtain as the equivalent of the iterative procedure (4.82) in terms of $\Delta \bar{a}_n^\dagger$, using the value (4.113) of h_n ,

$$(4.115) \quad \Delta \bar{a}_1^\dagger = \Delta \bar{a}_0^\dagger - \frac{\Delta \bar{a}_0^\dagger \cdot \Delta \bar{a}'_0^\dagger}{\Delta \bar{a}_0^\dagger \cdot L^\dagger \cdot \Delta \bar{a}'_0^\dagger} \cdot \Delta \bar{a}_0^\dagger \cdot L^\dagger + \dots$$

Postmultiplication with $\Delta \bar{a}'_0^\dagger$ shows that successive adjustments $\Delta \bar{a}_n^\dagger$ of \bar{a}_n^\dagger are, to a first approximation, orthogonal to each other. Figure 4.3.3.10 demonstrates the first-order properties of (4.115). It is seen from (4.115) that first-order terms in \mathcal{P}_{h_n} are homogeneous of degree zero in the characteristic values l_q of L^\dagger . It fol-

¹Taking $n = 0$ and 1, respectively.

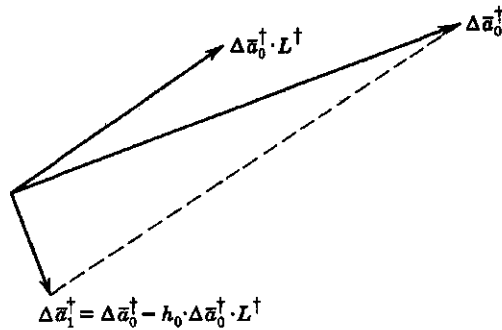


Figure 4.3.3.10

lows that only the ratios of the values l_q affect the asymptotic speed of convergence of \mathcal{P}_h^n . This circumstance appears in the figure, in that only the direction, not the length, of the vector $\Delta\bar{a}_0^\dagger L^\dagger$ determines the vector $\Delta\bar{a}_1^\dagger$. This shows that characteristic values l_q near to zero, indicating proximity to a situation of incomplete identification using all restrictions, are a more fundamental difficulty in computing maximum-likelihood estimates than characteristic values l_q near to or equal -2 , indicating "almost" or "altogether" lack of identification using restrictions on A alone. Components corresponding to the latter characteristic values can always be reduced by a suitable value of h .

4.3.3.11. *Convergence properties of \mathcal{P}_h and \mathcal{P}_h^n under incomplete identification.* It should be added that if any of the processes so far discussed is applied in cases where identification is incomplete (using all restrictions), and where therefore at least one characteristic value l_q vanishes, this circumstance is not revealed by the convergence properties of the process. Corresponding components $\Delta\bar{a}_0^\dagger$ in the initial displacement \bar{A}_0 are preserved through iterative multiplication by a factor unity. In other words, while reasonably fast convergence may be obtained if no other characteristic values near to zero exist, the limiting value $A_\infty = A + \bar{A}_\infty$ now depends on the initial components $\Delta\bar{a}_0^\dagger$ concerned.

4.3.3.12. *Comparisons between \mathcal{P}_h and \mathcal{P}_h^n .* Experience with

actual data is needed to determine whether the greater speed of convergence that might be expected of ρ_{h^n} as compared with ρ_h for any constant h justifies the greater cost of computation per iteration. One would expect ρ_{h^n} to be particularly economical if the ratios of the characteristic values l_q are close to unity. In any case, ρ_{h^n} would seem to be less subject to as yet unknown risks connected with the distribution of the characteristic values l_q , with the selection of the initial approximation A_0 , and with the sampling variations of M_{xx} around its expected value (4.89), which leads in general to depression of the restricted maximum of the likelihood function.

*4.3.3.13. *Effect of the \tilde{A} -restrictions on the characteristic values l_q .* On the basis of (4.99) we decompose \tilde{A} uniquely in terms of characteristic "vectors" of four types:

$$\begin{aligned}
 \tilde{A} &= [\tilde{B} \quad \tilde{C}] \\
 &= [\tilde{B}_{\text{dia}} \quad 0] + [\tilde{B}_{\text{sym}} \quad 0] + [\tilde{B}_{\text{ant}} \quad 0] + [0 \quad \tilde{C}] \\
 &= \tilde{B}_{\text{dia}} I_{[K_y K_x]} + \tilde{B}_{\text{sym}} I_{[K_y K_x]} + \tilde{B}_{\text{ant}} I_{[K_y K_x]} + \tilde{C} I'_{[K_x K_x]} \quad ,
 \end{aligned}
 \tag{4.116}$$

\tilde{B}_{dia} is diagonal,

$\tilde{B}_{\text{sym}} = \tilde{B}'_{\text{sym}}$ is symmetric with vanishing diagonal elements,

$\tilde{B}_{\text{ant}} = -\tilde{B}'_{\text{ant}}$ is antisymmetric.

In the absence of any restrictions, (4.116) gives the components of \tilde{A} according to the subspaces corresponding to the three different characteristic values of L^\dagger . We have further decomposed the symmetric component of \tilde{B} into a diagonal component \tilde{B}_{dia} and a component \tilde{B}_{sym} with vanishing diagonal elements, because of the trivial nature of the components \tilde{B}_{dia} .

With a priori restrictions of the type here considered, whenever \tilde{A} satisfies the \tilde{A} -restrictions, the component $\tilde{B}_{\text{dia}} I_{[K_y K_x]}$ satisfies the same restrictions. In any study of the effect of these restric-

tions on the characteristic values $l_q, q = 1, \dots, Q$, therefore, the component $\tilde{B}_{\text{dia}} I_{[K_y K_x]}$ and the corresponding characteristic values $l_q, q = Q - K_y + 1, \dots, Q$, play no role.

If \tilde{A} -restrictions are imposed on \tilde{A} , the decomposition (4.116) is still unique but it does not in general represent \tilde{A} as a linear combination of characteristic "vectors." It will do so if and only if each of the last three components in (4.116) also satisfies the \tilde{A} -restrictions. Whenever an \tilde{A} -restricted \tilde{A} exists for which at least one of those three components fails to satisfy the \tilde{A} -restrictions, at least one new "intermediate" characteristic value, l_q has been introduced of which the (each) characteristic "vector" $\tilde{A}^{(i)}$ is a linear combination $\tilde{B}_{\text{sym}}^{(i)} I_{[K_y K_x]} + \tilde{B}_{\text{ant}}^{(i)} I_{[K_y K_x]} + \tilde{C}^{(i)} I_{[K_x K_x]}$ with at least two nonvanishing components. This is seen as follows: If, after imposing the \tilde{A} -restrictions, a complete set of new characteristic "vectors" could be chosen in such a manner that each of them consisted of one component (4.116) only, the unique expression of \tilde{A} as a linear combination of those new characteristic "vectors" [derived from (4.87) by transformation to the \tilde{A} -space] would coincide with the unique decomposition (4.116). This would contradict the assumption that at least one component (4.116) of \tilde{A} fails to satisfy the \tilde{A} -restrictions.

*4.3.3.14. *Case where all relevant characteristic values coincide at $l = -2$.* We have met already with one example where the a priori restrictions imply \tilde{A} -restrictions that are preserved in the decomposition (4.116). This is the case, discussed in section 4.3.2.3 in connection with $\mathbb{P}_{1/2}$, where the only a priori restrictions are

$$(4.117) \quad C = 0, \quad \text{and hence} \quad \bar{C} = 0.$$

From (4.93) and (4.95) the corresponding \tilde{A} -restrictions are seen to be

$$(4.118) \quad \tilde{C} = 0.$$

This wipes out all characteristic values $l = -1$, while leaving unaffected the values $l = 0$ and $l = -2$ and the corresponding vectors. It follows that one application of $\mathbb{P}_{1/2}$ will reduce all components corresponding to $l = -2$ to second-order magnitude.

This accounts for the high speed of convergence of $\mathcal{P}_{\frac{1}{2}}$ found in this special case. It is seen as follows, however, that a single additional restriction on B would introduce a characteristic value $-2 < l < 0$ with probability one in the sample space. Let the additional restriction be denoted $\text{tr } \tilde{A} \Phi' = 0$. The \tilde{A} -restriction $\text{tr } \tilde{B} A \Phi' = 0$ will only then be identically satisfied by $\tilde{B}_{\text{sym}} I_{[K_y, K_x]}$ if it also implies $\text{tr } \tilde{B}' A \Phi' = \text{tr } \Phi A' \tilde{B} = \text{tr } \tilde{B} \Phi A' = 0$, which requires $A \Phi' = \Phi A'$, an occurrence of probability zero.

*4.3.3.15. *Case where all characteristic values coincide at $l = -1$.* It is of interest to inquire whether there exists a counterpart to the foregoing case, in which the nontrivial component $\tilde{B}_{\text{sym}} I_{[K_y, K_x]}$ with the characteristic value $l = -2$ is wiped out by restrictions without introducing any intermediate characteristic values. Such a case can be found easily if we also require complete identification, i.e., wiping out of the component $\tilde{B}_{\text{ant}} I_{[K_y, K_x]}$ with characteristic value $l = 0$. For in that case the restrictions must imply that

$$(4.119) \quad \begin{aligned} \tilde{a}_{gh} &= \bar{a}(g) H^{-1} \nu'(h) = \bar{a}(g) \begin{bmatrix} B^{-1} \\ 0 \end{bmatrix} \nu'(h) \\ &= \bar{a}^g X^g B^{-1} \nu'(h) = 0 \quad \text{for } g \neq h, \quad g, h \leq G, \end{aligned}$$

whatever \bar{a}^g , or that

$$(4.120) \quad X^g B^{-1} \nu'(h) = 0 \quad \text{for } g \neq h.$$

Here X^g is a submatrix of Φ^g as defined in (4.18). This again requires the existence of column vectors λ'^g such that

$$(4.121) \quad X^g = \lambda'^g b(g),$$

so that the matrix B is restricted by

$$(4.122) \quad \beta(g) = \alpha^g X^g = \alpha^g \lambda'^g b(g) \equiv \lambda b(g),$$

where λ is the scalar quantity $\alpha^g \lambda^g$. Thus, the ratios of all elements of any row of B must be prescribed by the restrictions, and B is the product $\Lambda_{yy} B^\oplus$ of an unknown nonsingular diagonal matrix Λ_{yy} with a known nonsingular matrix B^\oplus . Hence a known nonsingular transformation,

$$(4.123) \quad y' = B^\oplus y'^{\oplus},$$

of the dependent variables alone will then bring the system of structural equations into the form

$$(4.124) \quad \Lambda_{yy} y'^{\oplus} + \Gamma z' = u',$$

which differs from the reduced form only by the principle of normalization. In this form, therefore, maximum-likelihood estimation is equivalent to the least-squares principle applied to each equation separately.

Thus, the exclusion of both components $\tilde{B}_{sym} I_{[K_y K_x]}$ and $\tilde{B}_{ant} I_{[K_y K_x]}$ leads to a trivial case which can be treated by more elementary methods. In fact, one single application of \mathcal{P}_1 , which is the optimal procedure in the present case, is the equivalent to the least-squares procedure and leads to the exact solution of the maximum-likelihood equation in one iteration. However, if only a single restriction on B is relaxed, nondiagonal values of \tilde{B} are made possible that are in general neither symmetric nor antisymmetric, so that new intermediate characteristic values are introduced. A simple example is obtained if in (4.124) we assume only $K_y = 2$ dependent variables y^\oplus and leave λ_{12} unrestricted while λ_{21} is still required to vanish.

4.3.3.16. *Asymptotic speed of convergence of $\mathcal{P}_h = \mathcal{P}_{h_n}$ in the foregoing case.* This discussion shows that cases where all relevant characteristic values coincide are rare and, from the point of view of the problem of measuring economic relationships, either trivial or accidental. In general, different characteristic values are present, and the discrepancy $\bar{A}_{n+1} = A_{n+1} - A$ between the result of the $(n+1)$ th iteration and the solution of the maximum-likelihood equations is even asymptotically a nonvanishing fraction of the corresponding difference \bar{A}_n computed from the result of the n th iteration with any of the methods so far discussed.

As to speed of convergence per iteration, therefore, the present methods fall short of the well-known Newton method, in which by definition \bar{A}_{n+1} is of second-order magnitude compared with \bar{A}_n .

4.3.4. Arrangement of computations for $\mathcal{P}_1, \mathcal{P}_h, \mathcal{P}_{h_n}$

4.3.4.1. *A constructed example for numerical illustration of $\mathcal{P}_h, \mathcal{P}_{h_n}$ involving only single-parameter restrictions.* Before exploring the application of the Newton method to our present problem, we shall give a few numerical illustrations of the procedures so far discussed, and give further comments of a practical character with regard to computational procedure.

We have constructed an example of a three-equation system characterized by the matrices

$$A = \begin{bmatrix} 0 & 1 & 4 & 1 & 0 & 0 \\ 1 & 0 & -3 & 0 & 1 & 0 \\ -2 & 1 & 0 & 0 & 0 & 1 \end{bmatrix},$$

$$(4.125) \quad \Sigma = \begin{bmatrix} 0.2 & 0.1 & 0.0 \\ 0.1 & 0.2 & 0.1 \\ 0.0 & 0.1 & 0.3 \end{bmatrix},$$

$$M_{zz} = I,$$

in which the variables z are regarded as fixed in repeated samples. The matrix M_{yx} then fluctuates from sample to sample, but, in order to abstract from the effect of sampling fluctuations, we have assumed that the observed moment matrix in the sample imagined to be drawn is equal to its expected value, which is easily computed from (4.125):

$$(4.126) \quad M_{yx} = M_{yx} = \frac{1}{T} \sum_{t=1}^T \mathcal{E} y'(t) x(t)$$

$$\begin{aligned}
&= B^{-1} \Sigma B'^{-1} I_{[K_y K_x]} - B^{-1} \Gamma M_{zz} [-\Gamma' B'^{-1} \quad I_{[K_z]}] \\
&= \begin{bmatrix} 0.417 & 0.484 & -0.021 & -0.300 & -0.400 & 0.300 \\ 0.484 & 1.568 & -0.192 & -0.600 & -0.800 & -0.400 \\ -0.021 & -0.192 & 0.073 & -0.100 & 0.200 & 0.100 \end{bmatrix}.
\end{aligned}$$

In the present section 4.3 the diagonality restriction is imposed on Σ , although the example has been constructed with a non-diagonal Σ . Comparison of the "maximum-likelihood estimates" of A so obtained with the true values will illustrate the effect of the diagonality assumption regarding Σ in a case where it is incorrect.

The restrictions on A are that those elements that are zero in (4.125) are known and required to be zero. In that case it is profitable to state the definition (4.57) of ΔA_0 in matrix form

$$\begin{aligned}
(4.127) \quad \mathcal{L}(\Delta A_0 M_{xx}) &= \mathcal{L}(B_0'^{-1} I_{[K_y K_x]} - A_0 M_{xx}), \\
\mathcal{L} \Delta A_0 &= \Delta A_0,
\end{aligned}$$

because the operator \mathcal{L} now consists in replacing by zero the elements numbered (1,1), (1,5), (1,6), (2,2), (2,4), (2,6), (3,3), (3,4), (3,5), and can according to (4.62) also be applied to submatrices. Alternatively, the restrictions can be expressed through the set of basic matrices

$$\begin{aligned}
(4.128) \quad \Phi^1 &= \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}, \\
\Phi^2 &= \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix},
\end{aligned}$$

$$\Phi^3 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

4.3.4.2. *Elimination of the unknowns C_n under single-parameter restrictions.* If \mathcal{P}_1 is applied,

$$(4.129) \quad \mathcal{Q}(A_n M_{xz}) = 0, \quad \mathcal{Q}(\Delta A_n M_{xz}) = 0$$

will automatically be satisfied for $n \geq 1$. The condition (4.129) will be recognized as the condition (4.5) for maximization of the likelihood function with respect to the parameters Γ only. For other choices of h_n it is possible and desirable to insure the validity of (4.129) for all n by imposing it as a condition on the initial value A_0 . Any initial value A_0 derived from any set of single-equation least-squares estimates in which one of the variables y_i , $i = 1, \dots, K_y$, is selected as "dependent" variable for each value of g , satisfies that condition. Since (4.129) also holds for A , a similar condition is satisfied by $\bar{A}_n = \tilde{A}_n H$. Let the submatrices of H'^{-1} partitioned similarly to (4.93) be denoted by subscript combinations yy , yz , etc. Then, since $(H'^{-1})_{yz} = 0$ and $(H'^{-1})_{zz}$ is nonsingular, substitution of (4.94) in (4.129) shows that the latter condition secures the identical vanishing of $\mathcal{Q}\tilde{C}_n$ rather than its general reduction in successive iterations. It was already recognized in section 4.1.5 in connection with \mathcal{P}_1 that this cuts down the number of unknowns that need to be determined in each iteration. The unknowns participating in the iterations are the unrestricted elements of $\mathcal{Q}B_n$ corresponding, through $\bar{B}_n = \tilde{B}_n B$, to the first three components of \tilde{A}_n in the decomposition (4.116).

4.3.4.3. *Arrangement of computations under single-parameter restrictions.* Because the restrictions imposed are of the single-parameter type, the notation α^g simply indicates that the elements in $\alpha(g)$ prescribed to be zero are deleted. Solving for $\mathcal{Q}\Delta C_n$ and $\mathcal{Q}C_n$ from (4.129) and substituting in (4.127), we have as the defi-

notation of $\mathcal{L} \Delta B_0$ in that notation,

$$(4.130) \quad \Delta b_0^g \cdot {}^z M_{yy}^g = \iota(g) \cdot B_0'^{-1} \cdot \Phi'^g - b_0^g \cdot {}^z M_{yy}^g,$$

from which we solve for Δb_0^g for computational purposes,

$$(4.131) \quad \Delta b_0^g = \iota(g) B_0'^{-1} \Phi'^g ({}^z M_{yy}^g)^{-1} - b_0^g.$$

If \mathcal{P}_{h_n} is applied, the value (4.69) of h is obtained most conveniently from

$$(4.132) \quad h_0 = \frac{\sum_g \Delta b_0^g \cdot {}^z M_{yy}^g \cdot \Delta b_0'^g}{\text{tr}\{(B_0'^{-1} \cdot \Delta B_0'^{-1})^2 + \sum_g \Delta b_0^g \cdot {}^z M_{yy}^g \cdot \Delta b_0'^g\}}$$

in which the first term in the denominator is not transformed to vector coordinates.

First the K_y matrices ${}^z M_{yy}^g$ and their inverses are prepared from M_{xx} and the restrictions. Then an initial value A_0 (a set of vectors b_0^g) is selected. Each row of $({}^z M_{yy}^g)^{-1}$ represents a least-squares regression as a possible choice of initial vector b_0^g . The initial vectors are then normalized by

$$(4.133) \quad b_0^g \cdot {}^z M_{yy}^g \cdot b_0'^g = 1,$$

and $B_0'^{-1}$ is computed by any suitable method of inversion. Then either Δb_0^g or $b_0^g + \Delta b_0^g$ is determined from (4.131), from which b_1^g is obtained, in the case of \mathcal{P}_{h_n} with the help of (4.132), and the next iteration can proceed.

If the process is terminated after the n th iteration, $\mathcal{L}C_n$ is obtained from $\mathcal{L}B_n$ by (4.5) or

$$(4.134) \quad c_n^g = -b_n^g \cdot M_{yz}^g \cdot (M_{zz}^g)^{-1}.$$

4.3.4.4. Application of \mathcal{P}_1 , $\mathcal{P}_{3/4}$, \mathcal{P}_{h_n} to the constructed example. Table 4.3.4.4 gives the results of applying \mathcal{P}_1 , $\mathcal{P}_{3/4}$,

TABLE 4.3.4.4. Numerical illustration of \mathcal{P}_1 ,

Method	Iteration $n =$	Ratios of elements of $b_n(g)$			Scales: $m\{a(g)\} \equiv \{a_n(g) M_{xx} a'_n(g)\}^{1/2}$		
		$-\frac{b_{12}^n}{b_{13}^n}$	$-\frac{b_{21}^n}{b_{23}^n}$	$-\frac{b_{32}^n}{b_{31}^n}$	$m\{a_1(g)\}$	$m\{a_2(g)\}$	$m\{a_3(g)\}$
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
\mathcal{P}_1	0	0.25000	0.22960	0.42887	1.00000	0.99805	1.00000
	1	0.23673	0.32274	0.46075	1.00787	1.05568	1.01041
	2	0.23277	0.33141	0.47046	0.99301	0.94762	0.99060
	3	0.23230	0.33408	0.47160	1.00709	1.05530	1.00858
	4	0.23219	0.33440	0.47187	0.99292	0.94757	0.99087
	5	0.23217	0.33448	0.47191	1.00716	1.05368	1.00947
	6	0.23217	0.33448	0.47192	0.99286	0.94758	0.99091
$\mathcal{P}_{3/4}$	0	0.25000	0.22960	0.42887	1.00000	0.99805	1.00000
	1	0.23991	0.29946	0.45275	1.00443	1.03172	1.00592
	2	0.23524	0.32130	0.46427	0.99884	0.98748	0.99837
	3	0.23336	0.32965	0.46895	1.00073	1.00677	1.00105
	4	0.23263	0.33265	0.47078	0.99639	0.99669	0.99950
	5	0.23235	0.33379	0.47148	1.00021	1.00025	1.00025

$\rho_{3/4}$, and ρ_{h_n} when Σ is required to be diagonal.

ρ_{h_n} $h_n = 1.06540$ 0.57158 -0.87201 0.59751 0.85300 0.60588	0 1 2 3 4 5 6	0.25000 0.23591 0.23399 0.23259 0.23238 0.23223 0.23220	0.22961 0.32884 0.33053 0.33291 0.33368 0.33434 0.33436	0.42887 0.46283 0.46747 0.47089 0.47140 0.47176 0.47184	1.00000 1.00902 0.99893 1.00088 0.99983 1.00011 0.99995	0.99805 1.06286 0.99330 1.00521 0.99898 1.00075 0.99986	1.00000 1.01181 0.99860 1.00118 0.99971 1.00016 0.99997
Newton	0 1 2 3	0.25000 0.23197 0.23224 0.23218	0.22960 0.33566 0.33463 0.33432	0.42887 0.47251 0.47192 0.47185	1.00000 1.00283 0.99924 0.99988	0.99805 1.00458 0.99989 1.00022	1.00000 1.00542 1.00044 1.00032
Solutions		$\frac{b_{12}}{b_{13}}$	$-\frac{b_{21}}{b_{23}}$	$-\frac{b_{32}}{b_{31}}$	$\frac{b_{14}}{b_{13}}$	$-\frac{b_{25}}{b_{23}}$	$-\frac{b_{36}}{b_{31}}$
Σ diagonal (from ρ_1 above)		0.23217	0.33448	0.47192	0.23931	0.33796	0.48876
Σ nondiagonal (true values)		0.25000	0.33333	0.50000	0.25000	0.33333	0.50000

and ρ_{h_n} respectively in the examples given. It also gives the application, to the same data, of the Newton method, to be discussed below. Initial values were the same in all cases and were based on the rows of $({}^z M_{yy}^g)^{-1}$ numbered 3, 1, and 2 for $g = 1, 2, 3$, respectively.

Comparison of the three methods shows that in the present case ρ_1 is superior to $\rho_{3/4}$, and even slightly better than ρ_{h_n} , with respect to the most essential property: the speed of convergence of the ratios of elements of each $a_n(g)$. ρ_1 shows the characteristic alternation in successive values $m\{a_n(g)\}$, $n = 0, 1, \dots$. These values converge gradually in $\rho_{3/4}$ and slightly faster in ρ_{h_n} . The wide variation in successive values of h_n is remarkable. The apparent alternation of these values is peculiar to the present example. Results substantially similar to those shown here were obtained in another constructed example of only two equations, but the sequence of values h_n was found to be more irregular in that case.¹

4.3.4.5. Arrangement of computations under more general restrictions. Under the simple type of restrictions imposed in the foregoing example the vectors $a(g)$ differ from the corresponding vectors a_n^g only in that they contain in addition certain vanishing elements. With more general basic matrices Φ^g it is necessary to decide whether the quantities a_n^g are actually to be computed as a separate step in each iteration. We shall show that this can be avoided, with a resulting saving of computational labor.

Using successively (4.40), the second relation of (4.30), (4.37), and (4.72), we can rewrite the definition (4.57) of ΔA_0 in the form

$$(4.135) \quad \Delta a_0^* M^{**} = \text{vec}[B_0'^{-1} \quad 0] \Phi'^* - a_0^* M^{**}.$$

Even with orthogonalization of Φ^* , which we do not now require, the nonsingular matrix $(\Phi^* \Phi'^*)^{-1}$ appearing in (4.40) can be and has been omitted from (4.135) because it appears originally as postmultiplicand to all terms.

¹This example was discussed as "case I" in another article by one of the present authors [Koopmans, 1945]. Initial values were the least-squares regression with x_1 as dependent variable in both structural equations. Successive values of h_n were 0.394, 0.839, 1.133, 0.582, 0.795,

From (4.135) we derive

$$(4.136) \quad \Delta a_0^* = \text{vec}[B_0'^{-1} \quad 0] \Phi'^* (M^{**})^{-1} - a_0^*,$$

to which we again apply (4.37) to obtain

$$(4.137) \quad \Delta a_0 = \text{vec}[B_0'^{-1} \quad 0] \Phi'^* (M^{**})^{-1} \Phi^* - a_0.$$

Defining

$$(4.138) \quad P^{(g)} \equiv \Phi'^g (\Phi^g M \Phi'^g)^{-1} \Phi^g \equiv P_{xx}^{(g)} \equiv \begin{bmatrix} P_{yy}^{(g)} & P_{yz}^{(g)} \\ P_{zy}^{(g)} & P_{zz}^{(g)} \end{bmatrix},$$

we have

$$(4.139) \quad P \equiv \Phi'^* (M^{**})^{-1} \Phi^* = \begin{bmatrix} P^{(1)} & 0 & \dots & 0 \\ 0 & P^{(2)} & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & P^{(G)} \end{bmatrix}.$$

Therefore, (4.137) is equivalent to

$$(4.140) \quad \Delta a_0(g) = \iota(g) B_0'^{-1} P_{yx}^{(g)} - a_0(g), \quad g = 1, \dots, G.$$

The preparations for any of the iterative procedures based on (4.140) now consist in the evaluation of the G matrices $M^g = \Phi^g M \Phi'^g$ from (4.72), their inversion, and the transformation of the inverses $(M^g)^{-1}$ back to the x -coordinates by (4.138), to obtain the $P_{yx}^{(g)}$.

4.3.4.6. Canonical form of the basic matrices Φ^g .* In many cases a further saving, similar to that obtained under simpler restrictions, can be secured by choosing an initial value satisfying (4.129). It will be noted that $P_{zz}^{(g)}$ does not occur in (4.140), and that $P_{yz}^{(g)}$ is not needed to obtain B by an iterative process based on (4.140). This suggests that the inversion of M^{} can profitably be replaced by successive inversions of lower-order matrices.

In order to obtain the full benefit of this consideration, we must analyze the basic matrices Φ^g with the help of the following lemma:

LEMMA 4.3.4.6. *If $\bar{\Phi}^g$ is such that its number of rows equals its rank,*

$$(4.141) \quad r(\bar{\Phi}^g) = \rho(\bar{\Phi}^g) = Q_g,$$

and if $Q_g^I, Q_g^{II}, Q_g^{III}$ are defined by

$$(4.142) \quad Q_g \equiv \rho(\bar{X}^g) + Q_g^{III} \equiv \rho(\bar{\Psi}^g) + Q_g^I \equiv Q_g^I + Q_g^{II} + Q_g^{III},$$

then there exists a nonsingular transformation matrix¹ Ω such that

$$(4.143) \quad \Omega \bar{\Phi}^g = \Phi^g = [X^g \quad \Psi^g] = \begin{bmatrix} X_I^g & 0 \\ X_{II}^g & \Psi_{II}^g \\ 0 & \Psi_{III}^g \end{bmatrix},$$

with

$$(4.144) \quad \begin{aligned} \rho(X_I^g) &= r(X_I^g) = Q_g^I, \\ \rho(\Psi_{III}^g) &= r(\Psi_{III}^g) = Q_g^{III}, \\ \rho(\Phi_{II}^g) &= r(\Phi_{II}^g) = Q_g^{II}, \end{aligned}$$

and there exists no nonsingular transformation matrix Ω such that in (4.143) X_I^g or Ψ_{III}^g or both have higher ranks than given by (4.144).

Proof: Owing to (4.142) there exist matrices Ω_I, Ω_{III} such that

$$(4.145) \quad \rho(\Omega_I) = r(\Omega_I) = Q_g^I, \quad \rho(\Omega_{III}) = r(\Omega_{III}) = Q_g^{III},$$

¹For brevity no subscript g is appended to Ω .

and

$$(4.146) \quad \Omega_{\text{III}} \bar{X}^g = 0, \quad \Omega_{\text{I}} \bar{\Psi}^g = 0.$$

We must have

$$(4.147) \quad \rho \begin{pmatrix} \Omega_{\text{I}} \\ \Omega_{\text{III}} \end{pmatrix} = \rho(\Omega_{\text{I}}) + \rho(\Omega_{\text{III}})$$

because, if the right-hand member in (4.147) exceeded the left-hand member, there would, according to Lemma 2.3.2, exist nonvanishing vectors $\lambda_{\text{I}}, -\lambda_{\text{III}}$ such that

$$(4.148) \quad \mu = \lambda_{\text{I}} \Omega_{\text{I}} = \lambda_{\text{III}} \Omega_{\text{III}} \neq 0.$$

In that case (4.146) would imply

$$(4.149) \quad \mu [\bar{X}^g \quad \bar{\Psi}^g] = \mu \bar{\Phi}^g = 0, \quad \mu \neq 0,$$

contrary to the assumption (4.141) about the rank of $\bar{\Phi}^g$.

Owing to (4.147) we can find a matrix Ω_{II} with

$$(4.150) \quad \rho(\Omega_{\text{II}}) = r(\Omega_{\text{II}}) = Q_g^{\text{II}} \geq 0$$

such that the matrix Ω defined by

$$(4.151) \quad \Omega \equiv \begin{bmatrix} \Omega_{\text{I}} \\ \Omega_{\text{II}} \\ \Omega_{\text{III}} \end{bmatrix}$$

is nonsingular (rank Q_g). This matrix Ω produces the partitioning required in (4.143) with submatrices having the number of rows required in (4.144). Since $\bar{\Phi}^g$ has the same rank (4.141) as Φ^g , the ranks in (4.144) cannot fall below the corresponding number of rows.

Finally, if a nonsingular matrix Ω existed such that the rank and hence the number of rows of X_{I}^g , say, exceeded Q_g^{I} , we should have, using successively (4.141), (4.143), and (4.142),

$$(4.152) \quad Q_g - Q_g^I > \rho(\Phi^g) - \rho(X_I^g) = \rho(\Psi^g) = Q_g - Q_g^I,$$

an obvious contradiction. This completes the proof of Lemma 4.3.4.6. The form (4.143) of Φ^g will be called its canonical form. It is worth stressing that this form neither requires nor precludes orthogonalization of Φ^g according to (4.24).

*4.3.4.7. *Elimination of the unknowns $\alpha_{n, III}^g$.* Let us now assume that the basic matrices Φ^g are already in the canonical form. Then the expression (4.22) for $\alpha(g)$ in terms of a set of unrestricted parameters partitions as follows:

$$(4.153) \quad \begin{aligned} \beta(g) &= \alpha_I^g X_I^g + \alpha_{II}^g X_{II}^g, \\ \gamma(g) &= \alpha_{II}^g \Psi_{II}^g + \alpha_{III}^g \Psi_{III}^g. \end{aligned}$$

Thus the parameters α_{III}^g do not enter into the Jacobian B. Iterative processes involving B_n only can therefore be constructed on the basis of the parameters

$$(4.154) \quad {}_{III}\alpha^g \equiv [\alpha_I^g \quad \alpha_{II}^g]$$

alone. With each approximation ${}_{III}\alpha_n^g$ to ${}_{III}\alpha^g$, there are associated "silent" values $\alpha_{n, III}^g$ which are those linear functions of ${}_{III}\alpha_n^g$ that maximize the likelihood function with ${}_{III}\alpha_n^g$ inserted for ${}_{III}\alpha^g$. Only at the termination of iterations do these values need to be determined explicitly or implicitly. In the computational arrangement of the Newton method discussed below, an explicit determination from equation (4.185) involves no extra cost. Since we have for the present decided against explicit evaluation of any part of α_n^g , we shall operate equivalently from (4.140) and (4.138) on the basis of properties of inverses of partitioned matrices.

We define

$$(4.155) \quad (M^g)^{-1} \equiv N^g \equiv \begin{bmatrix} N_{I I}^g & N_{I II}^g & N_{I III}^g \\ N_{II I}^g & N_{II II}^g & N_{II III}^g \\ N_{III I}^g & N_{III II}^g & N_{III III}^g \end{bmatrix}.$$

As before, the postsubscripts I, II can be subsumed in the presubscript III, the postsubscripts II, III in the presubscript I. For the iterations in (4.140) we need only

$$(4.156) \quad P_{yy}^g = X'^g N^g X^g = {}_{III}X'^g \cdot {}_{III III}N^g \cdot {}_{III}X^g,$$

where ${}_{III III}N^g$ is obtained from N^g by deleting the rows and columns intersecting in $N_{III III}^g$. If the ultimate evaluation of C_n is based on (4.140), we need in addition

$$(4.157) \quad P_{yz}^g = X'^g N^g \Psi^g = {}_{III}X'^g \cdot {}_{III I}N^g \cdot {}_I\Psi^g.$$

Thus $N_{III III}^g$ is not needed. ${}_{III III}N^g$ is computed from¹

$$(4.158) \quad {}_{III III}N^g = \{ {}_{III III}M^g - {}_{III}M_{III}^g (M_{III III}^g)^{-1} {}_{III}M'_{III}^g \}^{-1}$$

and the submatrix ${}_{III}N_{III}^g$ needed in addition for (4.157) is obtainable from¹

$$(4.159) \quad {}_{III}N_{III}^g = - {}_{III III}N^g \cdot {}_{III}M_{III}^g \cdot (M_{III III}^g)^{-1},$$

using once more a matrix ${}_{III}M_{III}^g (M_{III III}^g)^{-1}$ already computed for (4.158).

These formulae show that the most important saving from the use of the canonical form (4.143) - avoiding the calculation of $N_{III III}^g$ - is due to the separation of Φ_{III}^g from ${}_{III}\Phi^g$. The further separation of Φ_I^g from Φ_{II}^g leads to a minor additional saving by reducing the

¹[Hotelling, 1943-A, p.4].

number of elements involved in the second matrix multiplication in (4.157).

In general the ranks of P_{yy}^g and P_{yz}^g are lower than the maximum compatible with the number of rows and columns. This expresses the fact that the elements of the vectors $a_n(g)$ depend linearly on a smaller number of parameters a_n^g .

4.3.4.8. *The final evaluation of the a_{III}^g .* The return from B_n to C_n on the basis of (4.140) requires in addition the inversion of B_n , which then serves for all values of g ,

$$(4.160) \quad c_n(g) = v(g) B_0'^{-1} P_{yz}^g.$$

Depending on the circumstances, an alternative formula for $c_n(g)$ may be more economical. This is based on the expression for the "silent" values,

$$(4.161) \quad a_{III}^g = -_{III} a_n^g \cdot_{III} M_{III}^g \cdot (M_{III \ III}^g)^{-1},$$

which is the equivalent of (4.134) under the present form of the Φ^g . From (4.22) and (4.143) we have

$$(4.162) \quad b_n(g) =_{III} a_n^g \cdot_{III} X^g, \quad c_n(g) = a_{n,II}^g \cdot \Psi_{II}^g + a_{n,III}^g \cdot \Psi_{III}^g.$$

The first of these relations is solved for $_{III} a_n^g$ by

$$(4.163) \quad_{III} a_n^g = b_n(g) \Xi,$$

where the relation

$$(4.164) \quad_{III} X^g \cdot \Xi =_{III \ III} I$$

defines all that needs to be defined concerning Ξ . The condition (4.164) may in simple cases offer more ready ways of finding a suitable value of Ξ than the explicit calculation of the particular solution

$$(4.165) \quad \Xi =_{III} X'^g \cdot (_{III} X^g \cdot_{III} X'^g)^{-1}.$$

Combining the second relation (4.162) with (4.161) and (4.163), we have

$$(4.166) \quad \begin{aligned} c_n(g) &= b_n(g) \Xi \left\{ \text{III} \Psi^g - \text{III} M_{\text{III}}^g (M_{\text{III}}^g \text{III})^{-1} \Psi_{\text{III}}^g \right\} \\ &\equiv b_n(g) J(g), \end{aligned}$$

say. The application of this formula requires the evaluation of G matrices $J(g)$, involving in principle G inversions (4.165) of orders equal to the respective values of $(Q_g^I + Q_g^{II})$.

4.3.4.9. *A formula for the computation of h_n .* In the application of \mathbb{P}_{h_n} , the formula (4.166) is preferable to (4.160) because it also holds if $a_n(g)$ is replaced by the scalar multiple $\Delta a_n(g)$ of the difference between two successive approximations. It can therefore be used to derive from (4.69) the formula

$$(4.167) \quad h_0 = \frac{\sum_{g=1}^G \Delta b_0(g) \cdot M_{yy}(g) \cdot \Delta b'_0(g)}{\text{tr}(B_0'^{-1} \cdot \Delta B'_0) + \sum_{g=1}^G \Delta b_0(g) \cdot M_{yy}(g) \cdot \Delta b'_0(g)}$$

for the evaluation of h_0 , in which

$$(4.168) \quad M_{yy}(g) = M_{yy} + J(g) M_{zy} + M_{yz} J'(g) + J(g) M_{zz} J'(g).$$

4.3.5. The Newton method

4.3.5.1. *The principle of the Newton method.* Unlike the methods discussed so far, the principle of the Newton method has no connection with the particular form of the likelihood function. Its application to our problem proceeds as follows. If we write

$$(4.169) \quad A_1 \equiv A_0 + \Delta A_0,$$

the Taylor expansion (4.68) of the likelihood function $L(A_1)$ in terms of ΔA_0 can be written

$$\begin{aligned}
 L(A_1) - L(A_0) &= \text{tr}\{(B_0'^{-1} I_{[K_y K_x]} - A_0 M_{xx}) \Delta A_0'\} \\
 (4.170) \quad &+ \frac{1}{2} \text{tr}\{-(B_0'^{-1} \Delta B_0')^2 - \Delta A_0 M_{xx} \Delta A_0'\} + \dots
 \end{aligned}$$

This expansion, transformed to the unrestricted parameters Δa_0^* defined as in (4.37), may be denoted

$$(4.171) \quad L(a_1^*) - L(a_0^*) = l_0^* \Delta a_0'^* + \frac{1}{2} \Delta a_0^* L_0^* \Delta a_0'^* + \dots$$

The vector l_0^* and symmetric matrix L_0^* so defined depend, of course, on the initial value A_0 , as distinct from the vector $l^* = 0$ and the matrix L^* defined by (4.77) on the basis of the solution A .

The Newton method determines Δa_0^* from the requirement that the first two terms in the expansion

$$(4.172) \quad \frac{dL(a_1^*)}{d\Delta a_0^*} = l_0^* + \Delta a_0^* L_0^* + \dots$$

shall cancel out. This leads to

$$(4.173) \quad a_1^* L_0^* = (a_0^* + \Delta a_0^*) L_0^* = a_0^* L_0^* - l_0^*$$

as the formula defining a_1^* .

4.3.5.2. Comparisons between the Newton method and $\mathcal{P}_h, \mathcal{P}_{h_n}$.

The following differences between this method and the procedures based on the earlier choice (4.57) of ΔA_0 deserve discussion.

The Newton method seeks any stationary point of the likelihood function to which the initial value A_0 is sufficiently near. The earlier methods converge only to maxima. This establishes a presumption that the Newton method requires for convergence a closer proximity of the initial value to the maximum sought. It also means that after a stationary point has been obtained, second-order conditions must now be investigated to determine whether the point found is actually a maximum. Saddle points, maxima, and minima can

be distinguished through indefiniteness, negative definiteness, and positive definiteness, respectively, of L^* .

The Newton method breaks down in the case of incomplete identification because then L^* is singular, and hence cannot be inverted.

If A_0 is sufficiently close to the desired maximum A , the speed of convergence per iteration in the Newton method is superior. Writing again $a_n^* = a^* + \bar{a}_n^*$, it is seen as follows that \bar{a}_1^* is quadratic in \bar{a}_0^* - a property which was found to be present in the earlier methods only in the rare case that all relevant characteristic values of L^\dagger coincide.

From the expansion of the likelihood function $L(\alpha^*)$ with respect to $\bar{\alpha}^* = \alpha^* - a^*$,

$$(4.174) \quad L(\alpha^*) - L(a^*) = \frac{1}{2} \bar{\alpha}^* L^* \bar{\alpha}^{\prime*} + \dots,$$

we have

$$(4.175) \quad l_0^* \equiv \left(\frac{dL(\alpha^*)}{d\bar{\alpha}^*} \right)_{\bar{\alpha}^* = \bar{a}_0^*} = \bar{a}_0^* L^* + \dots$$

From (4.173) and (4.175) we have the relation

$$(4.176) \quad \bar{a}_1^* = \bar{a}_0^* - l_0^* (L_0^*)^{-1} = \bar{a}_0^* \{ I - L^* (L_0^*)^{-1} + \dots \},$$

from which it is easily seen, by expanding L_0^* in terms of \bar{a}_0^* , that the first-order term in \bar{a}_0^* in this expression vanishes, and that the quadratic term in \bar{a}_0^* depends on the third derivatives of the likelihood function in the point a^* .

Against the superior speed of convergence per iteration in the Newton method must be set the greatly increased computational labor per iteration. Disregarding for a moment the saving arising in both methods from the canonical form of the basic matrices, the Newton method requires for each iteration afresh the calculation of the matrix L_n^* and the solution of the linear equations (4.173) in Q unknowns a_1^* . In contrast, the matrix M^{**} occurring in the left-hand member $\Delta a_0^* M^{**}$ of (4.57) remains the same for all iterations, so that its inversion paves the way for evaluation of successive values Δa_n^* by matrix multiplication only. Finally, under

the type of restrictions here considered, M^{**} partitions into diagonal blocks M^g , and its inversion therefore requires only the inversion of G matrices of orders Q_1, \dots, Q_G . The matrix L^* is not similarly partitionable, although it has other regularities of which a smaller advantage could probably be taken.

*4.3.5.3. *Computational procedure in the Newton method.* Since the matrix L_n^* depends on n , there is no incentive in the Newton method to avoid the explicit appearance of the vectors a_n^* in the computations. We shall therefore develop the formulae largely in terms of $*$ -coordinates.

As in the other methods, there is a possibility of saving computational work whenever certain elements of α^* do not enter the Jacobian B . Assume, therefore, that Φ^* is in canonical form. There will be a further computational advantage in assuming that at least the submatrices $\Phi_I^g, \Phi_{II}^g, \Phi_{III}^g$ are made mutually orthogonal by suitable choice of the Ω in (4.143). To simplify the formulae, we shall assume row-by-row orthogonality

$$(4.177) \quad \Phi^* \Phi'^* = I$$

of Φ^* , although in actual computations it need not be economical to go that far.

We shall now relate l_0^* and L_0^* in (4.173) to the initial vector a_0^* . Using (4.177), we evaluate l_0^* from (4.46) and (4.49) as

$$(4.178) \quad l_0^* \equiv \left(\frac{dL(\alpha^*)}{d\alpha^*} \right)_{\alpha^* = a_0^*} = \text{vec}^* (B_0'^{-1} I_{[K_y K_x]} - A_0 M_{xx}).$$

On the other hand we have, for any $A = \text{mat}^* \alpha^*$, from the definition (4.171) of L_0^* ,

$$(4.179) \quad \begin{aligned} \alpha^* L_0 &= \frac{1}{2} \frac{d}{d\alpha^*} (\alpha^* L_0^* \alpha'^*) \\ &= \text{vec}^* (-B_0'^{-1} B' B_0'^{-1} I_{[K_y K_x]} - A M_{xx}), \end{aligned}$$

and, hence, in particular,

$$(4.180) \quad \alpha_0^* L_0^* = \text{vec}^* (-B_0'^{-1} I_{[K_y K_x]} - A_0 M_{xx}).$$

When (4.178) and (4.179) are inserted in (4.173), the terms containing $A_0 M_{xx}$ cancel. Using (4.179) again with α_1^* substituted for α^* , we see that (4.173) is equivalent to

$$(4.181) \quad \begin{aligned} -\alpha_1^* L_0^* &\equiv \text{vec}^* (B_0'^{-1} B_1' B_0'^{-1} I_{[K_y K_x]} + A_1 M_{xx}) \\ &= 2 \text{vec}^* (B_0'^{-1} I_{[K_y K_x]}). \end{aligned}$$

Computation can conveniently be based on this formulation of the Newton method or on an equivalent formulation in terms of ΔA_0 instead of A_1 .

We shall use postsubscripts and presubscripts I, II, III to denote submatrices of Φ^* , L_n^* , and subvectors of α^* , etc., corresponding to the canonical form of Φ^* . For instance

$$(4.182) \quad \text{III } X_I^* \equiv \begin{bmatrix} \text{III } X^1 & 0 & \dots & 0 \\ 0 & \text{III } X^2 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \text{III } X^G \end{bmatrix}.$$

It is seen from (4.181) [or from the definition (4.171) of L_0^*] that

$$(4.183) \quad L_0^* = \begin{bmatrix} \text{III } \text{III } L_0^* & -\text{III } M_{\text{III}}^{**} \\ -\text{III } M_{\text{III}}^{***} & -M_{\text{III } \text{III}}^{**} \end{bmatrix},$$

because the first term in the second member of (4.181) does not give rise to any III-components of L_0^* . Similarly, the last member in (4.181) has vanishing III-components. It follows that

$$(4.184) \quad a_1^* \begin{bmatrix} - {}_{III} M_{III}^{**} \\ - M_{III III}^{**} \end{bmatrix} = - {}_{III} a_1^* \cdot {}_{III} M_{III}^{**} - a_{1,III}^* \cdot M_{III III}^{**} = 0,$$

from which $a_{1,III}^*$ can be solved according to (4.161) and inserted in (4.181) to give

$$(4.185) \quad {}_{III} a_1^* \cdot {}_{III} L_0^* = {}_{III} \text{vec}^* [B_0'^{-1} \quad 0],$$

with

$$(4.186) \quad {}_{III} L_0^* \equiv {}_{III III} L_0^* + {}_{III} M_{III}^{**} \cdot (M_{III III}^{**})^{-1} \cdot {}_{III} M_{III}^{**}.$$

The following steps arise in the application of (4.185). Initial "overhead" work consists in orthogonalization of Φ^* , determination of M^{**} from (4.72), and calculation of the last term in (4.186), which remains the same in all iterations. The inversion of $M_{III III}^{**}$ can be carried out for each of its diagonal blocks $M_{III III}^{**g}$ separately. Then one chooses B_0 , calculates its inverse and uses it in calculating both

$$(4.187) \quad {}_{III} \text{vec}^* [B_0'^{-1} \quad 0] = \text{vec} [B_0'^{-1} \quad 0] \cdot {}_{III} \Phi'^*$$

and

$$(4.188) \quad {}_{III III} L_0^* = - {}_{III} X^* \cdot K_0 \cdot {}_{III} X'^* - {}_{III III} M^{**}.$$

The matrix K_0 is defined by

$$(4.189) \quad \text{tr}(B_0'^{-1} B')^2 \equiv {}_{III} a^* \cdot {}_{III} X^* \cdot K_0 \cdot {}_{III} X'^* \cdot {}_{III} a'^*,$$

and has as elements

$$(4.190) \quad k_{g i, h j} = \iota(i) \cdot B_0^{-1} \cdot \iota'(h) \cdot \iota(j) \cdot B_0^{-1} \cdot \iota'(g),$$

arranged (with $G = K_y$) according to

$$(4.191) \quad \begin{bmatrix} k_{11,11} & \cdots & k_{11,1G} & \cdots & k_{11,G1} & \cdots & k_{11,GG} \\ \cdot & \cdots & \cdot & \cdots & \cdot & \cdots & \cdot \\ k_{1G,11} & \cdots & k_{1G,1G} & \cdots & k_{1G,G1} & \cdots & k_{1G,GG} \\ \cdot & \cdots & \cdot & \cdots & \cdot & \cdots & \cdot \\ k_{G1,11} & \cdots & k_{G1,1G} & \cdots & k_{G1,G1} & \cdots & k_{G1,GG} \\ \cdot & \cdots & \cdot & \cdots & \cdot & \cdots & \cdot \\ k_{GG,11} & \cdots & k_{GG,1G} & \cdots & k_{GG,G1} & \cdots & k_{GG,GG} \end{bmatrix} \cdot$$

Finally, ${}_{III}L_0^*$ is put together from (4.186), and ${}_{III}a_1^*$ solved from (4.185), leading to

$$(4.192) \quad B_1 = \text{mat}({}_{III}a_1^* \cdot {}_{III}X^*).$$

At the termination of iterations, C_n is obtained from

$$(4.193) \quad \begin{aligned} C_n &= \text{mat}({}_{III}a_n^* \cdot {}_{III}\Psi^* + a_n^* \cdot {}_{III}\Psi_{III}^*) \\ &= \text{mat} [{}_{III}a_n^* \{ {}_{III}\Psi^* - {}_{III}M_{III}^{**} \cdot (M_{III}^{**})^{-1} \cdot {}_{III}\Psi_{III}^* \}]. \end{aligned}$$

It appears from (4.188) that the matrix ${}_{III}L_0^*$ has considerable regularity in its make-up. The problem of how to best utilize those regularities for the inversion of ${}_{III}L_0^*$ or for the solution of (4.185) has not been investigated by us.

4.3.5.4. *Numerical illustration of the Newton method.* This method has likewise been applied to the constructed example already discussed in which the basic matrices have the simple form (4.128). The superior speed of convergence of the Newton method comes out clearly in the results shown in Table 4.3.4.4. More experience with actual data is required to determine whether and in what circumstances the greater speed of convergence is adequate compensation for the greatly increased labor per iteration.

4.3.5.5. *Estimated sampling variances and covariances of the*

estimated coefficients a_{gk} . Even if another method is used to obtain a satisfactory approximation A_n to A , it is still advisable to make one final iteration with the Newton method in order to obtain the matrix of estimated sampling variances and covariances

$$(4.194) \quad \text{est } \mathcal{E}\{(a'^* - \alpha'^*)(a^* - \alpha^*)\} = -\frac{1}{T} \left[\frac{\partial^2 L(\alpha^*)}{\partial \alpha'^* \partial \alpha^*} \right]_{\alpha^* = a^*}^{-1} \\ = -\frac{1}{T} L^{*-1}$$

of the estimated parameters a^* as a by-product. It was shown in Theorem 3.3.10 that the estimates (4.194) are consistent. A suitable method of obtaining L^{*-1} in the present circumstances is the partitioning method whereby ${}_{III} {}_{III} (L^{*-1})$ is obtained as the inverse

$$(4.195) \quad {}_{III} {}_{III} (L^{*-1}) = ({}_{III} L^*)^{-1}$$

of ${}_{III} L^*$ as a step in solving ${}_{III} a_1^*$ from (4.185), and $(L^{*-1})_{III} {}_{III}$ and $(L^{*-1})_{III}$ are found from similar formulae [Hotelling, 1943 A, p. 4], quoted and used before.

Because of the normalization rule (4.25) here employed, the sampling variances (4.194) are not in the form in which they are normally expressed. One will usually regard as final parameters the ratios

$$(4.196) \quad \frac{\alpha^{g \cdot \nu'}(q)}{\alpha^{g \cdot \nu'}(1)}, \quad q = 2, \dots, Q_g, \quad g = 1, \dots, G,$$

of the elements of each α^g . Since the estimates (4.194) themselves are first-order approximations that become only asymptotically exact as the sample size T tends to infinity, sampling variances and covariances of the estimates $\alpha^{g \cdot \nu'}(q) / \alpha^{g \cdot \nu'}(1)$ of the parameters (4.196), of an equal order of approximation, can be found by Taylor expansions in which only terms linear in the estimates (4.194) are retained.

Alternatively, one may normalize on the $\alpha^{g \cdot \nu'}(1)$ by (4.26) and treat the diagonal elements σ_{gg} of Σ as unknown parameters. This procedure may lead to a further saving in computational labor because the parameters σ_{gg} so introduced fall in the same category as

the parameters α_{III}^g : for any $III a_n^g$ the corresponding maximizing values of α_{III}^g and σ_{gg} are easily found. The order of the inversion (4.195) can therefore be further reduced by the number G of parameters σ_{gg} . It is not necessary to go into the details of this procedure since the application of the normalization (4.26) will be demonstrated in section 4.4 in the case where Σ is entirely unrestricted.

From the estimated sampling variances and covariances of the $\alpha^{g \cdot v}(q) / \alpha^{g \cdot v}(1)$ we may revert to the singular matrix of estimated sampling variances and covariances of the estimates a_{gk} of the structural coefficients α_{gk} through the transformations (4.22).

4.4. The Case of Unrestricted Correlations between the Disturbances

4.4.1. *No restrictions on Σ .* We shall now study the case in which no a priori restrictions are imposed on the matrix Σ of variances and covariances of the disturbances in the structural equations (1.1) except the symmetry and positive-definiteness conditions arising from its definition. The discussion can be brief in those aspects of the problem that are also found in the case of uncorrelated disturbances just discussed. The main emphasis will be on points of difference between the two cases.

4.4.2. *Normalization.* With the nondiagonal elements of Σ unrestricted in any case, it is not convenient to impose normalization through the diagonal elements σ_{gg} of Σ . We shall either impose no normalization at all or normalize through one element of each vector α^g , for which we may conveniently take the first element $\alpha^{g \cdot v}(1) = 1$. In the latter case we shall employ the notation

$$(4.197) \quad \begin{aligned} \alpha_{[1]}^* &\equiv [\alpha^{1 \cdot v}(1) \quad \alpha^{2 \cdot v}(1) \quad \dots \quad \alpha^{G \cdot v}(1)] \\ &= [1 \quad 1 \quad \dots \quad 1] \end{aligned}$$

to express the normalization rule, and introduce similar notations

$$\Phi^g = \begin{bmatrix} \alpha^g \\ \Phi_{[1]}^g \\ [1] \Phi^g \end{bmatrix},$$

$$(4.198) \quad \Phi_{[1]}^* \equiv \begin{bmatrix} \Phi_{[1]}^1 & 0 & \dots & 0 \\ 0 & \Phi_{[1]}^2 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & \Phi_{[1]}^G \end{bmatrix},$$

$$[1] \Phi^* \equiv \begin{bmatrix} \Phi^1 & 0 & \dots & 0 \\ [1] \Phi^2 & \dots & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & [1] \Phi^G \end{bmatrix},$$

for the corresponding partitioning of the basic matrices.

4.4.3. *Elimination of the parameters Σ .* We shall first maximize the likelihood function

$$(4.199) \quad L(\mathbf{A}, \Sigma) = \text{const} + \log \det \mathbf{B} - \frac{1}{2} \log \det \Sigma - \frac{1}{2} \text{tr}(\Sigma^{-1} \mathbf{A} M_{xx} \mathbf{A}')$$

with respect to the unrestricted parameters Σ while the parameters \mathbf{A} are kept constant. From (3.17) and (3.18) it is easily seen that the first derivatives $\partial L / \partial \sigma_{gh}$, $g, h = 1, \dots, G$, vanish if

$$(4.200) \quad \Sigma = \hat{\Sigma} \equiv \mathbf{A} M_{xx} \mathbf{A}' = \hat{\Sigma}'.$$

The derivation of (4.200) must take account of the required symmetry of Σ but the result is not affected thereby. It follows from (3.35) that (4.200) indicates the unique and absolute maximum of the function (4.199) with respect to Σ . Upon inserting (4.200) in (4.199), we obtain

$$(4.201) \quad L(A) = \text{const} + \log \det B - \frac{1}{2} \log \det(A M_{xx} A')$$

as the likelihood function after maximization with respect to Σ . This function is homogeneous of degree 0 in each row of A , i.e., it is invariant for changes in normalization of A through premultiplication with a nonsingular diagonal matrix Υ . This is easily verified directly or can be seen as a consequence of the invariance of (4.201) under the wider group of nonsingular transformations implied in Theorem 2.1.3.5.

The maximum of (4.201) in the absence of restrictions on A has already been studied in the analysis leading to Theorem 3.1.10.

4.4.4. *The maximum-likelihood equations.* We shall continue to use the symbol $\hat{\Sigma}$ as an abbreviation for the expression (4.200) in terms of A . Similarly, we shall use the abbreviations

$$(4.202) \quad S_n = A_n M_{xx} A'_n, \quad S = A M_{xx} A'.$$

Again writing $A = A_0 + \Delta A_0$, where A_0 is a trial value, we have, using (3.16) and (3.17),

$$(4.203) \quad \begin{aligned} L(A) - L(A_0) &= \text{tr}(B_0'^{-1} \Delta B_0') \\ &- \frac{1}{2} \text{tr}\{S_0^{-1} (A_0 M_{xx} \Delta A_0' + \Delta A_0 M_{xx} A_0')\} + \dots \\ &= \text{tr}\{B_0'^{-1} I_{[K_y K_x]} - S_0^{-1} A_0 M_{xx}\} \Delta A_0' + \dots \end{aligned}$$

The restrictions are

$$(4.204) \quad \alpha = \alpha^* \Phi^*, \quad a_0 = a_0^* \Phi^*, \quad \text{vec } \Delta A_0 = (\text{vec }^* \Delta A_0) \Phi^*.$$

In the absence of normalizing restrictions on A , A_0 can coincide with a restricted maximum A of the likelihood function only if the linear term in (4.203) vanishes for all values of $\text{vec }^* \Delta A_0$. The first-order maximum-likelihood conditions in this case are therefore, owing to Lemma 4.2.4,

$$(4.205) \quad \text{vec}^* (B'^{-1} I_{\begin{bmatrix} K_y & K_x \end{bmatrix}} - S^{-1} A M_{xx}) = 0,$$

with S again depending on A according to (4.202). As before, these conditions permit premultiplication of A by a diagonal matrix, through which we may satisfy the normalization (4.197) if desired.

4.4.5. *The processes \mathfrak{P}_h and \mathfrak{P}_{h_n} .* For the generalization of these processes to the present case, we shall for the time being not impose a normalization rule on A . Given an initial value A_0 , the following extension of the definition (4.57) of the direction matrix ΔA_0 seems natural as well as the simplest possible:

$$(4.206) \quad \text{vec}^*(S_0^{-1} \Delta A_0 M_{xx}) = \text{vec}^*(B_0'^{-1} I_{\begin{bmatrix} K_y & K_x \end{bmatrix}} - S_0^{-1} A_0 M_{xx}).$$

Comparison with (4.203) shows that a property similar to (4.58) in the previous case again holds: If $A_1 = A_0 + h \Delta A_0$, a sufficiently small value of h will always lead to $L(A_1) > L(A_0)$ if a stationary value of $L(A)$ is not already reached in A_0 .

One can again choose a suitable constant value of h , or a value h_0 determined from the principle underlying \mathfrak{P}_{h_n} . To obtain the latter we write as an extension of (4.203), using also (3.19),

$$\begin{aligned} L(A_1) - L(A_0) &= \\ &= h \text{tr} \{ (B_0'^{-1} I_{\begin{bmatrix} K_y & K_x \end{bmatrix}} - S_0^{-1} A_0 M_{xx}) \Delta A_0' \} \\ (4.207) \quad &+ \frac{1}{2} h^2 \text{tr} \{ - (B_0'^{-1} \Delta B_0')^2 + S_0^{-1} (A_0 M_{xx} \Delta A_0' \\ &+ \Delta A_0' M_{xx} A_0') S_0^{-1} A_0 M_{xx} \Delta A_0' - S_0^{-1} \Delta A_0 M_{xx} \Delta A_0' \} \\ &+ \dots \end{aligned}$$

Using (4.206) we find that the sum of the two terms shown in (4.207) is maximized if h is given the value

$$(4.208) \quad h_0 = \frac{\text{tr}(S_0^{-1} \Delta A_0 M_{xx} \Delta A_0')}{\text{tr}\{(B_0'^{-1} \Delta B_0') - S_0^{-1}(A_0 M_{xx} \Delta A_0' + \Delta A_0 M_{xx} A_0) S_0^{-1} A_0 M_{xx} \Delta A_0' + S_0^{-1} \Delta A_0 M_{xx} \Delta A_0'\}}$$

We define matrices M_0^* and L_0^* such that the numerator and denominator in (4.208) are identical with those in

$$(4.209) \quad h_0 = \frac{\Delta a_0^* M_0^* \Delta a_0'^*}{-\Delta a_0^* L_0^* \Delta a_0'^*},$$

postponing their explicit evaluation until the discussion of computational arrangements below.

4.4.6. *Asymptotic properties of $\hat{\rho}_h$ and $\hat{\rho}_{h_n}$.* If we write as before $A_n = A + \bar{A}_n$, the expansion of (4.206) in terms of \bar{A}_0 is, by use of (4.205),

$$(4.210) \quad \text{vec}^*(S^{-1} \Delta \bar{A}_0 M_{xx}) + \dots \\ = \text{vec}^*\{-B'^{-1} \bar{B}_0' B'^{-1} I_{[K_y K_x]} \\ + S^{-1}(A M_{xx} \bar{A}_0' + \bar{A}_0 M_{xx} A') S^{-1} A M_{xx} - S^{-1} \bar{A}_0 M_{xx}\}.$$

Omitting bars from Δa_n and ΔA_n (see p.174), this can be written as

$$(4.211) \quad \Delta a_0^* M^* + \dots = \bar{a}_0^* L^* + \dots,$$

with suitable definitions of the matrices M^* and L^* , which are now the same functions of A and M_{xx} as the matrices M_0^* and L_0^* , respectively, are of A_0 and M_{xx} .

The study of (4.211) is exactly similar to that of (4.79) in the previous case of uncorrelated disturbances. Formulae (4.80) through (4.88) and the discussion connected therewith remain valid

with the new definitions of M^* and L^* . Limits on the characteristic values l_q , $q = 1, \dots, Q$, of L^* can again be determined in the case that the maximum of the likelihood function is not depressed by the restrictions, as follows: Retaining for that case the definition (4.93) of H , we have instead of (4.94)

$$(4.212) \quad H M_{xx} H' = \begin{bmatrix} S & 0 \\ 0 & I_{[K_z]} \end{bmatrix} = T,$$

say. Writing for the moment, instead of (4.95),

$$(4.213) \quad \bar{A} = \tilde{A}^\oplus H,$$

we have from (4.212) and (4.202)

$$(4.214) \quad \begin{aligned} L_{(2)} &\equiv \text{tr}\{- (B'^{-1} \bar{B}')^2 + S^{-1}(A \bar{A}' + \bar{A} M_{xx} A') S^{-1} A M_{xx} \bar{A}' - S^{-1} \bar{A} M_{xx} \bar{A}'\} \\ &= \text{tr}\{- (\tilde{B}'^\oplus)^2 + (\tilde{B}'^\oplus + S^{-1} \tilde{B}^\oplus S) \tilde{B}'^\oplus - S^{-1} (\tilde{B}^\oplus S \tilde{B}'^\oplus + \tilde{C}^\oplus \tilde{C}'^\oplus)\} \\ &= -\text{tr}(S^{-1} \tilde{C}^\oplus \tilde{C}'^\oplus) \end{aligned}$$

and

$$(4.215) \quad M_{(2)} \equiv \text{tr}(S^{-1} \bar{A} M_{xx} \bar{A}') = \text{tr}\{S^{-1} (\tilde{B}^\oplus S \tilde{B}'^\oplus + \tilde{C}^\oplus \tilde{C}'^\oplus)\}.$$

Through a further transformation

$$(4.216) \quad S = U U', \quad U^{-1} \tilde{A}^\oplus \begin{bmatrix} U & 0 \\ 0 & I_{[K_z]} \end{bmatrix} = \tilde{A},$$

it is seen that, in the absence of restrictions on A , the characteristic values l_q are the stationary values of the quadratic form

$$(4.217) \quad L_{(2)} = -\text{tr}(\tilde{C} \tilde{C}') = -\sum_{g=1}^{K_y} \sum_{h=K_y+1}^{K_x} \tilde{\alpha}_g^2 h$$

under the restrictions

$$(4.218) \quad M_{(2)} = \text{tr}(\tilde{A} \tilde{A}') = \sum_{g=1}^{K_y} \sum_{h=1}^{K_x} \tilde{a}_{gh}^2 = 1.$$

We record in one formula combining (4.213) and (4.216) the transformation

$$(4.219) \quad \bar{A} = U \tilde{A} \begin{bmatrix} U^{-1} & 0 \\ 0 & I_{[K_x]} \end{bmatrix} H = U \tilde{B} U^{-1} A + U \tilde{C} F I_{[K_x K_x]}$$

through which the forms (4.217) and (4.218) have been derived. These forms lead to the following complete table of characteristic values and vectors.

	(a)	(b)	(c)	(d)
(4.220)	Value of l	Value of $k = 1 + h l$	Multiplicity	Characteristic "vectors" \tilde{A} satisfy
	0	1	$K_y K_z$	$\tilde{C} = 0$
	-1	$1 - h$	$(K_y)^2$	$\tilde{B} = 0$

Therefore, under any a priori restrictions that permit the likelihood function to attain its absolute maximum,

$$(4.221) \quad -1 \leq l_q \leq 0, \quad 1 - h \leq k_q \leq 1, \quad q = 1, \dots, Q.$$

Furthermore, under any such restrictions that in addition ensure (as here supposed) complete identification of each structural equation,

$$(4.222) \quad -1 \leq l_q < 0, \quad 1 - h \leq k_q < 1, \quad q = 1, \dots, Q - K_y,$$

from which we exclude the K_y characteristic values $l = 0$, connected with the freedom of normalization of A (choice of diagonal elements of $\tilde{B}^\oplus = U \tilde{B} U^{-1}$) and unaffected by homogeneous restrictions (4.204).

4.4.7. *Considerations in choosing a constant value of h .* It follows that, among processes \mathcal{P}_h with a constant value of h , \mathcal{P}_1 does

not have the excellence it possesses in one important case with uncorrelated disturbances. In large samples under valid restrictions, \mathcal{P}_1 confines the characteristic values $k_q = 1 + h l_q$ approximately to the interval $0 \leq k_q < 1$, so that, unless all k_q vanish, $\max |k_q|$ can be decreased by taking $h > 1$.

As a guide in determining how far above 1 to choose h , it is of interest to ask what type of restrictions will exclude the characteristic value $l = -1$. This value will remain present as long as the restrictions permit an addition to A of the type

$$(4.223) \quad \bar{A} = U \tilde{C} F I_{[K_z K_x]} = \bar{C} I_{[K_z K_x]}$$

containing only the second term of (4.219). It follows that, for the exclusion of the characteristic value $l = -1$, it is necessary and sufficient that in the canonical form of the basic matrix Φ^* the submatrix Φ_{III}^* be absent. If this is the case, the same reasoning from ignorance that previously favored \mathcal{P}_1 , now leads to the recommendation of \mathcal{P}_2 : the relevant values k_q are then confined to the interval $-1 < k_q < 1$. However, if Φ_{III}^* is present, any constant value of h should be chosen below 2, and the nearer to 2, the nearer the highest of the values l_q , $q = 1, \dots, Q - K_y$, is suspected of being to zero.

*4.4.8. *Problems in the arrangement of computations for \mathcal{P}_h , \mathcal{P}_h^n .* We shall now write (4.206), using (4.40) and the orthogonalization (4.177) of Φ^* , in the form

$$(4.224) \quad \Delta a_n^* M_n^* = \text{vec}(B_n'^{-1} I_{[K_y K_x]}) \cdot \Phi'^* - a_n^* M_n^*,$$

$n = 0, 1, \dots,$

where M_n^* is defined by

$$(4.225) \quad M_n^* \equiv \Phi^* M_n \Phi'^*,$$

$$(4.226) \quad M_n \equiv S_n^{-1} \otimes M_{xx} \equiv \begin{bmatrix} s_n^{11} M_{xx} & s_n^{12} M_{xx} & \dots & s_n^{1G} M_{xx} \\ s_n^{21} M_{xx} & s_n^{22} M_{xx} & \dots & s_n^{2G} M_{xx} \\ \cdot & \cdot & \dots & \cdot \\ s_n^{G1} M_{xx} & s_n^{G2} M_{xx} & \dots & s_n^{GG} M_{xx} \end{bmatrix}.$$

with s_n^{gh} denoting the elements of

$$(4.227) \quad S_n^{-1} \equiv [s_n^{gh}],$$

which are again functions of A_n . Combining (4.225) and (4.226), we can alternatively write for M_n^* , using (4.31),

$$(4.228) \quad M_n^* = \begin{bmatrix} M_n^{11} & \dots & M_n^{1G} \\ \cdot & \dots & \cdot \\ M_n^{G1} & \dots & M_n^{GG} \end{bmatrix},$$

with the further definition

$$(4.229) \quad M_n^{gh} \equiv s_n^{gh} \Phi^g M_{xx} \Phi'^h \equiv s_n^{gh} V^{gh},$$

say. The symbol V^{gh} is merely an abbreviation for the matrix product it represents. There is no meaning, in the present context, in putting the matrices V^{gh} together to form a larger matrix V^* . In the special case that $S_n = I_{[K_y]}$, M_n^* goes over into M^{**} as defined in (4.72), in which $M_n^{gh} = 0$ for $g \neq h$.

Since M_n^* changes from one iteration to the next, there is no advantage in avoiding explicit use of q -coordinates. Likewise, in solving for Δa_n^* from (4.224), there is no greater advantage from the use of the canonical form of the basic matrix Φ^* than there is in general from the use of any partitioning method for the inversion of a matrix or the solution of linear equations. The new element in the present situation, as compared with the application of Φ_h in the case of uncorrelated disturbances, is that now M_n^* does not partition into diagonal blocks. In principle, therefore, we now have one high-order inversion job instead of K_y lower-order inversions - a situation such as was already encountered in the Newton method in the case of uncorrelated disturbances (because L_n^* likewise does not partition). The main problem of computational economy now is to find an efficient method of solving for Δa_n^* from (4.224) which takes advantage of the special form of M_n^* . This problem again has not been systematically investigated by us. The following considerations seem relevant.

Of the matrices entering into the definition of M_n^* , those re-

maining the same through all iterations are M_{xx} and Φ^* . This suggests that it will be advantageous to go as far as possible toward the solution of Δa_n^* on the basis of these matrices alone before the matrix S_n specific to the n th iteration is brought into play. One possible procedure would be to start from the matrices V^{gh} as basic material, developing functions of these matrices that facilitate the solution of Δa_n^* from (4.224) for all n . This method would be similar to the partitioning method of matrix inversion, although a complete inversion of M_n^* may not be needed.

A perhaps more powerful method would be to utilize the common origin of all V^{gh} in M_{xx} on the basis of one initial inversion of M_{xx} used in

$$(4.230) \quad M_n^{-1} = S_n \otimes M_{xx}^{-1},$$

followed by

$$(4.231) \quad M_n^{-1}(\ast) = \Phi'^{-1}(\ast) M_n^{-1} \Phi^{-1}(\ast).$$

[The evaluation of (4.231) may be facilitated by orthogonalization of $\Phi(\ast)$.] This approach requires an economical method of finding $M_n^{\ast-1}$ if both M_n^* and $M_n^{-1}(\ast)$ are available.

*4.4.9. Processes \mathcal{P}_h and \mathcal{P}_n modified by normalization. In the derivation of the first-order maximum-likelihood conditions (4.205) from (4.203) we have not imposed any normalization on A_0 and A . If, alternatively, we had required that both A_0 and A satisfy the normalization rule (4.197), ΔA_0 would have been restricted by

$$(4.232) \quad \text{vec}_{[1]}^* \Delta A_0 = 0,$$

and the first-order condition for a maximum would be expressed by

$$(4.233) \quad \begin{cases} (4.233, -1) & [1] \text{vec}^*(B'^{-1} I_{[K_y K_x]} - S^{-1} A M_{xx}) = 0, \\ (4.233, +1) & a_{[1]}^* = [1 \ 1 \ \cdots \ 1]. \end{cases}$$

It follows from the homogeneity properties of the likelihood function that (4.233) is equivalent to (4.205). However, the anal-

ogous iterative procedures using ΔA_0 defined by

$$(4.234) \left\{ \begin{array}{l} (4.234, -1) \quad [1] \text{vec}^*(S_0^{-1} \Delta A_0 M_{xx}) \\ \qquad \qquad \qquad = [1] \text{vec}(B_0'^{-1} I_{[K_y K_x]} - S_0^{-1} A_0 M_{xx}), \\ (4.234, +1) \quad \Delta a_0^*, [1] = 0, \quad \text{or} \\ \qquad \qquad \qquad \Delta a_0 = (\Delta [1] a_0^*) [1] \Phi^* , \end{array} \right.$$

are essentially different from those based on (4.206). For, even if A_0 satisfies the normalization rule (4.197), the solution ΔA_0 of (4.206) cannot satisfy (4.234, +1) for all possible choices of that row of each Φ^{β} on which normalization is based. For if that were so, Δa_0^* and therewith ΔA_0 would vanish identically. In general, therefore, the substitution of (4.234, +1) for an equal number K_y of the equations (4.206) leads to nonproportional changes in the elements of the solution ΔA_0 .

It follows that the convergence properties of the modified processes ρ_h , ρ_{h_n} based on (4.234) differ from those derived from (4.206) and depend on what rows of Φ^* have been selected for normalization purposes. We have not investigated the effect of this application of the normalization rule (4.197) on the asymptotic convergence properties. There is reason to believe that the effect is not a radical modification. For, whatever value of h is chosen, the characteristic value $l = 0$ corresponding to the diagonal elements of $\tilde{B}^{\oplus} = U \tilde{B} U^{-1}$ in (4.220) connected with the scales of the rows of A leads to $k = 1$. There is therefore no alternation or other unsteadiness in scales in the application of (4.206). Thus, in the first approximation, the elements of Δa_0^* determined from (4.206) differ from zero only to an extent required for improving the ratios of the elements of a_0^* in the next approximation a_1^* . Hence the fixing of certain elements of Δa_0^* at the value zero while relaxing an equal number of the conditions (4.206) might somewhat retard, but need not destroy, convergence.

This point is important because the modification of ρ_h through normalization reduces by K_y the number of unknowns in Δa_n^* to be

determined in each iteration. The resulting saving of labor per iteration may be considerable except possibly in methods based on the initial inversion of M_{xx} for use in (4.230).

The computational arrangement for the modified processes \mathcal{P}_h , \mathcal{P}_{h_n} differs from that described earlier only in easily perceived details.

4.4.10. *Numerical illustrations of \mathcal{P}_1 , $\mathcal{P}_{5/4}$.* We have applied \mathcal{P}_1 and $\mathcal{P}_{5/4}$ without normalization, and \mathcal{P}_1 modified by normalization, to the constructed example discussed before. As initial values we have taken the result of the 6th iteration with \mathcal{P}_1 in the case where Σ was assumed to be diagonal. The lower cost per iteration under that assumption is a good reason to apply it initially, albeit only to get closer to the maximizing value A without restrictions on Σ . The results are shown in Table 4.4.10.

*4.4.11. *The Newton method.* All properties previously derived from the general formulation (4.173) of the Newton method carry over, of course, to the present case. All that is here required is to derive new expressions for the first and second derivatives of the likelihood function in the initial point A_0 , generalizing formulae (4.178) and (4.179) of the previous case.

It will be remembered that the Newton method requires the matrix of second derivatives of the likelihood function to be nonsingular. As before, complete identification, and, as a new requirement, use of the normalization rule (4.197), are therefore now indispensable. Instead of the previous formulation (4.172) we thus obtain as the definition of the Newton method

$$(4.235) \quad \begin{aligned} a_1^* &= a_0^* + \Delta a_0^*, \\ (\Delta_{[1]} a_0^*)_{[1]} L_0^* &= - [1] l_0, \quad \Delta a_{0,[1]}^* = 0. \end{aligned}$$

Assuming again complete orthogonality (4.177) of Φ^* (which is compatible with any choice of one row of each Φ^g for normalization purposes), we have from (4.203)

$$(4.236) \quad [1] l_0 \equiv \left(\frac{dL(\alpha^*)}{d[1]\alpha^*} \right)_{\alpha^* = \alpha_0^*} = [1] \text{vec}^* \{ B_0'^{-1} I_{[K_y K_x]} - S_0^{-1} A_0 M_{xx} \}.$$

TABLE 4.4.10
 Numerical illustrations of \mathcal{P}_1 and $\mathcal{P}_{5/4}$ without restrictions on Σ .

Method	Iteration $n =$	Row $g =$	Matrices B_n			Matrices S_n			Scale Factors ² $a_{n;g,g+3}$
			$k = 1$	2	3	$h = 1$	2	3	
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
\mathcal{P}_1 without normalization ¹	0	1	0.00000	0.97020	4.17885	0.19385	0.08535	-0.02277	
		2	1.00207	0.00000	-2.99581	0.08535	0.20000	0.08700	
		3	-2.04599	0.96553	0.00000	-0.02277	0.08700	0.29509	
	1	1	0.00000	0.99521	4.04308	0.19834	0.09630	-0.00515	
		2	0.99144	0.00000	-2.97534	0.09630	0.19674	0.09717	
		3	-2.01349	0.99600	0.00000	-0.00515	0.09717	0.29990	
	2	1	0.00000	0.99892	4.00812	0.19958	0.09919	-0.00103	
		2	0.99789	0.00000	-2.99534	0.09919	0.19929	0.09948	
		3	-2.00275	0.99930	0.00000	-0.00103	0.09948	0.30000	
	3	1	0.00000	0.99976	4.00157	0.19990	0.09982	-0.00021	
		2	0.99951	0.00000	-2.99918	0.09982	0.19986	0.09991	
		3	-2.00058	0.99989	0.00000	-0.00021	0.09991	0.30001	
True Values		1	0.00000	1.00000	4.00000	0.20000	0.10000	0.00000	
		2	1.00000	0.00000	-3.00000	0.10000	0.20000	0.10000	
		3	-2.00000	1.00000	0.00000	-0.00000	0.10000	0.30000	

¹As described in section 4.4.8, no normalization has been imposed as part of the computation of each iteration result from the preceding iteration result. See, however, note 2.

²For the purpose of comparison of successive iteration results, each iteration result has been re-normalized by $a_{14} = a_{25} = a_{36} = 1$ before being entered in this table. Column (10) gives values of $a_{n,14}$, $a_{n,25}$, $a_{n,36}$ obtained, before such re-normalization, by one application of \mathcal{P}_1 (without normalization) to the matrices B_{n-1} , S_{n-1} as stated in the table.

TABLE 4.4.10
(Continued)

Method	Iteration $n =$	Row $g =$	Matrices B_n			Matrices S_n			Scale ² Factors
			$k = 1$	2	3	$h = 1$	2	3	$a_{n;g,g+3}$
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
$P_{5/4}$ without normalization ¹	0	1	0.00000	0.97020	4.17885	0.19385	0.08535	-0.02277	1.00000
		2	1.00207	0.00000	-2.99581	0.08535	0.20000	0.08700	1.00000
		3	-2.04599	0.96553	0.00000	-0.02277	0.08700	0.29509	1.00000
	1	1	0.00000	1.00146	4.00914	0.20059	0.09890	0.00002	1.00140
		2	0.98878	0.00000	-2.97022	0.09890	0.19595	0.09957	1.05436
		3	-2.00537	1.00362	0.00000	0.00002	0.09957	0.30200	0.99498
	2	1	0.00000	0.99944	3.99793	0.19978	0.10020	0.00000	1.00215
		2	1.00258	0.00000	-3.00778	0.10020	0.20104	0.10020	0.99025
		3	-1.99864	0.99936	0.00000	0.00000	0.10020	0.29961	1.00302
	3	1	0.00000	1.00013	4.00054	0.20005	0.09995	0.00000	0.99947
		2	0.99935	0.00000	-2.99806	0.09995	0.19974	0.09995	1.00263
		3	-2.00034	1.00017	0.00000	0.00000	0.09995	0.30010	0.99934
True Values		1	0.00000	1.00000	4.00000	0.20000	0.10000	0.00000	
		2	1.00000	0.00000	-3.00000	0.10000	0.20000	0.10000	
		3	-2.00000	1.00000	0.00000	0.00000	0.10000	0.30000	

¹As described in section 4.4.8, no normalization has been imposed as part of the computation of each iteration result from the preceding iteration result. See, however, note 2.

²For the purpose of comparison of successive iteration results, each iteration result has been re-normalized by $a_{14} = a_{25} = a_{36} = 1$ before being entered in this table. Column (10) gives values of $a_{n,14}$, $a_{n,25}$, $a_{n,36}$ obtained, before such re-normalization, by one application of $P_{5/4}$ (without normalization) to the matrices B_{n-1} , S_{n-1} as stated in the table.

TABLE 4.4.10
(Continued)

Method	Iteration $n =$	Row $g =$	Matrices B_n			Matrices S_n			
			$k = 1$	2	3	$h = 1$	2	3	
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	
P_1 modified by normalization ¹	0	1	0.00000	0.97020	4.17885	0.19385	0.08535	-0.02277	
		2	1.00207	0.00000	-2.99581	0.08535	0.20000	0.08700	
		3	-2.04599	0.96553	0.00000	-0.02277	0.08700	0.29509	
	1	1	1	0.00000	0.99331	4.03780	0.19760	0.09666	-0.00542
		2	2	0.99910	0.00000	-3.00006	0.09666	0.19986	0.09740
		3	3	-2.00830	0.99236	0.00000	-0.00542	0.09740	0.29796
	2	1	1	0.00000	0.99772	4.00567	0.19911	0.09911	-0.00126
		2	2	0.99923	0.00000	-2.99997	0.09911	0.19987	0.09938
		3	3	-2.00071	0.99768	0.00000	-0.00126	0.09938	0.29915
	3	1	1	0.00000	0.99931	4.00083	0.19973	0.09976	-0.00031
		2	2	0.99973	0.00000	-3.00002	0.09976	0.19996	-0.09984
		3	3	-1.99991	0.99931	0.00000	-0.00031	0.09984	0.29972
True Values		1	0.00000	1.00000	4.00000	0.20000	0.10000	0.00000	
		2	1.00000	0.00000	-3.00000	0.10000	0.20000	0.10000	
		3	-2.00000	1.00000	0.00000	0.00000	0.10000	0.30000	

¹As described in section 4.4.9.

TABLE 4.4.10
 (Continued)

Method	Iteration $n =$	Row $g =$	Matrices B_n			Matrices S_n			
			$k = 1$	2	3	$h = 1$	2	3	
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	
$\rho_{5/4}$ modified by normalization ¹	0	1	0.00000	0.97020	4.17885	0.19385	0.08535	-0.02277	
		2	1.00207	0.00000	-2.99581	0.08535	0.20000	0.08700	
		3	-2.04599	0.96553	0.00000	-0.02277	0.08700	0.29509	
	1	1	1	0.00000	0.99909	4.00254	0.19964	0.09950	-0.00038
		2	2	0.99835	0.00000	-3.00112	0.09950	0.19983	0.09999
		3	3	-1.99888	0.99907	0.00000	-0.00038	0.09999	0.29952
	2	1	1	0.00000	0.99972	4.00037	0.19995	0.09997	-0.00003
		2	2	0.99963	0.00000	-3.00031	0.09997	0.20000	0.09999
		3	3	-1.99973	0.99975	0.00000	-0.00003	0.09999	0.29996
	3	1	1	0.00000	0.99997	3.99998	0.20001	0.10002	0.00002
		2	2	0.99999	0.00000	-3.00002	0.10002	0.20001	0.10000
		3	3	-1.99998	0.99998	0.00000	0.00002	0.10000	0.30002
True Values		1	0.00000	1.00000	4.00000	0.20000	0.10000	0.00000	
		2	1.00000	0.00000	-3.00000	0.10000	0.20000	0.10000	
		3	-2.00000	1.00000	0.00000	0.00000	0.10000	0.30000	

¹As described in section 4.4.9.

Furthermore, because $\Delta a_{0,[1]}^* = 0$, we obtain

$$\begin{aligned}
 (\Delta_{[1]} a_0^*)_{[1]} L_0^* &= \frac{1}{2} \frac{d}{d(\Delta_{[1]} a_0^*)} \{ (\Delta_{[1]} a_0^*)_{[1]} L_0^* (\Delta_{[1]} a_0'^*) \} \\
 (4.237) \qquad \qquad \qquad &= \frac{1}{2} \frac{d}{d(\Delta_{[1]} a_0^*)} (\Delta a_0^* L_0^* \Delta a_0'^*).
 \end{aligned}$$

Combining (4.235), (4.236), (4.237), and the definition (4.208) - (4.209) of L_0^* , we obtain

$$\begin{aligned}
 (\Delta_{[1]} a_0^*)_{[1]} L_0^* &= [1] \text{vec}^* \{ - [B_0'^{-1} \Delta B_0' B_0'^{-1} \quad 0] \\
 &\quad + S_0^{-1} (A_0 M_{xxx} \Delta A_0 + \Delta A_0 M_{xxx} A_0') S_0^{-1} A_0 M_{xxx} \\
 (4.238) \qquad \qquad \qquad &\quad - S_0^{-1} \Delta A_0 M_{xxx} \} \\
 &= [1] \text{vec}^* (- B_0'^{-1} I_{[K_y K_x]} + S_0^{-1} A_0 M_{xxx}).
 \end{aligned}$$

as a formulation of the Newton method adapted to computational use. The middle member of (4.238) serves to define $[1] L_0^*$ through

$$(4.239) \qquad \Delta A_0 \equiv \text{mat}^* \{ (\Delta_{[1]} a_0^*)_{[1]} \Phi^* \}.$$

The repeated evaluation of $[1] L_n^*$ from (4.238) and (4.239) and its inversion (or other method of solving for $\Delta_{[1]} a_n^*$) are laborious. The problem of how to take advantage of the regularities in $[1] L_n^*$ for its inversion also appears more formidable than in the case of \mathcal{P}_h where only M_n^* as defined in (4.225) needs to be inverted.

*4.4.12. *Numerical experiment with the Newton method.* A numerical experiment in which (4.238) was applied to the constructed example discussed earlier with the (least-squares) initial value A_0

used in Table 4.3.4.4 did not produce convergent results, in contrast with the modified $\hat{\mathcal{P}}_1$, as defined by (4.234, -1), which led to convergent iterations from the same initial value. We have not repeated the experiment with the (closer) initial value used in Table 4.4.10.

*4.4.13. *Estimated sampling variances and covariances of the estimated coefficients a_{gk} .* It follows from Theorem 3.3.10 that

$$(4.240) \quad \text{est } \mathcal{E} \{ ([1]a'^* - [1]\alpha'^*) ([1]a^* - [1]\alpha^*) \} = -\frac{1}{T} ([1]L^*)^{-1}$$

defines consistent estimates of the sampling variances and covariances of the estimates $[1]a^*$ of the parameters $[1]\alpha^*$. Their evaluation requires inversion of the matrix $[1]L_n^*$ computed from the final value A_n with which iterations are terminated. If this is done from a value A_n obtained by a method other than the Newton method, a check on A_n is obtained at small extra cost by employing $([1]L_n^*)^{-1}$ for one more iteration by the Newton method.

If estimated sampling variances and covariances of the estimates $[1]a^*$, S of all parameters $[1]\alpha^*$, Σ are desired, it is necessary to operate with the second-derivative matrix of the original likelihood function $L(A, \Sigma)$ defined by (4.199). Denoting by

$$(4.241) \quad \begin{aligned} \sigma &= [\sigma_{11} \quad \sigma_{12} \quad \cdots \quad \sigma_{1G} \quad \sigma_{21} \quad \cdots \quad \sigma_{2G} \quad \cdots \quad \sigma_{GG}] \\ &= \text{vec } \Sigma \end{aligned}$$

a row vector containing all independent elements of Σ , we have for any direction $\delta \Sigma = \delta \Sigma'$ of variation of Σ

$$(4.242) \quad \frac{\partial L(A, \Sigma)}{\partial \sigma} \delta \sigma' = -\frac{1}{2} \text{tr} \{ \Sigma^{-1} \delta \Sigma (I - \Sigma^{-1} A M_{xx} A') \}.$$

We introduce the notational definition

$$\left(\begin{array}{c} \left[\begin{array}{cc} \frac{\partial^2}{\partial_{[1]\alpha'^*} \partial_{[1]\alpha^*}} & \frac{\partial^2}{\partial_{[1]\alpha'^*} \partial \sigma} \\ \frac{\partial^2}{\partial \sigma' \partial \alpha^*} & \frac{\partial^2}{\partial \sigma' \partial \sigma} \end{array} \right] L(A, \Sigma) \\ \begin{array}{l} A = A \\ \Sigma = S \end{array} \end{array} \right)$$

(4.243)

$$= \begin{bmatrix} [11] \hat{L}^* & [1] \hat{L}_{\sigma}^* \\ [1] \hat{L}'_{\sigma} & \hat{L}_{\sigma \sigma}^* \end{bmatrix} = \hat{L}.$$

Then, if $\delta a_{[1]}^* = 0$, we have from (4.242) and (4.202), after using (3.16) and (3.19),

$$\begin{aligned} \delta_{[1]\alpha'^*} [11] \hat{L}^* \delta_{[1]\alpha'^*} &= \text{tr} \{ -(B'^{-1} \delta A')^2 - S^{-1} \delta A M_{xx} \delta A' \}, \\ \delta_{[1]\alpha^*} [1] \hat{L}_{\sigma}^* \delta s' &= \text{tr} \{ S^{-1} \delta \Sigma S^{-1} A M_{xx} \delta A' \}, \\ (4.244) \quad \delta s \hat{L}_{\sigma \sigma}^* \delta s' &= -\frac{1}{2} \text{tr} \{ S^{-1} \delta \Sigma \}^2, \end{aligned}$$

where s is defined analogously to σ in (4.241). These formulae serve to evaluate \hat{L} . The desired sampling covariances are now obtained from

$$\begin{aligned} (4.245) \quad \text{est } \mathcal{E} \begin{bmatrix} ([1]\alpha'^* - [1]\alpha'^*) \\ (s' - \sigma') \end{bmatrix} \begin{bmatrix} ([1]\alpha^* - [1]\alpha^*) & (s - \sigma) \end{bmatrix} \\ = -\frac{1}{T} \hat{L}^{-1}, \end{aligned}$$

possibly by the partitioning method of inversion already quoted

[Hotelling, 1943 -A, p. 4]. Of course, the matrix $[_1]L^{*-1}$ inverse to the matrix $[_1]L^*$ defined by (4.238) is a principal submatrix of \hat{L}^{-1} located in the upper left corner.

4.5. Concluding Remarks

4.5.1. *Nature of the concluding remarks.* In this concluding section 4.5 we shall first make rough comparisons between the costs of computation in the various methods discussed. These comparisons will give occasion to recall certain problems of matrix computation which have not been investigated by us and to make some remarks on suitable methods for the various matrix inversions required. Secondly, we shall indicate a possible generalization of the restrictions on A. Finally, we shall draw attention to important problems connected with the number and nature of different maxima of the likelihood function which require solution before full reliance can be placed in the methods developed.

4.5.2. *Uncertainty in computation costs.* A good measurement of computation cost requires counts of the number of operations involved (initially and per iteration), distinguishing additions, multiplications, and divisions, and indicating the number of decimals required in intermediate steps for a given decimal accuracy in the result. If such measurements were available, cost comparisons would still depend on insufficiently known relative speeds of convergence per iteration. However, even initial cost and cost per iteration cannot be measured by the counts indicated because in applications to economic equation systems so much depends on the precise form of the basic matrix Φ^* . In addition, there is still considerable uncertainty about the most economical method of inversion or solution of linear equations in cases where Φ^* is already specified.

4.5.3. *Cost comparisons between various methods.* For these reasons we shall confine ourselves to setting out in Table 4.5.3 in a comparative fashion the main features of each method affecting cost of computation. In reading this table, which requires reference to earlier formulae for detailed comparisons, it must be remembered that the inversion of a matrix of order N involves a number of operations proportional to N^3 and that the multiplication of an $N_1 \times N_2$ -matrix into an $N_2 \times N_3$ -matrix requires $N_1 N_2 N_3$ multiplications and an almost equal number of additions.

Comparison of Various Methods

TABLE 4.5.3

Method	I. Case of uncorrelated disturbances (Σ diagonal)	II. Case of unrestricted variances and covariances Σ of disturbances
<p>A. Methods ρ_h and ρ_{h_n} based on M_n^*</p>	<ol style="list-style-type: none"> $M_n^* = M^{**} = \Phi^*(I_{[X_y]} \otimes M_{xxx}) \Phi'^*$ is constant in iterations. Hence (a) α_{III}^* can be eliminated, and (b) the required extent of the inversion of M^{**}, and the transformations $B = \text{mat}_{(III \alpha^*, III \Phi^*)}$, can be carried out for all iterations at once. M^{**} partitions into diagonal blocks, reducing the number of operations in its inversion in the ratio¹ $(\Sigma III Q_g)^3$ to $\Sigma(III Q_g)^3$. Each iteration requires the inversion of B_n. 	<ol style="list-style-type: none"> $M_n^* = \Phi^*(S_n^{-1} \otimes M_{xxx}) \Phi'^*$ is not constant in iterations. Its inversion (or other processing in solving for α_{n+1}^*) must be repeated for each iteration. There are probable advantages from the regularities in M^* for its inversion or other processing, perhaps dispensing with inversion of S_n for each iteration. Each iteration requires the inversion of B_n.

¹ $III Q_g$ is defined as $Q_g^I + Q_g^{II}$.

TABLE 4.5.3

(Continued)

<p>B. Newton method based on L_n^*</p>	<ol style="list-style-type: none"> 1. L_n^* contains only two terms. 2. The first term requires for its evaluation the inversion of B_n for each iteration. This term vanishes outside $III III L_n^*$. 3. The remaining term $-M_n^*$ is constant in iterations, so that use of $III L_n^*$ permits the elimination of α_{III}^*. The required inversion of $M_{III III}^{**}$ is facilitated by partitioning into diagonal blocks. 4. The inversion of $III L_n^*$ (or other processing in solving for $III a_{n+1}^*$) must be repeated for each iteration. 5. There may be advantages in the regularities of $III L_n^*$ for its inversion or other processing. 	<ol style="list-style-type: none"> 1. $[1] L_n^*$ contains four terms. 2. Its evaluation requires the inversion of B_n and S_n for each iteration. 3. No submatrix of $[1] L_n^*$ remains constant in iterations. Its inversion or other processing must be repeated for each iteration. 4. The advantage from regularities in $[1] L_n^*$ for its inversion or other processing are highly uncertain.
---	---	---

The relative cheapness of \mathcal{P}_h and \mathcal{P}_n in the case of uncorrelated disturbances stands out clearly from this table. Not only can M^{**} be inverted (to the extent required) once for all iterations, but its partitioning into diagonal blocks greatly reduces the amount of work involved in that inversion. Precise comparisons between the remaining three methods (entries BI, A II, B II of the table) are made difficult by the uncertainties already mentioned. The general inference can be made that each transition, either from method A to method B within the same case, or from case I to case II within the same method, leads to a considerable increase in cost of computation.

4.5.4. *Methods of inversion.* If we include the computation of sampling variances and covariances of estimated parameters a^* , the inverse of each of the matrices B_n , S_n , III III M^g , III L_n^* , M_n^* , $[1]L_n^*$ is required in at least one of the methods or cases. The problems encountered in taking advantage of the regularities in the definitions (4.186), (4.225) and (4.208) - (4.209) of the last three of these matrices have already been mentioned. Here we only point to the importance of these problems for the methods discussed, and to their inherent mathematical interest. Our further remarks are directed to those inversions where such "advantages" are not present or are not real advantages in the sense that their exploitation is not worth the cost. This will often be the case in computing the relevant submatrices of $(M^g)^{-1}$, which could all be derived from one larger-order inverse M_{xx}^{-1} with the help of orthogonalized basic matrices Φ^g .

Five of the seven inversion jobs listed have to be repeated in successive iterations. This places a premium on iterative methods of inversion since, for instance, S_{n-1}^{-1} can serve as initial value for the iterative inversion of S_n . Iterative methods for inverting matrices have been discussed by Hotelling [1943-A, especially paragraphs 5, 7, 9, 10]. When such a method is applied, the approximation to S_n^{-1} must not be pushed beyond a certain level corresponding to the degree of approximation to A expected to be reached by A_{n+1} . A certain balance between the frequency of iterations in the various parts of the whole calculation should thus be preserved.

The inversion of B_n may offer special opportunities for economies because usually many of its elements are prescribed to vanish. In such cases it is advisable to permute the rows and the first K_y

columns of A (i.e., structural equations and dependent variables, respectively) so as to bring B as nearly as possible into triangular form. Partitioning of B according to (2.82) will be noticed as a by-product of such analysis. If the corresponding partitioning (2.82) of Σ is not assumed, the partitioning (2.82) of B still facilitates the inversion of B_n . In cases where permutation of rows and columns can only produce a compact block of zeros in the southwest corner of B that does not extend to the main diagonal, considerable savings will still be encountered in any of the variants of the Doolittle method applied to the inversion of B_n .

4.5.5. *Generalization of the restrictions on A.* The formulae for all methods discussed admit, without serious complications, of a generalization of the restrictions on A which has already been mentioned in earlier sections. This is the restriction (2.73 α) requiring two pairs of coefficients occurring in different structural equations to have the same unknown ratio. In combination with suitably chosen normalization rules, such a restriction can be given the linear form

$$(4.246) \quad \begin{cases} (4.246k) & \alpha_{g_1 k_1} = \alpha_{g_2 k_2} = 1, \\ (4.246l) & \alpha_{g_1 l_1} - \alpha_{g_2 l_2} = 0, \end{cases}$$

which differs from restrictions considered earlier only in that elements of different rows of A enter into the same restriction (4.246l). The treatment of restrictions of the type (4.246k) has already been demonstrated above. A restriction of the type (4.246l) can be introduced into the various iterative procedures by incorporating in Φ^* , as defined by (4.31), the row

$$(4.247) \quad \left[\begin{array}{cccccccc} 0(1) & \dots & 0(g_1 - 1) & \varphi(g_1) & 0(g_1 + 1) & & & \\ & & & & & \dots & 0(K_y) & \end{array} \right]$$

with

$$(4.248) \quad \begin{aligned} \varphi(g_1) &\equiv \left[\begin{array}{cccccc} 0_1 & \dots & 0_{l_1-1} & 1 & 0_{l_1+1} & \dots & 0_{K_x} \end{array} \right], \\ \varphi(g_2) &\equiv \left[\begin{array}{cccccc} 0_1 & \dots & 0_{l_2-1} & 1 & 0_{l_2+1} & \dots & 0_{K_x} \end{array} \right], \end{aligned}$$

where the zeros in (4.247) represent vectors of order K_x , whereas the zeros in (4.248) represent scalar components. This new row makes Φ^* different from all the Φ^* previously considered. Previously Φ^* was a matrix with vanishing elements except in the diagonal blocks occupied by $\Phi^1, \Phi^2, \dots, \Phi^G$, where Φ^i expressed restrictions on the parameters of the i th equation only. It will be clear that the number of restrictions (4.246) that can be expressed in this manner is limited by the fact that only one normalization rule can be imposed on each structural equation. The only computational complication arising from the presence of rows like (4.247) in Φ^* is that M^* partitions into fewer and larger diagonal blocks.

4.5.6. *Unsolved problems in distinguishing the highest maximum of the likelihood function.* An important class of unsolved problems, presumably requiring methods quite different from those here employed, is connected with the question of how to make sure that any maximum of the likelihood function found is the highest maximum or, if possible, of how to ensure by choice of initial values that iterations will converge to the highest maximum. Of course, the theory of the asymptotic properties of the maximum-likelihood estimates a^*, S has approximative value only if the highest maximum is well above the next highest. But how will proximity of the two highest maxima be recognized? Will it necessarily be revealed by high sampling variances of the estimates?

The condition

$$(4.249) \quad \det B = 0$$

divides the space of the elements α_{gh} , $g, h = 1, \dots, K_y$, into two connected regions. The logarithmic likelihood function

$$(4.250) \quad L(A) = \text{const} + \log \det B - \frac{1}{2} \text{tr } A M_{xx} A'$$

in the case of uncorrelated disturbances approaches $-\infty$ whenever B approaches the boundary (4.249) of the two regions. It follows that, whatever the linear restrictions on A , there are at least two maxima of the likelihood function (4.250). Under linear restrictions that are more than adequate in number and variety for complete identification of the structural equations, many more maxima can be expected to arise: there will be at least one maximum for each connected part of the restricted-parameter space cut out by the condition (4.249).

In the case where Σ is unrestricted, no difficulties arise in the subspace of the parameters Σ because the positive definiteness of Σ precludes the passing of a border analogous to (4.249), and for a given A only one maximum (4.200) with respect to Σ exists. Further discussion can therefore be based on the function (4.201) which we rewrite as

$$\begin{aligned}
 L(A) &= \text{const} + \log \det B - \frac{1}{2} \log \det A M_{xx} A' \\
 (4.251) \quad &= \text{const} - \frac{1}{2} \log \det B^{-1} A M_{xx} A' B^{-1} \\
 &= \text{const} - \frac{1}{2} \log \det \begin{bmatrix} -I_{[K_y]} & \Pi_{[K_y, K_z]} \end{bmatrix} M_{xx} \begin{bmatrix} -I_{[K_y]} \\ \Pi_{[K_y, K_z]} \end{bmatrix}.
 \end{aligned}$$

This function will still approach $-\infty$ if B approaches a point B_0 on the boundary (4.249), provided Γ does not simultaneously approach a point Γ_0 such that $\Pi_{[K_y, K_z]}$ has a finite limit. It is easily seen that, if the point A_0 approached by A is finite, $\Pi_{[K_y, K_z]}$ can remain finite only if A_0 is of rank $K_y - 1$. Points A_0 of this character form a "bridge" across the boundary (4.249) which may complicate the analysis of the number of maxima under linear restrictions on A .

These remarks may suffice to indicate a class of difficult problems, the solution of which is vital to the computation methods here developed. Pending a systematic attack on these problems, the best one can do is to accumulate "practical" experience by trying out various alternative initial values in order to learn from what range of plausible initial values the same maximum is approached.

4.5.7. *Choice of initial values A_0 .* The single-equation least-squares estimates for various choices of "dependent variables" in each equation, obtained anyhow as a by-product of the preparations for the simpler ones of the iterative processes discussed, would seem to be suitable material for such experimenting. If divergence of iterations or convergence to a different maximum for different least-squares initial values occurs frequently, or even occasionally, it will be an important problem to find initial values as near as possible to the highest maximum of the likelihood

function. In the case of unrestricted Σ , probably the best possible initial values for that purpose are obtained by the reduced-form method developed by Anderson and Rubin¹. While more costly than the least-squares estimates, the reduced-form estimates are superior in that they are consistent estimates. They are, moreover, maximum-likelihood estimates under sacrifice of an amount of a priori information that is perhaps in some sense the minimum sacrifice consistent in general with direct (i.e., noniterative) methods of computation. If so, these estimates are in a sense the nearest one can get to the highest maximum of the likelihood function by direct methods. They may, however, be less economical than least-squares estimates in cases where no doubt exists as to the identity of the highest maximum of the likelihood function. An intermediate choice is given by the maximum-likelihood estimates with diagonally restricted Σ , using all a priori information relating to A , and determined iteratively. These estimates are, of course, not consistent if Σ is actually nondiagonal.

¹See [IX].

III. NOTE ON THE IDENTIFICATION OF ECONOMIC RELATIONS

BY A. WALD

	Page
1. Definitions and Formulation of the Problem	238
2. Two Lemmas	240
3. Necessary and Sufficient Conditions for the Identification of a Coordinate θ_p of a Parameter Point θ	243

T. C. Koopmans and H. Rubin have discussed the problem of identification of economic relations in [II - 2], and have obtained a number of very interesting results. In this note the problem treated by Koopmans and Rubin is somewhat generalized and a different approach to its solution is briefly discussed.

1. Definitions and Formulation of the Problem

Let x_1, \dots, x_K be a set of K variables¹ and let $A = [\alpha_{gk}]$ ($g = 1, \dots, G; k = 1, \dots, K$) be a given matrix of rank G . Denote the linear form $\sum_{k=1}^K \alpha_{gk} x_k$ by l_g ($g = 1, \dots, G$) and let $\Sigma = [\sigma_{gh}]$ ($g, h = 1, \dots, G$) be a given symmetric and positive definite matrix. Furthermore, let

$$(1.1) \quad \varphi_r(\alpha_{11}, \alpha_{12}, \dots, \alpha_{GK}, \sigma_{11}, \sigma_{12}, \dots, \sigma_{GG}) = 0$$

($r = 1, \dots, R$)

be a given system of equations, called a priori restrictions, that are satisfied by the quantities α_{gk} and σ_{gh} . For any nonsingular matrix $Y = [v_{gh}]$ ($g, h = 1, \dots, G$) we shall denote the matrix YA by $A(Y)$ and the elements of $A(Y)$ by $\alpha_{gk}(Y)$ ($g = 1, \dots, G; k = 1, \dots, K$). Thus, $\alpha_{gk}(Y) = \sum_{h=1}^G v_{gh} \alpha_{hk}$. Furthermore, we shall

¹The integer K corresponds to what was denoted K_x in [II].

denote the matrix $\Upsilon \Sigma \Upsilon'$ (Υ' is the transpose of Υ) by $\Sigma(\Upsilon)$ and the elements of $\Sigma(\Upsilon)$ by $\sigma_{gh}(\Upsilon)$ ($g, h = 1, \dots, G$). Finally, the linear form $\sum_{h=1}^G v_{gh} l_h$ will be denoted by $l_g(\Upsilon)$.

DEFINITION 1.1. A nonsingular matrix $\Upsilon = [v_{gh}]$ ($g, h = 1, \dots, G$) is said to be an admissible transformation if and only if the equations

$$(1.2) \quad \varphi_r \{ \alpha_{11}(\Upsilon), \alpha_{12}(\Upsilon), \dots, \alpha_{GK}(\Upsilon), \\ \sigma_{11}(\Upsilon), \sigma_{12}(\Upsilon), \dots, \sigma_{GG}(\Upsilon) \} = 0 \\ (r = 1, \dots, R)$$

are fulfilled.

DEFINITION 1.2. An element α_{gk} of the matrix A is said to be identifiable¹ if $\alpha_{gk}(\Upsilon)$ takes only a finite number of different values over the domain of all admissible transformations Υ . Similarly, an element σ_{gh} of Σ is said to be identifiable if $\sigma_{gh}(\Upsilon)$ takes only a finite number of different values over the domain of all admissible transformations Υ .

DEFINITION 1.3. The linear form l_g is said to be identifiable if the coefficients $\alpha_{g1}, \dots, \alpha_{gK}$ are identifiable.

The matrix A has GK elements and the matrix Σ has $(G^2 + G)/2$ elements. Thus, the total number of elements in the two matrices A and Σ is equal to $GK + (G^2 + G)/2 = P$ (say). Consider the elements of A and Σ arranged in an ordered sequence and denote them by $\theta_1, \dots, \theta_P$, respectively. The set $\theta = (\theta_1, \dots, \theta_P)$ can be represented by a point in the P -dimensional space, called parameter space. For any nonsingular transformation Υ we shall denote the point $(\theta_1(\Upsilon), \dots, \theta_P(\Upsilon))$ by $\theta(\Upsilon)$.

DEFINITION 1.4. A coordinate θ_p of a point θ will be said to be locally identifiable if there exists an open set ω containing θ such that for any admissible transformation Υ either $\theta_p(\Upsilon) = \theta_p$ or $\theta(\Upsilon)$ lies outside ω .

¹This concept corresponds to what was called multiple identifiability in [II-2.4.4].

The problem considered in this note is to formulate conditions under which a coordinate θ_p of a point θ of the parameter space is identifiable or is locally identifiable.

2. Two Lemmas

In this section we shall prove two lemmas which will then be used for deriving necessary and sufficient conditions for the identification of θ_p .

Consider the quadratic form

$$(2.1) \quad X = \sum_{h=1}^G \sum_{g=1}^G \sigma^{gh} l_g l_h,$$

where $[\sigma^{gh}]$ is the inverse of $[\sigma_{gh}]$. Let ξ_{kl} ($k, l = 1, \dots, K$) denote the coefficient of $x_k x_l$ in X . For any nonsingular transformation Υ we shall put

$$(2.2) \quad X(\Upsilon) = \sum_{h=1}^G \sum_{g=1}^G \sigma^{gh}(\Upsilon) l_g(\Upsilon) l_h(\Upsilon),$$

where $[\sigma^{gh}(\Upsilon)]$ denotes the inverse of $[\sigma_{gh}(\Upsilon)]$. We shall denote by $\xi_{kl}(\Upsilon)$ the coefficient of $x_k x_l$ in $X(\Upsilon)$.

LEMMA 2.1. *For any nonsingular transformation Υ we have $\xi_{kl}(\Upsilon) = \xi_{kl}$ ($k, l = 1, \dots, K$).*

Proof: Denote by l the row vector $[l_1 \dots l_G]$. Using matrix notation we can write

$$(2.3) \quad X = l \Sigma^{-1} l'$$

and

$$(2.4) \quad X(\Upsilon) = l(\Upsilon) \Sigma^{-1}(\Upsilon) l'(\Upsilon),$$

where l' is the transpose of l and Σ^{-1} is the inverse of Σ . We have

$$(2.5) \quad l'(\Upsilon) = \Upsilon l',$$

$$(2.6) \quad l(\Upsilon) = l \Upsilon',$$

and

$$(2.7) \quad \Sigma^{-1}(\Upsilon) = \Upsilon'^{-1} \Sigma^{-1} \Upsilon^{-1}.$$

Hence, from (2.4) - (2.7) we obtain

$$(2.8) \quad X(\Upsilon) = l \Upsilon' \Upsilon'^{-1} \Sigma^{-1} \Upsilon^{-1} \Upsilon l = l \Sigma^{-1} l' = X,$$

and Lemma 2.1 is proved.

Let $\theta^* = (\theta_1^*, \dots, \theta_p^*)$ be a parameter point different from θ and denote by l_g^* , X^* , and ξ_{kl}^* the expressions we obtain from l_g , X , and ξ_{kl} , respectively, by substituting θ^* for θ . Now we shall prove the following lemma.

LEMMA 2.2. *If θ^* is a point such that $\xi_{kl}^* = \xi_{kl}$ ($k, l = 1, \dots, K$), then there exists a nonsingular transformation Υ such that*

$$(2.9) \quad \xi_{kl}^*(\Upsilon) = \xi_{kl}(\Upsilon) \quad (k, l = 1, \dots, K).$$

Proof: From $\xi_{kl}^* = \xi_{kl}$ it follows that $X^* = X$ identically in x_1, \dots, x_K . Thus we have

$$(2.10) \quad X^* = \sum \sum \sigma^{gh} l_g^* l_h^* = \sum \sum \sigma^{gh} l_g l_h = X.$$

Since l_1, \dots, l_G are independent linear forms and since $[\sigma_{gh}]$ is nonsingular, the rank of the quadratic form X is equal to G . Hence, the rank of X^* is also equal to G and, therefore, l_1^*, \dots, l_G^* are independent linear forms. From this and (2.10) it follows that each linear form l_g^* is a linear combination of the forms l_1, \dots, l_G . Hence there exists exactly one nonsingular transformation Υ such that

$$(2.11) \quad l_g(\Upsilon) = l_g^* \quad (g = 1, \dots, G).$$

From Lemma 2.1 it follows that

$$(2.12) \quad \sum \sum \sigma^{gh}(\Upsilon) l_g(\Upsilon) l_h(\Upsilon) = \sum \sum \sigma^{gh} l_g l_h.$$

From (2.10), (2.11), and (2.12) we obtain

$$(2.13) \quad \sum \sum \sigma^{gh}(\Upsilon) l_g^* l_h^* = \sum \sum \sigma^{*gh} l_g^* l_h^*.$$

Hence

$$(2.14) \quad \sigma^{gh}(\Upsilon) = \sigma^{*gh} \quad (g, h = 1, \dots, G),$$

and, therefore,

$$(2.15) \quad \sigma_{gh}(\Upsilon) = \sigma_{gh}^* \quad (g, h = 1, \dots, G). \quad 1$$

Since (2.11) implies that $\alpha_{gk}(\Upsilon) = \alpha_{gh}^*$, Lemma 2.2 is proved.

The coefficients ξ_{kl} ($k, l = 1, \dots, K$) depend, of course, on the parameter point θ . To make this evident, we shall occasionally replace ξ_{kl} by $\xi_{kl}(\theta)$, and $\xi_{kl}(\Upsilon)$ by $\xi_{kl}\{\theta(\Upsilon)\}$. Since $\xi_{kl}\{\theta(\Upsilon)\} = \xi_{kl}(\theta)$, we shall say that the functions $\xi_{kl}(\theta)$ are invariant under nonsingular transformations Υ .

Let $F(\theta)$ be a function of θ . We shall say that $F(\theta)$ is invariant under nonsingular transformations if for any nonsingular transformation Υ we have $F\{\theta(\Upsilon)\} = F(\theta)$. Clearly, if $F(\theta)$ is a function of $\xi_{11}(\theta)$, $\xi_{12}(\theta)$, \dots , $\xi_{KK}(\theta)$, then $F(\theta)$ is invariant under nonsingular transformations. We shall show that the converse is also true. Let $F(\theta)$ be a function such that $F\{\theta(\Upsilon)\} = F(\theta)$ for all nonsingular transformations Υ . Suppose that $F(\theta)$ is not a function of $\xi_{11}(\theta)$, $\xi_{12}(\theta)$, \dots , $\xi_{KK}(\theta)$. Then there exist two points θ'' and θ''' such that

$$(2.16) \quad \xi_{kl}(\theta'') = \xi_{kl}(\theta''') \quad (k, l = 1, \dots, K)$$

and

$$(2.17) \quad F(\theta'') \neq F(\theta''').$$

From Lemma 2.2 and (2.16) it follows that there exists a nonsingular transformation Υ such that

$$(2.18) \quad \theta''(\Upsilon) = \theta'''.$$

But then

$$(2.19) \quad F(\theta'') \neq F\{\theta''(\Upsilon)\},$$

which contradicts our assumption about $F(\theta)$. Hence $F(\theta)$ must be a

function of $\xi_{11}(\theta)$, $\xi_{12}(\theta)$, ..., $\xi_{KK}(\theta)$. Thus, the functions $\xi_{11}(\theta)$, $\xi_{12}(\theta)$, ..., $\xi_{KK}(\theta)$ form a fundamental set of invariants.

*3. Necessary and Sufficient Conditions for the
Identification of a Coordinate θ_p of a Parameter Point θ*

Let θ be a parameter point satisfying the a priori conditions (1.1). And further, let $\theta^* = (\theta_1^*, \dots, \theta_p^*)$ be an unknown parameter point and consider the equations in θ^* :

$$(3.1) \quad \xi_{kl}(\theta^*) = \xi_{kl}(\theta) \quad (k, l = 1, \dots, K)$$

and

$$(3.2) \quad \varphi_r(\theta^*) = 0 \quad (r = 1, \dots, R) \text{ (a priori conditions).}$$

The following two theorems are immediate consequences of Lemmas 2.1 and 2.2.

THEOREM 3.1. *A necessary and sufficient condition that θ_p be identifiable is that the equations (3.1) and (3.2) in the unknowns θ_1^* , ..., θ_p^* should admit of only a finite number of solutions for θ_p^* .*

THEOREM 3.2. *A necessary and sufficient condition that θ_p be locally identifiable is that there exists a finite neighborhood ω of θ such that for any solution θ^* in ω of the equations (3.1) and (3.2) we have $\theta_p^* = \theta_p$.*

In what follows in this section we shall assume that the R equations (3.2) have unique solutions in R unknowns, i.e., in R coordinates of θ^* . We may assume without loss of generality that these R coordinates are the last ones, i.e., θ_{p-R+1}^* , ..., θ_p^* . Thus, equations (3.2) can be written as

$$(3.3) \quad \theta_p^* = \psi_p(\theta_1^*, \dots, \theta_{p-R}^*) \quad (p = p - R + 1, \dots, p).$$

We shall assume that the functions ψ_p admit of continuous first-order partial derivatives. For any parameter point $\theta = (\theta_1, \dots, \theta_p)$ in the p -dimensional parameter space we shall denote by $\bar{\theta}$ the parameter point in the $(p-R)$ -dimensional space obtained from θ by omitting the last R coordinates, i.e., $\bar{\theta} = (\theta_1, \dots, \theta_{p-R})$.

Denote by $\bar{\xi}_{kl}(\bar{\theta})$ the function we obtain from $\xi_{kl}(\theta)$ by substituting $\psi_p(\bar{\theta})$ for θ_p ($p = P - R + 1, \dots, P$). Then the system of equations (3.1) and (3.2) is equivalent to the system

$$(3.4) \quad \bar{\xi}_{kl}(\bar{\theta}^*) = \bar{\xi}_{kl}(\bar{\theta}) \quad (k, l = 1, \dots, K)$$

and

$$(3.5) \quad \theta_p^* = \psi_p(\bar{\theta}^*) \quad (p = P - R + 1, \dots, P).$$

Denote the $(P - R)$ -dimensional parameter space by $\bar{\Omega}$. For any point $\bar{\theta}$ of $\bar{\Omega}$ we shall denote by $\Delta(\bar{\theta})$ the Jacobian of the functions $\bar{\xi}_{11}(\bar{\theta}), \bar{\xi}_{12}(\bar{\theta}), \dots, \bar{\xi}_{KK}(\bar{\theta})$ taken at the point $\bar{\theta}$. A point $\bar{\theta}$ of $\bar{\Omega}$ will be called regular if the following condition is satisfied: Any minor of the Jacobian of the $K^2 + P - R$ functions $\bar{\xi}_{11}(\bar{\theta}), \bar{\xi}_{12}(\bar{\theta}), \dots, \bar{\xi}_{KK}(\bar{\theta}), d_p(\bar{\theta}) = \theta_p$ ($p = 1, \dots, P - R$) is either unequal to zero at $\bar{\theta}$ or is *identically* zero in some finite neighborhood of $\bar{\theta}$.

THEOREM 3.3. *Let $\bar{\theta}^0$ be a regular point and denote by $\delta_p(\bar{\theta})$ the Jacobian of the $K^2 + 1$ functions $\bar{\xi}_{11}(\bar{\theta}), \bar{\xi}_{12}(\bar{\theta}), \dots, \bar{\xi}_{KK}(\bar{\theta}), d_p(\bar{\theta}) = \theta_p$ for any value of p satisfying $p \leq P - R$. A necessary and sufficient condition that θ_p be locally identifiable for any point $\bar{\theta}$ in a finite neighborhood $\bar{\theta}^0$ is that the rank of $\Delta(\bar{\theta}^0)$ be equal to the rank of $\delta_p(\bar{\theta}^0)$.*

Proof: Since $\bar{\theta}^0$ is a regular point, a necessary and sufficient condition that $\bar{\theta}_p$ be a single valued function of $\bar{\xi}_{11}(\bar{\theta}), \bar{\xi}_{12}(\bar{\theta}), \dots, \bar{\xi}_{KK}(\bar{\theta})$ in a finite neighborhood of $\bar{\theta}^0$ is that the rank of $\Delta(\bar{\theta}^0)$ be equal to that of $\delta_p(\bar{\theta}^0)$. Theorem 3.3 follows from this and the fact that the functions $\bar{\xi}_{11}(\bar{\theta}), \bar{\xi}_{12}(\bar{\theta}), \dots, \bar{\xi}_{KK}(\bar{\theta})$ form a fundamental set of invariants.

IV. GENERALIZATION OF THE CONCEPT OF IDENTIFICATION¹

BY LEONID HURWICZ

	Page
1. Introduction	245
2. Structures and Models	246
3. Structural Estimation and Identification	248
4. Fundamental Definition of Identification Power of a Model	249
5. Models Identifying over Certain Submodels	249
6. Identification Power with Regard to a Criterion	253
7. Summary	255
8. General Remarks	257

1. Introduction

1.0. A model² may or may not possess the property of being *structure-identifying* (or, for briefness, *identifying*). When the model is *not* structure-identifying, no amount of empirical information will make it possible to determine the *structure*³ of the system under investigation. Thus it is desirable to deal with identifying models when the knowledge of the structure is needed; this need frequently arises when forecasts for policy purposes are to be made.

There are several kinds and "degrees" of identification power.⁴ *Totally unique identification power* ("totally" is usually omitted) implies that if an infinite sample of observations were available at most one structure consistent with the model can be made to

¹Part of the work on this paper was done in 1945-46 during the author's tenure of the Guggenheim Memorial Fellowship. Thanks are due to Tjalling Koopmans and Herman Rubin for valuable suggestions.

²Say a particular theory of the "business cycle."

³I.e., the set of numerical values of the parameters characterizing the model; e.g., the numerical values of the marginal propensity to consume or of the effects that profits may have on investment, etc.

⁴An identifying model is said to possess identification power.

"fit" a given set of observational data; in a finite sample only one "best (in some given sense) fit" (estimate) consistent with the model is possible. When more than one structure is compatible with the observed data, it will, in general, so happen that all structures compatible with a given set of data will have a certain property in common.¹ Then the model is said to be *uniquely identifying with regard to that property*. In such a case we speak of *partially* unique identification power of the model, unless the model is uniquely identified with regard to *all* its properties so that *totally* unique identification power, as already defined, is present.

A model lacking unique identifying power with regard to a given property may nevertheless possess *complete* identifying power with regard to that property if the set of structures compatible with any given data is finite or, at worst, denumerably infinite.

In what follows these concepts are defined in a more rigorous manner; a partial summary is given in section 7. The definitions given differ from those found elsewhere² in that they are independent of the nature of the models and the structures involved. In particular, linearity of the "structural relations" and the normality of the disturbance distribution are not assumed in the present paper.³ It is not necessary that either the structural relations or the disturbance distributions should have a parametric form. In fact, the definitions given would apply to situations where structures are defined in terms other than those of "structural relations" and disturbance distributions.

One may hope that these generalizations will be of some help in clarifying the logic of the identification problems; they are not a substitute for the study of conditions under which specified models of practical importance possess one of the various types of identification power.

2. Structures and Models

2.1.1. The cumulative probability distribution G of the ob-

¹E.g., just the consumption functions may be uniquely determined, but not the investment function; it might even conceivably happen that the effect of income on consumption should be identified, but not the effect on consumption of the assets held.

²E.g., [II-2].

³Illustrative examples have, however, been provided which are based on normal linear, hence parametric, systems. This may facilitate the task of relating theory and results already in existence to the generalizations of this note.

served variate (vector) x may be regarded as produced by an operation (usually a transformation) \mathfrak{J} performed on the distribution F of the (nonobservable) disturbance (vector) u , i.e.,

$$(1) \quad G = \mathfrak{J} \rightarrow F.$$

2.1.2. As an illustration, consider the usually treated parametric case where \mathfrak{J} is a linear transformation and F (hence also G) is normal. Adopting two simplifying assumptions¹, viz., that the mean of u is zero and that the predetermined variables are absent, one may in this case represent F and G by their respective covariance matrices Σ^u and Σ^x while \mathfrak{J} is represented by the inverse A^{-1} of the structural coefficient matrix² A . We then have, as a counterpart of (1), the equation

$$(1') \quad \Sigma^x = A^{-1} \Sigma^u (A^{-1})',$$

where Σ^x is produced by performing the operation $A^{-1}\{\dots\}(A^{-1})'$ on Σ^u . [The operation consists in premultiplying by A^{-1} and postmultiplying by $(A^{-1})'$.]

2.2. It will be noted that (1) does not require a parametric representation for F , G , or \mathfrak{J} . It is sufficient that F and G should be functions of u and x , respectively, and that these functions should belong to the class of distribution functions.

\mathfrak{J} may be any well-defined operation carrying one distribution function into another.

u and x need not be of the same dimensionality³.

2.3. We shall refer to

$$(2) \quad S = (F, \mathfrak{J})$$

as a structure. An a priori postulated class \mathfrak{S}_0 of structures S that is

¹The theory developed in this note is not subject to these restrictions. In particular, it is valid for systems containing fixed variates as well lagged endogenous variables.

²In terms of operations on variables (rather than on distributions) we have $Ax = u$, or, if the transformation need not be linear, $\Phi(x) = u$. In the latter case \mathfrak{J} is represented by Φ^{-1} . This is the symbol used in [1].

³When the dimensionality of u exceeds that of x we are, except for trivial cases, dealing with *nonadditive disturbances*. Cf. [XVIII].

a proper subset of the class \mathcal{S} of all structures is called a model¹.

When (1) holds, the structure S , as defined in (2), is said to *generate* the distribution G (written $S :. G$); G is then said to be *generated by* S (written $G :. S$).

3. Structural Estimation and Identification

3.1. The process of *structural estimation* consists in estimating G and then obtaining S generating that particular G . (The two stages may be combined into one computational process.)

In general, more than one structure generating a given G can be found. Hence structural estimation would lead to a class of structures, say \mathcal{S}_0 , containing all structures S that could have generated the given G . However, when certain additional assumptions are made with regard to S - these are the "identifying restrictions" - the class \mathcal{S}_0 may be narrowed down to a proper subset \mathcal{S}_{00} of \mathcal{S}_0 . (If the identifying restrictions exclude a priori all structures not in \mathcal{S}_1 , say, we have $\mathcal{S}_{00} = \mathcal{S}_0 \mathcal{S}_1$.) When \mathcal{S}_{00} is denumerably infinite or contains a finite number N of elements, we speak of *complete identification power of the model*: *unique identification power* if $N = 1$, *multiple* if $N > 1$; when \mathcal{S}_{00} is nondenumerable, we speak of *incomplete identification power*.

3.2.1. The nature of the additional (identifying) assumptions mentioned in section 3.1 will now be stated and a more explicit definition of identification power provided.

3.2.2. The identifying assumptions postulate that the structure S which has generated a given G belongs to a certain class \mathcal{S}_1 of structures².

The elements of \mathcal{S}_1 may be distinguished from those of $\text{comp } \mathcal{S}_1 \equiv \mathcal{S} - \mathcal{S}_1$ (where \mathcal{S} is, as before, the class of all structures) through restrictions on F , \mathfrak{J} , both F and \mathfrak{J} , or, more generally, through a relation to be satisfied by F and \mathfrak{J} . In the parametric

¹Thus a structure is defined by one completely specified distribution F and one completely specified operation \mathfrak{J} . In the parametric case treated above, section 2.1.2, the structure is given if both matrices Σ^u and A are numerically given. On the other hand, postulating, say, zeros in certain parts of either matrix without restricting other elements of these matrices defines a whole class of structures, i.e., a model.

²Thus, by a previous definition, \mathcal{S}_1 is a model.

case used for illustrative purposes (cf. section 2.1.2) restrictions are imposed on A , Σ^u , or both.

4. Fundamental Definition of Identification Power of a Model

We shall now define the class \mathcal{G}_1 of all G 's generated by the elements of \mathcal{G}_1 (the class of all structures not excluded by the identifying restrictions as defined in the preceding section), so that symbolically

$$(3.1) \quad G_1 \in \mathcal{G}_1 : \supset: \exists S_1 \in \mathcal{G}_1 \ni S_1 :. G_1$$

(for every element G_1 of \mathcal{G}_1 there exists a structure S_1 in \mathcal{G}_1 such that S_1 generates G_1) and

$$(3.2) \quad S_1 \in \mathcal{G}_1 : \supset: \exists! G_1 \in \mathcal{G}_1 \ni S_1 :. G_1$$

(for every element S_1 of \mathcal{G}_1 there exists *one and only one* G_1 in \mathcal{G}_1 such that G_1 is generated by S_1). It may so happen that \mathcal{G}_1 has the property that for any element G_1 in \mathcal{G}_1 there exists *one and only one* S_1 in \mathcal{G}_1 such that S_1 generates G_1 , i.e.,

$$(4) \quad G_1 \in \mathcal{G}_1 : \supset: \exists! S_1 \in \mathcal{G}_1 \ni S_1 :. G_1.$$

When (4) holds we say that the model \mathcal{G}_1 is *uniquely identifying*.¹

If for any G_1 in \mathcal{G}_1 the set \mathcal{G}_{1,G_1} of all S_1 in \mathcal{G}_1 which generate G_1 is finite or denumerably infinite, the model \mathcal{G}_1 is said to be *completely identifying*. When \mathcal{G}_{1,G_1} has a finite number of elements for all G_1 in \mathcal{G}_1 , \mathcal{G}_1 is said to be *multiply or uniquely identifying* depending on whether or not $N > 1$ for some G_1 . (This definition of unique identification power is, of course, equivalent to that given earlier.)

5. Models Identifying over Certain Submodels

5.1. The above definition will now be generalized so as to take care of a situation frequently arising in applications. Sup-

¹It will be noted that this definition of identification could have been stated in identical terms for definitions of S and of the relation $S :. G$ entirely different from those given in section 2.3.

pose one is willing to postulate that the unknown structure is an element of the model \mathfrak{G}_1 . We write

$$(5) \quad \mathfrak{G}_1 = \mathfrak{G}_{11} + \mathfrak{G}_{12} ,$$

where \mathfrak{G}_{11} is defined as the set of *all*¹ elements of \mathfrak{G}_1 that possess the following property: if S_{11} is an element of \mathfrak{G}_{11} and if S_{11} generates the distribution G_{11} in \mathfrak{G}_1 (where \mathfrak{G}_1 is defined as before), then there exists no other element S_1 of \mathfrak{G}_1 that also generates G_{11} . It follows that if S_{12} is an element of \mathfrak{G}_{12} and if S_{12} generates G_{12} , there exists at least one other element S_{12} of \mathfrak{G}_{12} that also generates G_{12} ; however, G_{12} cannot be generated by any elements of \mathfrak{G}_{11} . In such a case we shall say that the (sub-)model \mathfrak{G}_{11} is *uniquely identifying in the model \mathfrak{G}_1* . Alternatively, it will be said that \mathfrak{G}_1 is *uniquely identifying over \mathfrak{G}_{11}* . The definition of identification power given in section 4 is consistent with the one just given. In section 4, \mathfrak{G}_1 was identifying over (or in) itself, so that in the present notation $\mathfrak{G}_1 = \mathfrak{G}_{11}$ and \mathfrak{G}_{12} is empty.

5.2. One may find it helpful to use a graph depicting the situation described in the previous paragraph. In Figure 1 the lower axis corresponds to the set \mathfrak{G} of all G 's, while the upper one corresponds to the set \mathfrak{S} of all S 's. A point on the lower axis corresponds to a given G , a point on the upper axis to a given S . A

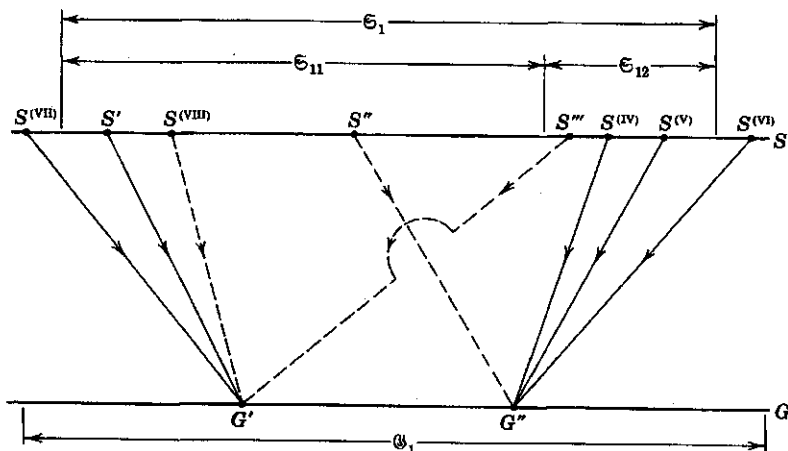


Figure 1.

¹It is important to note the word "all" in this definition.

solid line (with an arrow pointing downward) connects a structure S with the distribution G which that particular structure generates. Broken lines (with arrows) indicate "prohibited" generating relationships in the case where \mathcal{G}_1 is uniquely identifying over \mathcal{G}_{11} .

Thus it can be seen that, when \mathcal{G}_1 is uniquely identifying over \mathcal{G}_{11} , G' which is assumed to be generated by S' cannot be generated by S''' or $S^{(viii)}$, although it can be generated by $S^{(vii)}$. On the other hand, G'' can be generated by both $S^{(iv)}$ and $S^{(v)}$ as well as $S^{(vi)}$, although it cannot be generated by S'' .

Similar diagrams could be constructed for other types of identification.

5.3. As an illustration of the preceding concepts, consider the parametric case $Ax = u$, where $x = (x_1 \ x_2 \ x_3 \ x_4)$, $u = (u_1 \ u_2 \ u_3 \ u_4)$, and $A = [\alpha_{ij}]$, $i, j = 1, 2, 3, 4$. Now let \mathcal{G}_1 be defined by the following zeros in Σ^u :

$$\begin{bmatrix} \times & \times & 0 & 0 \\ \times & \times & 0 & 0 \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \end{bmatrix}$$

and by the following zeros in A :

$$\begin{bmatrix} 1 & \times & \times & 0 \\ \times & 1 & 0 & \times \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

where \times indicates an unrestricted element and 1 indicates a nonvanishing element which has been equated to unity by the normalization rule chosen. Then \mathcal{G}_{12} is defined by zeros for α_{13} or α_{24} or both, so that, for S in \mathcal{G}_{12} , A is one of the following:

$$\begin{bmatrix} 1 & \times & 0 & 0 \\ \times & 1 & 0 & \times \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \begin{bmatrix} 1 & \times & \times & 0 \\ \times & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \begin{bmatrix} 1 & \times & 0 & 0 \\ \times & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

The requirement that the structure should belong to \mathfrak{G}_{11} is in this case the "rank condition" treated in [II-2.2.2 and 2.2.3].

5.4.1. \mathfrak{G}_{11} in section 5.1 is said to be the *maximal set* of \mathfrak{G}_1 of multiplicity 1.¹

\mathfrak{G}_1 will be said to be *multiply identifying* (of multiplicity N) over \mathfrak{G}_{11} if 1) for every G_{11} in \mathfrak{G}_{11} , the subset $\mathfrak{G}_{1, G_{11}}$ of all structures S_1 in \mathfrak{G}_{11} which generate G_{11} contains a finite number N of elements (possibly less than N for some, but not all, G_{11}), and 2) \mathfrak{G}_{11} contains all structures generating G 's with property 1). \mathfrak{G}_{11} is then called the maximal set of \mathfrak{G}_1 of multiplicity N .

Similarly, \mathfrak{G}_1 is said to be *denumerably identifying* when the above definition holds with denumerable infinity substituted for N . Thus the maximal set whose multiplicity is denumerable infinity is defined as in the case of finite multiplicity.²

Now let $\{\mathfrak{G}_1^{(M_N)}\}$, $N = 2, 3, \dots$, be the (finite or denumerably infinite) sequence of the maximal sets of \mathfrak{G}_1 of multiplicities N , and write $\mathfrak{G}_1^{(1)}$, $\mathfrak{G}_1^{(D)}$ for the maximal sets of \mathfrak{G}_1 of multiplicities 1 and denumerable infinity, respectively. Then \mathfrak{G}_1 is said to be *multiply identifying* over $\mathfrak{G}_1^{(M)} = \sum_{1 < N < \infty} \mathfrak{G}_1^{(M_N)}$, *finitely identifying* over $\mathfrak{G}_1^{(F)} = \mathfrak{G}_1^{(1)} + \mathfrak{G}_1^{(M)}$, *completely identifying* over $\mathfrak{G}_1^{(C)} = \mathfrak{G}_1^{(F)} + \mathfrak{G}_1^{(D)}$, and *incompletely identifying* over $\mathfrak{G}_1^{(I)} = \mathfrak{G}_1 - \mathfrak{G}_1^{(C)}$. (Nondenumerable infinity is defined as the multiplicity in the case of incomplete identification power.)

5.4.2. It is convenient to have a term for the property of being a subset of a given maximal set.

Let \mathfrak{G}'_1 be any subset of $\mathfrak{G}_1^{(1)}$. We then say that \mathfrak{G}'_1 is *uniquely identifying over and possibly beyond* \mathfrak{G}'_1 . (If it is known that \mathfrak{G}'_1 is a proper subset of $\mathfrak{G}_1^{(1)}$, we omit "possibly." If it is known that $\mathfrak{G}'_1 = \mathfrak{G}_1^{(1)}$, we omit "and possibly beyond.") Clearly, if \mathfrak{G}_1 is

¹There exist smaller sets of structures with no two elements of the set generating the same G_1 , but \mathfrak{G}_{11} contains all such sets; hence the term *maximal*.

²It is important to note that the choice of \mathfrak{G}_1 determines all the maximal sets. Some of the maximal sets may, of course, happen to be empty.

uniquely identifying over and possibly beyond \mathcal{G}'_1 , it is so identifying over any subset of \mathcal{G}'_1 . Similar language may be used for multiple identification power and identification power of other multiplicities.

This terminology helps formulate the correspondence between the language of this note and that of [II]. When [II] says that a given set of equations is uniquely identified over some subset Ω_{11} of the parameter space Ω , this can be expressed in the language of the present note as saying that the model is uniquely identifying *over and possibly beyond* the set \mathcal{G}_{11} corresponding to Ω_{11} .

The language of [II] does not specify whether Ω_{11} is a maximal set.

6. Identification Power with Regard to a Criterion

6.1. Since often only some of the properties of the structure are of interest, it is desirable to broaden the concept of identification so as to cover situations where some properties of S can be determined uniquely from the knowledge of G while other properties of S perhaps cannot.

6.2. Consider a criterion C which establishes a partition of the class \mathcal{G} of all structures into a system of nonoverlapping subclasses $\mathcal{G}^{(i)}$ ($\sum_i \mathcal{G}^{(i)} = \mathcal{G}$); this partition, in general, need not be finite or denumerable. By definition, if both S_1 and S_2 belong to $\mathcal{G}^{(i)}$ they are indistinguishable with regard to C .

Let there be given a model \mathcal{G}_1 and its not necessarily proper subset $\mathcal{G}_{11} \subseteq \mathcal{G}_1$. Let \mathcal{G}_1 be the class of all G 's generated by the elements of \mathcal{G}_1 and let G_1 in \mathcal{G}_1 be generated by S_{11} in \mathcal{G}_{11} . Consider the set \mathcal{G}_{1, G_1} of all structures that generate G_1 . Then if, for every G_1 in \mathcal{G}_1 , all the elements of \mathcal{G}_{1, G_1} belong to the same subclass, say $\mathcal{G}^{(i_0)}$, it is said that \mathcal{G}_{11} is *uniquely identifying* in \mathcal{G}_1 with regard to C .

6.3. It may again be helpful to present the matter diagrammatically. In Figure 2 solid generating lines are the permissible ones while the broken ones are prohibited in the case of unique (partial) identification with regard to the criterion on which the partition is based.

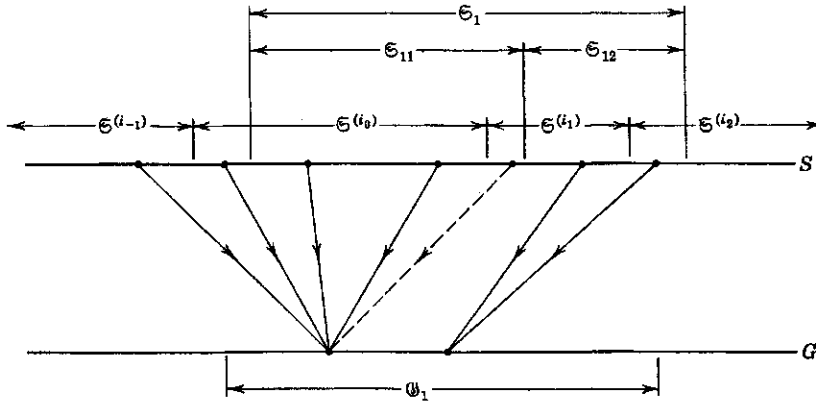


Figure 2.

6.4. As an illustrative example we use the one given in section 5.3 except that for S in \mathfrak{E}_1 we now have A with the following zeros:

$$\begin{bmatrix} 1 & \times & 0 & 0 \\ \times & 1 & 0 & \times \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

\mathfrak{E}_{12} is here given by a matrix A with $\alpha_{24} = 0$, i.e.,

$$\begin{bmatrix} 1 & \times & 0 & 0 \\ \times & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix};$$

here as criterion C we have the values of the elements of the first row of A . This is the familiar case of only one of the equations of the system being identified, to use the language of [II].

With other identifying restrictions, the value of a single element of A could have been chosen as C , cf. [III].

6.5. Multiplicities other than 1 of identification power with regard to a given criterion are defined by analogy with earlier sections of this note.

Now let \mathcal{G}_1 possess identification power of a given multiplicity over \mathcal{G}_{11} with regard to *all* criteria \mathcal{C} . Then \mathcal{G}_1 is said to possess *total* identification power over \mathcal{G}_{11} of the appropriate multiplicity. ("Total" may be omitted where no danger of confusion exists.)

On the other hand, let \mathcal{G}_1 possess identification power over \mathcal{G}_{11} of varying multiplicities with regard to different criteria. Then we say that \mathcal{G}_1 has *total* identification power over \mathcal{G}_{11} of the *highest* multiplicity and *partial* of all other multiplicities.

Thus, for instance, if we say that \mathcal{G}_1 is *partially* denumerably identifying over \mathcal{G}_{11} , the following is implied: over \mathcal{G}_{11} , \mathcal{G}_1 is denumerably identifying with regard to some criteria and *incompletely* with regard to some others.¹ (There may or may not be some criterion with regard to which \mathcal{G}_1 is uniquely or multiply identifying over \mathcal{G}_{11} .)

As in section 5.4.1, \mathcal{G}_1 can be divided into the maximal sets $\mathcal{G}_1^{(1)}$, $\mathcal{G}_1^{(M)}$, ..., with regard to a given criterion \mathcal{C} . This process is too obvious to require detailed description.

7. Summary

At this stage it may be desirable to provide a verbal summary of the more important concepts introduced.

A model is said to be *uniquely identifying with respect to a given criterion* if all structures compatible with the probability distribution of the observations are indistinguishable on the basis of that criterion.

A model is said to be *multiply identifying with respect to a given criterion* if all structures compatible with the probability distribution of the observations can be grouped into a *finite* number of classes such that within each class structures are indistinguishable on the basis of that criterion. The number N of such groups is the *multiplicity* of identification power of the model.

A model is said to be *denumerably identifying with respect to a given criterion* if all structures can be grouped into a *denumer-*

¹If it were completely identifying with regard to all criteria, and denumerably with regard to some, it would be (totally) denumerably identifying.

ably infinite number of such classes.

A model uniquely or multiply identifying with respect to a given criterion is said to be *finitely identifying with respect to that criterion*.

A model *finitely* or *denumerably* identifying with respect to a given criterion is said to be *completely identifying with respect to that criterion*.

A model *not* completely identifying with regard to a given criterion is said to be *incompletely identifying with regard to that criterion*.

The multiplicity of identification power of a model with regard to a given criterion is, in general, defined as 1, N , denumerable infinity, or nondenumerable infinity when the identification power is unique, multiple (of multiplicity N as defined above), denumerable, or incomplete, respectively. Multiplicity N is *higher* than 1, etc.¹

If a model possesses identification power of a given multiplicity with respect to *all* criteria, it is said to possess *total* identification power of that multiplicity. Thus, if a model is uniquely identifying with respect to all criteria, it is said to be *totally-unique* identifying, or - more simply - *uniquely* identifying.

When a model is uniquely identifying, there exists only one structure compatible with the probability distribution of the observations.

If a model possesses identification power of different multiplicities with respect to different criteria, it is said to possess *total* identification power of the *highest* multiplicity and *partial* identification power of all other multiplicities.

Thus, if a model is uniquely identifying with respect to one set of properties, multiply with respect to another, and denumerably with respect to a third set of properties, then the model is said to be *totally denumerably* identifying as well as *partially multiply* and *partially uniquely* identifying.²

¹Thus the higher the multiplicity, the more difficult it is to determine the structure from the probability distribution of the observations.

²When the model is said to be *partially* uniquely identifying with regard to a given criterion, "partially" serves as a warning that there exist other criteria with regard to which the model possesses identification power of a higher multiplicity, i.e., multiple, denumerable, or even incomplete, identification power.

8. General Remarks

8.1. In practice we are not often certain of the proper choice of identifying assumptions. In fact, if the choice of \mathfrak{G}_1 is made before G is known, it may happen that the structures S that generate G are not among the elements of \mathfrak{G}_1 . In this case we obtain no S whatever.

Also, a given G will, in general, yield different structures, say S'_1, S''_1, \dots , corresponding to different models $\mathfrak{G}'_1, \mathfrak{G}''_1, \dots$, even if each of them is uniquely identifying. If no a priori information permits us to choose one of these models, or, equivalently, when the model is not completely identifying, general principles of making decisions (e.g., minimization of maximum risk) may be applied. However, a solution will not always exist¹.

8.2. Identification power, as here defined, is a property of a model (with reference to specified criteria). This in itself makes it clear why the problem of identification does not arise in prediction under unchanged structure². Such prediction involves only certain properties of G , and the structure does not even enter the picture. On the other hand, for prediction under changed structure it is, in general, desirable to know S ; to do this requires, in general, an identifying model.

8.3. Since the concept of identification refers to structure and to the distributions thus generated, it should be clear that the problem of identification is entirely independent of any sampling aspects of structural estimation. The latter can only refer to the relationship of the estimate of G obtained from the sample to the true value of G . The former refers to the possibility of ascribing a unique S to a given G .

¹Cf. [Hurwicz, 1946].

²Cf. [VI].

V. REMARKS ON FRISCH'S CONFLUENCE ANALYSIS AND ITS USE IN ECONOMETRICS

BY TRYGVE HAAVELMO

	Page
1. Confluence Analysis and the Markoff Theorem on Least Squares . . .	258
2. Confluence Analysis and Econometrics	261

Confluence Analysis [Frisch, 1934] was written in part as a protest against a mechanical and uncritical use of the classical least-squares method to estimate demand functions and other economic relations. Frisch pointed out that economic data, in general, do not satisfy the conditions required to justify the use of the classical method of least-squares. It is of course true that the strict conditions of a theoretical model perhaps never are exactly fulfilled in any observational material. And good theoretical models should be able to absorb moderate discrepancies between model and facts without the inference drawn becoming valueless or nonsensical. In many respects the classical model of the least-squares method fulfills this requirement. One important exception, however, is the case where the "independent variables" are themselves highly intercorrelated while, at the same time, they are subject to "errors" which tend to hide such intercorrelation. Then the application of the classical method of least-squares might give unreliable or even nonsensical results. I shall try to explain this a little more explicitly.

1. *Confluence Analysis and the Markoff Theorem on Least-Squares*

Let x'_{it} , $i = 1, 2, \dots, n$, $t = 1, 2, \dots, T$, denote a system of values of nT variables, fixed in repeated samples, and interconnected by the linear relation

$$(1.1) \quad x'_{1t} = \alpha_2 x'_{2t} + \alpha_3 x'_{3t} + \dots + \alpha_n x'_{nt},$$

$t = 1, 2, \dots, T.$

(The case of a separate constant term is covered by setting, e.g., x'_{nt} equal to 1, identically, for all values of t .) We make the

following assumptions [David and Neyman, 1938]:

ASSUMPTION I. The $(n-1)T$ quantities $x'_{2t}, x'_{3t}, \dots, x'_{nt}$, $t = 1, 2, \dots, T$, are known observations while, instead of the values of x'_{1t} , we can only observe the values of x_{1t} defined by

$$(1.2) \quad x_{1t} = x'_{1t} + x''_{1t}, \quad t = 1, 2, \dots, T,$$

where x''_{1t} for each value of t is a random variable that cannot be observed separately..

ASSUMPTION II. The T random variables x''_{1t} , $t = 1, 2, \dots, T$, are independent random variables with $\mathcal{E}(x''_{1t}) = 0$ and $\mathcal{E}(x''_{1t}{}^2) = \sigma^2$ for all values of t .

ASSUMPTION III. The $(n-1)$ -rowed and T -columned matrix $[x'_{it}]$, $i = 2, 3, \dots, n$, $t = 1, 2, \dots, T$, is of rank $(n-1)$.

If these assumptions are fulfilled, and if $x_{11}, x_{12}, \dots, x_{1T}$, is a sample of the T variables x_{1t} , $t = 1, 2, \dots, T$, Markoff's theorem on least-squares states that of all the unbiased estimates of α_i that are linear in the variables x_{1t} the one that has the smallest variance is given by the value of a_i which minimizes the

$$\text{sum} \sum_{t=1}^{t=T} (x_{1t} - a_2 x'_{2t} - a_3 x'_{3t} - \dots - a_n x'_{nt})^2. \text{ Further,}$$

the variance of this estimate is given by

$$(1.3) \quad \text{variance of } a_i = \sigma^2 \frac{M'_{ii}}{M'}, \quad i = 2, 3, \dots, n,$$

where M' is the $(n-1)$ th order determinant of the matrix

$$\left[\sum_{t=1}^{t=T} x'_{kt} x'_{jt} \right], \quad k, j = 2, 3, \dots, n, \text{ while } M'_{ii} \text{ is the } (n-2)\text{th}$$

order principal minor obtained from this determinant by omitting the i th row and column.

If the conditions of the Markoff theorem are met by the data to which the theorem is applied, the accuracy of the estimation of the α 's is shown by their variances (1.3). For example, even if Assumption III is very near to being violated, the least-squares method will still give unbiased estimates of the α 's, but their

variances might then become very large. The estimate, s^2 , say, of σ^2 , and hence that of the variance of a_i , is not impaired by the fact that the independent variables are highly intercorrelated. Therefore, if the fact that Assumption III is near to being violated leads to very uncertain estimates of the α 's, this will be shown by large values of the estimated variances of the a 's. It is only when M' is exactly equal to zero that the method breaks down.

Frisch pointed out, however, that in practice, and in particular when dealing with economic statistics, it is hardly ever correct to assume that the expected values of one variable are linearly related to the *observable* values of the other variables. It is more realistic to assume that every one of the variables contains some random elements having no functional or stochastic connection with the other variables. Instead of (1.2) above, Frisch therefore assumed that all we can observe are variables x_{it} defined by

$$(1.4) \quad x_{it} = x'_{it} + x''_{it},$$

$$i = 1, 2, \dots, n, \quad t = 1, 2, \dots, T,$$

where all the variables x''_{it} are assumed to be uncorrelated and to have stochastic properties similar to those of the variables x''_{1t} in (1.2). As Frisch did not use a probabilistic approach, his "language" in describing the model differs considerably from that employed here. The stochastic interpretation of Frisch's model as outlined above is due to Koopmans [1937].

Suppose now that, by analogy, we attempt to estimate α_i in (1.1) by the value of a_i^* that satisfies

$$(1.5) \quad \sum_{t=1}^{t=T} (x_{1t} - a_2^* x_{2t} - a_3^* x_{3t} - \dots - a_n^* x_{nt})^2 = \text{minimum.}$$

Clearly, the Markoff theorem cannot be called upon as a justification for this procedure as the conditions of the theorem are no longer met. And the resulting "estimates" of the α 's obtained from (1.5) will have the following properties:

1. No matter whether Assumption III is fulfilled or not, and no matter how many observations we have, the statistics a_i^* will no longer be unbiased estimates of the α_i 's. The bias will depend on

the variances of the disturbances x''_{it} , $i = 2, 3, \dots, n$. Even if Assumption III is not fulfilled and, therefore, the α 's could not possibly be determined uniquely, the procedure (1.5) will yield *some* apparently unique values of the quantities a_i^* , depending mainly on the variances of the disturbances x''_{it} .

2. If, by analogy, the variance of a_i^* is taken as

$$(1.6) \quad \text{variance of } a_i^* = \sigma^2 \frac{M_{ii}}{M}, \quad i = 2, 3, \dots, n,$$

where the determinants M are obtained from the determinants M' in (1.3) by replacing x' by x , this "variance" will no longer control the accuracy with which the α 's are being estimated by (1.5), because (1.6) is simply not the variance of an unbiased estimate of α_i .

Frisch was mainly interested in deriving a method of control which would reveal the particularly nonsensical results that this procedure might lead to if Assumption III of the Markoff theorem is not fulfilled. The result was his now well-known "bunch-map" method, which is a graphical method of exposing the manner in which the T sample points $(x_{1t}, x_{2t}, \dots, x_{nt})$, $t = 1, 2, \dots, T$, are clustered in the n -dimensional scatter diagram of the n observable variables. The most important objective of the "bunch analysis," therefore, is to reveal a possible lack of fulfillment of Assumption III of the Markoff theorem when this fact - contrary to the assumptions of the theorem - is hidden by errors in all the variables.

The method of "bunch analysis" has been much discussed, and, because of its partly subjective nature, it is no doubt open to criticism. We shall not add to this discussion here as we shall be interested in another and perhaps more fundamental question, namely the question of whether or not the *statistical model* discussed above actually represents a workable model for the type of relationships we may expect between economic variables.

2. Confluence Analysis and Econometrics

If the set of "systematic parts" x'_{it} defined above satisfies one or more independent linear relations, *in addition* to (1.1), the set $(x'_{1t}, x'_{2t}, \dots, x'_{nt})$, in Frisch's terminology, is said to

be multicollinear. The occurrence of multicollinearity might sometimes be due to pure accident. An alternative is that the data under observation are ruled, simultaneously, by two or more linear relations that are permanent and inherent in the economic structure that produces the data observed.

If, however, Frisch's statistical model (or, as a special case, that of Markoff) is accepted, I think it will be very difficult to construct any realistic economic theory that would imply such multicollinearity in the strict sense. In fact, I believe that the model and the method based upon it are tenable only in connection with a single-equation approach where other relations, if any, between the variables are accidental, or "given from outside." I shall try to make this clear.

Let $x_{1t}, x_{2t}, \dots, x_{nt}$ be n observable series of economic variables. Usually, no exact functional relationship holds between such observable variables. It is much more likely that they would be interconnected through some sort of stochastic relationship. And the stochastic elements involved would usually be not only "errors of measurements" but random variables of a more fundamental nature that are characteristic elements of economic actions and decisions. If, therefore, the disturbances x''_{it} in (1.4) are thought of only as errors of measurement, the "true" variables x'_{it} would hardly ever satisfy an exact relation such as (1.1); nor could there be any case of multicollinearity in the strict sense, except by accident. In other words, if the "true" variables are thought of as economic variables which we might be able to observe if we had better statistics, then any assumption of exact linear relations between these variables is unrealistic.¹

If, nevertheless, we would require such relations to be exact, one or more of the variables x'_{it} would have to be considered as theoretical constructions for the purpose of building a model. The remainder, x''_{it} , would then no longer consist only of errors of measurements. Moreover, the behavior and, indeed, the economic meaning of these theoretical variables would then *depend radically on the whole network of stochastic economic relations between the variables studied* rather than on the particular relation (1.1) under investigation. Now, the real reason for splitting the variables into "systematic parts" and "disturbances" would seem to be, first, the idea that the systematic parts of the variables observed are the would-be "true economic variables" which according to eco-

¹For a more extensive discussion of this subject reference is made to [Haavelmo, 1943, 1944].

conomic theory would fit exactly the theoretical relation considered, second, that this relation would be autonomous in the sense that it would hold regardless of whether or not other economic relations were fulfilled, and, third, that there might be some hope of eventually being able to experience and observe the "true" variables so that the exact relation between them would have value for prediction purposes. But a set of "systematic parts" derived from a simultaneous system of stochastic equations will in general not meet these conditions.

It might be useful to illustrate these points by a simple example. Let x_{1t} , x_{2t} , x_{3t} , $t = 1, 2, \dots, T$, be three observable economic series, each observation x_{it} containing some error of measurement x''_{it} , the true parts $x'_{it} = x_{it} - x''_{it}$ being the observations we should actually make if there were no errors of measurement. Then any relationship assumed to exist between the systematic parts x'_{1t} , x'_{2t} , and x'_{3t} (except in the case of a bookkeeping identity or a similar, "uninteresting," relation) would have to be of a stochastic type rather than an exact functional relation. Suppose now that economic theory has led us to a relation

$$(2.1) \quad x'_{1t} = \alpha_2 x'_{2t} + \alpha_3 x'_{3t} + u_{1t}, \quad t = 1, 2, \dots, T,$$

and that we want to estimate the parameters α_2 and α_3 , the variables u_{1t} being nonobservable random elements. Suppose further that, perhaps without the knowledge of the investigator, another, similar relation

$$(2.2) \quad x'_{1t} = \beta_2 x'_{2t} + \beta_3 x'_{3t} + u_{2t}, \quad t = 1, 2, \dots, T,$$

is also fulfilled. Let us further assume that the series x'_{3t} are "exogenous variables" having a fixed value for each value of t . Both x'_{1t} and x'_{2t} will then have to be stochastic variables, their stochastic properties being determined by those of the u 's and the transformations (2.1) and (2.2). Let us for simplicity assume that all the $2T$ u 's are normally and independently distributed and that $\mathcal{E}(u_{1t}) = \mathcal{E}(u_{2t}) = 0$ and $\mathcal{E}(u_{1t}^2) = \sigma_1^2$, $\mathcal{E}(u_{2t}^2) = \sigma_2^2$, for all values of t . We assume that the σ 's are unknown parameters.

On the basis of the definitions above it is possible to split the variables x_{1t} , x_{2t} , x_{3t} , artificially, into systematic parts and disturbances as follows:

$$(2.3) \quad x_{1t} = \left(\frac{\alpha_2 \beta_3 - \alpha_3 \beta_2}{\alpha_2 - \beta_2} x'_{3t} \right) + \left(\frac{\alpha_2 u_2 - \beta_2 u_1}{\alpha_2 - \beta_2} + x''_{1t} \right),$$

$$(2.4) \quad x_{2t} = \left(\frac{\beta_3 - \alpha_3}{\alpha_2 - \beta_2} x'_{3t} \right) + \left(\frac{u_2 - u_1}{\alpha_2 - \beta_2} + x''_{2t} \right),$$

$$(2.5) \quad x_{3t} = (x'_{3t}) + (x''_{3t}).$$

If the first term of the right-hand side in each of these equations is taken as the systematic part of x_{1t} , x_{2t} , x_{3t} , respectively, any two among these three systematic parts are linearly related. But the artificially constructed "disturbances" in the second term of (2.3) and (2.4) would here, in general, be stochastically dependent, which is another expression for the fact that the "true" systematic parts x'_{1t} , x'_{2t} , x'_{3t} are not multicollinear according to Frisch's definition.

It is well known that in the example discussed above *no unique estimate* exists for any of the parameters α_2 , α_3 , β_2 , β_3 , unless some additional a priori knowledge is available. In Koopmans' and Rubin's terminology [II-1.11] neither (2.1) nor (2.2) can be identified. But there exist two well-defined and identifiable relationships among the three variables x'_{1t} , x'_{2t} , x'_{3t} , namely the expected value of x'_{1t} , given x'_{2t} and x'_{3t} , and the expected value of x'_{2t} , given x'_{1t} and x'_{3t} . These relations are

$$(2.6) \quad \begin{aligned} & \mathcal{E}(x'_{1t} \mid x'_{2t}, x'_{3t}) \\ &= \frac{\alpha_2 \sigma_2^2 + \beta_2 \sigma_1^2}{\sigma_1^2 + \sigma_2^2} x'_{2t} + \frac{\alpha_3 \sigma_2^2 + \beta_3 \sigma_1^2}{\sigma_1^2 + \sigma_2^2} x'_{3t}, \end{aligned}$$

$$(2.7) \quad \begin{aligned} & \mathcal{E}(x'_{2t} \mid x'_{1t}, x'_{3t}) \\ &= \frac{\alpha_2 \sigma_2^2 + \beta_2 \sigma_1^2}{\alpha_2^2 \sigma_2^2 + \beta_2^2 \sigma_1^2} x'_{1t} - \frac{\alpha_2 \alpha_3 \sigma_2^2 + \beta_2 \beta_3 \sigma_1^2}{\alpha_2^2 \sigma_2^2 + \beta_2^2 \sigma_1^2} x'_{3t}. \end{aligned}$$

Apart from bias caused by the presence of the disturbances x''_{1t} , x''_{2t} , and x''_{3t} , these two prediction equations could be estimated from the observations x_{1t} , x_{2t} , x_{3t} . Now it is easy to see that for certain values of the unknown parameters the two planes (2.6) and (2.7) might not differ very much, in which case the method of "bunch analysis," by its very nature, would not show any sign of multicollinearity. This, however, does not mean that the α 's in the original equation (2.1) could be estimated. Thus, the apparent absence of multicollinearity would be no guarantee that the equation which we are seeking could actually be identified.

The purely geometric properties of a set of points in the sample space are insufficient as a basis for statistical inference. In fact, a sample of observations is just a set of cold, uninteresting numbers unless we have a theory concerning the stochastic mechanism that has produced them. To know the meaning of the results that a certain statistical method will yield we shall have to adopt a well-defined stochastic model. Then we can usually also find more efficient statistical tools than those that are not based on a probabilistic approach.

VI. PREDICTION AND LEAST SQUARES¹

BY LEONID HURWICZ

	Page
0. Introduction	266
0.1. Prediction and Structural Changes	266
0.2. Regression: A Special Case of Predictive Estimation . . .	271
0.3. Outline of the Remainder of this Paper	273
1. Notation and Definitions: Distributions of Observed Variates; Structure and Structural Changes; Predictive Estimation	274
2. Prediction for Policy Purposes	278
3. Optimal Properties and Mutual Relationships of Alternative Methods (Maximum-Likelihood, Markoff, Least-Squares) of Estimating the Regression Functions	278
3.0. Regression and Its Estimation; a Special Case of Predictive Estimation	278
3.1. Equivalence of Best Unbiased Linear and Least-Squares Estimates of the Regression Coefficients.	280
3.2. Equivalence of Maximum-Likelihood and Least-Squares Estimates of the Regression Coefficients	293
3.3. Regression Estimates under A Priori Restrictions	299

0. INTRODUCTION

0.1. *Prediction and Structural Changes*

0.1.1. *Predictive estimation.* Among the problems confronting the statistician is that of making predictions. As an example, one may be asked to forecast next year's national income on the basis of this year's income and government spending planned for the next year. Next year's income is the *predictand*, this year's

¹Part of the work on this paper was done in 1945-46 during the author's tenure of the Guggenheim Memorial Fellowship. Some of the problems considered arose in connection with the author's research at the Institute of Meteorology at the University of Chicago in 1944.

income and next year's planned government spending are the *predictors*. The *prediction* is a statement concerning, say, the most likely value of the predictand given the predictors; or the prediction may be formulated as an interval within which the predictand is likely to fall given the predictor values. Thus, in general, *prediction* is a statement about the (*conditional*) *probability distribution of the predictand given the predictors*.¹

In order to make the prediction, the statistician must obtain this conditional distribution on the basis of past observations. Thus, one may *estimate*, on the basis of past observations, the conditional distribution of a given year's income, with income in the preceding year and government spending planned for the given year specified.

However, the past observations can only tell what relationship between the values of the predictors and probability distribution of the predictand existed during the (past) observation period. Will the same relationship hold in the *future* period for which the prediction is to be made?

An affirmative answer implies that *no structural changes* have occurred, or are expected to occur, between the observation period and the period for which prediction is made. A *structural change*² is a change in any of the not directly observable ("structural") properties of the system. Thus structural changes occur if people's tastes shift or industrial productivity increases, provided neither tastes nor productivity are observed directly; should they have been observed, they could simply have been included among the variables of the system.

Hence, if *no structural changes* are expected to occur between

¹Actually, there are three steps in the prediction process: first, a *rule* is derived which enables us to construct the probability distribution of the predictand for any given values of the predictors (the regression equation is an example of such a rule); second, the relevant predictor values are found; third, with the help of the above-mentioned rule, the given predictor values are used to construct the relevant predictand distribution.

²A structural change transforms one *structure* into another. An example may be given as follows: Let an equation system consist of a supply equation and a demand equation for some commodity. Then the structure is given by the parameters of these two equations ("the structural parameters") and the parameters of the distribution of the disturbances. (More rigorous definitions are given below in section 1; they are also discussed in connection with the identification concept in [IV].) For brevity we often refer to the structure before and after the change as the "old" and "new" structures, respectively.

the observation period and the period for which prediction is to be made, all the statistician is required to do for prediction purposes is to provide a "best" (in some well-defined sense) estimate of the relationship that existed during the observation period between the probability distribution of the predictand and the values of the predictors. Once the relationship has been estimated, the forecast is made by specifying the predictor values to be substituted in that relationship and the latter will yield the appropriate (best estimate of the) probability distribution of the predictand.

The situation is fundamentally different when *structural changes are expected* to occur (or are already known to have taken place) between the observation period and the period for which prediction is to be made. Here it is not legitimate to apply to the future the past relationship between (the probability distribution of) the predictand and (the values of) the predictors: the structural changes, in general, will modify that relationship and it is this modified relationship that has to be determined. Once the modified relationship has been determined, a forecast can again be made by substituting into it the values of predictors; the relationship will then again yield (the best estimate of the probability distribution of) the predictand.

The problem is thus reduced to the *determination of the modified predictor-predictand relationship* after structural changes have taken place. How can this be accomplished? To begin with, it should be evident that if nothing is known about the nature of the expected structural changes, the new predictor-predictand relationship cannot be determined. Thus the nature of these changes must first be stated.

A structural change, by definition, transforms the "old" structure into the "new" (modified) structure. The nature of the *structural change* must be *completely specified* in the following sense: enough must be known about the change to make possible the derivation of the "new" structure if the "old" one were known. One example of such specification is that of the structural change consisting in, say, a productivity coefficient increasing by 10 per cent; this statement in itself does not tell us either the old or the new value of the coefficient, but it would enable us to find the new value if the old one were known.

Assuming that the structural change is completely specified, the "new" predictor-predictand relationship that is being sought may be determined in the following manner. From the available observations we estimate the ("old") joint probability of the ob-

served variates. Then, provided the model is *structure-identifying*¹ we can determine the "old" structure, i.e., the structure during the observation period. Next, we derive the "new" structure from the "old" one. (This can be done since the structural change is completely specified.)

Then from this "new" structure we derive the "new" distribution that the observed variates will have after the structural change has occurred.² The predictor-predictand relationship is "embedded" in the distribution of observed variates, so that from the distribution the relationship can always be obtained. (The converse is not true.) This procedure - obtaining an estimate of the modified predictor-predictand relationship given past observations, a completely specified structural change, and, possibly, an identifying (i.e., structure-identifying) model³ - is called *predictive estimation under changed structure*.

It may be found helpful to have a diagrammatic presentation of this procedure as given in Figure 1. As indicated on the diagram, the link B (from 2'' to 3) requires that the model be identifying. Hence, in general, *prediction under changed structure* requires that the identification power³ of the model be established so as to make *structural estimation* possible. (Structural estimation consists of links A'' and B, but computationally one can go directly from 1 to 3; this, of course, does not eliminate the need for identification power.⁴)

The point to be emphasized is that, in general, one cannot get to 6 directly from 2': it is necessary to go through 2'', 3, 4, and 5. However, if no structural change occurs, the "new" predictor-predictand relationship 6 is identical with the "old" one 2'. Hence there is no need for going all around from 1 through the stages 2'', 3, 4, and 5, to 6. Instead one may go directly from 1 to 2', using the link A'⁵. Here the link B is not used and hence there is *no need for identification power* in the case of *prediction under unchanged structure*.

¹I.e., provided enough a priori knowledge is available to make the structure determinable from the distribution of observed variates; cf. [IV].

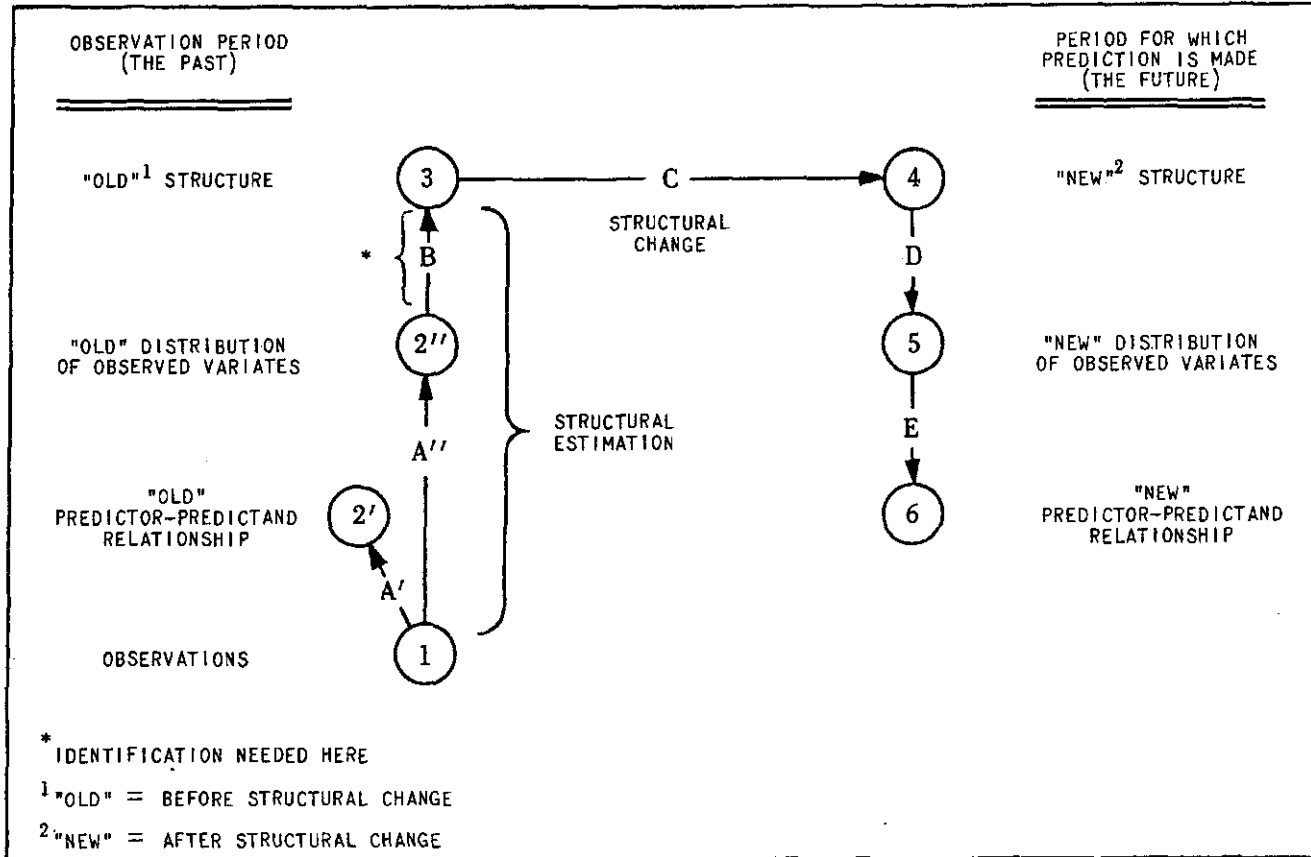
²The distribution of the observed variates can *always* be found when the structure is known; the converse is true only in structure-identifying models.

³A model that is structure-identifying is said to possess identification power.

⁴Cf. footnote on page 271.

⁵One could also go from 1 to 2'' and from there to 2'. This would be analogous to using the link E.

Figure 1. Predictive Estimation



0.1.2. *Prediction for structural policies.* Structural changes may be due either to causes consciously controlled by man ("policies") or to other factors. (However, not all policies are of this *structural* type. Some policies are carried out through certain observable variables, say the quantity of money, and need not imply any structural changes.) If several courses of action are available, the policy decision will be made on the basis of the prediction of the likely effects of the different policies. If the policy is of the type that implies a structural change, the decision requires that prediction be made of behavior of certain economic variables under the different possible structural changes. Thus one is forced to go through the procedure of predictive estimation under changed structure, a phase of which is structural estimation, and one must have an identifying model¹.

0.1.3. *Prediction concept generalized.* The term *prediction* may be given a somewhat broader meaning than that of making probability statements about the future on the basis of the past. Sometimes it is desired to estimate the values of a variable given a contemporary or later value of another variable. The need for such a procedure may be due to the fact that the variable to be "predicted" from subsequent values of another is not directly available or, perhaps, is available only with a considerable delay.

Thus, the essence of the prediction concept lies not in the temporal sequence of the predictor and the predictand, but rather in the conditional nature of the (probability) statement about some set of variables given another such set.

Also, prediction need not apply to one variable only. One may, for instance, wish to know the future correlation of some two variables given the past values of a third one, etc.

0.2. *Regression: A Special Case of Predictive Estimation*

0.2.1. *Prediction of a single variable under unchanged structure.* When the structure is unchanged, there is, as we have seen, no need for estimation of the structure of the system; our objec-

¹There are exceptional situations where some short-cuts can be made: it might even happen that *partial identification power* (see [IV]) is sufficient, but in general this will not be the case. In certain case where prediction of policy effects does require *total identification power* (*ibid.*), it is nevertheless possible to decide on a preferable course of action when such identification power (*ibid.*, 8.1) does not exist.

tive is 2' in Figure 1 and this can be reached directly from 1 through the link A'. This link consists in obtaining the conditional distribution of the (univariate) predictand. The distribution may be regarded as given by its moments. If the distribution is normal, the first two moments (say the mean and the variance) suffice to determine the distribution; in general, the higher moments are also needed, but they will not be considered here.

The relationship between the mean of the predictand and (the values of) the predictors is the *regression* function; the relationship between the variance of the predictand and the predictors is the *scedastic* function.

0.2.2. *Methods of estimating the regression functions.* The problem of statistical inference is now reduced to finding *optimal methods of estimating* these functions. We shall see that, under rather general conditions, in models with normally distributed disturbances, the *maximum-likelihood* estimates of the (true) regression coefficients will be given by the application of the *least-squares* procedure. Under even more general conditions, the least-squares method will yield *consistent*, but probably somewhat inefficient, estimates of the regression coefficients.

Hence, if prediction under unchanged structure is intended, and the sample is sufficiently large, the method of least squares is a safe one (in the sense of consistency) to use, although not necessarily the best one (in the sense of efficiency).

But for small and medium samples, the asymptotic properties of consistency and efficiency are not adequate safeguards against serious sampling errors. It is therefore of interest to investigate the small-sample properties of the least-squares estimates both for the case where the least-squares estimates satisfy the maximum-likelihood criterion and the case where the maximum-likelihood estimates differ from the least-squares estimates.

It is found that in both cases the least-squares estimates are strongly *biased*¹ when the series is of the *autoregressive* type, i.e., where the variables satisfy a system of stochastic difference equations with *lagged endogenous* variables present.² On the other hand,

¹An estimate is biased if its expectation, for any given size sample, is not identically equal to the true value of the parameter.

²The example given in section 0.1.1 was of the autoregressive type since national income, which is an endogenous variable, appeared both as a predictand and, with a lag of one year, as a predictor.

in *lagless* (nonautoregressive) systems, the least-squares estimates of the regression coefficients do possess an important optimal property: they are the best¹ unbiased linear estimates of those coefficients. It is important to note that this property is not asymptotic, but holds for any finite-sized sample.²

0.2.3. *Regression coefficients and structural coefficients.*

The question of the extent to which the least-squares estimates are "good" estimates of the regression coefficients, treated in the preceding paragraph, should not be confused with the following one: under what conditions do the structural coefficients of an equation in a simultaneous stochastic system coincide with a set of regression coefficients in some predictor-predictand relationship?

When such a coincidence does occur, it is clear that the problem of structural estimation (for that particular equation) is exactly identical with the problem of predictive estimation under unchanged structure, when the predictors and the predictand are suitably chosen. Hence in such a case the least-squares estimates are "good" or "best" estimates, not only of the regression coefficients, but also of the structural coefficients. Just when this does happen is discussed elsewhere in the volume.³

0.3. *Outline of the Remainder of This Paper*

The remainder of this paper is divided into three sections. Section 1 provides definitions of the concepts used as well as the principles of notation. The very brief section 2 is devoted to

¹I.e., with smallest sampling variance among all the unbiased linear estimates.

²It is conceivable that for certain populations there exist better unbiased linear estimates than the least-squares estimates, but there are no linear estimates unbiased for all populations and better than the least-squares estimates.

³[I], [II]. One case where structural and regression coefficients coincide is when in a given equation all variables but one are either lagged or exogenous. In general, when a least-squares regression has a high multiple-correlation coefficient and the standard errors of the regression coefficients are small, this only proves that we have a good (and probably accurately estimated) predictive relationship, but not necessarily that this relationship coincides with one of the structural equations (even if the variables entering the predictive formula are exactly the same as those entering a given structural equation). See [I] and [II - 3.3.7] in the special case where the set *I* contains only one structural equation. See also 1.2.1 below.

policy aspects of the prediction problem. In section 3 the field of inquiry is narrowed down to that of regression estimation (see section 3.0): the optimal properties of the maximum-likelihood, Markoff (see 3.1.1.2), and least-squares regression estimates and the domain of their equivalence are investigated. A more detailed outline of section 3 is given in 3.0.3.

1. NOTATION AND DEFINITIONS:

DISTRIBUTIONS OF OBSERVED VARIATES; STRUCTURE AND STRUCTURAL CHANGES; PREDICTIVE ESTIMATION

1.1. Let the system involve a set (denoted by capitals)

$$(1) \quad Y = (y_0, y_1, \dots, y_{K_y})$$

of the $K_y + 1$ observed stochastic variates y_g and a set

$$(2) \quad Z = (z_1, z_2, \dots, z_{K_z})$$

of the K_z observed fixed¹ variates z_j . Let the cumulative distribution function of Y be

$$(3) \quad F \equiv F_y \equiv F_y(y; z),$$

where y and z are vectors whose components are the elements of the sets Y and Z . For any nonempty² proper³ subset Y^* of Y (i.e., $\Lambda \neq Y^* \subset Y$) we have the marginal cumulative distribution function

$$(4) \quad F_* \equiv F_{y^*} \equiv F_{y^*}(y^*; z).$$

¹A variable z is said to be a "fixed" variate if its t th observation has the cumulative distribution function given by

$$(2') \quad F(z_t) = \begin{cases} 0 & \text{for } z_t < z_{0t}, \\ 1 & \text{for } z_t \geq z_{0t}, \end{cases}$$

where the z_{0t} 's form an arbitrarily chosen sequence of constants.

²The empty set will be denoted by Λ .

³ $E' \supset E''$ means that E'' is a subset of E' but $E' \neq E''$. $E' \supset E''$ means that E'' is a subset of E' , but their equality is not excluded.

Given two sets of y^* and y^{**} such that¹

$$(5) \quad \Lambda \neq Y^* \subset Y$$

and

$$(6) \quad Y^{**} \equiv Y - Y^*,$$

we define the *conditional* marginal cumulative distribution function

$$(7) \quad F_{*|**} \equiv F_{y^*|y^{**}} \equiv F_{y^*|y^{**}}(y^* | y^{**}; z).$$

We shall refer to F in (3) as the *complete* distribution, while the marginal and conditional distributions in (4) and (7) will be called *derivates* of F , or derived distributions.

1.2. Let the structure² of the system be denoted by S . Then S determines F , which may be written as

$$(8) \quad F = \varphi(S).$$

Now $F_{*|**}$ is a derivate of F , hence

$$(9) \quad F_{*|**} = \varphi_{*|**}(S).$$

Suppose a *structural change* takes place, so that the "old" structure $S^{(0)}$ is replaced by the "new" one $S^{(1)}$. We have

$$(10) \quad S^{(1)} = \mathfrak{J}(S^{(0)}),$$

where \mathfrak{J} must be so specified that when $S^{(0)}$ is given, $S^{(1)}$ can be uniquely determined. Then

$$(11) \quad F^{(1)} = \varphi(S^{(1)}) = \varphi\{\mathfrak{J}(S^{(0)})\} = \varphi(S^{(0)}; \mathfrak{J})$$

and

¹It will be noted that Y^{**} is not required to precede Y^* in time. The opposite may well be true.

²For general definitions of structure, model, and identification, see [IV]. In the usually treated parametric normal linear case described by the matrix equation $Ax = u$ (where A is the structural coefficient matrix, u the nonobservable disturbance, and x the observed variate), the structure is defined by A and the disturbance covariance matrix Σ^u so that $S = (A, \Sigma^u)$.

$$(12) \quad F_{**|**}^{(1)} = \varphi_{**|**}(S^{(1)}) = \varphi_{**|**}\{\mathfrak{D}(S^{(0)})\} \equiv \varphi_{**|**}(S^{(0)}, \mathfrak{D})$$

while

$$(13) \quad F_{**|**}^{(0)} = \varphi_{**|**}(S^{(0)}).$$

In general, there does *not* exist a relationship

$$(14) \quad F_{**|**}^{(1)} = \psi_{**|**}\{F_{**|**}^{(0)}; \mathfrak{D}\}.$$

That is, apart from special cases, the knowledge of $S^{(0)}$ is necessary in order that $F_{**|**}^{(1)}$ should be obtained.

1.2.1. The situation discussed in the preceding section may be illustrated by the following example. Let our model consist of equations

$$y_0 = \beta_1 y_1 + \gamma_1 z_1 + u_1, \quad y_1 = \beta_2 y_0 + \gamma_2 z_2 + u_2,$$

where the z 's are exogenous with a moment matrix $[\zeta_{ij}]$ and the u 's are the disturbances with zero means and a covariance matrix $[\sigma_{ij}]$.

If we are to predict y_0 given y_1 and z_1 , we must have the regression function

$$\mathcal{E}(y_0 | y_1, z_1) = \chi_1^y y_1 + \chi_1^z z_1.$$

The χ 's can be expressed in terms of the structural parameters of the system. Thus, for instance,

$$\begin{aligned} \chi_1^y - \beta_1 &= (1 - \beta_1 \beta_2)(\sigma_{12} + \beta_2 \sigma_{11}) \\ &\times \left\{ \sigma_{22} + \beta_2^2 \sigma_{11} + 2\beta_2 \sigma_{12} + \gamma_2^2 \left(\frac{\zeta_{22}}{\zeta_{11}} - \frac{\zeta_{12}^2}{\zeta_{11}} \right) \right\}^{-1}. \end{aligned}$$

Unless $\sigma_{12} + \beta_2 \sigma_{11} = 0$, in which case the single-equation approach is applicable, we find that the difference $\Delta\beta_1$ between a "new" value $\beta_1^{(1)}$ and an "old" value $\beta_1^{(0)}$ is not sufficient to determine the corresponding increment $\Delta\chi_1^y$. Hence, in general, we need additional knowledge of the structural parameters in the right-hand member of the above expression.

It may be added that the expression for $\chi_1^y - \beta_1$ gives the magnitude of the bias involved in using the single-equation least-squares value χ_1^y as a measure of β_1 . This bias may be high even if the multiple correlation coefficient ρ^2 (y_0 on y_1 and z_1) is high and the (conventionally computed) standard errors σ_χ low. Thus, for instance, an example has been constructed by the author with $\beta_1 = 0$ while $\chi_1^y = .99$ despite the fact that $\rho^2 = .9901$ and $\sigma_{\chi_1^y}^2 = .005$.

1.3. Now let the observations (prior to the structural change \mathfrak{D}) be given by the matrix $X \equiv [X_{it}]$; ($i = 0, 1, \dots, K_y + K_z$; $t = 1, 2, \dots, T$).¹ We obtain an estimate, denoted by $\text{est} F^{(0)}(X)$, of $F^{(0)}$. Then, if we are dealing with an identifying model² we can obtain an estimate of $S^{(0)}$ from, say,³ $\text{est} S^{(0)} = \psi(\text{est} F^{(0)}; \mathfrak{G}_0)$. By virtue of (10), we have

$$(15) \quad \text{est} S^{(1)} = \mathfrak{D}(\text{est} S^{(0)})$$

and, consequently, we obtain

$$(16) \quad \text{est} F_{*|**}^{(1)} = \varphi_{*|**}(\text{est} S^{(1)}) \equiv \varphi_{*|**}[\mathfrak{D}\{\psi[\text{est} F^{(0)}(X); \mathfrak{G}_0]\}].$$

The process of obtaining $\text{est} F_{*|**}^{(1)}$ from X is called *predictive estimation*. Y^* is the *predictand*, Y^{**} and Z the *predictors*. In general, prediction involves the use of X , \mathfrak{D} , and \mathfrak{G}_0 . However, when $\mathfrak{D} = \mathfrak{I}$, where \mathfrak{I} is the identity transformation, we clearly have

$$(17) \quad \text{est} F_{*|**}^{(1)} = \text{est} F_{*|**}^{(0)}.$$

This is a special case where a relationship of the type (14) does exist, hence it is not necessary here to obtain $S^{(0)}$, etc. In particular, therefore, it is not necessary that the model be identifying.

It is with regard to this special case of *predictive estimation*

¹In lagless systems, t may, but need not necessarily, mean time.

²See [IV].

³ \mathfrak{G}_0 is the set of all structures permitted by the identifying restrictions.

under unchanged structure ($\mathfrak{D} = \mathfrak{A}$) that we examine the optimal properties of the least-squares method in section 3.

2. PREDICTION FOR POLICY PURPOSES

One of the most important cases of need for prediction under changed structure is that of certain types of economic policies. There are cases where structural changes are not due to policy decisions (e.g., spontaneous changes of tastes); nor do all policies imply structural change; some only involve manipulation of Z (non-structural policies). But when a change is to be made between two policies implying the respective structural changes \mathfrak{D}' and \mathfrak{D}'' , one cannot decide between them (no matter what the policy objective) without knowing just how they will affect certain variables. Here, since $\mathfrak{D} \neq \mathfrak{A}$ for one or both of \mathfrak{D}' and \mathfrak{D}'' , we deal with predictive estimation under structural change; hence in general we must resort to structural estimation, and it is necessary that the model possess a certain degree of identification power.

On the other hand, suppose that the policy contemplated involves *no structural change*, so that it is either one of inaction or manipulation of Z . Then, regardless of whether or not Z is to be manipulated, we have $\mathfrak{D} = \mathfrak{A}$, and the effects of such a policy (nonstructural policy or complete inaction) can be predicted without resort to structural estimation; hence it is not necessary to have an identifying model.

3. OPTIMAL PROPERTIES AND MUTUAL RELATIONSHIPS OF ALTERNATIVE METHODS (MAXIMUM-LIKELIHOOD, MARKOFF, LEAST-SQUARES) OF ESTIMATING THE REGRESSION FUNCTIONS

3.0. *Regression and its Estimation: A Special Case of Predictive Estimation*

3.0.1. The remainder of this paper is devoted to prediction under unchanged structure (see section 1.3 above).

The case considered is not the most general one even under unchanged structure. Three important specializing assumptions are made:

- (a) *The choice of the sets Y^* and Y^{**} .*

$$(18) \quad Y^* = (y_0),$$

so that Y^* is a one-element set.

$$(19) \quad Y^{**} = Y - Y^* = (y_1, y_2, \dots, y_{K_y}),$$

i.e., Y^{**} is the complement of Y^* in Y .¹

(b) *Admissible distributions.* It will be assumed that the regression of y_0 on the elements of Y^{**} is linear and homoscedastic,² so that³

$$(20) \quad \mathcal{E}(y_0 | Y^{**}; Z) = \sum_{i=1}^{K_y} \chi_i^y y_i + \sum_{j=1}^{K_z} \chi_j^z z_j + \chi_0$$

and

$$(21) \quad \sigma^2(y_0 | Y^{**}; Z) = \text{const.}$$

(c) *Properties to be estimated.* While the general definition of prediction calls for estimation of $F_{*|**}$, our study will be confined to the estimation of the first moment of $F_{*|**}$, i.e., to the estimation of the χ 's (*regression coefficients*) in (21).

3.0.2. *Estimation principles.* Three important principles have been applied in the past in estimating the regression coefficients. They are: (a) maximum-likelihood, (b) best unbiased linear (Mar-

¹(18) does narrow down the scope of the study, but (19) only seems to. For suppose that, contrary to (19), we have

$$(19') \quad Y - Y^* - Y^{**} \neq \Lambda;$$

then we form the marginal distribution of $Y^{***} \equiv Y^* + Y^{**}$, say F_{***} , and consider it as the parent distribution; obviously (19) now holds, provided Y is replaced by Y^{***} , and the nature of the problem is unchanged.

²The restriction of homoscedasticity is not difficult to relax, provided ratios of conditional variances are assumed known or even under more general conditions.

³ χ_0 in (21) may, without loss of generality, be assumed to vanish since it is always permissible to have one of the Z 's identically equal to unity.

koff), and (c) least-squares. It is well known that the three methods are related and that there are cases where all three, or at least two, are equivalent. But examples of divergence also exist. Therefore, it becomes important to investigate the domain of their equivalence in relation to the nature of the model.

3.0.3. *Outline of the remainder of this section.* Section 3.1 of this paper will be devoted to the study of the domain of equivalence of the best unbiased linear and least-squares estimates; section 3.2 to equivalence of maximum-likelihood and least-squares estimates. Each section is divided into two subsections; the first subsection of each section (i.e., 3.1.1 and 3.2.1) treats of lagless models.

The problem of additional restrictions on the parameters is treated in section 3.3 for both best unbiased linear and maximum-likelihood estimates, with or without lags.

3.1. *Equivalence of Best Unbiased Linear and Least-Squares Estimates of the Regression Coefficients*

3.1.1. *Lags absent.*

3.1.1.1. Consider a sample of size T where the t th observation is a $(1 + K_y + K_z)$ -dimensional vector (y_t, z_t) with

$$(22) \quad \begin{aligned} y_t &= (y_{0t}, y_{1t}, \dots, y_{K_y t}), \\ z_t &= (z_{1t}, z_{2t}, \dots, z_{K_z t}). \end{aligned}$$

It will be assumed that the observations for successive values of t are independent in the probability sense, but they need not come from the same universe. (This will be qualified later.)

Given the cumulative distribution function $F_t(y_t, z_t)$ of the universe from which the t th observation is drawn, the joint cumulative distribution function of the sample is

$$(23) \quad \prod_{t=1}^T F_t(y_t; z_t).$$

Now from each of the F_t it is possible to derive the conditional cumulative distribution function of y_0 given the remaining y 's and z 's. These conditional distributions, to be written as

$$(24) \quad F_t^G(y_{0t} | y_t^{(0)}; z_t),$$

where $y_t^{(0)} = (y_{1t}, \dots, y_{kt})$, will be assumed identical for all t so that the subscript t of F_t^G may be dropped. Moreover, the first two moments of this conditional distribution F^G are assumed to satisfy equations (20) and (21), so that the regression of y_0 on the other y 's and z 's is linear and homoscedastic.

We now seek to form (out of observations of the sample described above) an estimate q for any one of the regression coefficients χ in (20). (For the moment the affixes of χ and q may be omitted.) These estimates are to be *linear* in the y_0 and can thus be written as

$$(25) \quad q = \sum_{t=1}^T \varphi_t y_{0t},$$

where

$$(26) \quad \varphi_t = \varphi_t(y_1^{(0)}, y_2^{(0)}, \dots, y_T^{(0)}; z_1, z_2, \dots, z_T).$$

The expectation of q is given by the Stieltjes integral

$$(27) \quad \mathcal{E}(q) = \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} \sum_{t=1}^T \varphi_t y_{0t} \prod_{t=1}^T dF_t(y_t; z_t)$$

and it may or may not equal χ . Whether this is the case depends a) on the choice of the φ_t 's, and b) on the nature of the function in (23).

For a given choice of φ_t 's, it may happen that q would equal χ only provided the function in (23) belongs to a specified family \mathfrak{F}_0 of such functions. In such a case we say that the estimate q corresponding to that particular choice of the φ_t 's is (conditionally) *unbiased with regard to \mathfrak{F}_0* . If the family \mathfrak{F}_0 is the family \mathfrak{F} of all distribution functions, the estimate q will be said to be *absolutely unbiased*.

We shall now derive a *necessary and sufficient condition* to be satisfied by the φ_t 's if q is to be an *absolutely unbiased estimate* of χ . The regression is assumed to be linear, that is, (20) holds.

Using (20) and the well-known properties of conditional distri-

butions, we obtain¹ from (27)

$$\begin{aligned}
 \mathcal{E}(q) &= \int \cdots \int \sum_t \varphi_t y_{0t} \prod_t dF^G(y_{0t} | y_t^{(0)}; z_t) \prod_t dF_t^M(y_t^{(0)}; z_t) \\
 (28) \quad &= \int \cdots \int \sum_t \varphi_t \mathcal{E}(y_{0t} | y_t^{(0)}; z_t) \prod_t dF_t^M(y_t^{(0)}; z_t) \\
 &= \int \cdots \int \sum_t \varphi_t \left(\sum_i \chi_i^y y_{it} + \sum_j \chi_j^z z_{jt} \right) \prod_t dF_t^M(y_t^{(0)}; z_t),
 \end{aligned}$$

or, rearranging,

$$\begin{aligned}
 \mathcal{E}(q) &= \sum_i \chi_i^y \int \cdots \int \left(\sum_t \varphi_t y_{it} \right) \prod_t dF_t^M \\
 (29) \quad &+ \sum_j \chi_j^z \int \cdots \int \left(\sum_t \varphi_t z_{jt} \right) \prod_t dF_t^M,
 \end{aligned}$$

where F^G is a conditional, and F_t^M a marginal, cumulative distribution function.

So far neither q nor the φ_t 's have had any affixes referring to the χ which is being estimated. Now, for definiteness, we shall assume that it is χ_1^y that is being estimated; its estimate will be written q_1^y ; the corresponding φ_t 's will be denoted by $\varphi_t^{y_1}$, but the superscripts may be omitted where no danger of confusion exists.

Now $\mathcal{E}(q_1^y)$ will be given by (29) with proper affixes attached to q and the φ_t 's, thus yielding

$$\begin{aligned}
 \mathcal{E}(q_1^y) &= \sum_i \chi_i^y \int \cdots \int \sum_t \varphi_t^{y_1} y_{it} \prod_t dF_t^M \\
 (30) \quad &+ \sum_j \chi_j^z \int \cdots \int \sum_t \varphi_t^{y_1} z_{jt} \prod_t dF_t^M.
 \end{aligned}$$

If q_1^y is to be an absolutely unbiased estimate of χ_1^y , the right-hand member of (30) must equal χ_1^y for any choice of $\prod_t dF_t^M$. This

¹The integrals are between $-\infty$ and $+\infty$ unless otherwise indicated.

will be the case if and only if

$$\begin{aligned}
 & \int \cdots \int \left(\sum_t \varphi_t^{y_1} y_{it} - 1 \right) \prod_t dF_t^M \equiv 0, \\
 (31) \quad & \int \cdots \int \sum_t \varphi_t^{y_1} y_{it} \prod_t dF_t^M \equiv 0, \quad \text{for all } i \neq 1, \\
 & \int \cdots \int \sum_t \varphi_t^{y_1} z_{jt} \prod_t dF_t^M \equiv 0, \quad \text{for all } j.
 \end{aligned}$$

It is easily seen that a set of *sufficient* conditions for (31) to hold is:

$$\begin{aligned}
 & \sum_t \varphi_t^{y_1} y_{it} - 1 \equiv 0, \\
 (32) \quad & \sum_t \varphi_t^{y_1} y_{it} \equiv 0, \quad \text{for all } i \neq 1, \\
 & \sum_t \varphi_t^{y_1} z_{jt} \equiv 0, \quad \text{for all } j.
 \end{aligned}$$

But it can be shown that the identities (32) are also *necessary* for q_1^y to be absolutely unbiased. To carry out the proof it is enough if a family \mathcal{F}^M of $\{F_t^M; (t = 1, 2, \dots, T)\}$ is found for which the vanishing of the integrals in (31) implies the vanishing of their integrands, as given in (32). One such family is given by making the y_{it} 's normally independently distributed with a common variance but different means¹. Many more could, no doubt, be found. Therefore, the identities (32) are the *necessary and sufficient* conditions for q_1^y to be absolutely unbiased. But if q_1^y were to be only conditionally unbiased, depending on the family with regard to which lack of bias is postulated, either (31) or (32) would have to be used.

3.1.1.2. We shall now proceed to define the *best* absolutely unbiased linear estimate of, say, χ_1^y . This estimate will be denoted by \bar{q}_1^y and will have the following properties:

¹[Wald, 1944].

- (a) linearity [i.e., (25) with (26)],
- (b) absolute unbiasedness [i.e., (32)],
- (c) minimum sampling variance (to be defined below).

The property (c) implies that the sampling variance $\sigma_{\bar{q}_1^y}^2$ of \bar{q}_1^y will not exceed the sampling variance of any other linear absolutely unbiased estimate of χ_1^y .

Mathematically, the problem of finding an estimate possessing the properties (a), (b), (c) (a "Markoff estimate") requires minimization of the integral giving $\sigma_{q_1^y}^2$ [where q_1^y has properties (a) and (b)] subject to the side relation (32).

Now

$$(33) \quad \sigma_{q_1^y}^2 = \int \cdots \int \left(\sum_t \varphi_t y_{0t} - \chi_1^y \right)^2 \prod_t dF_t(y_t; z_t)$$

or

$$(34) \quad \sigma_{q_1^y}^2 = \int \cdots \int \left(\sum_t \varphi_t y_{0t} \right)^2 \prod_t dF_t - (\chi_1^y)^2.$$

$\mathcal{E}(q_1^y)^2$ is given by the integral in the right-hand member of (34).

Now, using the properties of homoscedasticity and applying transformations analogous to those in (28), we have

$$(35) \quad \mathcal{E}(q_1^y)^2 = \int \cdots \int \sum_{t', t''} \varphi_{t'} \varphi_{t''} y_{0t'} y_{0t''} \prod_t dF_t^C \prod_t dF_t^M$$

which, after simplification gives

$$(36) \quad \sigma_{q_1^y}^2 = \sigma^2 \int \cdots \int \sum_t \varphi_t^2 \prod_t dF_t^M.$$

Hence the problem is reduced to that of minimizing the integral (36) subject to the restrictions (32). Formally, this is a Calculus of Variations problem, but actually, since no derivatives enter the integrand, we may simply minimize the integrand $\sum_t \varphi_t^2$

subject to (32).

At this stage the mathematical problem is exactly the same as that treated in the older version of the Markoff Theorem.¹ Therefore, the estimate \tilde{q}_1^y , obtained by minimizing the sampling variance $\sigma_{\tilde{q}_1^y}^2$ with regard to the φ 's subject to conditions (a) and (b) listed above, is identically equal to the least-squares estimate \tilde{q}_1 , obtained by minimizing the expression

$$(37) \quad S \equiv \sum_t (y_{0t} - \sum_i q_i^y y_{it} - \sum_j q_j^z z_{jt})^2$$

with regard to the q 's. Thus "best" (in the Markoff sense) *absolutely unbiased linear* estimates are *identical with the least-squares estimates*.

3.1.1.3. In the preceding section we considered Markoff estimates that were *absolutely unbiased*. They turned out to be identical with the least-squares estimates.

In general, this is not the case for Markoff estimates *conditionally unbiased* with regard to a family of distributions for which (31) does not imply (32).²

Since this is a proposition of a negative nature, it can be proved by demonstration of a special case. We shall assume that only χ_1^y (to be written simply as χ) is different from zero, while all the other χ 's vanish.

Here we again proceed to minimize σ_q^2 subject to

$$(38) \quad \int \cdots \int \left(\sum_t \varphi_t y_{1t} - 1 \right) \prod_t dF_t^M \equiv 0$$

which is all that remains of (31).

Now in this case (35) gives

$$(39) \quad \mathcal{E}(q^2) = \int \cdots \int \left\{ \sigma^2 \sum_t \varphi_t^2 + \chi^2 \left(\sum_t \varphi_t y_{1t} \right)^2 \right\} \prod_t dF_t^M$$

or

¹[David and Neyman].

²If the estimates are to be conditionally unbiased with regard to a family of distributions for which (31) does imply (32), then such conditionally unbiased Markoff estimates would, of course, still be identical with the least-squares estimates.

or

$$(40) \quad \mathcal{E}(q^2) = \sigma^2 \int \cdots \int \left\{ \sum_t \varphi_t^2 + \frac{\chi^2}{\sigma^2} \left(\sum_t \varphi_t y_{1t} \right)^2 \right\} \prod_t dF_t^M.$$

It is evident that we may minimize the integral in (40) instead of σ_q^2 . Writing

$$(41) \quad v^2 = \frac{\chi^2}{\sigma^2},$$

we may therefore reformulate the problem as that of minimizing the integral

$$(42) \quad M = \int \cdots \int \left\{ \sum_t \varphi_t^2 + v^2 \left(\sum_t \varphi_t y_{1t} \right)^2 \right\} \prod_t dF_t^M$$

subject to (38).

Now

$$(43) \quad M = \int \cdots \int \sum_{t', t''} \varphi_{t'} \varphi_{t''} s_{t' t''} \prod_t dF_t^M,$$

where

$$(44) \quad s_{t' t''} = \delta_{t' t''} + v^2 y_{1t'} y_{1t''},$$

and where $\delta_{t' t''}$ is the Kronecker symbol.

Now the integrand in (43) may be written in matrix form as

$$(45) \quad \varphi' S \varphi$$

where¹

$$(46) \quad \varphi = \{\varphi_1, \dots, \varphi_T\},$$

$$S = [s_{t' t''}],$$

and φ' is the transpose of φ .

¹The notation used here is the customary one (φ is a column vector, φ' a row vector) rather than that used in other parts of this volume.

Similarly, the integrand of (38) is

$$(47) \quad \varphi' y_1 - 1.$$

We now minimize (45) with respect to φ' subject to the vanishing of (47) and denote the Lagrange multiplier by μ (a scalar).

Differentiating the expression

$$(48) \quad \varphi' S \varphi - \mu \varphi' y_1$$

with respect to φ' , we obtain

$$(49) \quad \frac{\partial}{\partial \varphi'} \{ \varphi' S \varphi - \mu \varphi' y_1 \} = S \varphi - \mu y_1 = 0,$$

and, hence,

$$(50) \quad \varphi = \mu S^{-1} y_1.$$

To evaluate

$$(51) \quad S^{-1} y_1$$

we observe that

$$(52) \quad S = I + v^2 y_1 y_1'.$$

Hence

$$(53) \quad \begin{aligned} S y_1 &= (I + v^2 y_1 y_1') y_1 = I y_1 + v^2 (y_1' y_1) y_1 \\ &= I \{ 1 + v^2 (y_1' y_1) \} y_1 \end{aligned}$$

and

$$(54) \quad S^{-1} S y_1 = \{ 1 + v^2 (y_1' y_1) \} S^{-1} y_1,$$

so that the desired expression is

$$(55) \quad S^{-1} y_1 = \frac{1}{1 + v^2 (y_1' y_1)} y_1.$$

Thus, substituting (55) into (56), we have

$$(56) \quad \varphi = \mu \frac{1}{1 + v^2(y_1' y_1)} y_1$$

Only the multiplier μ remains to be evaluated.

Premultiplying (56) with y_1' and integrating, we find that

$$(57) \quad \int \cdots \int y_1' \varphi \prod_t dF_t^M = \mu \int \cdots \int \frac{y_1' y_1}{1 + v^2(y_1' y_1)} \prod_t dF_t^M,$$

where the left-hand member equals unity by (42). Therefore,

$$(58) \quad \mu^{-1} = \int \cdots \int \frac{y_1' y_1}{1 + v^2(y_1' y_1)} \prod_t dF_t^M,$$

and, finally,

$$(59) \quad \varphi = \frac{y_1}{\{1 + v^2(y_1' y_1)\} \int \cdots \int \frac{y_1' y_1}{1 + v^2(y_1' y_1)} \prod_t dF_t^M}.$$

Thus, if matrix notation is abandoned, \bar{q} , the Markoff "estimate"¹ of the regression coefficient χ is

$$(60) \quad \bar{q} = \frac{\sum_t y_{1t} y_{0t}}{(1 + v^2 \sum_t y_{1t}^2) \int \cdots \int \frac{\sum_t y_{1t}^2}{1 + v^2 \sum_t y_{1t}^2} \prod_t dF_t^M},$$

¹ q depends not only on the observations, but also on v^2 which is a population parameter and, in general, is not known; an analogous "estimate" was obtained in [Johnson], where the knowledge of the coefficient of variation is needed in order to obtain a least variance (about true value, not necessarily unbiased) estimate of the mean in a normal population.

In order to avoid this type of solution one should subject the above Calculus of Variations procedure to the additional restrictions $\partial \varphi_t / \partial \chi =$

while the least-squares estimate \tilde{q} of χ would have been

$$(61) \quad \tilde{q} = \frac{\sum_t y_{1t} y_{0t}}{\sum_t y_{1t}^2} .$$

3.1.1.4. Three observations should be made concerning the foregoing result.

(a) \bar{q} is actually a better "estimate" than \tilde{q} , i.e.,

$$(62) \quad \sigma_{\bar{q}}^2 < \sigma_{\tilde{q}}^2 .$$

To see this we shall calculate both variances. Writing

$$(63) \quad \begin{aligned} y_1' y_1 &= p , \\ 1 + v^2 p &= k , \end{aligned}$$

we may, using (56), write the integrand $Q = \varphi' S \varphi$ of M in (43) as

$$(64) \quad Q = \mu^2 (p/k)$$

where μ , given by (58), may be written as

$$(65) \quad \mu = \left\{ \int \cdots \int (p/k) \prod_t dF_t^M \right\}^{-1} .$$

Hence the integral M in (43) is

$$(66) \quad \begin{aligned} M &= \int \cdots \int Q \prod_t dF_t^M = \int \cdots \int \mu^2 \frac{p}{k} \prod_t dF_t^M \\ &= \mu^2 \int \cdots \int \frac{p}{k} \prod_t dF_t^M = \mu^2 \mu^{-1} = \mu , \end{aligned}$$

0, $\partial \varphi_t / \partial \sigma = 0$. This has not yet been done, but it would be of interest to see whether these additional restrictions would imply $\bar{q} = \tilde{q}$. (Asymptotically, at least, this would seem to be the case because of the efficiency property of \tilde{q} as a maximum-likelihood estimate of χ ; see below, section 3.2. Exceptions might conceivably arise if linear unbiased estimates of χ exist whose asymptotic distributions are not normal).

and, using (65) and (66),

$$(67) \quad \sigma_{\bar{q}}^2 = \sigma^2 M - \chi^2 = \sigma^2 \left\{ \int \cdots \int \frac{p}{k} \prod_t dF_t^M \right\}^{-1} - \chi^2.$$

Having obtained the variance of the Markoff "estimate" \bar{q} of χ , we shall now evaluate the variance of the least-squares estimate \tilde{q} of χ .

Since

$$(68) \quad \tilde{q} = \frac{y_1' y_0}{y_1' y_1}$$

is also unbiased,¹ we have

$$(69) \quad \sigma_{\tilde{q}}^2 = \mathcal{E}(\tilde{q}^2) - \chi^2,$$

with

$$\begin{aligned} \mathcal{E}(\tilde{q}^2) &= \int \cdots \int \frac{1}{p} \sum_{t', t''} y_{1t'} y_{1t''} \int \cdots \int y_{0t'} y_{0t''} \\ &\quad \times \prod_t dF^C(y_{0t} | y_{1t}) \prod_t dF_t^M(y_{1t}) \\ (70) \quad &= \sigma^2 \int \cdots \int \frac{1}{p^2} y_1' S y_1 \prod_t dF_t^M \\ &= \sigma^2 \int \cdots \int \frac{1}{p^2} (p + v^2 p^2) \prod_t dF_t^M \\ &= \sigma^2 \int \cdots \int \frac{k}{p} \prod_t dF_t^M, \end{aligned}$$

and, hence,

$$(71) \quad \sigma_{\tilde{q}}^2 = \sigma^2 \int \cdots \int \frac{k}{p} \prod_t dF_t^M - \chi^2.$$

¹In fact it was shown (cf. section 3.1.1.2) to be the best *absolutely unbiased* estimate.

Comparing this with (67) we see that in order to prove (62) we must show that

$$(72) \quad \int \cdots \int \frac{k}{p} \prod_t dF_t^M > \left\{ \int \cdots \int \frac{p}{k} \prod_t dF_t^M \right\}^{-1}.$$

But (72) follows from a multivariate analogue of a well-known inequality.¹ Thus q is actually better than \tilde{q} .

(b) From a practical viewpoint, (62) offers little consolation since q cannot be computed unless the F_t^M 's and v^2 are known, and v^2 is not likely to be known.

Unless an unbiased linear estimate not requiring the knowledge of v^2 can be found with a variance between σ_q^2 and $\sigma_{\tilde{q}}^2$, the least-squares estimate \tilde{q} is still the best practical solution, although it does not utilize whatever knowledge we may possess about the F_t^M 's.

(c) However, it is of considerable value to have the expression for σ_q^2 since a lower bound is obtained for the variance of all the linear unbiased estimates of χ . If the use of a priori information concerning v^2 is permitted, σ_q^2 is obviously the greatest lower bound. If the use of v^2 is not permitted, σ_q^2 may or may not be the greatest lower bound, but we can say at least that the greatest lower bound does not exceed² σ_q^2 and is not less than σ_q^2 . It would, of course, be of interest to find such a greatest lower bound for the class of all (linear and nonlinear) unbiased estimates.

3.1.2. (*Best unbiased linear and least-squares estimates.*)
Lags present. In (1) some of the y 's may be lagged values of other y 's. Then the right-hand member of (20) might contain lagged values of y_0 .

We shall consider the simplest case of this type, where

$$(73) \quad y_{1t} \equiv y_{0,t-1}$$

¹[Hardy, Littlewood, Pólya; pp. 150-151, Theorem 204]. (72) would become an equality only if $(1 + \rho^2 p)/p = \text{const.}$

²It would seem that asymptotically σ_q^2 is the greatest lower bound because of the maximum-likelihood properties of \tilde{q} (see section 3.2), and their thereby implied efficiency. Exceptions might be due to the existence of estimates whose asymptotic distribution is nonnormal.

and there are no other y 's or z 's present. In this case (20) and (21) are reduced to

$$(74) \quad \begin{aligned} \mathcal{E}(y_t | y_{t-1}) &\equiv \mathcal{E}(Y_t | Y_{t-1} = y_{t-1}) = \chi y_{t-1}, \\ \sigma^2(y_t | y_{t-1}) &= \text{const.}, \end{aligned}$$

where y_t is written for y_{0t} and the affixes of χ are dropped.

The sample available is that of T observations on y , say

$$(75) \quad (y_1, \dots, y_T).$$

When χ is being estimated by the least-squares method under the assumption (74), it is considered that $T - 1$ pairs of observations of the type (y_t, y_{t-1}) are given, with the first element in the pair as the "dependent" variable. The estimate \tilde{q} of χ obtained from the least-squares principle is

$$(76) \quad \tilde{q} = \frac{\sum_{t=2}^T y_t y_{t-1}}{\sum_{t=2}^T y_{t-1}^2}.$$

The question to be answered is whether \tilde{q} is a best (conditionally or absolutely) unbiased estimate. The complete solution of this problem is not known. So far there is no known class of cumulative distribution functions of (y_1, \dots, y_T) satisfying (74) with regard to which \tilde{q} is (conditionally) unbiased. On the other hand, an important class of cumulative distribution functions for which (74) holds but where \tilde{q} is biased has been found¹. Thus it is certain the \tilde{q} is not in any case an *absolutely unbiased estimate* of χ , and a fortiori not a *best* absolutely unbiased estimate.

Several questions of interest remain unanswered. Are there cumulative distribution functions satisfying (74) with regard to which \tilde{q} is (conditionally) unbiased? What is the best (linear or nonlinear) unbiased² estimate of χ ? (Asymptotically, \tilde{q} probably is a best unbiased estimate due to its maximum-likelihood properties.)

¹However, this bias tends to zero, though not very rapidly, as the sample size increases. Cf. [XV].

²It is shown in [XV] that an unbiased estimate of χ does exist.

3.2. Equivalence of Maximum-Likelihood and Least-Squares Estimates of the Regression Coefficients

3.2.1. Lags absent.

3.2.1.1. The results given in this section are not in any way new, but there does not seem to be an easily accessible source treating the general case and giving the necessary proofs.

As is well-known, the maximum-likelihood method leads to least-squares estimates only when a normal universe is assumed and this assumption will be implicit in all that follows.

Thus, consider a set Y of jointly normally distributed variates [Y given in (1)] with the following properties:

$$(77) \quad \begin{aligned} \mathcal{E}(y_{gt}) &= \sum_{j=1}^{K_z} \pi_{gj} z_{jt} & (g = 0, 1, \dots, K_y), \\ \text{cov}(y_{gt'}, y_{ht''}) &= \omega_{gh} \delta_{t't''} & (g, h = 0, 1, \dots, K_y), \end{aligned}$$

where the z 's are (observed) fixed variates, the π 's and ω 's are unknown parameters¹, and $\delta_{t't''}$ the Kronecker symbol. Then, clearly,

$$(78) \quad \mathcal{E}(Y_{0t} | Y^{**}; Z) = \sum_{i=1}^{K_y} \chi_i^y y_{it} + \sum_{j=1}^{K_z} \chi_j^z z_{jt}$$

(which is identical with (20) where $\chi_0 = 0$) and the regression, as is well-known, is homoscedastic.

The regression coefficients χ are related to the parameters in (77) by

$$(79) \quad \chi_i^y = - \frac{\omega^{0i}}{\omega^{00}},$$

$$(80) \quad \chi_j^z = \sum_g \pi_{gj} \chi_g^y.$$

¹The problem of possible restrictions on the π 's and ω 's is discussed below in section 3.3.

3.2.1.2. It is now to be shown that the least-squares estimates \tilde{q} and the maximum-likelihood estimates $\hat{\chi}$ of the χ 's are identical.¹ The proof will be divided into two parts: first it will be shown that $\tilde{q}^y = \hat{\chi}^y$, and then that $\tilde{q}^z = \hat{\chi}^z$. In each case we first derive the least-squares estimates and the the maximum-likelihood estimates.

I. To show that $\tilde{q}^y = \hat{\chi}^y$, we first obtain the least-squares estimates by minimizing with respect to the q 's the expression S given in (37).

Setting

$$(81) \quad q_0^y = -1,$$

we may rewrite S as

$$(82) \quad S = \sum_{t=1}^T \left(\sum_{g=0}^{K_y} q_g^y y_{gt} + \sum_{j=0}^{K_z} q_j^z z_{jt} \right)^2$$

or, for short

$$(83) \quad S = \sum_t \left(\sum^y q^y y + \sum^z q^z z \right)^2$$

where $\sum^y \equiv \sum_0^{K_y}$, and $\sum^z \equiv \sum_1^{K_z}$.

Differentiating S with respect to the q 's we obtain

$$(84) \quad \frac{1}{z} \frac{\partial S}{\partial q_i^y} = \sum_t \left[\left(\sum_h^y q_h^y y_h + \sum_j^z q_j^z z_j \right) y_i \right] = 0$$

($i = 1, 2, \dots, K_y$)

and

$$(85) \quad \frac{1}{z} \frac{\partial S}{\partial q_k^z} = \sum_t \left[\left(\sum_h^y q_h^y y_h + \sum_j^z q_j^z z_l \right) z_l \right] = 0$$

($l = 1, 2, \dots, K_z$).

¹The proof given in the text is a reproduction of the one originally presented at the Cowles Commission Conference in February, 1945. Following this presentation Koopmans suggested an alternative proof [VII], which is much simpler, more elegant, and more easily generalized.

Using notation similar to that of [II] for sample moments, we write

$$\begin{aligned}
 m_{y_g y_h} &= \frac{1}{T} \sum_t y_{gt} y_{ht} , \\
 (86) \quad m_{y_g z_j} &= \frac{1}{T} \sum_t y_{gt} z_{jt} , \\
 m_{z_j z_l} &= \frac{1}{T} \sum_t z_{jt} z_{lt} ,
 \end{aligned}$$

where T is the size of the sample. We may then rewrite equations (84) and (85), thus obtaining

$$(87) \quad \sum_h^y m_{y_i y_h} q_h^y + \sum_j^z m_{y_i z_j} q_j^z = 0 \quad (i = 1, 2, \dots, K_y),$$

$$(88) \quad \sum_h^y m_{z_l y_h} q_h^y + \sum_j^z m_{z_l z_j} q_j^z = 0 \quad (l = 1, 2, \dots, K_z).$$

Now write

$$(89) \quad \sum_h^y m_{y_0 y_h} q_h^y + \sum_j^z m_{y_0 z_j} q_j^z = \varphi,$$

where the value of φ is immaterial, and consider (89) as the first equation of a system consisting of (87) and (89). Then (87) with (89) may be written in matrix form as

$$(90) \quad M_{yy} q^y + M_{yz} q^z = \varphi e ,$$

where $e = [1 \ 0 \ \dots \ 0]$. Similarly, (88) may be written

$$(91) \quad M_{zy} q^y + M_{zz} q^z = 0 ,$$

where $M_{zy} = M_{yz}'$ is the transpose of M_{yz} , and $0 = [0 \ 0 \ \dots \ 0]$.

It follows from (91) that

$$(92) \quad q^z = -M_{zz}^{-1} M_{zy} q^y,$$

and substituting this value into (90) we obtain

$$(93) \quad M_{yy} q^y - M_{yz} M_{zz}^{-1} M_{zy} q^z = \varphi e,$$

or, finally,

$$(94) \quad q^y = \tilde{q}^y = (M_{yy} - M_{yz} M_{zz}^{-1} M_{zy})^{-1} \varphi e.$$

Now, in obtaining the maximum-likelihood estimates, we see from the elementary properties of determinants that

$$(95) \quad \sum_h^y w_{gh} w^{0h} = \delta_{0g} |W| \quad (g = 0, 1, \dots, K_y)$$

where w_{gh} is the maximum-likelihood estimate of w_{gh} . Dividing by w^{00} , we have

$$(96) \quad \sum_h^y w_{gh} = \delta_{0g} \frac{|W|}{w^{00}}.$$

Thus, by virtue of (80)¹ and setting $q_0^y = -1$, we have

$$(97) \quad \sum_h^y w_{gh} \hat{\chi}_h^y = -\delta_{0g} \frac{|W|}{w^{00}},$$

which may be written in matrix form as

$$(98) \quad W \hat{\chi}^y = \psi e,$$

where $\psi = -\frac{|W|}{w^{00}}$ and $e = [1 \ 0 \ \dots \ 0]$ as before.

It will be noted that the value of ψ does not affect $\hat{\chi}^y$. Now, as has been shown in [II],

¹This operation is justified since (79) and (80) hold for the maximum-likelihood estimates as well as for true values of the parameters.

$$(99) \quad W = M_{yy} - M_{yz} M_{zz}^{-1} M_{zy}$$

so that (98) becomes

$$(100) \quad (M_{yy} - M_{yz} M_{zz}^{-1} M_{zy}) \hat{\chi}^z = \psi e$$

or

$$(101) \quad \hat{\chi}^y = (M_{yy} - M_{yz} M_{zz}^{-1} M_{zy})^{-1} \psi e,$$

which implies

$$(102) \quad \tilde{q}^y = \hat{\chi}^y$$

since the expressions in (101) and (94) are identical (except possibly for φ and ψ which are without effect on \tilde{q}^y and $\hat{\chi}^y$).

II. To show that $\tilde{q}^z = \hat{\chi}^z$, we first note, as stated in (92), that

$$(103) \quad \tilde{q}^z = -M_{zz} M_{zy} \tilde{q}^y.$$

Since, as is shown in [II], the maximum-likelihood estimate P of π is given by

$$(104) \quad P = -M_{yz} M_{zz}^{-1},$$

the maximum-likelihood estimate $\hat{\chi}^z$ of χ^z is obtained by substituting P from (104) and $\hat{\chi}^y$ from (101) into the sample counterpart of (80), thus obtaining

$$(105) \quad \hat{\chi}^z = -M_{zz}^{-1} M_{zy} \hat{\chi}^y,$$

with the symmetry of M_{zz} taken into account.

Hence, by virtue of (102),

$$(106) \quad \hat{\chi}^z = -M_{zz}^{-1} M_{zy} \tilde{q}^y,$$

where the right-hand member is identical with that in (103). Thus

$$(107) \quad \tilde{q}^z = \hat{\chi}^z,$$

and the second part of the proof is completed.

3.2.2. *Lags present.* No distinction is made in the derivation of the maximum-likelihood (or least-squares) estimates between "fixed variates" and lagged values of the endogenous variables Y . Hence we may regard the z 's in section 3.2 as being "predetermined"¹ (i.e., including lagged y 's). The stochastic difference equations are covered by the above proof.

The only word of warning concerns the nature of the *initial observation* or observations. An example will clarify this point. Let there be only one variable y and a sample

$$(108) \quad (y_1, \dots, y_T).$$

We assume that y_t satisfies a stochastic difference equation of, say, first order:

$$(109) \quad y_t = \chi y_{t-1} + u_t,$$

where the u 's are normally independently distributed with a zero mean and a constant variance.

To specify the problem completely, we must still provide information concerning the (marginal) distribution of y_1 (i.e., the initial observation). Now if y_1 is assumed to be a fixed variate, the maximum-likelihood estimate of χ is identical with the least-squares estimate. But if the (marginal) distribution of y_1 equals the (marginal) distribution of any other y_t (i.e., y_1 has a zero mean and a variance $\sigma_{y_1}^2 = \sigma_u^2 / (1 - \chi^2)$), the resulting maximum-likelihood estimate of χ is different from that obtained from least-squares.² This discrepancy will, in general, occur whenever the distribution of y_1 depends on χ .

In most economic problems it is the writer's opinion that it is not realistic to consider y_1 as a fixed variate; the only merit this assumption seems to have is that of expediency in that it leads to least-squares estimates.³

¹See section 1.1.

²[Koopmans, 1942].

³As the sample size increases, the maximum-likelihood estimate obtained for the case where y_1 is not a fixed variate will tend stochastically to the least-squares estimate.

3.3. Regression Estimates under A Priori Restrictions

The proofs of the optimal properties of the least-squares estimates \tilde{q} of the regression coefficients χ were implicitly based on the following assumption: the a priori restrictions (if any) on the nature of the structure¹ S do not imply any restriction on the nature of the distribution F of the observed variates.²

Now consider the case where the a priori restrictions on the nature of the structure S are such that F itself is subject to restrictions.

$$(110) \quad F \equiv F_{*|**}^C(y^* | y^{**}; z; \chi, \sigma^2) F_{**}^M(y^{**}; z; \theta^M).$$

(The regression in $F_{*|**}$ is assumed linear in y^{**} and z , with regression coefficients χ , and homoscedastic with variance σ^2 .)

We may then distinguish three cases of restrictions³ on F (implied by the a priori restrictions on S):

The restrictions on F do not involve χ :

$$(111) \quad \psi(\chi)(\sigma^2; \theta^M) = 0.$$

The optimal properties of \tilde{q} are not affected.

The restrictions on F involve χ only:

$$(112) \quad \psi_\chi(\chi) = 0.$$

Then it is easily shown that the least-squares method, where minimization is carried out subject to (112), gives results equivalent to both best unbiased linear and maximum-likelihood estimates when the latter also take (112) into account (whenever the equivalence

¹In the linear normal case $S = (A, \Sigma^u)$ where A is the structural coefficient matrix and Σ^u the covariance matrix of the disturbances. As an example of restrictions on S one may mention the specification of zeros in A or Σ^u .

²In the linear normal case, normality of F is postulated, but the first two moments remain unrestricted.

³The discussion in the remainder of this section treats restrictions in the form of equalities only; the extension to inequalities is straightforward.

would exist in the absence of restrictions).

Restrictions involve both the regression coefficients, and some other parameters of F (σ^2 or θ^M or both); they may then be written as

$$(113) \quad \psi(\chi, \sigma^2; \theta^M) = 0.$$

In this case the equivalence of least-squares with best unbiased linear and maximum-likelihood estimates, in general, no longer holds.

It may be remarked that in the latter two cases, if the existing restrictions are ignored, the unbiasedness or consistency is not removed, but the variance of the estimates is no longer minimal.

VII. THE EQUIVALENCE OF MAXIMUM-LIKELIHOOD AND LEAST-SQUARES ESTIMATES OF REGRESSION COEFFICIENTS

BY T. C. KOOPMANS

The equivalence of maximum-likelihood and least-squares estimates, under the conditions assumed by Hurwicz [VI-3.2], can also be established on the basis of concepts and propositions developed in other contributions to this volume. The equation¹ {(78)} of which we desire to estimate the coefficients may be combined with the equations {(77)} with $g = 1, \dots, K_y$ (i.e., omitting the first equation with $g = 0$) to give a system

(1)

$$(1.0) \quad y_0 - \chi_1^y y_{1t} - \dots - \chi_{K_y}^y y_{K_y t} - \chi_1^z z_{1t} - \dots - \chi_{K_z}^z z_{K_z t} = v_{0t}^*$$

$$(1.1) \quad y_{1t} - \pi_{11} z_{1t} - \dots - \pi_{1K_z} z_{K_z t} = v_{1t},$$

⋮
⋮
⋮

$$(1.k') \quad y_{K_y t} - \pi_{K_y 1} z_{1t} - \dots - \pi_{K_y K_z} z_{K_z t} = v_{K_y t},$$

in which the z_1, \dots, z_{K_z} are fixed variables, and the $v_{0t}^*, v_{1t}, \dots, v_{K_y t}$ have the same joint normal distribution for each value of t . In particular, since the first equation (1.0) in (1) represents the regression {(78)} of one variable (y_0) in a normal system on the others (y_1, \dots, y_{K_y}), the quantity

$$(2) \quad v_0^* = y_0 - \mathcal{E}(y_0 | y_1, \dots, y_{K_y}; z_1, \dots, z_{K_z})$$

is distributed independently of y_1, \dots, y_{K_y} , and therefore also

¹Equation numbers in braces { } refer to formulae in article [VI] by Hurwicz.

independently of v_1, \dots, v_{K_y} ,

$$(3) \quad \mathcal{E} v_0^* v_k = 0, \quad k = 1, \dots, K_y.$$

In fact, this condition alone uniquely identifies¹ the equation (1.0) within the system (1), given the fact that each of the remaining equations (1.1), ..., (1.k') is uniquely identified by the ("reduced"²) form prescribed for it.

Because of the independence of v_0^* from v_1, \dots, v_{K_y} , and because of the special form of the matrix of coefficients of y_0, y_1, \dots, y_{K_y} in (1), the system (1) meets the conditions, stated elsewhere,³ for a factorization

$$(4) \quad F = F_1(\chi_i^y, \chi_j^z, \sigma_{v_0^*}; y_{it}, z_{jt}) \\ \times F_2\{\pi_{gk} (g \geq 1), \omega_{gh} (g, h \geq 1), y_{it}, z_{it}\}$$

of the likelihood function.⁴ Here each factor depends, besides on the observed variables, on a different set of parameters (as indicated), relating to the equation (1.0) and to the set of equations (1.1), ..., (1.K_y), respectively. The two factors can therefore be maximized separately. In particular, the variables y_1, \dots, y_{K_y} can in equation (1.0) be regarded as "predetermined"⁵ by the equations (1.1), ..., (1.K_y). That is, for purposes of maximizing the first factor

$$(5) \quad F_1 = \text{constant} \times \sigma_{v_0^*}^{-1} \exp - \frac{1}{2 \sigma_{v_0^*}^2} \sum_t (v_{0t}^*)^2$$

of the likelihood function, where v_{0t}^* is given by (1.0), these variables can be treated in (1.0) as if they were fixed variables. In

¹See [II, Definition 2.4.4.2].

²See [II-3.1.6].

³See [XVII-3 & 6].

⁴ F here denotes the probability density in the sample space, not the cumulative distribution function.

⁵See [XVII-6].

view of the form of (5), this establishes the equivalence of maximum-likelihood and least-squares procedures with regard to the estimation of the parameters of (1.0).

To complete the link with other articles in this volume, the question should be raised as to what extent the equivalence just established remains valid if the system (1) is obtained by linear transformation from a system of structural equations as discussed in [II], the parameters of which are subject to a number of a priori restrictions, and if the likelihood function is maximized subject to those restrictions.

It is seen from formula (3.13) in [II] that the parameters of the second factor F_2 in the likelihood function (4) depend on the structural parameters through

$$(6) \quad \Pi = -B^{-1}\Gamma, \quad \Omega = (B^{-1}\Sigma B'^{-1})^{-1}.$$

It should be noted that the elements in the zero row of Π and those in the zero row (and column) of Ω do not enter in F_2 . On the other hand, the parameters in F_1 , according to {(77)} and (1), satisfy

$$(7) \quad \chi_i^y = -\frac{\omega^{0i}}{\omega^{00}}, \quad \chi_j^z = -\sum_{g=0}^{K_y} \chi_g^y \pi_{gj},$$

$$\sigma_{v_0}^2 = \sum_{g,h=0}^{K_y} \chi_g^y \omega_{gh} \chi_h^y,$$

where $\chi_0^y = -1$. It follows that the parameters in F_1 can be given any arbitrarily chosen values by suitable choice of the elements in the zero rows in Π and Ω . The parameters in F_1 and those in F_2 are, therefore, truly independent of each other.

This remains true under such restrictions on the structural parameters A, Σ as do not entail any restrictions on the possible values of Π and Ω . In the terminology of [II-3.2], the equivalence of maximum-likelihood and least-squares procedures remains valid under a priori restrictions that do not depress the likelihood function below its absolute maximum, because both the likelihood function¹ and the parameters of (1.0) depend on A, Σ only through the parameters Π and Ω which in that case are unrestricted.

¹See either (4) and (7) above or formula (3.15) of [II].

On the other hand, maximization of the likelihood function under restrictions on A, Σ that do not permit the likelihood function to attain its absolute maximum, can only "by accident" be equivalent to the least-squares method, depending on how the vectors χ^y, χ^z of regression coefficients are affected by those restrictions (as explained in [VI-3.3]).

VIII. REMARKS ON THE ESTIMATION OF UNKNOWN PARAMETERS IN INCOMPLETE SYSTEMS OF EQUATIONS

BY A. WALD

	Page
1. Formulation of the Problem	305
2. Existence of Consistent Estimates	306
3. Construction of Confidence Regions for the Parameter Point α	307
4. Some Examples	308

1. Formulation of the Problem

Let x_1, \dots, x_H be N variables and denote by x_{nt} the value of x_n at the time point t where t can take any integral value. It is assumed that these variables satisfy a system of H ($H < N$) stochastic equations

$$\begin{aligned}
 & f_h(x_{1t}, x_{1, t-1}, \dots, x_{1, t-\tau^\square}; \dots; \\
 (1) \quad & x_{Ht}, x_{H, t-1}, \dots, x_{H, t-\tau^\square}; \alpha_1, \dots, \alpha_P) = u_{ht}, \\
 & \qquad \qquad \qquad (h = 1, \dots, H)
 \end{aligned}$$

where $\alpha_1, \dots, \alpha_P$ are unknown parameters, f_h is a given function of the variables $x_{n, t-\tau}$ ($n = 1, \dots, N; \tau = 0, 1, \dots, \tau^\square$) and of the parameters α_ρ ($\rho = 1, \dots, P$), and u_{1t}, \dots, u_{Ht} are random variables. It is assumed that the distribution of the random vector $u_t = (u_{1t}, \dots, u_{Ht})$ is independent of t and that the vectors u_1, u_2, \dots , are independently distributed. It is also assumed that the unknown distribution function of u_t is known to be an element of a given finite-parameter family Ω of cumulative distribution functions. Denote by $\theta_1, \dots, \theta_Q$ the unknown parameters involved in the distribution function of u_t .

The problem considered here is to estimate all or some of the unknown parameters $\alpha_1, \dots, \alpha_p; \theta_1, \dots, \theta_\rho$ on the basis of the observed values x_{nt} ($n = 1, \dots, N; t = 1, \dots, T$).

Since $H < N$, the distribution of the random vectors $u_t = (u_{1t}, \dots, u_{Ht})$, $t = 1, \dots, T$, does not determine the distribution of the variables x_{nt} , $n = 1, \dots, N$, but only imposes a restriction on the latter distribution. For this reason, the title of this note refers to incomplete systems of equations. It will appear from what follows that the estimation problems in incomplete systems are essentially different from those in complete systems discussed in other contributions to this volume.

2. Existence of Consistent Estimates

A basic problem is, of course, the question of the existence of consistent estimates of the unknown parameters. We shall say that a parameter point $\alpha = (\alpha_1, \dots, \alpha_p)$ satisfies condition C if the following conditions are simultaneously fulfilled: (1) the distribution of the vector

$$u_t(\alpha) = \left[f_1(x_{1t}, \dots, x_{1,t-\tau^{\square}}; \dots; x_{Nt}, \dots, x_{N,t-\tau^{\square}}; \alpha_1, \dots, \alpha_p), \right. \\ \left. \dots, f_H(x_{1t}, \dots, x_{1,t-\tau^{\square}}; \dots; x_{Nt}, \dots, x_{N,t-\tau^{\square}}) \right]$$

is independent of t ; (2) the vectors $u_1(\alpha), u_2(\alpha), \dots$, are independently distributed; (3) the distribution of $u_t(\alpha)$ is an element of Ω . It is clear that if there exist two parameter points α' and α'' such that both satisfy condition C , no consistent estimate of the parameter point α exists. On the other hand, if there exists one and only one parameter point α that satisfies condition C , then, under some further mild restrictions that we do not propose to discuss here, a consistent estimate of α will exist.

Whether or not there exist several parameter points α satisfying condition C depends on the joint probability distribution of the observable variables x_{nt} ($n = 1, \dots, N; t = 1, \dots, T$). Thus, the existence of consistent estimates depends on the joint probability distribution of the observable variables. Since this distribution is usually unknown a priori, we cannot be sure that a consistent estimate exists. This difficulty, however, is not as serious as it would appear at first sight. In fact, instead of

point estimates, we are usually more interested in constructing a confidence region for the unknown parameters corresponding to a given confidence coefficient. We shall see in the next section that a confidence region can be obtained irrespective of the existence of consistent estimates. The only effect of the nonexistence of consistent estimates on the confidence region is that the diameter of the confidence region (maximum distance between two points of the confidence region) will not approach zero as the number of observations approaches infinity.

3. Construction of Confidence Regions for the Parameter Point α

To construct a confidence region for the parameter point α we may proceed as follows: Suppose that the prescribed confidence coefficient is δ . For any given parameter point $\alpha^0 = (\alpha_1^0, \dots, \alpha_p^0)$ we construct a critical region $W(\alpha^0)$ of size $1 - \delta$ for testing the hypothesis that $\alpha = \alpha^0$. The critical region $W(\alpha^0)$ is a subset in the HT -dimensional space of the variables u_{ht} ($h = 1, \dots, H$; $t = 1, \dots, T$) and we reject the hypothesis that $\alpha = \alpha^0$ if and only if the observed point $u_{11}(\alpha^0), u_{12}(\alpha^0), \dots, u_{Ht}(\alpha^0)$ falls in $W(\alpha^0)$. Of course, $W(\alpha^0)$ is to be constructed in such a way that the probability measure of $W(\alpha^0)$ is equal to $1 - \delta$ under any distribution of u_t that is an element of Ω . Usually such a region $W(\alpha^0)$ can be constructed without any difficulty. Now consider all possible parameter points α^0 and the corresponding family of critical regions $W(\alpha^0)$. The set of all parameter points α^0 that are not rejected on the basis of this test procedure will form a confidence set with confidence coefficient δ .

Thus, we see that there is no serious difficulty in obtaining a confidence region for α . There will be, in general, infinitely many possible confidence regions, since $W(\alpha^0)$ can usually be chosen in infinitely many different ways. The problem of proper choice of $W(\alpha^0)$ is not yet solved. The theory of confidence intervals as developed by J. Neyman [1937] does not apply, since Neyman deals with the parametric case, while in our problem the class of all possible distribution functions of the observable variables x_{nt} ($n = 1, \dots, N$; $t = 1, \dots, T$) compatible with the relations (1) cannot be described by a finite number of parameters owing to the fact that $H < N$.

4. Some Examples

In this section we shall give a few examples to illustrate the procedure for the construction of confidence regions. It should be emphasized, however, that the confidence regions described later in this section are by no means "best." As a matter of fact they may be very inefficient under certain conditions and it is not our intention to recommend them for practical use.

Consider the following problem: Let x_1 , x_2 , and x_3 be three observable variables and denote by x_{nt} ($n = 1, 2, 3; t = 1, 2, \dots$) the value of x_n at the time point t . Suppose that x_{1t} , x_{2t} , and x_{3t} satisfy the relation

$$(2) \quad x_{1t} + \alpha_2 x_{2t} + \alpha_3 x_{3t} + \alpha_0 = u_t,$$

where α_2 , α_3 , α_0 are unknown constants and u_1, u_2, \dots are independently and normally distributed random variables with zero means and a common variance σ^2 . We shall construct a confidence region for the parameters α_2 , α_3 , and α_0 on the basis of the observed values x_{nt} ($n = 1, 2, 3; t = 1, \dots, T$). For simplicity we shall assume that $T = 3V$, where V is a positive integer. Denote by \bar{x}_n^m the arithmetic mean of the observations on x_n in the m th group of V observations, i.e.,

$$(3) \quad \bar{x}_n^m = \frac{1}{V} \sum_{t=(m-1)V+1}^{mV} x_{nt}, \quad m = 1, 2, 3.$$

Furthermore, let

$$(4) \quad s_m^2 = \frac{1}{V-1} \sum_{t=(m-1)V+1}^{mV} \left[(x_{1t} - \bar{x}_1^m) + \alpha_2(x_{2t} - \bar{x}_2^m) + \alpha_3(x_{3t} - \bar{x}_3^m) \right]^2, \quad m = 1, 2, 3.$$

It is clear that $(V-1)s_1^2/\sigma^2$, $(V-1)s_2^2/\sigma^2$, and $(V-1)s_3^2/\sigma^2$ are independently distributed, each having the χ^2 -distribution with $V-1$ degrees of freedom. Thus, each of the expressions

$$(5) \quad t_m = \frac{(\bar{x}_1^m + \alpha_2 \bar{x}_2^m + \alpha_3 \bar{x}_3^m + \alpha_0) V^{\frac{1}{2}}}{s_m}, \quad m = 1, 2, 3,$$

has the t -distribution with $V - 1$ degrees of freedom, and t_1 , t_2 , and t_3 are independently distributed.

A confidence region for $(\alpha_2, \alpha_3, \alpha_0)$ with confidence coefficient δ can be constructed as follows: Let λ_δ be the value for which the probability that $-\lambda_\delta \leq t \leq \lambda_\delta$ is equal to $\delta^{1/3}$. Here t denotes a random variable that has the t -distribution with $V - 1$ degrees of freedom. The set of all parameter points $(\alpha_2, \alpha_3, \alpha_0)$ that satisfy simultaneously the following three inequalities,

$$(6) \quad |t_m| \leq \lambda_\delta, \quad m = 1, 2, 3,$$

forms a confidence region with confidence coefficient δ .

The confidence region given by (6) is certainly far from being the best possible one. There is a loss of efficiency in taking three different estimates s_1^2 , s_2^2 , s_3^2 for σ^2 . A better procedure would be to combine these three estimates into a single one given by $s^2 = \frac{1}{3}(s_1^2 + s_2^2 + s_3^2)$. This, however, would make the variables t_1 , t_2 , and t_3 dependent and would complicate the derivation of a confidence region.

If the determinant

$$(7) \quad \begin{vmatrix} \bar{x}_2^1 & \bar{x}_3^1 & 1 \\ \bar{x}_2^2 & \bar{x}_3^2 & 1 \\ \bar{x}_2^3 & \bar{x}_3^3 & 1 \end{vmatrix}$$

is bounded away from zero as $T \rightarrow \infty$, and if s_1^2 , s_2^2 , and s_3^2 are bounded functions of T , the diameter of the confidence region given by (6) will approach zero as $T \rightarrow \infty$. If the distribution of the vector $v_t = (x_{1t}, x_{2t}, x_{3t})$ is independent of t , if v_1, v_2, \dots are independently distributed, and if the first two moments of x_{nt} exist, then the determinant (7) will converge stochastically to zero and the probability is one that the diameter of the confidence

set given by (6) will not approach zero as $T \rightarrow \infty$.

Another possible procedure for obtaining a confidence region is the following: Consider the expression

$$(8) \quad F = \frac{V \sum_{m=1}^3 (\bar{x}_1^m + \alpha_2 \bar{x}_2^m + \alpha_3 \bar{x}_3^m + \alpha_0)^2}{(V-1)(s_1^2 + s_2^2 + s_3^2)} \frac{T-3}{3}.$$

This has the F -distribution with 3 and $T-3$ degrees of freedom. Let F_δ be a positive value for which the probability that $F \leq F_\delta$ is δ . Then the set of all parameter points $\alpha_2, \alpha_3, \alpha_0$ that satisfy the inequality

$$(9) \quad F \leq F_\delta$$

will form a confidence set with confidence coefficient δ .

A confidence region for $\alpha_2, \alpha_3, \alpha_0$ can also be obtained as follows: Consider the serial correlation

$$(10) \quad r = \frac{\sum_{t=1}^{T-1} \left[(x_{1t} + \alpha_2 x_{2t} + \alpha_3 x_{3t} + \alpha_0) \times (x_{1,t+1} + \alpha_2 x_{2,t+1} + \alpha_3 x_{3,t+1} + \alpha_0) \right]}{\sum_{t=1}^T (x_{1t} + \alpha_2 x_{2t} + \alpha_3 x_{3t} + \alpha_0)^2}$$

and let r_δ be a value such that the probability that $|r| \leq r_\delta$ is equal to δ . Then the set of all parameter points $(\alpha_2, \alpha_3, \alpha_0)$ for which $|r| \leq r_\delta$ will form a confidence set with confidence coefficient δ .

IX. ESTIMATION OF THE PARAMETERS OF A SINGLE
EQUATION BY THE LIMITED-INFORMATION
MAXIMUM-LIKELIHOOD METHOD

BY T. W. ANDERSON

	Page
1. The Purpose of the Method	311
2. Derivation of the Coefficients from the Population Reduced Form .	312
3. The Maximum-Likelihood Estimates	314
4. Asymptotic Properties of the Estimates	
4.1. Consistency	316
4.2. Asymptotic Normality	317
5. Small-Sample Confidence Regions	317
6. Statistical Tests of the Assumption of Zero Coefficients	
6.1. The Likelihood-Ratio Test	320
6.2. A Test Based on Small-Sample Theory	320
7. Advantages and Disadvantages of the Limited-Information Maximum-Likelihood Method	321

1. The Purpose of the Method

In many economic studies the investigator wishes to study only a single equation out of a complete system. For example, an economist may be interested in the demand function for a given product. This interest would involve only one structural relation out of a set of such relations, among which might be a supply equation, an income equation, etc. If a particular equation is to be studied, a complete model must exist that explains the formation of all the variables in this equation that are considered as endogenous. It is possible to set up the complete model explicitly and, if all the equations are linear, estimate all of the coefficients by the methods described elsewhere in this volume.

To obviate this often complicated process a method has been developed which gives consistent estimates of the coefficients of a single equation without requiring estimates of the complete system. This method is known under two names, the "limited-information maximum-likelihood method" or the "reduced form method," each of which refers to an important aspect of the method.

The only information about the system that is required for this method is the knowledge of all the predetermined variables. The restrictions on the system that are imposed in the estimation procedure are that the coefficients of certain variables be zero in the equation estimated. The estimates are then derived on the basis of maximum likelihood. Since there is in general a larger number of effective a priori restrictions (i.e., depressing the likelihood function¹) on the system than the restrictions on a single equation, this method of estimation is not as efficient as the maximum-likelihood method using all a priori restrictions on the complete system.

This article describes the estimates and gives some of their statistical properties. A full mathematical treatment is given in [T. W. Anderson and Rubin].

2. Derivation of the Coefficients from the Population Reduced Form

Let the equation in which we are interested be

$$(1) \quad \beta y_t' + \gamma z_t' = u_t,$$

where y_t is a row vector of H ($\leq G$) jointly dependent variables², z_t is a row vector of all the predetermined variables in the complete system (K in number), and u_t is the random disturbance associated with this equation. Suppose that there are K^{**} coordinates of γ that are known and prescribed to be zero. Then (1) can be written as

$$(2) \quad \beta y_t' + \gamma^* z_t^* = u_t,$$

where z_t^* has K^* ($= K - K^{**}$) of the coordinates of z_t . Here γ^* is the vector of the K^* coefficients of the coordinates of z_t^* . We shall assume that (2) is identified by the zero coefficients prescribed through the omission of z_t^{**} and possibly certain jointly dependent variables which are not coordinates of y_t (if $H < G$).

¹See [II-3.2.1].

²For the definition of *jointly dependent* and *predetermined* variables see [II-1.8].

This condition implies that¹ $K^{**} \geq H - 1$.

The reduced form of the complete system is the system written in the form of a regression on the predetermined variables. The part of the reduced form² of the complete system that involves the H coordinates of y_t in (2) is

$$(3) \quad y_t' = \Pi^* z_t^{*'} + \Pi^{**} z_t^{**'} + v_t',$$

where v_t is a row vector of disturbances. If we premultiply (3) by β we obtain

$$(4) \quad \beta y_t' = \beta \Pi^* z_t^{*'} + \beta \Pi^{**} z_t^{**'} + \beta v_t'.$$

Since (4) must be identical with (2), we have

$$(5) \quad \gamma^* = -\beta \Pi^*,$$

$$(6) \quad 0 = \beta \Pi^{**}.$$

If β is unknown, but Π^{**} is known, β can be found from (6) except for a factor of proportionality. In turn γ^* can be deduced from (5). The condition³ that (2) be identified by zeros is that (6) can be solved for β . In terms of Π^{**} the condition is that its rank be $H - 1$.

If Ω , the matrix of variances and covariances of the v_t' of (3), is known, the formula

$$(7) \quad \sigma^2 = \beta \Omega \beta'$$

gives the variance of u_t . The normalization of β can be written as

$$(8) \quad \beta \Xi \beta' = 1,$$

where Ξ may be Ω or may be some known constant matrix.

The basic idea of the "reduced-form method" is that if the rel-

¹See the corollary to Theorem 2.2.1 in [II].

²See equation (3.11) in [II]. If $H < G$, the present equations (3) form a subset of equations (3.11) in [II].

³See Theorem 2.2.1 in [II].

evant part of the reduced form (3) is known together with the matrix Ω , the coefficients β and γ^* and the variance σ^2 can be deduced from (5), (6), (7), and (8).

If $K^{**} = H - 1$, then we can solve for β from the sample equivalent of (6). However, when $K^{**} \geq H - 1$, the (unrestricted) estimate of Π^{**} will be of rank H , and then the sample form of (6) cannot be solved for a nontrivial β . In the next section this difficulty is met by obtaining a maximum-likelihood estimate of β (under disregard of a part of the a priori information) which also estimates Π^{**} with rank $H - 1$.

3. The Maximum-Likelihood Estimates

To derive estimates of the vectors β and γ^* on the basis of maximum likelihood we assume that all the disturbances in the system are normally distributed with mean zero and covariance matrix independent of t and of the exogenous variables and that the disturbances are serially independent. Knowledge of all the predetermined variables in the system is assumed. The only restriction on the system that is taken into account is that the coefficients of z^{**} and possibly of some jointly dependent variables (if $H < G$) in the specified structural equation shall be zero. The likelihood function is

$$(9) \quad L = (2\pi)^{-\frac{1}{2}TH} (\det \Omega)^{-\frac{1}{2}T} \\ \times \exp\left\{-\frac{1}{2} \sum_{t=1}^T (y_t - z_t^* \Pi^{*'} - z_t^{**} \Pi^{**'}) \Omega^{-1} (y_t - \Pi^* z_t^{*'} - \Pi^{**} z_t^{**'})\right\}.$$

The restrictions on the parameters are

$$(10) \quad \beta \Pi^{**} = 0,$$

$$(11) \quad \beta' \Xi \beta = 1.$$

To define the estimates of β , γ^* , and σ^2 we need the following statistics:

$$(12) \quad \frac{1}{T} \sum_{t=1}^T y_t' z_t = M_{yz} = [M_{yz^*} \quad M_{yz^{**}}],$$

the matrix of mean cross-moments of y and z , partitioned according to whether or not the coordinates of z are in the equation (2);

$$(13) \quad \frac{1}{T} \sum_{t=1}^T z_t' z_t = M_{zz} = \begin{bmatrix} M_{z^* z^*} & M_{z^* z^{**}} \\ M_{z^{**} z^*} & M_{z^{**} z^{**}} \end{bmatrix},$$

the matrix of mean second-order moments of z ;

$$(14) \quad P = M_{yz} M_{zz}^{-1},$$

the matrix of sample regression coefficients which is partitioned as

$$(15) \quad P = [P^* \quad P^{**}];$$

$$(16) \quad \tilde{P}^* = M_{yz^*} M_{z^* z^*}^{-1},$$

the matrix of sample regressions of y on z^* ;

$$(17) \quad M_{z^{\circ} z^{\circ}} = M_{z^{**} z^{**}} - M_{z^{**} z^*} M_{z^* z^*}^{-1} M_{z^* z^{**}},$$

the mean moments of the residuals (z°) calculated from the regression of z^{**} on z^* ;

$$(18) \quad W = M_{yy} - M_{yz} M_{zz}^{-1} M_{zy},$$

which is T^*/T times the least-squares estimate of Ω (where $T^* = T - K$). Any of these quantities will be underlined to denote multiplication by T (for example $\underline{M}_{z^{\circ} z^{\circ}} = T M_{z^{\circ} z^{\circ}}$).

The second step in obtaining the estimates is to find the smallest root of the determinantal equation

$$(19) \quad | P^{**} M_{z^{\circ} z^{\circ}} P'^{**} - v W | = 0$$

and the corresponding characteristic vector b defined by

$$(20) \quad (P^{**} M_{z^{\circ} z^{\circ}} P'^{**} - v W) b = 0,$$

$$(21) \quad b W b' = 1.$$

The population quantities corresponding to v and b are 0 and $\beta(\beta \Xi \beta') / (\beta \Omega \beta')$.

The estimates of β , γ^* , and σ^2 are

$$(22) \quad \hat{\beta} = \frac{b}{(b \Xi b')^{1/2}},$$

$$(23) \quad \hat{\gamma}^* = -\hat{\beta} \tilde{p}^*,$$

$$(24) \quad \hat{\sigma}^2 = \frac{(1+v)}{(b \Xi b')}.$$

For normalization $\beta \Omega \beta' = 1$, the above equation can be simplified. In this case σ^2 is unity and (22) is

$$(25) \quad \hat{\beta} = \frac{b}{(1+v)^{1/2}}.$$

4. Asymptotic Properties of the Estimates

4.1. *Consistency.* Although the usual theorems concerning maximum-likelihood estimates cannot be applied in general to the estimates of section 3 because successive observations may not be independent, most of the usual asymptotic properties hold. In particular when the assumptions of sections 2 and 3 are true (i.e., assumptions for deriving maximum-likelihood estimates), the estimates are consistent and are asymptotically normally distributed. However, stronger results can be given. For certain weaker conditions on the system, asymptotic normality of the estimates holds, and under still weaker conditions, the estimates are consistent. We shall consider consistency under three sets of conditions where neither normality nor a constant covariance matrix is required.

(i) The estimates are consistent if (a) the stochastic process is stable, (b) the limit in probability of the matrix M_{zz} is non-singular, (c) the ratio of the largest characteristic root of W_{yy}

to the smallest is bounded in probability, and (d)

$$(26) \quad \text{plim}_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T u_t' z_t = 0,$$

where u_t is the row vector of disturbances.

(ii) The estimates of β and γ^* may be consistent even though not all of z^{**} in the system is used. Of course, at least $H-1$ coordinates must be used.

(iii) In certain nonlinear systems, a linear equation (1) can be consistently estimated by this method. The important condition is that the part of the reduced form that is not a linear function of z^* has a population regression on some z^{**} that is of rank $H-1$. There are, of course, certain identified nonlinear systems for which this does not hold.

4.2. *Asymptotic normality.* The requirements for asymptotic normality are stronger. Here, we require of a linear, stable system that the $(4 + \epsilon)$ th moments (for any $\epsilon > 0$) of the disturbances be uniformly bounded in t . In addition there are certain requirements on z .

It should be pointed out that in most cases for the kind of linear systems that an economist is likely to postulate for his problems, these theorems are adequate for the desired asymptotic properties.

5. *Small-Sample Confidence Regions*

If all of the predetermined variables in the system are exogenous (that is, for purposes of obtaining distributions they can be treated as fixed in repeated samples), then confidence regions for the structural coefficients can be found from small-sample theory based on the normality of the disturbance in the specified equation. The regions are based on the idea that given a sequence of observations y_t, z_t^*, z_t^{**} ($t = 1, 2, \dots, T$) from a complete system containing (1), the numbers

$$(27) \quad b y_t' + c z_t^{*'} = w_t$$

will behave like observations from a normal distribution if

$$(28) \quad b = \beta,$$

$$(29) \quad c = \gamma^*.$$

Roughly speaking, if b differs from β and c from γ^* by small amounts, then the sequence $\{w_t\}$ will not be much different from $\{u_t\}$. The regressions of w_t on z_t ,

$$(30) \quad b M_{yz} M_{zz}^{-1} + c (M_{z^*z}) M_{zz}^{-1},$$

are normally distributed if (28) and (29) hold true. Then

$$(31) \quad b \underline{M}_{yz} \underline{M}_{z^{**}z^{**}}^{-1} \underline{M}_{z^0z^0} \underline{M}_{z^{**}z^{**}}^{-1} \underline{M}_{z^{**}y} b' = X_{K^{**}}$$

has the χ^2 -distribution with K^{**} degrees of freedom if

$$(32) \quad \beta \Omega \beta' = 1.$$

Similarly,

$$(33) \quad (b \underline{M}_{yz} + c \underline{M}_{z^*z}) \underline{M}_{zz}^{-1} (b \underline{M}_{yz} + c \underline{M}_{z^*z})' = X_K$$

has the χ^2 -distribution with K degrees of freedom and

$$(34) \quad b \underline{M}_{yy} b' - b \underline{M}_{yz} \underline{M}_{zz}^{-1} \underline{M}_{zy} b' = X_{T^*}$$

has the χ^2 -distribution with T^* degrees of freedom and is independent of either (31) or (33). If the normalization of β is

$$(35) \quad \beta \Xi \beta' = 1$$

for fixed Ξ , the three quantities given above must be multiplied by the reciprocal of $\beta \Omega \beta'$ in order to obtain statistics with χ^2 -distributions.

For normalization (32), a confidence region for β of confidence $\epsilon_1 \epsilon_2$ is given by the vectors b for which

$$(36) \quad X_{K^{**}} \leq \chi_{K^{**}}^2(\epsilon_1),$$

$$(37) \quad \underline{\chi}_{T^*}^2(\varepsilon_2) \leq X_{T^*} \leq \bar{\chi}_{T^*}^2(\varepsilon_2),$$

where the limits are chosen so that the probability of (36) is ε_1 , and that of (37) is ε_2 and

$$(38) \quad \underline{\chi}_{T^*}^2(\varepsilon_2) \leq T^* \leq \bar{\chi}_{T^*}^2(\varepsilon_2).$$

To derive a region for β and γ^* we substitute

$$(39) \quad X_K \leq \chi_K^2(\varepsilon_1)$$

for (36).

If the normalization is (35) a region of size ε for β is (35) and

$$(40) \quad \frac{T^*}{K^{**}} \frac{X_{K^{**}}}{X_{T^*}} \leq F_{K^{**} T^*}(\varepsilon),$$

where $F_{K^{**} T^*}(\varepsilon)$ is chosen from Snedecor's F -table [Snedecor, pp. 88 - 91] so that the probability of (40) is ε . A region for β and γ^* simultaneously is (35) and

$$(41) \quad \frac{T^*}{K} \frac{X_K}{X_{T^*}} \leq F_{K T^*}(\varepsilon).$$

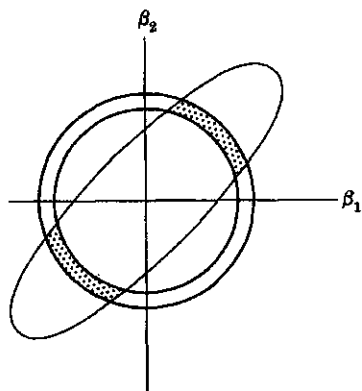


Figure 1.

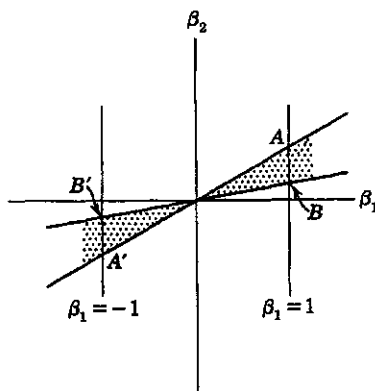


Figure 2.

The first-mentioned region for β of two coordinates is illustrated in Figure 1. The region for β of two coordinates when the normalization is

$$(42) \quad \beta_1^2 = 1$$

is given in Figure 2 as the two segments $A'B'$ and AB .

The regions suggested converge in probability to the parameter points as $T \rightarrow \infty$ in very general systems. One type of condition is that the covariance matrix of successive disturbance vectors be such that the average of the matrices converges.

6. Statistical Tests of the Assumption of Zero Coefficients

6.1. *The likelihood-ratio test.* The hypothesis that the coefficient vector of z^{**} is zero is fundamental for the proposed method. If there are more than $H - 1$ zero coefficients prescribed, we can test the hypothesis that all of these coefficients actually are zero. In a sense we are testing information beyond that which is necessary for identification.

If the coefficient vector of z^{**} is zero in the specified equation, the rank of Π^{**} cannot be greater than $H - 1$. Hence, we shall test the hypothesis that the rank of Π^{**} is $H - 1$ against the set of alternative hypotheses that the rank is H . Of course, this test is unnecessary if $K^{**} < H$.

The likelihood-ratio criterion for this hypothesis is

$$(43) \quad (1 + v)^{-\frac{1}{2}T},$$

where v is the smallest root of (19). The asymptotic distribution of

$$(44) \quad T \log(1 + v)$$

is the χ^2 -distribution with $K^{**} - H + 1$ degrees of freedom if the conditions for asymptotic normality of the estimates $\hat{\beta}$, $\hat{\gamma}^*$ are met.

6.2. *A test based on small-sample theory.* If $\beta \Pi^{**}$ is not zero, the quantities given in section 5 do not have the χ^2 - and F -distributions as stated. For example, the inequality (40) is satisfied with probability less than ϵ if $\beta \Pi^{**}$ is different from zero. It can be shown that if the number of observations is made large enough the probability of (40) can be made as small as we please.

The minimum of the ratio $(X_{K^{**}}/K^{**}) / (X_{T^*}/T^*)$ for variations of b is given when b is set equal to a solution of (20). Then the value of the ratio is the smallest root of

$$(45) \quad \left| \frac{P^{**} M_{z^0 z^0} P'^{**}}{K^{**}} - \lambda \frac{W}{T^*} \right| = 0.$$

If the inequality (40) is not satisfied by the minimum value of the ratio, it is not satisfied by the ratio for any choice of b . The probability, say δ , that

$$(46) \quad \lambda \geq F_{K^{**} T^*}(\epsilon)$$

is greater than $1 - \epsilon$ for $\beta \Pi^{**}$ equal to zero. The inequality (46) constitutes a test for the hypothesis that Π^{**} is of rank $H - 1$. However, the significance level, δ , is not exactly known. The inequality $\delta \geq 1 - \epsilon$ shows that if the test is used as if it were at significance level $1 - \epsilon$ it will be conservative.

If one attempts to compute a confidence region when there is disagreement with (46), one will find the confidence region given by (46) imaginary.

7. Advantages and Disadvantages of the Limited-Information Maximum-Likelihood Method

The significant advantage of the limited-information or reduced form method is the saving in computation when the coefficients in only one equation (or several, but not all, equations) out of a system are to be estimated. In this case there would be a greater computational expense to estimate the many more coefficients in the entire system. Moreover, simultaneous consideration of all restrictions in the system leads to more complicated formulas. It is possible to estimate the complete system, equation by equation, with the limited-information method. Then the laborious methods of [II-4] are not used.

In systems in which there are restrictions depressing the likelihood function in addition to the condition that Π^{**} be of rank $H - 1$, the limited-information method does not take into account all of the restrictions on the system while estimating the coefficients of one equation. This fact leads us to believe that the

limited-information method is not as efficient as the maximum-likelihood method, which utilizes restrictions on the entire system.

Some systems are restricted by bilinear restrictions (i.e., restrictions connecting two equations). If bilinear restrictions are imposed on an equation, then it may not be possible to estimate it by the limited-information method. On the other hand, all linear restrictions can easily be put into the form of prescribed zero coefficients.

X. SOME COMPUTATIONAL DEVICES

BY HAROLD HOTELLING

In solving normal equations it should be understood that it is highly desirable first to obtain the inverse matrix. This fact is often overlooked, but the necessity of computing standard errors and the frequent desirability of other uses of the inverse matrix - for example, adding to or taking from the set of predictors - are cogent arguments for adopting a standing rule that no large system of linear equations should ever be solved excepting by first finding the inverse of the matrix of their coefficients.

Direct calculation by methods such as those of Doolittle and Dwyer requires labor of the order of p^3 , where p is the number of rows or of unknowns. Thus inverting a matrix of 50 rows requires only about one-eighth as much work as inverting a matrix of 100 rows. This consideration contributes interest to methods of inverting a matrix by partitioning it and inverting submatrices. A method of doing this has been set forth¹, involving four inversions of matrices of half the order of the given one. An improvement given by Waugh involves only two inversions of submatrices. In each case there are also multiplications and additions of submatrices, but these are more straightforward operations than inversion, even though the labor of multiplying two matrices of order p is also of the order of p^3 . In partitioning a matrix for this purpose there is an advantage in dividing the rows into *equally numerous* groups, since when the sum of two positive numbers is fixed, the sum of their cubes is a minimum when they are equal.

In special cases there are further advantages in partitioning. Thus there may be whole blocks of zeros, or there may be triangles of zeros that make it easy to invert particular submatrices.

Another important method in matrix inversion is iteration. Several iterative methods are discussed in the paper cited [Hotelling, 1943-1]. The method recommended when a fairly good approximation C_0 to A^{-1} has been reached is to use

$$C_{m+1} = C_m(2 - A C_m)$$

¹[Hotelling, 1943-A]. Further results on matrix calculation are given in [Hotelling, 1943-B and 1949] and [Ullman, 1944].

repeatedly to compute successively C_1, C_2, \dots , while considering at each state the matrix of errors

$$D_m = I - A C_m.$$

The *norm* of a real matrix is the square root of the sum of the squares of the elements. It is an extremely useful means of setting an upper bound for the errors. Let the norm of D_0 , denoted by $N(D_0)$, be k . Then if $k < 1$, we have the following upper bound for the errors in the approximations:

$$N(C_m - A^{-1}) \leq N(C_0) \frac{k^{2^m}}{1 - k}.$$

This formula shows two advantages of this particular method of iteration. One is that the limit of error is expressed only in terms of known numbers, without involving things calculated with a degree of accuracy not previously determined. The other advantage is that this limit of error decreases with great speed as m increases, because of the exponential of the exponential appearing in it. If k happens to be greater than unity, the method will sometimes converge nevertheless, and a new $k < 1$ will emerge at a later stage. If the method diverges, a better first approximation must be found in some other way.

Both iteration and partitioning are useful in least-squares problems with constraints. These problems are of two main kinds. The first is typified in surveying, where the Euclidean value for the sum of the angles of a triangle provides a set of side conditions that modify the solution.

The second kind is represented by the regression treatment of the analysis of variance with disproportionate class frequencies. Here the side conditions are introduced solely for convenience, to make definite the selection of a particular solution among an infinity of possible solutions. When this fact is realized, the arbitrary character of the side conditions becomes evident, and it will be recognized that essentially the same final results will be reached even if the usual simple conditions are considerably altered.

It is customary to take side conditions asserting that the simple sums of certain of the regression coefficients are zero. The labor of calculation is cut down if these are altered by intro-

ducing the marginal total frequencies as coefficients.

The whole system of normal equations and side conditions in such cases is treated in a new paper by the author with the help of partitioned matrices, which clarify many of the relationships. The matrix of the whole system has an inverse, of which a particular submatrix plays a part in this theory closely analogous to that of the inverse of the matrix of the normal equations in the ordinary nonsingular case. This analogy includes both the computational and the probability aspects. Furthermore, the matrix of the whole system can be inverted with the help of the iterative method set forth above. This will be advantageous whenever a good first approximation is available. A suitable source to which we may look for such a first approximation is the familiar treatment of proportionate frequencies.

PART TWO

PROBLEMS SPECIFIC TO TIME SERIES

XI. VARIABLE PARAMETERS IN STOCHASTIC PROCESSES: TREND AND SEASONALITY

BY LEONID HURWICZ¹

1. When an economic time series is of any considerable length, the existence of *trend* is often suspected. The customary fashion of dealing with such a situation is either to eliminate the trend component by one of the several well-known methods or to treat *time* as a fixed variate entering the stochastic equations (or the regression equations) in the form of some given function, usually a polynomial.

On the other hand, when the data cover periods shorter than a year (days, months, quarters, etc.), *seasonal* behavior is often suspected.

In meteorological series there also exists, in addition to trends and seasonals, the diurnal component when, e.g., hourly data are available. The methods used in this case are analogous to those in treating trend - elimination (by subtraction or division) and fixed-variate regression.

These methods, presupposing as they do an additive or multiplicative effect of trend and seasonality,² are bound to fail in all but the simplest cases.

The experience in dealing with many time series, especially those in meteorology, indicates the need for a more flexible and general approach. In the present note the essentials of such an approach are sketched. The note is incomplete in many respects: among its deficiencies is the lack of proof of the consistency of the maximum-likelihood estimates used. For the sake of simplicity the discussion is restricted to time series in one variable only.

¹Part of the work on this paper was done in 1945 - 46 during the author's tenure of the Guggenheim Memorial Fellowship. Some of the problems considered arose in connection with the author's research at the Institute of Meteorology at the University of Chicago.

²Diurnal effects are fully analogous to seasonals and will not be discussed separately.

The generalization to several variables is straightforward.

2. Let there be an infinite time series

$$(2.1) \quad x_{\infty} = (\dots, x_{t_0-1}, x_{t_0}, x_{t_0+1}, \dots)$$

such that every finite segment x of x_{∞} has the probability density function $f(x)$. The function f is at this stage assumed to be *factorable*¹ (this assumption will later be relaxed); thus there exists ϕ such that

$$(2.2) \quad f(x) = \prod_{\tau} f_{\tau}(x_{\tau}^{(\phi)}), \quad \tau = \dots, t_0 - 1, t_0, t_0 + 1, \dots,$$

where

$$(2.3) \quad x_{\tau}^{(\phi)} = (x_{\tau}, x_{\tau+1}, \dots, x_{\tau+\phi}), \quad 1 \leq \phi < \infty.$$

3. We shall now indicate the possibility of imposing restrictions on the nature of

$$(3.1) \quad f = (\dots, f_{t_0-1}, f_{t_0}, f_{t_0+1}, \dots).$$

In the very simplest case we might have

$$(3.2) \quad f_{\tau} = f_0 \quad \text{for all } \tau.$$

But the case of interest is when f_{τ} changes with τ . If the change is strictly *periodic*, so that

$$(3.3) \quad f_{\tau} = f_{\tau+k} \quad \text{for all } \tau,$$

where k is some integer, we shall speak of *seasonal* fluctuations of f_{τ} .² Clearly, (3.2) is a special case of (3.3) with $k = 1$.

On the other hand, there may be *nonperiodic* variations in f_{τ} ; these may be referred to as *trend*, although, in the narrow sense

¹This implies nonautocorrelated disturbances.

²More generally, f_{τ} might satisfy some given functional equation, say a difference equation of finite order.

of the word, "trend" usually implies some uniform change in the value of $E(x_t)$.

4. In what follows we shall be concerned with the parametric case and assume that

$$(4.1) \quad f_{\tau} \equiv f(x_{\tau}^{(\phi)} \mid \theta_{\tau}^{(s)})$$

where $\theta_{\tau}^{(s)} \equiv (\theta_{\tau 1}, \dots, \theta_{\tau s})$ is a parameter vector variable in time. It is assumed that the form of f does not change with time. [$x_{\tau}^{(\phi)}$ was defined in (2.3).]

5. Let the sample $O_T = (x_1, \dots, x_T) = x_1^{(T-1)}$, where $T > \phi$, be drawn from the universe given by (2.2) and (4.1) after the initial $\phi - 1$ values of x have been fixed.¹

Then the likelihood function Φ_T of the sample is

$$(5.1) \quad \Pr(O_T) = \prod_{\tau=0}^{T-\phi} f(x_{\tau+\phi}^{(\phi)} \mid \theta_{\tau}^{(s)}; x_1^{(\phi-2)})$$

with $x_1^{(\phi-2)}$ fixed.

Before we proceed with estimation of the $\theta_{\tau}^{(s)}$, we shall introduce assumptions concerning the behavior of $\theta_{\tau}^{(s)}$ in time. (The superscript s will be omitted where no danger of confusion exists.)

In the *seasonal* (periodic) case, corresponding to (3.3) we would have

$$(5.2) \quad \theta_{\tau} \equiv \theta_{\tau+k},$$

which implies that all θ 's can be expressed in terms of the first k θ 's, i.e., $(\theta_{\tau_0}, \theta_{\tau_0+1}, \dots, \theta_{\tau_0+k-1})$, or, say, $\alpha \equiv (\alpha_1, \dots, \alpha_q)$, where there is a one-to-one correspondence between the components of α and of the first k θ 's. Thus the problem is reduced to that of estimating the components of α .

Similarly in the case of *trend* (in the broad sense given to this term earlier) it is assumed that the components of θ_{τ} are

¹Alternatively, the initial value vector $x_1^{(\phi-2)}$ might be considered stochastic and assigned some given probability density function. This course, while more realistic, is less expedient since it does not lead to least-squares estimates when the maximum-likelihood criterion is used. Cf. [VI].

given functions of time (e.g., polynomials), say,

$$(5.3) \quad \theta_{\tau} \equiv \theta(\tau; \alpha),$$

and the problem is again reduced to that of estimating α .

Thus both in the case of seasonals (5.2) and trend (5.3), the likelihood function Φ_{τ}^* given in (5.1) may be expressed in terms of α and τ and written

$$(5.4) \quad \Pr(O_{\tau}) = \prod_{\tau=0}^{T-p} f(x_{\tau+p}^{(p)} | \tau; \alpha; x_1^{(p-2)}),$$

and, if the maximum-likelihood criterion is used, $\Pr(O_{\tau})$ is to be maximized with respect to α .¹

6. As an example, consider the stochastic difference equation

$$(6.1) \quad x_t = \sum_{j=1}^p \theta_{jt} x_{t-j} + \theta_{0t} + \theta_{p+1,t} u_t,$$

where the u 's are independently and normally distributed [in view of the presence of θ_{0t} and $\theta_{p+1,t}$, u_t may, without any loss of generality, be assumed to have the normal distribution $N(0,1)$ with zero mean and unit variance]; the behavior of the θ 's is given by

$$(6.2) \quad \theta_{it} = \theta_{i,t+k}, \quad i = 0, 1, \dots, p+1,$$

where k may be either a year or a day (24 hours).

This type of equation is strongly suggested by meteorological data which exhibit seasonal (and diurnal) fluctuations not only in their "normal (i.e., mean) values," but also in their variances and lag-correlation coefficients.

Thus, for instance, in a series consisting of many years of monthly data it is found that the lag correlation of, say, January with February may differ from that of February with March, etc. Moreover, even for the 12-month lag the correlation of, say, the Januaries of two successive years differs from the corresponding 12-month lag correlation for February. This suggests an equation of the type (6.1) with $p > 1$. (For $p = 1$ the 12-month lag corre-

¹The alternative procedure would be to apply Lagrange's method and to maximize $\Pr(O_{\tau})$ subject to restrictive side relations.

lations do not vary from month to month.)

7. The maximum-likelihood equations can be set up easily, provided the variable coefficients θ_τ are linear in the α 's. We shall show this in two examples: one of trend and one of seasonal fluctuations.

7.1. *Trend.* Let the behavior of the time series be described by a first-order stochastic difference equation

$$(7.1) \quad x_t = (\alpha_0 + \alpha_1 t)x_{t-1} + \alpha_2 + u_t,$$

where the u 's are independently normally distributed with a zero mean and a common variance. (It will be observed that in this case the trend affects lag regression and not the expected value of x_t .) This may be considered as a special case of a one-equation stochastic system with one dependent variable x_t , and x_{t-1} , $t x_{t-1}$, and 1 as "predetermined" variables.¹

Thus let

$$(7.2) \quad \begin{aligned} z_{1t} &\equiv x_{t-1}, & \alpha_0 &= \gamma_1, \\ z_{2t} &\equiv t x_{t-1}, & \alpha_1 &= \gamma_2, \\ z_{3t} &\equiv 1, & \alpha_2 &= \gamma_3. \end{aligned}$$

Then the likelihood function of the sample is

$$(7.3) \quad \Pr(O_T) = \prod_{t=2}^T f(x_t - (\gamma_1 z_{1t} + \gamma_2 z_{2t} + \gamma_3 z_{3t})),$$

where f is normal with a variance σ^2 and a zero mean. Hence

$$(7.4) \quad \Pr(O_T) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^T \exp \left\{ - \frac{1}{2\sigma^2} \sum_t (x_t - \sum_{i=1}^3 \gamma_i z_{it})^2 \right\},$$

and the maximum-likelihood estimates $\hat{\gamma}$ of γ are obtained by minimizing the sum of squares

¹Cf. [II-1.8].

$$(7.5) \quad S = \sum_t (x_t - \sum_i \gamma_i z_{it})^2.$$

This is an example of equivalence of the maximum-likelihood and least-squares method. (The equivalence would have been absent if x_1 had not been assumed fixed!)

7.2. *Seasonal.* A strictly analogous seasonal situation is given by

$$(7.6) \quad x_t = (\alpha_0 + \alpha_1 \sin t) x_{t-1} + \alpha_2 + u_t,$$

where the period is 2π .

Here let

$$(7.7) \quad \begin{aligned} z_{1t} &\equiv x_{t-1}, & \alpha_0 &= \gamma_1, \\ z_{2t} &\equiv (\sin t) x_{t-1}, & \alpha_1 &= \gamma_2, \\ z_{3t} &\equiv 1, & \alpha_2 &= \gamma_3, \end{aligned}$$

and the procedure is exactly the same as in the preceding example: the maximum-likelihood criterion again leads to the least-squares procedure.

Clearly no obstacles will be encountered as long as the equations remain *linear in the unknown coefficients*.

One problem which may arise is that of extending the available proofs of consistency of the maximum-likelihood estimates of the α 's so as to cover the "explosive" (or, in general, undamped) situations. For a very special case this has been accomplished in [XIV].

8. In this section we consider a stochastic difference equation of order p whose coefficients are periodic with period k , and the same is true of the variance of the normally distributed "disturbance." This is the case described by (6.1).

Assuming for simplicity that $T = nk$, where n is an integer, we obtain the likelihood function

$$(8.1) \quad \Pr(O_T) = (2\pi)^{-1/2 T} \prod_{j=1}^k \sigma_j^{-n} \exp\left\{-\frac{1}{2} \sum_{j=1}^k q_j\right\},$$

where

$$(8.2) \quad q_j = \sigma_j^{-2} \sum_{r=0}^{n-1} (x_{j+rk} - \sum_{i=1}^p \alpha_{ji} x_{j+r k-i})^2 .$$

In the case of monthly observations $k = 12$. Let $j = 1$ correspond to January, $j = 2$ to February, etc. Then the *estimates of the α 's* in (8.2) would be the least-squares estimates of the regression coefficients in the regression with, say, the i th month as the "dependent" variable and some or all of the remaining eleven months as the "independent" variables.

Thus denoting the January observation in the μ th year of the sample by $y_{1\mu}$, that for February by $y_{2\mu}$, etc., we proceed as follows: Minimizing with regard to the α 's the j th sum of squares

$$(8.3) \quad s_j^2 = \sum_{\mu=1}^n (y_{j\mu} - \sum_{s=1}^p \alpha_{js} y_{j-s,\mu})^2, \quad y_{j-s,\mu} \equiv y_{j-s+12,\mu},$$

we obtain the maximum-likelihood estimates $\hat{\alpha}_{js}$ for the given j and $s = 1, 2, \dots, p$. When the $\hat{\alpha}_{js}$ are substituted in (8.3) we obtain $n \hat{\sigma}_j^2$.

This procedure is then simply that of estimating the regression of each month on its predecessors, and the 12 such regressions can be computed quite independently. This has been done in the past, although it is not easy to find justification for the methods used.

9. In practice it frequently occurs that the economist (or the meteorologist) has at his disposal data extending over a period of, say, 15 or 20 years and the number of parameters to be estimated is 3 or more. The significance of the results obtained from the use of annual figures cannot then be too reassuring. Therefore it becomes natural to attempt to utilize monthly (or quarterly) observations. Clearly, this means that seasonal phenomena (absent from annual figures) must now be taken into account. The simplest way to do this is to take the stochastic equations of the model used for annual figures and endow their parameters with seasonal variations.

Let, for instance, the "annual model" be

$$(9.1) \quad x_t = \alpha_1 x_{t-k} + \alpha_0 + \alpha_2 u_t ,$$

where x_t is the annual average of x for the t th year, the time unit is one year, and the properties of u_t are given in section 6 [following equation (6.1)].

Then the simplest way of setting up a "seasonal model" would be to assume periodic behavior of the α 's; thus

$$(9.2) \quad x_t = \alpha_{1t} x_{t-q} + \alpha_{0t} + \alpha_{2t} u_t,$$

where q is, say, one-twelfth of the time unit, i.e., a month, x_t is the monthly average, and u_t has the same properties as before.

If this were the case, it would be possible to apply the theory developed in earlier sections and no serious difficulties would arise. Unfortunately, however, such a passage from the annual to the seasonal model would not, in general, be justified. The crucial assumption, which might hold in (9.1)¹ but not in (9.2), is that of independence (or *nonautocorrelation*)² of the successive disturbances u_t .

To see why this should be the case it will suffice to remember that the disturbance itself is a stochastic process and it is not unreasonable to assume that its properties can be described by a stochastic difference equation (for simplicity, of first-order)

$$(9.3) \quad u_t = \beta u_{t-r} + v_t,$$

where the time unit r is small compared with q - say, a day. v_t , might be called the *secondary* disturbance (while u_t is *primary*), is a stochastic variable with a constant mean and standard deviation and is nonautocorrelated.

It can easily be seen that the autocorrelation $\rho(u_t, u_{t-\theta})$ of u_t tends to zero as the lag θ increases. It might still be appreciable for lags of one month but negligible for one-year lags.

Hence it may be quite legitimate to consider u_t as nonautocorrelated in the annual model (9.1), but the same assumption might be highly unrealistic for the (monthly) seasonal model (9.2).

It then becomes necessary to generalize the theory of seasonal models as developed in the earlier sections of this note by considering disturbances of the type (9.3) instead of the nonautocorrelated ones.

¹Although even here it would not be very realistic. Cf. [Hurwicz, 1944].

²As long as we deal with normally distributed disturbances, the two terms are synonymous.

10.1. Before investigating the properties of models of the type (9.2) with u_t satisfying (9.3), we shall first consider the nonseasonal case, viz.,

$$(10.1) \quad x_t = \sum_{\tau=1}^{\phi} \alpha_{\tau} x_{t-\tau} + u_t,$$

where x_t is now the *instantaneous* ("daily") value of x at the time t , and u_t the *instantaneous* value of u at time t .

10.2. It is important to see that autocorrelation of the disturbance does not necessarily imply the loss of identification, provided the nature of the autocorrelation is properly specified. (If the autocorrelation function of the disturbance is unrestricted, no identification is, in general, possible.) We distinguish two cases, depending on whether the autocorrelation of the disturbance is of the autoregressive or moving-average type [Wold].

10.2.1. *Autoregressive disturbance.* Let

$$(10.2) \quad x_t + \alpha x_{t-1} = u_t,$$

and

$$(10.3) \quad u_t + \rho u_{t-1} = v_t,$$

where v_t is random (i.e., nonautocorrelated) and has a constant variance. Then we have

$$(10.4) \quad x_t + (\alpha + \rho)x_{t-1} + \alpha\rho x_{t-2} = v_t.$$

This may also be written as

$$(10.5) \quad x_t + \beta_1 x_{t-1} + \beta_2 x_{t-2} = v_t,$$

where

$$(10.6) \quad \beta_1 = \alpha + \rho,$$

$$\beta_2 = \alpha\rho.$$

Now since the β 's are functionally independent, it is possible to obtain α and ρ from them.¹ But (10.5) is an ordinary difference

¹This is a case of *multiple* identification; cf. [IV-5.4.1] and [II-2.4.4].

equation with a random disturbance; hence there is no difficulty in obtaining the β 's as well as the variance of v_t . Thus (10.2) is identifiable¹ in spite of the autocorrelation in u_t .

More generally, we have

$$(10.7) \quad \sum_{j=0}^J \alpha_j x_{t-j} = u_t, \quad \alpha_0 = 1,$$

and

$$(10.8) \quad \sum_{i=0}^I \rho_i u_{t-i} = v_t, \quad \rho_0 = 1,$$

where v_t has the same properties as before.

Then we obtain

$$(10.9) \quad \sum_i \rho_i \sum_j \alpha_j x_{t-i-j} = v_t,$$

or

$$(10.10) \quad \sum_{s=0}^S \beta_s x_{t-s} = v_t, \quad \beta_0 = 1, \quad S = I + J,$$

where

¹Identification would not have been lost even if in addition to the "disturbance in the equation" u we also had had to deal with a "disturbance in the variable" (e.g., error of observation) w . Thus, let x_t in the equation above denote the "true" value while x_t^o is the observed value, so that

$$(a) \quad x_t^o = x_t + w_t$$

where the w_t 's are nonautocorrelated. Then

$$(b) \quad x_\tau^o = w_\tau + v_\tau + \gamma_1 v_{\tau-1} + \gamma_2 v_{\tau-2} \dots;$$

also from (10.5) we have

$$(c) \quad x_t^o + \beta_1 x_{t-1}^o + \beta_2 x_{t-2}^o = v_t + w_t + \beta_1 w_{t-1} + \beta_2 w_{t-2}.$$

By multiplying (b) and (c) and taking the expectations of both sides, enough independent equations are available to determine all the unknown parameters in terms of the lag moments of x_t^o . Cf. [T. W. Anderson and Hurwicz].

$$(10.11) \quad \beta_s = \varphi_s(\rho_1, \dots, \rho_I; \alpha_1, \dots, \alpha_J), \quad s = 1, 2, \dots, S.$$

Explicitly,

$$(10.12) \quad \begin{aligned} \beta_0 &= \alpha_0 \rho_0 = 1, \\ \beta_1 &= \rho_0 \alpha_1 + \rho_1 \alpha_0, \\ \beta_2 &= \rho_0 \alpha_2 + \rho_1 \alpha_1 + \rho_2 \alpha_0, \\ &\vdots \end{aligned}$$

Since there are as many β 's as there are ρ 's and α 's together, a necessary and sufficient condition for the identification of (10.7) is that the Jacobian of (10.12) should not vanish identically. This was seen to be the case for $I = J = 1$, where the β 's were obviously independent, but it can also be proved more generally.

10.2.2. *Moving-average disturbance.* In the previous two sections of this note the disturbance u was assumed to possess autocorrelation of the *autoregressive* type, expressed in (10.8). Another case of considerable interest is autocorrelation of the *moving-average* type where

$$(10.13) \quad u_t = \sum_{\tau=0}^R \gamma_\tau v_{t-\tau}, \quad \gamma_0 = 1,$$

with random v 's. In this case we have

$$(10.14) \quad \text{cov}(u_t, u_{t-R-k}) = 0, \quad k \geq 1,$$

where $\text{cov}(x, y)$ denotes the covariance of x and y .

Using the fact that

$$(10.15) \quad x_t = u_t + \delta_1 u_{t-1} + \dots,$$

where the δ 's are easily obtainable functions of the α 's, and writing

$$(10.16) \quad \text{cov}(x_t, x_{t-p}) = \mu_p,$$

we easily find [with the help of (10.14)] that

$$(10.17) \quad \sum \alpha_j \mu_{p-j} = 0, \quad p \geq p_0,$$

for p_0 sufficiently large in relation to R .

Since the μ 's are given if the distribution of the x 's is known, the α 's can always be determined. Similarly the lower μ 's serve to evaluate the weights γ in (10.13), although a proof of functional independence should again be supplied to make this statement rigorous. Hence the moving-average type of autocorrelated disturbance does not destroy identifiability.

So far we have confined ourselves to single equations in one unknown. However, the situation does not seem to be fundamentally different in *simultaneous systems*. As an example we use the system given by Koopmans and Rubin [II-2.5.6].

Simplifying the notation, we have

$$(10.18) \quad \begin{aligned} L'_t &\equiv x'_t + \alpha_1 x'_{t-1} + \alpha_2 x''_{t-1} = u'_t, \\ L''_t &\equiv \beta x'_t + x''_t = u''_t \end{aligned}$$

which is undoubtedly identifiable when the u 's are random.

Now let

$$(10.19) \quad u_t^{(i)} + \rho^{(i)} u_{t-1}^{(i)} = v_t^{(i)}, \quad i = ', ''$$

(v 's random) and

$$(10.20) \quad L_t^{(i)*} = L_t^{(i)} + \rho^{(i)} L_{t-1}^{(i)} = v_t^{(i)}.$$

In terms of the original parameters we have

$$(10.21) \quad \begin{aligned} x'_t + (\alpha_1 + \rho') x'_{t-1} + \rho' \alpha_1 x'_{t-2} + \rho' \alpha_2 x''_{t-2} &= v'_t, \\ \beta x'_t + x''_t + \rho' \beta x'_{t-1} &= v''_t, \end{aligned}$$

or

$$(10.22) \quad x'_t + \epsilon_1 x'_{t-1} + \epsilon_2 x''_{t-1} + \epsilon_3 x'_{t-2} + \epsilon_4 x''_{t-2} = v'_t,$$

$$\eta_1 x'_t + x''_t + \eta_2 x'_{t-1} + \eta_3 x''_{t-1} = v''_t.$$

Clearly, the system (10.22) is identified, even without taking into account the additional restriction on the ϵ 's and η 's which follow from (10.21).

It should be emphasized that both (10.18) and (10.19) are very special cases, but it is not unreasonable to conjecture that in a large class of cases identification is not destroyed by introducing autocorrelated disturbances.

10.3. Now define the "monthly" average of x_t

$$(10.23) \quad X_t = \frac{1}{H} \sum_{\tau=t-H+1}^t x_\tau$$

and of u_t

$$(10.24) \quad U_t = \frac{1}{H} \sum_{\tau=t-H+1}^t u_\tau,$$

where H is the number of "days" in the "month."

We assume, as is usually the case in practice, that only the X_t but not the x_t are known from observation. These are given for non-overlapping time intervals of H "days" each:

$$(10.25) \quad \dots, X_{t_0}, X_{t_0+H}, X_{t_0+2H}, \dots,$$

which may be denoted by

$$(10.26) \quad \dots, Y_{t_0}, Y_{t_0+1}, Y_{t_0+2}, \dots,$$

respectively.

Now in order to estimate the unknown parameters (the mean and variance of v_t , β , and the α 's) by the maximum-likelihood method, we must obtain the likelihood function for the Y 's.¹ Since in the

¹In general, of course, the system may be unidentified. This can be remedied by assuming a priori restrictions on the unknown parameters or by introducing, say, a fixed variate in the right-hand member of (10.1). In

case of normal v_t all the other variables are also normal, the likelihood function is specified if the means and the lag covariance matrix of Y_t are obtained. Making the means of v_t zero, we are left only with the problem of finding the lag covariance matrix of Y_t . (The case of a nonzero mean may be treated in an analogous manner.)

11. We first observe that

$$(11.1) \quad X_t = \sum_{\tau=1}^p \alpha_{\tau} X_{t-\tau} + U_t$$

so that

$$(11.2) \quad X_t = \sum_{\tau=0}^{\infty} \gamma_{\tau} U_{t-\tau},$$

where the γ 's are functions of the α 's and may be obtained from the identity

$$(11.3) \quad \sum_{\tau=0}^{\infty} \gamma_{\tau} U_{t-\tau} \equiv \sum_{\tau=1}^p \alpha_{\tau} \sum_{\tau'=0}^{\infty} \gamma_{\tau'} U_{t-\tau-\tau'}.$$

Thus

$$(11.4) \quad \begin{aligned} \mathcal{E}(X_t X_{t-k}) &= \mathcal{E} \left[\left(\sum_{\tau=1}^{\infty} \gamma_{\tau} U_{t-\tau} \right) \left(\sum_{\tau'=0}^{\infty} \gamma_{\tau'} U_{t-k-\tau'} \right) \right] \\ &= \mathcal{E} \left[\sum_{\tau, \tau'}^{\infty} \gamma_{\tau} \gamma_{\tau'} U_{t-\tau} U_{t-k-\tau'} \right] \\ &= \sum_{\tau, \tau'}^{\infty} \gamma_{\tau} \gamma_{\tau'} \mu_{t-\tau, t-k-\tau'}, \end{aligned}$$

where

$$(11.5) \quad \mu_{t-\tau, t-k-\tau'} = \mathcal{E}(U_{t-\tau} U_{t-k-\tau'}).$$

The μ 's depend on σ^2 of v and β (as well as H) and may easily be evaluated with the help of (9.3) and (10.24).

the latter case the procedure would be essentially of the same type as, though more complex algebraically than, that given in the text.

It remains to obtain the lag moments of the Y 's, that is, of every H th X_t . One will note that the Y 's satisfy a stochastic difference equation

$$(11.6) \quad Y_t = \sum_{\tau=1}^p \alpha'_\tau Y_{t-H\tau} + U'_t,$$

where U' is a linear combination of lagged values¹ of the U_t . The coefficients of this linear combination are functions of the α_τ and so are the α'_τ in (11.6).

These facts make it possible to obtain the lag moments of Y_t and, because of normality and the zero-means assumptions, the likelihood function of the observed variates is obtained.

12. The procedure carried out in sections 10 and 11 was based on the assumption that the parameters were constant. Now, in order to take account of *seasonal fluctuations*, we may make σ^2 , β , or the α 's vary with time t but subject to the restriction of periodicity. This will change the likelihood function, but its derivation follows the same pattern. Once the likelihood function has been obtained the unknown parameters may be estimated provided the system is identified.

13. The models of stochastic processes used so far have all been of the discrete type. It is not difficult to construct continuous models with properties analogous to those presented above. However, even for the nonseasonal case the study of estimation in continuous systems presents serious difficulties. A treatment of this problem is not yet available.²

Nevertheless, for the sake of completeness, an example will be given which is the continuous counterpart of (10.1) and its seasonal version.

This model may be written as an integral equation

$$(13.1) \quad x(t) = \int_{-\infty}^t \varphi(t - \tau) x(\tau) d\tau + u(t),$$

where $u(t)$ has an autocorrelation function tending to zero as the

¹Hence (11.6) is a difference equation with an autocorrelated disturbance of the moving-average type. See above, section 10.2.

²Cf. [XVI].

lag increases.

To introduce seasonality we use a modified kernel, thus obtaining

$$(13.2) \quad x(t) = \int_{-\infty}^t \psi(t, t - \tau) x(\tau) d\tau + u(t),$$

where

$$(13.3) \quad \psi(t, s) = \psi(t + k, s).$$

To make the mathematics of (13.3) more manageable one might assume

$$(13.4) \quad \psi(t, s) = \psi_1(t) \psi_2(s).$$

XII. NONPARAMETRIC TESTS AGAINST TREND¹

BY HENRY B. MANN

In testing against trend we are testing the hypothesis that the members of a certain sequence of random variables x_1, \dots, x_n, \dots are distributed independently of each other, each with the same distribution, which we shall throughout this paper assume to be continuous. Although the null hypothesis is thus well defined, we are not in the same happy position with respect to the alternatives that we wish to admit. We have some rather vague ideas of a sequence where the variables tend to decrease, but no clear notion of how this decrease should be expressed in terms of distribution functions has as yet been advanced. It seems advantageous to eliminate from such a definition any tendency to decrease which has its roots in the dependence between successive observations. We, therefore, propose the following definition for a downward trend:

A sequence of random variables x_1, x_2, \dots will be said to have a downward trend if the variables are independently distributed so that x_i has the cumulative distribution function f_i and $f_i(x) < f_j(x)$ for every $i < j$ and every x .

An upward trend is similarly defined with $f_i(x) > f_j(x)$ for $i < j$.

If we are justified in restricting the null hypothesis to the form

$$(1) \quad f(x) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^x e^{-\frac{(y-\mu)^2}{2\sigma^2}} dy$$

and in restricting the alternatives to expressions of the type

$$(2) \quad f_i(x) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^x e^{-\frac{(y-\mu_i)^2}{2\sigma^2}} dy,$$

¹This is an abstract of a paper published in *Econometrica* under the same title [Mann].

where $\mu_i = P(ih)$ and P is a fixed polynomial of known degree with monotonically decreasing values $P(ih)$ for $i > 0$, $h > 0$, then the test against trend may be carried out in the usual fashion. That is to say, we can test the significance of the regression coefficients and this test is known to possess certain desirable properties. It is, for instance, in the case where P is a linear function, a uniformly most powerful test with respect to the alternatives considered. However, if we are not justified in assuming the null hypothesis to be of the form (1), then the customary procedure will not give the correct size of the critical region. Even if we are justified in assuming the null hypothesis to be of the form (1), but wish to admit as alternatives all downward trends as defined before, then we can not state that the usual procedure of testing the significance of the regression coefficients is in any respect superior to other procedures. In fact, we shall describe a test, based on ranks, which has maximum power with respect to a class of alternatives just as likely to occur in practical work as the alternatives (2). Since we very rarely have a knowledge of the actual distribution, it is of importance to discuss tests for which the size of the critical region does not depend on a priori assumptions about the distributions concerned.

A fairly simple test against trend is the T -test first introduced by Kendall [1938] for the testing of independence in a bivariate distribution. The statistic T counts, in a sample x_1, \dots, x_n , the number of inequalities $x_i < x_j$ for $i < j$. One such inequality will be termed a reverse arrangement. The distribution of T is readily computed from the recursion

$$(3) \quad P_n(T) = P_{n-1}(T) + P_{n-1}(T-1) + \dots + P_{n-1}(T-n+1),$$

where $P_n(T)$ is the number of permutations of n different values x_1, \dots, x_n , which possess exactly T reverse arrangements. The recursion (3) was derived by M. G. Kendall. He also derived the variance of T and proved ultimate normality. The approach to normality is remarkably rapid.

The odd moments of T around its mean value $n(n-1)/4$ are 0. From the recursion (2) it is possible to derive the following difference equation for the even moments of T :

$$(4) \quad \mathcal{E}_n(x^{2i}) = \mathcal{E}_{n-1}(x^{2i}) + \binom{2i}{2} B_n^{(2)} \mathcal{E}_{n-1}(x^{2i-2}) + \dots + B_n^{(2i)},$$

where $\mathcal{E}_n(x^{2i})$ is the $2i$ th moment of T around its mean in permutations of n variables and

$$B_n^{(2j)} = \frac{1}{n} \sum_{K=0}^{n-1} \left(\frac{n-1}{2} - K \right)^{2j}.$$

From (4) we find

$$\sigma_n^2(T) = \frac{2n^3 + 3n^2 - 5n}{72}$$

and

$$\mathcal{E}_n(x^4) = \frac{100n^6 + 228n^5 - 455n^4 - 870n^3 + 625n^2 + 372n}{43.200}.$$

(4) can also be utilized to give a simple proof for the ultimate normality of T . This can be done in showing by mathematical induction that $\mathcal{E}_n(x^i) / \{\sigma_n(T)\}^i$ converges to the moments of a normal distribution.

To give lower bounds for the power of the T -test, we introduce a parameter λ_n defined by the equation

$$\mathcal{E}_n(T|H) - \frac{n(n-1)}{4} = \lambda_n \frac{n(n-1)}{2},$$

where $\mathcal{E}_n(T|H)$ is the expectation of T under the alternative H . Denote by $\sigma_n(T|H)$ the variance under the alternative H , by σ_{0n} the variance under the null hypothesis, by $P(T \leq \bar{T} | H)$ the probability that $T \leq \bar{T}$ if H is the true situation. If H is a trend, then the following inequalities can be derived:

$$(5) \quad \sigma_n^2(T|H) \leq \frac{n(n-1)}{4} + \frac{n(n-1)(n-2)}{6} - \frac{\lambda_n^2 n^2 (n-1)}{4},$$

$$(6) \quad P(T \leq \bar{T} | H) \geq 1 - \frac{\sigma_n^2(T|H)}{\left(\lambda_n \frac{n(n-1)}{2} + t_n \sigma_{0n} \right)^2},$$

for $\lambda_n \{n(n-1)/2\} \leq -t_n \sigma_{0n}$, where $\bar{T} = \{n(n-1)/4\} - t_n \sigma_{0n}$

defines t_n . From the ultimate normality of T it follows that t_n converges to a constant t if the size of the critical region is fixed.

For small values of n another inequality gives a better lower bound for the power. It can be shown that

$$(7) \quad P(T \leq \bar{T} \mid H) \geq \frac{-2\lambda_n n(n-1) - 4t_n \sigma_{0n}}{n(n-1) - 4t_n \sigma_{0n}}.$$

From (6) it follows that the T -test is consistent with respect to alternatives for which λ_n is negative and of order larger than $1/\sqrt{n}$. (6) and (7) can also be utilized to derive sufficient conditions for the unbiasedness of the T -test.

An example of a class of alternatives with respect to which the T -test is a most powerful test based on ranks is the following class C of alternatives H .

Let p_{ij} be the probability that x_i is the j th largest of the variables x_i, x_{i+1}, \dots, x_n . If under an alternative H , $p_{ij} = a_i p^{j-1}$ ($p < 1$), independent of the ranks of x_1, \dots, x_{i-1} , then the T -test will be most powerful with respect to H . It is not known, however, whether there are such alternatives among the trends as defined in the introduction.

Another test that the author considers worth serious consideration is the K -test. This test is carried out as follows. We determine the smallest K for which

$$(8) \quad x_i > x_{i+j+K} \quad \text{for } j \geq 1, i = 1, \dots, n-K.$$

We then fix \bar{K} so that $P(K \leq \bar{K} \mid H_0)$ equals the size of the critical region. If then $K \leq \bar{K}$ we proceed on the hypothesis of a trend. If $K > \bar{K}$ we proceed as if the null hypothesis were true.

$P(K \leq \bar{K} \mid H_0)$ may be found for $K \leq 3$ from the following relations:

$$(9) \quad \begin{aligned} Q_n(1) &= 1, & Q_n(2) &= Q_{n-1}(2) + Q_{n-2}(2), \\ Q_n(3) &= Q_{n-1}(3) + Q_{n-2}(3) + 3Q_{n-3}(3) + Q_{n-4}(3), \end{aligned}$$

where $Q_n(\bar{K})$ is the number of permutations x_{i_1}, \dots, x_{i_n} satisfying (8).

It is also easy to prove the relations

$$(10) \quad P_n(n - \bar{K}) = P_{2\bar{K}}(\bar{K}) \quad \text{for } n \geq 2\bar{K},$$

where $P_n(\bar{K}) = P(K \leq \bar{K})$ under the null hypothesis in samples of size n . It may further be shown that

$$(11) \quad P_n(n - \bar{K}) = \sum \{ [i_1 + 1] [\max(i_1, i_2) + 2] \cdots [\max(i_1, i_2, \dots, i_K) + \bar{K}] \}^{-1},$$

where Σ denotes summation over all permutations i_1, \dots, i_K of $1, 2, \dots, \bar{K}$.

Also

$$(12) \quad P_9(4) = \sum' \{ [\bar{\max}(i_1) + 1] [\bar{\max}(i_1, i_2) + 2] \cdots [\bar{\max}(i_1, \dots, i_5) + 5] \}^{-1},$$

where $\bar{\max}(i_1, \dots, i_e) = \min[\max(i_1, \dots, i_e), 4]$ and Σ' denotes summation over all permutations i_1, \dots, i_5 of $1, 2, \dots, 5$ for which 1 precedes 5.

The relations (9), (10), (11), and (12) permit tabulation of $P_n(\bar{K})$ for $n \leq 9$. For higher values of n the following critical regions are available:

$$(13) \quad \begin{aligned} P_n(n - 5) &= P_{10}(5) = 0.0098, & \text{for } n \geq 10, \\ P_n(n - 4) &= P_8(4) = 0.0284 & \text{for } n \geq 8 \\ P_n(n - 3) &= P_6(3) = 0.0792 & \text{for } n \geq 6, \\ P_n(n - 2) &= 0.2083 & \text{for } n \geq 4, \\ P_n(n - 1) &= P_2(1) = 0.5 & \text{for } n \geq 2. \end{aligned}$$

These are all the critical regions between 0.0098 and 0.5 possible for the K -test for $n \geq 10$. The K -test has the disadvantage

TABLE 1*
PROBABILITY OF OBTAINING A PERMUTATION WITH $T \leq \bar{T}$ IN PERMUTATIONS OF n VARIABLES.

$n \backslash \bar{T}$	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
3	167	500	833																			
4	042	167	375	625	833	958																
5	008	042	117	242	408	592	758	883	958	992												
6	001	008	028	068	136	235	360	500	500	640	765	864	932	972	992	999						
7	000	001	005	015	035	068	119	191	281	386	500	500	614	719	809	881	932	965	985	995	999	
8	000	000	001	003	007	016	031	054	089	138	199	274	360	452								
9	000	000	000	000	001	003	006	012	022	038	060	090	130	179	238	306	381	460				
10	000	000	000	000	000	000	001	002	005	008	014	023	036	054	078	108	146	190	242	300	364	431
$P(c)$	000	000	000	000	001	001	002	004	006	010	016	025	037	054	076	105	142	186	237	296	360	429

* Tabular values should be divided by 1000.

$$P(c) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-c} e^{-x^2/2} dx, \quad c = \left(\frac{n(n-1)}{4} - T - \frac{1}{2} \right) / \sqrt{\frac{2n^3 + 3n^2 - 5n}{72}}, \quad n = 10.$$

TABLE 2
PROBABILITY OF OBTAINING A PERMUTATION WITH $K \leq \bar{K}$ IN PERMUTATIONS OF n VARIABLES.

$n \backslash \bar{K}$	1	2	3	4	5	6	7
3	0.1667	0.5000					
4	0.0417	0.2083	0.5000				
5	0.0083	0.0667	0.2083	0.5000			
6	0.0014	0.0181	0.0792	0.2083	0.5000		
7	0.0002	0.0042	0.0246	0.0792	0.2083	0.5000	
8	0.0000	0.0008	0.0066	0.0284	0.0792	0.2083	0.5000
9	0.0000	0.0002	0.0016	0.0086	0.0284	0.0792	0.2083

that the choice of critical regions is rather limited.

It is easy to see that there are trend alternatives for which $P(x_i > x_{i+j+\bar{K}}) = 1$ for $j \geq 0$ ($i = 1, \dots, n - K$). The K -test has the power 1 with respect to such alternatives, while other tests, for instance the T -test or the significance tests ordinarily used, may have a considerably lower power. This seems remarkable in view of the fact that the K -test is a test based on ranks and in fact often uses only a fraction of the sample.

XIII. TESTS OF SIGNIFICANCE IN TIME-SERIES ANALYSIS

BY R. L. ANDERSON

1. Suppose we have data on the monthly prices of some product in a given locality over a period of n years. We assume that the prices P_{ij} in the i th month and j th year ($i = 1, 2, \dots, 12$; $j = 1, 2, \dots, n$) are independently distributed normal variables with a common variance σ^2 and means μ_{ij} where

$$(1) \quad \mu_{ij} = \mu + \alpha_i + \beta_j$$

and

$$(2) \quad \sum_i \alpha_i = \sum_j \beta_j = 0,$$

so that the likelihood function for the sample is

$$(3) \quad (\sqrt{2\pi}\sigma)^{-12n} \exp\left\{-\frac{1}{2\sigma^2} \sum_{j,i} (P_{ij} - \mu_{ij})^2\right\}.$$

2. We test the hypothesis H_0 given by

$$(4) \quad \alpha_i = 0, \quad i = 1, 2, \dots, 12.$$

Thus H_0 specifies that there is no monthly ("seasonal") pattern in the variations of the P_{ij} . The likelihood-ratio criterion leads to the use of the statistic

$$(5) \quad F = \frac{[(12-1)(n-1)]}{(12-1)} \frac{\hat{\sigma}'^2 - \hat{\sigma}^2}{\hat{\sigma}^2} = (n-1) \frac{\hat{\sigma}'^2 - \hat{\sigma}^2}{\hat{\sigma}^2},$$

which has Snedecor's distribution with 11 and $(n-1)$ degrees of freedom when H_0 is true.¹ The symbols $\hat{\sigma}'^2$ and $\hat{\sigma}^2$ have the following

¹For proof, cf., for instance, [Wilks, p.178 ff.].

meaning: $\hat{\sigma}^2$ is the maximum-likelihood estimate of σ^2 obtained by maximizing [subject to (2)] the expression in (3) with regard to the α_i , β_j , and to σ^2 and is given by

$$(6) \quad \hat{\sigma}^2 = \frac{1}{12n} \sum_{i,j} (P_{ij} - m_{ij})^2,$$

where

$$(7) \quad \begin{aligned} m_{ij} &= \hat{\mu}_{ij} = m + a_i + b_j, \\ m &= \hat{\mu} = \frac{1}{12n} \sum_{i,j} P_{ij}, \\ a_i &= \hat{\alpha}_i = \frac{1}{n} \sum_j P_{ij} - m, \\ b_j &= \hat{\beta}_j = \frac{1}{12} \sum_i P_{ij} - m. \end{aligned}$$

$\hat{\sigma}'^2$ is the maximum likelihood estimate of σ^2 on the assumption that H_0 is true, and can thus be obtained by maximizing [subject to (2)] the expression in (3) with regard to the β_j and to μ and σ^2 after all the α_i have been set equal to zero.

We then obtain

$$(8) \quad \hat{\sigma}'^2 = \frac{1}{12n} \sum_{i,j} (P_{ij} - m - b_j)^2,$$

with the values of m and b_j given by (7)

It can then be shown that

$$(9) \quad \hat{\sigma}'^2 - \hat{\sigma}^2 = \frac{1}{12n} \sum_{i,j} a_i^2 = \frac{1}{12} \sum_i a_i^2.$$

Thus (5) becomes¹

$$(10) \quad F = (n-1) \frac{\sum_i a_i^2}{\sum_{i,j} (P_{ij} - m_{ij})^2}.$$

¹The form of the denominator can be simplified for computational purposes.

The analysis of variance table is then as follows:

Source of variation	Degrees of freedom	Sums of squares	Estimate of variation ¹
Months	11	$n \sum_i a_i^2$	$11(\sigma^2 + n \sigma_m^2)$
Years	$n - 1$	$12 \sum_j b_j^2$	$(n - 1)(\sigma^2 + 12 \sigma_y^2)$
Residual	$11(n - 1)$	$12n \hat{\sigma}^2$	$11(n - 1) \sigma^2$

¹ $\hat{\sigma}^2$ is a biased estimate of the true error variance σ^2 . The unbiased estimate is $12n \hat{\sigma}^2 / 11(n - 1)$. σ_m^2 and σ_y^2 are the true variances between months and years.

3. The usefulness of the above analysis of variance depends on the extent to which the data fulfill the assumptions made. Some of the problems that merit consideration are:

(a). The simplicity of the analysis of variance depends upon the assumption of no missing prices. Simple methods have been devised for handling a two-way-classification analysis with missing observations, but these methods become more complicated for more than two classifications, especially if there are several missing observations. If many of the observations are missing, or if there is a different number of observations for each category under consideration, the analysis of variance may become quite complicated.

(b). A system of stochastic (difference) equations has been suggested as a proper model for economic phenomena [Haavelmo, 1943, p.1 ff.]. Our problem should be reconsidered from this viewpoint to find whether the likelihood function in (3) is a realistic description of observed phenomena. In particular, equation (1) might include other (economic) variables besides the year and the month; and the presence of lagged values in the equations would invalidate the assumption of independence of successive observations inherent in (3). If it is desired to test for such independence from the residuals, a way should be found to allow for the serial correlation introduced into the residuals by the fact that, as a result of

(7), they add up to zero in all directions.

(c). The assumption of normality usually can be severely violated without seriously affecting the general conclusions; however, the assumption of the same variance σ^2 for every observation cannot be violated to any great extent. A somewhat arbitrary rule often applied in biological experiments is to permit the assumption of constant variance only if no two means differ by more than 50 per cent.

(d). Finally, we should state that more thought should be devoted to the problem of designing social experiments, especially in setting up sampling methods. The analysis of variance set-up has been made much more useful for biological and agricultural research by connecting it with many designs of experiments other than the simple randomized block design, of which our month \times year analysis is an example.

XIV. CONSISTENCY OF MAXIMUM-LIKELIHOOD ESTIMATES IN THE EXPLOSIVE CASE

BY HERMAN RUBIN

	Page
1. Introduction	356
2. Transformation of the Problem	357
3. System Stable Without Disturbances	358
4. System Explosive Without Disturbances	362

1. Introduction

It was shown by Mann and Wald [1943] that if a temporal stochastic process described by linear difference equations is damped without disturbances, it is stable (the expectations of the squares of all variables are uniformly bounded) with disturbances and the "maximum-likelihood" estimates of the parameters involved are consistent. However, a system which is stable or explosive (the expectations of the squares of some variables being unbounded) without disturbances is explosive with them. Here we prove consistency in a simple example of the explosive case.

Let us consider the equation

$$(1) \quad x_t = \rho x_{t-1} + u_t, \quad |\rho| \geq 1 \quad (t = 1, 2, \dots),$$

where ρ is a real number, the u_t are real stochastic variables independently distributed with mean 0 and variance σ^2 , and x_0 is a given real number. (The results derived here hold equally well if ρ , u_t , and x_0 are complex numbers, quaternions, or Cayley numbers.) The maximum-likelihood estimate $\hat{\rho}$ of the parameter ρ is defined as

$$(2) \quad \hat{\rho} = \frac{\sum x_t x_{t-1}}{\sum x_{t-1}^2},$$

where

$$(3) \quad \sum = \sum_{t=1}^T.$$

[In the case of complex numbers, quaternions, or Cayley numbers, replace (2) by

$$(2') \quad \hat{\rho} = \frac{\sum x_t \bar{x}_{t-1}}{\sum x_{t-1} \bar{x}_{t-1}},$$

where \bar{x} denotes the conjugate of x .] We shall show that

$$(4) \quad \text{plim}_{T \rightarrow \infty} \hat{\rho} = \rho.$$

2. Transformation of the Problem

We can write (2) in the form

$$(5) \quad \hat{\rho} = \frac{\sum (\rho x_{t-1} + u_t) x_{t-1}}{\sum x_{t-1}^2} = \frac{\sum \rho x_{t-1}^2}{\sum x_{t-1}^2} + \frac{\sum u_t x_{t-1}}{\sum x_{t-1}^2} \\ = \rho + \frac{\sum u_t x_{t-1}}{\sum x_{t-1}^2}.$$

By the Cauchy-Schwarz inequality,

$$(6) \quad \left| \sum u_t x_{t-1} \right| \leq (\sum u_t^2)^{1/2} (\sum x_{t-1}^2)^{1/2}.$$

Hence

$$(7) \quad \left| \frac{\sum u_t x_{t-1}}{\sum x_{t-1}^2} \right| \leq \frac{(\sum u_t^2)^{1/2}}{(\sum x_{t-1}^2)^{1/2}}.$$

But the u_t^2 are real variables independently distributed with identical distribution functions and means σ^2 . Hence [Cramér, Theorem 15]

$$(8) \quad \text{plim}_{T \rightarrow \infty} \frac{1}{T} \sum u_t^2 = \sigma^2.$$

Let us consider the quantities x_t . By (1), we have

$$(9) \quad x_t = \rho^t x_0 + \rho^{t-1} u_1 + \cdots + \rho u_{t-1} + u_t.$$

Dividing (9) by ρ^t , we obtain

$$(10) \quad \rho^{-t} x_t = x_0 + \sum_{\tau=1}^t \rho^{-\tau} u_{\tau}.$$

3. System Stable Without Disturbances

In this section let us assume $|\rho| = 1$. Then

$$(11) \quad |\rho^{-\tau} u_{\tau}| = |u_{\tau}|.$$

We observe that for any $\delta > 0$ there exists a number $A(\delta)$ such that if $B > A(\delta)$ then

$$(12) \quad \int_{|u_{\tau}| > B} |\rho^{-\tau} u_{\tau}|^2 dF(u_{\tau}) = \int_{|u_{\tau}| > B} |u_{\tau}|^2 dF(u_{\tau}) < \delta,$$

since $\mathcal{E}(|u_{\tau}|^2) = \mathcal{E}(u_{\tau}^2) = \sigma^2 < \infty$. Therefore, if $t > A^2(\delta)/\epsilon^2$, we have

$$(13) \quad \frac{1}{t} \sum_{\tau=1}^t \int_{|u_{\tau}| > \epsilon\sqrt{t}} |\rho^{-\tau} u_{\tau}|^2 dF(u_{\tau}) < \frac{1}{t} \sum_{\tau=1}^t \delta = \delta.$$

Hence, for every $\epsilon > 0$,

$$(14) \quad \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=1}^t \int_{|u_{\tau}| > \epsilon\sqrt{t}} |\rho^{-\tau} u_{\tau}|^2 dF(u_{\tau}) = 0.$$

By the central limit theorem [Cramér, Theorem 21a],

$$x = \frac{1}{\sqrt{t}} \sum_{\tau=1}^t \rho^{-\tau} u_{\tau}$$

is asymptotically normally distributed with mean 0 and variance σ^2 . Therefore, its cumulative distribution function $G_t(x)$ converges

uniformly¹ to the normal distribution $N(x, 0, \sigma^2)$. Let us now prove the following lemma:

LEMMA 3. *There exists a number λ such that for every $\varepsilon > 0$ there is a t_0 such that if $t > t_0$ then for every b , $P_t(|x-b| < \varepsilon) < \lambda \varepsilon$.*

Since $G_t(x)$ converges uniformly to $N(x, 0, \sigma^2)$, we see that for $\delta\varepsilon > 0$ there is a t_0 such that if $t > t_0$ then

$$(15) \quad |G_t(x) - N(x, 0, \sigma^2)| < \delta\varepsilon.$$

Then

$$(16) \quad G_t(b+\varepsilon) - G_t(b-\varepsilon) < N(b+\varepsilon, 0, \sigma^2) - N(b-\varepsilon, 0, \sigma^2) + 2\delta\varepsilon.$$

But $N(b+\varepsilon, 0, \sigma^2) - N(b-\varepsilon, 0, \sigma^2) < 2\varepsilon/\sqrt{2\pi}\sigma$, since $1/\sqrt{2\pi}\sigma$ is the maximum of dN/dx . Therefore

$$(17) \quad P_t(|x-b| < \varepsilon) = G_t(b+\varepsilon) - G_t(b-\varepsilon) < \varepsilon\left(\frac{2}{\sqrt{2\pi}\sigma} + 2\delta\right)$$

for $t > t_0$. Take b in the foregoing lemma to be $-x_0/\sqrt{t}$. Then we obtain from (10)

$$(18) \quad P\left(\left|\frac{\rho^{-t}x_t}{\sqrt{t}}\right| < \varepsilon\right) < \lambda\varepsilon,$$

or, since $|\rho| = 1$,

$$(19) \quad P(|x_t| < \varepsilon\sqrt{t}) < \lambda\varepsilon.$$

From (1) we see that

$$(20) \quad \left||x_{t+1}| - |x_t|\right| = \left||x_{t+1}| - |\rho x_t|\right| \leq |x_{t+1} - \rho x_t| = |u_{t+1}|.$$

In the Tchebycheff inequality,

$$(21) \quad P(|u| > \theta\sigma) < \frac{1}{\theta^2},$$

²Let F_n be a sequence of distribution functions converging to a continuous distribution function F . Then F_n converges uniformly to F .

let us take

$$(22) \quad \theta = \frac{\varepsilon\sqrt{t}}{2\mu\sigma}.$$

Then

$$(23) \quad P(|u_\tau| > \frac{\varepsilon\sqrt{t}}{2\mu}) < \frac{4\mu^2\sigma^2}{\varepsilon^2 t}.$$

Therefore

$$(24) \quad P(\{|u_{t+1}| + |u_{t+2}| + \dots + |u_{t+\mu}|\} > \frac{\varepsilon\sqrt{t}}{2}) < \frac{4\mu^3\sigma^2}{\varepsilon^2 t},$$

since if $|u_{t+1}| + \dots + |u_{t+\mu}| > \varepsilon\sqrt{t}/2$, at least one $|u_{t+i}|$ is greater than $\varepsilon\sqrt{t}/2\mu$, and

$$(25) \quad P(\text{at least one } u_{t+i} > \frac{\varepsilon\sqrt{t}}{2\mu}) \leq \sum_{i=1}^{\mu} P(|u_{t+i}| > \frac{\varepsilon\sqrt{t}}{2\mu}).$$

On the other hand it follows that if $k \leq \mu$,

$$(26) \quad |x_{t+k}| \geq |x_t| - \sum_{\tau=1}^k |u_{t+\tau}| \geq |x_t| - \sum_{\tau=1}^{\mu} |u_{t+\tau}|.$$

Therefore

$$(27) \quad P(\min\{|x_t|, |x_{t+1}|, \dots, |x_{t+\mu}|\} < \frac{\varepsilon\sqrt{t}}{2}) < \lambda\varepsilon + \frac{4\mu^3\sigma^2}{\varepsilon^2 t}.$$

Hence we see that

$$(28) \quad P(\sum_{\tau=t}^{t+\mu} x_\tau^2 < \frac{\varepsilon(\mu+1)t}{4}) < \lambda\varepsilon + \frac{4\mu^3\sigma^2}{\varepsilon^2 t}.$$

But

$$(29) \quad \sum_{\tau=0}^{t+\mu} x_\tau^2 \geq \sum_{\tau=t}^{t+\mu} x_\tau^2.$$

Take $\mu = [4Av/\varepsilon^2]$, $v > 1$, where $[x]$ denotes the greatest integer less than or equal to x . Then

$$(30) \quad P\left(\sum_{\tau=0}^{t + [4Av/\epsilon^2]} x_{\tau} < Avt\right) < \lambda\epsilon + \frac{256\sigma^2 A^3 v^3}{\epsilon^8 t}.$$

And for any $t > ([4Av/\epsilon^2] + 1)/(v-1)$ we have

$$(31) \quad t + \left[\frac{4Av}{\epsilon^2}\right] < vt - 1.$$

Let $vt = T$. Since

$$(32) \quad \sum_{\tau=0}^{t + [4Av/\epsilon^2]} x_{\tau}^2 < \sum_{\tau=1}^T x_{\tau-1}^2,$$

it follows that

$$(33) \quad P\left(\sum x_{t-1}^2 < AT\right) < \lambda\epsilon + \frac{256\sigma^2 A^3 v^4}{\epsilon^8 T}$$

for $T > \max\{t_0, ([4Av/\epsilon^2] + 1)/(v-1)\}$, with t_0 defined as in the preceding lemma. Therefore

$$(34) \quad \text{plim}_{T \rightarrow \infty} \frac{1}{T} \sum x_{t-1}^2 \geq A$$

for every $A > 0$. Hence

$$(35) \quad \text{plim}_{T \rightarrow \infty} \frac{\sum u_t^2}{\sum x_{t-1}^2} = \text{plim}_{T \rightarrow \infty} \frac{\frac{1}{T} \sum u_t^2}{\frac{1}{T} \sum x_{t-1}^2} \leq \frac{\sigma^2}{A}.$$

Therefore

$$(36) \quad \text{plim}_{T \rightarrow \infty} \frac{(\sum u_t^2)^{1/2}}{(\sum x_{t-1}^2)^{1/2}} = 0,$$

and, from (5),

$$(37) \quad \text{plim}_{T \rightarrow \infty} \hat{\rho} = \rho.$$

4. System Explosive Without Disturbances

Let us assume now that $|\rho| > 1$. Denote by $F * G$ the convolution

$$(38) \quad F * G(x) = \int_{-\infty}^{+\infty} F(x-t) dG(t)$$

of the distributions F and G . If the distribution function of u_τ is $F(x)$, then the distribution function of $\rho^{-\tau} u_\tau$ is $F(\rho^\tau x)$. Since the mean of $\rho^{-\tau} u_\tau$ is 0 and the variance of $\rho^{-\tau} u_\tau$ is $\rho^{-2\tau} \sigma^2$, it follows [Wintner, Theorem 7.1] that if

$$(39) \quad F(x) * F(\rho x) * \cdots * F(\rho^n x) = G_n(x),$$

then $\lim_{n \rightarrow \infty} G_n(x)$ exists and is a distribution function $G(x)$.

Let us now prove the following lemma:

LEMMA 4. $G(x)$ is continuous.

We shall consider three cases, for although the proof for the third case applies also to the other two, it is the most complicated.

Case 1. $F(x)$ is continuous. Then [Cramér, p. 37, equation (40)] it follows that G is continuous.

For any distribution function $H(x)$ denote by the functional $J\{H(x)\}$ the sum of all jumps of $H(x)$. We observe that $J\{H(\lambda x)\} = J\{H(x)\}$, and that $J\{H(x)\} = 0$ if and only if $H(x)$ is continuous. If $J\{H(x)\} < 1$, we shall say $H(x)$ is partly continuous, and if $J\{H(x)\} = 1$, we shall say that $H(x)$ is discrete.

Case 2. $H(x)$ is partly continuous. Let $J\{H(x)\} = \delta < 1$. From (37) we have

$$(40) \quad G_n(x) = F(x) * G_{n-1}(\rho x).$$

But

$$(41) \quad J(H * K) = J(H) J(K).$$

Therefore

$$(42) \quad J\{G_n(x)\} = \delta^n.$$

We may derive from (37)

$$(43) \quad G(x) = G_n(x) * G(\rho^n x).$$

Using (39) and (40) we obtain

$$(44) \quad J\{G(x)\} = \delta^n J\{G(x)\} \leq \delta^n.$$

Therefore $J\{G(x)\} = 0$, and $G(x)$ is continuous.

Case 3. $H(x)$ is discrete. Denote by $L\{H(x)\}$ the maximum jump of $H(x)$. We see that $L\{H(\lambda x)\} = L\{H(x)\}$. We have $L\{F(x)\} = \varphi < 1$. For if $L\{F(x)\} = 1$, then

$$(45) \quad \begin{aligned} P(y \leq x) &= 0 & x < x^0, \\ P(y \leq x) &= 1 & x \geq x^0, \end{aligned}$$

and the variance of u_τ would be 0, which contradicts its being $\sigma^2 > 0$. Using the fact that [Wintner, Theorem 7.6] $J\{G(x)\} > 0$ if and only if

$$\prod_{k=1}^{\infty} L\{F(\rho^k x)\} > 0,$$

we see from

$$(46) \quad \prod_{k=1}^{\infty} L\{F(\rho^k x)\} = \prod_{k=1}^{\infty} \varphi = 0$$

that $G(x)$ is continuous.

Since the $G_n(x)$ approach $G(x)$, it follows that for any $\varepsilon > 0$, there are a δ and a t_0 such that if $n > t_0$ then

$$(47) \quad P_n(|x + x_0| < \delta) < \varepsilon.$$

We observe that $\rho^{-t}x_t - x_0$ has the distribution G_t . Therefore, for $t > t_0$, we have

$$(48) \quad P(|\rho^{-t}x_t| < \delta) < \varepsilon.$$

We further observe that for any $\delta > 0$, and for any K , there is a t_1 such that, if $t > t_1$, then

$$(49) \quad \frac{|\rho|^t \delta}{\sigma\sqrt{t}} > K.$$

From (46), we have for $t > t_0$,

$$(50) \quad P(|x_t| < |\rho|^t \delta) < \varepsilon.$$

Therefore, for $t > \max(t_0, t_1)$, we see that

$$(51) \quad P(|x_t| > K\sigma\sqrt{t}) > 1 - \varepsilon.$$

Then

$$(52) \quad \text{plim}_{T \rightarrow \infty} \frac{\sum u_t^2}{\sum x_{t-1}^2} \leq \text{plim}_{T \rightarrow \infty} \frac{\frac{1}{T} \sum u_t^2}{\frac{1}{T} x_{T-1}^2} \leq \lim_{T \rightarrow \infty} \frac{\sigma^2}{K^2 \sigma^2 \frac{T-1}{T}} = \frac{1}{K^2}.$$

Therefore

$$(53) \quad \text{plim}_{T \rightarrow \infty} \frac{(\sum u_t^2)^{1/2}}{(\sum x_{t-1}^2)^{1/2}} = \left(\text{plim}_{T \rightarrow \infty} \frac{\sum u_t^2}{\sum x_{t-1}^2} \right)^{1/2} = 0$$

and

$$(54) \quad \text{plim}_{T \rightarrow \infty} \hat{\rho} = \rho.$$

XV. LEAST-SQUARES BIAS IN TIME SERIES¹

BY LEONID HURWICZ

1.1. In this paper it is shown that there exist cases where the least-squares and the maximum-likelihood estimates of the regression and structural coefficients are biased² for any finite-sized sample³ drawn from a population defined by a noncircular stochastic difference-equation system. This bias is evaluated for certain special cases. It is found that for very small samples the bias may amount to as much as 25 per cent of the true value of the parameter⁴, while for medium sized samples (say of 20 observations) the bias is still almost 10 per cent (the numerical value of the expectation of the estimate always being below the true parameter value⁵). The relative bias seems to tend to zero, although rather slowly, as the damping of the system becomes weaker.

1.2. The initial objective of this paper was to prove that the least-squares method yields biased estimates of regression coeffi-

¹Part of the work on this paper was done in 1945-46 during the author's tenure of the Guggenheim Memorial Fellowship. Some of the problems considered arose in connection with the author's research at the Institute of Meteorology at the University of Chicago in 1944.

²I.e., the mathematical expectations of the estimates are not identically equal to the true parameter values. *Bias* is defined as the difference between the expectation of the estimate and the true parameter value. *Relative bias* is the ratio of this difference to the true parameter. The quantity N_T appearing in the formulae below is the relative bias plus 1 in samples of T observations; hence N_T equals the ratio of the expectation of the estimate to the true parameter, so that when N tends to 1 the relative bias tends to zero.

³Containing more than two observations.

⁴Highest relative bias, found in samples of four observations, is $26\frac{2}{3}$ per cent.

⁵I.e., the bias, where known, is always negative for positive parameter values and vice versa.

cients in autoregressive¹ noncircular time series. It was generally realized that the usual proofs of the Markoff Theorem [David and Neyman] were not valid for this case,² but the author does not know of any proof of the actual existence of a bias.³

It was possible to obtain explicit formulae for the bias in very small samples (three or four observations). This was done for two types of assumptions with regard to the initial value: (1) *fixed initial-value* case - where the initial observation may be regarded as a fixed variate (the initial value chosen being zero) see section 3.1.2 below; (2) *stochastic initial-value* case - where the initial value is a stochastic variable whose distribution is the marginal distribution of the later observations, see section 2.4.2 below. The author regards the latter assumption as being more realistic, but since the *exact* equivalence of least-squares and maximum-likelihood criteria of estimation of the regression coefficients applies only to the former (asymptotically to both), it is the existence of bias in the fixed initial-value case that constitutes a proof of *bias of the maximum-likelihood* criterion as well. Moreover, the case treated is one where the regression coefficient is identically equal to the corresponding structural coefficient. Hence, the *bias* of least-squares and maximum-likelihood methods exists in *structural* as well as *predictive* estimation.

Once the *existence* of the bias had been shown it was of interest to investigate its *magnitude* for samples of various sizes. The first three terms of a series expansion for the bias in a sample of arbitrary size were obtained; see equations (4.6) to (4.10) below.

Because of the equivalence of the structural and regression coefficients in the cases treated, the results described may be regarded as a first (and very modest) step in the small-sample theory of the maximum-likelihood estimates of the structural coefficients in (noncircular) stochastic difference-equation systems.

The results thus far obtained indicate the importance of the

¹I.e., generated by a system of stochastic difference equations, some or all of which contain *lagged* values of *endogenous* variables; autoregressive time series are to be distinguished from those composed *additively* of a given function of time (often a polynomial or a Fourier series) and a stochastic ("error") term.

²They are valid in the *lagless* case, as shown in [VI].

³In this note we are only concerned with "stable" systems, i.e., those satisfying Assumption IV₂ in [Mann and Wald, p. 192]; for the case here treated this assumption is equivalent to postulating $|\alpha| < 1$ in (2.1) below.

bias. There is urgent need for more intensive study of the small-sample properties of estimates in autoregressive time series.

2.1. Let a sample be given consisting of T observations on a stochastic variable X_t . It will be assumed that the joint cumulative distribution function of (X_1, \dots, X_T) , to be written as $F(X_1, \dots, X_T)$, has the following two properties:

$$(2.1) \quad \mathcal{E}(X_t | X_{t-1}) = \alpha X_{t-1}, \quad t = 2, 3, \dots, T,$$

$$(2.2) \quad \sigma^2(X_t | X_{t-1}) = \text{const.}$$

Thus the regression of any observation on its predecessor is linear¹ and homoscedastic.

The least-squares estimate \tilde{a}_T of α is in this case given by

$$(2.3) \quad \tilde{a}_T = \frac{\sum_{t=2}^T X_t X_{t-1}}{\sum_{t=2}^T X_{t-1}^2}.$$

2.2. An estimate a_T of α will be said to be (conditionally) unbiased with regard to a family \mathfrak{F} of cumulative distribution functions F if all F in \mathfrak{F} satisfy (2.1) and (2.2) and if

$$(2.4) \quad \mathcal{E}(a_T) \equiv \alpha \quad \text{for all } \alpha, \quad \text{for all } F \in \mathfrak{F}, \quad \text{and for all } T.$$

An estimate conditionally unbiased with regard to all cumulative distribution functions is said to be absolutely unbiased.²

Clearly, if a family \mathfrak{F}_0 can be found with regard to which a_T is not unbiased, a_T cannot be absolutely unbiased.

2.3. In this note we show that there exists a family³ \mathfrak{F}_0 , in

¹The condition (2.1) is stronger than linearity, but for the purposes of this paper no loss of generality is involved.

²These definitions are special cases of those given in [VI].

³This family, defined by (2.1) and (2.2) with the additional assumption of normality, is divided into two "branches" depending on whether the

fact one of considerable practical importance, with regard to which the least-squares estimate $\tilde{\alpha}_T$ of α is not unbiased. This shows that $\tilde{\alpha}_T$ is not absolutely unbiased. It is not known whether there exists a family, say \mathfrak{F}_1 , such that $\tilde{\alpha}_T$ is unbiased with regard to \mathfrak{F}_1 .

However, an unbiased estimate of α does exist. Thus in the sample (X_0, X_1, \dots, X_T) where X_0 is fixed and different from zero, the ratio X_1/X_0 is an unbiased, though in general very inefficient, estimate of α . For, assuming that

$$(2.5) \quad X_t = \alpha X_{t-1} + u_t, \quad t = 1, 2, \dots, T,$$

where the u 's have zero means, we have

$$(2.6) \quad \mathfrak{E}\left(\frac{X_1}{X_0}\right) = \mathfrak{E}\left(\frac{\alpha X_0 + u_1}{X_0}\right) = \alpha + \frac{1}{X_0} \mathfrak{E}(u_1) = \alpha.$$

It may be remarked that for $T = 2$, the ratio X_1/X_0 is a least-squares estimate. In fact, $T = 2$ is the only known case among finite-sized samples where the least-squares estimate $\tilde{\alpha}_T$ is unbiased.

On the other hand, in the sample (X_1, \dots, X_T) , where X_1 is stochastic and the likelihood function is given by (2.9), the expectation $\mathfrak{E}(X_2/X_1)$ exists only in the Cauchy principal-value sense. In that sense, however, X_2/X_1 is an unbiased estimate of α , since it has the Cauchy distribution with a mode at zero. In fact, with the u 's independent of each other, the mean

$$\frac{1}{T-1} \sum_{t=2}^T \frac{X_t}{X_{t-1}}$$

would also be unbiased in the Cauchy sense, although no more efficient than X_2/X_1 . One might conjecture that the median of the ratios X_t/X_{t-1} , $t = 2, \dots, T$, would be a more efficient estimate of α and perhaps an unbiased one.

2.4.1. In the following sections of the paper the proof of existence of the bias will be given and its magnitude evaluated for first-order stochastic difference equations with the initial value

initial value is assumed stochastic or fixed. The bias exists in both cases. The respective likelihood functions are given by (2.9) and (3.15).

stochastic and fixed.

For samples of three and four observations the size of relative bias¹ is given for all $|\alpha| < 1$ in Table 1 and Figure 1; these are based on the formulae (3.12), (3.19), and (3.37). The derivations and comments are to be found in section 3 of this paper. The case of fixed initial value (chosen as zero) is worked out for a sample of three observations only.

For larger samples, as well as the small ones, the limiting value of bias¹ as $|\alpha| \rightarrow 0$ is given in Table 2 and Figure 2; these are based on (4.4). The latter formula is valid for both the stochastic and fixed (zero) initial-value case.

Equations (4.7) to (4.10) give a more general result, viz., the first three terms of the Maclaurin expansion of the relative bias. The derivations are given in section 4. Sections 4.2 and 4.3 contain some conjectures with regard to the nature of the approximation provided by the expansion and with regard to the nature of relative bias in stochastic difference equations in general, as well as some suggestions for further research.

2.4.2. Define a stochastic process by

$$(2.7) \quad X_t = \alpha X_{t-1} + u_t, \quad t = 2, 3, \dots, T,$$

where the u 's are independently normally distributed with zero means and unit variances, and where $|\alpha| < 1$. Hence, if the process is stationary²,

$$(2.8) \quad \mathcal{E}(X_t) = 0, \quad \mathcal{E}(X_t^2) = \frac{1}{1 - \alpha^2}, \quad t = 1, 2, \dots, T.$$

Given a sample of size T , its likelihood function will be

$$(2.9) \quad (1 - \alpha^2)^{1/2} (2\pi)^{-T/2} \exp\left\{-\frac{1}{2} [(1 - \alpha^2)X_1^2 + \sum_{t=2}^T (X_t - \alpha X_{t-1})^2]\right\},$$

which is equivalent to equation (9) in [Koopmans, 1942]. It can be

¹The formulae, tables, and figures give the values of $N_T \equiv \mathcal{E}(\tilde{\alpha}_T)/\alpha$, i.e., the relative bias plus 1.

²This specifies the stochastic initial-value case. For fixed initial-value case see below, section 3.1.2.

seen that

$$(2.10) \quad \mathcal{E}(X_t | X_{t-1}) = \alpha X_{t-1}, \quad t = 2, 3, \dots, T,$$

so that (2.1) is satisfied. (2.2) is also satisfied with a unit variance.

3.1.1. To show the existence of bias it will suffice to find a T for which (2.4) does not hold. The following proof demonstrates that (2.4) does not hold for $T = 3$.

The proof consists in finding the expectation of \tilde{a}_T for $T = 3$, where

$$(3.1) \quad \tilde{a}_3 = \frac{X_1 X_2 + X_2 X_3}{X_1^2 + X_2^2}.$$

Following the procedure used in [Williams] and in [Dixon], we write

$$(3.2) \quad \tilde{a}_3 = \frac{A_3}{B_3}, \quad A_3 = X_1 X_2 + X_2 X_3, \quad B_3 = X_1^2 + X_2^2,$$

then find the characteristic function $\varphi_3(t_1, t_2)$ of A_3 and B_3 , and, finally, obtain the expectation of \tilde{a}_3 from

$$(3.3) \quad \mathcal{E}(\tilde{a}_3) = \int_{-\infty}^0 \left[\frac{\partial \varphi_3(t_1, t_2)}{\partial t_1} \right]_{t_1=0} dt_2.$$

We have

$$(3.4) \quad (1 - \alpha^2) \varphi_3^{-2}(t_1, t_2) = \begin{vmatrix} y & b & 0 \\ b & z & b \\ 0 & b & 1 \end{vmatrix} = \begin{vmatrix} y & f & 0 \\ 1 & z & f \\ 0 & 1 & 1 \end{vmatrix},$$

where

$$(3.5) \quad \begin{aligned} y &= -2t_2 + 1, \\ f &= b^2 = (t_1 + \alpha)^2, \\ z &= y + \alpha^2 = -2t_2 + 1 + \alpha^2. \end{aligned}$$

Denoting the determinant in (3.4) by $C^{(3)}$, we have

$$(3.6) \quad \mathcal{E}(\tilde{a}_3) = -\frac{\sqrt{1-\alpha^2}}{2} \int_{-\infty}^0 \frac{1}{C_0^{(3)} \sqrt{C_0^{(3)}}} \left. \frac{\partial C^{(3)}}{\partial t_1} \right|_{t_1=0} dt_2,$$

where

$$(3.7) \quad C_0^{(3)} \equiv C^{(3)} \Big|_{t_1=0} = 4(t_2^2 - t_2 + \frac{1-\alpha^2}{4})$$

and

$$(3.8) \quad \left. \frac{\partial C^{(3)}}{\partial t_1} \right|_{t_1=0} = 4\alpha(t_2 - 1).$$

The integral in (3.6) may now be evaluated [Pierce, formula 200] and we obtain

$$(3.9) \quad \mathcal{E}(\tilde{a}_3) = \frac{\alpha}{2} \left(1 + \frac{1 - \sqrt{1 - \alpha^2}}{\alpha^2} \right).$$

It can be seen that

$$(3.10) \quad \begin{aligned} \mathcal{E}(\tilde{a}_3) &\rightarrow \frac{3}{4} \alpha \quad \text{for } |\alpha| \rightarrow 0, \\ \mathcal{E}(\tilde{a}_3) &\rightarrow \alpha \quad \text{for } |\alpha| \rightarrow 1. \end{aligned}$$

Thus \tilde{a}_3 is a biased estimate of α .

Writing $\beta = \alpha^2$ and

$$(3.11) \quad N_T \equiv N_T(\beta) = \frac{\mathcal{E}(\tilde{a}_T)}{\alpha},$$

we may state the above results as

$$(3.12) \quad N_3(\beta) = \frac{1}{2} \left(1 + \frac{1 - \sqrt{1 - \beta}}{\beta} \right)$$

with $N_3(\beta)$ varying from¹ $N_3(0) = 0.75$ to $N_3(1) = 1$. The values of

¹ $N_T(0) \equiv \lim_{\beta \rightarrow 0} N_T(\beta); N_T(1) \equiv \lim_{\beta \rightarrow 1} N_T(\beta).$

$N_3(\beta)$ are shown below in Table 1 and plotted in Figure 1.

The convergence of N_3 to 1 as $|\alpha|$ tends to 1 is very slow. The relative bias is still 12½ per cent (i.e., $N_3 = 0.875$) for $|\alpha| = 0.94$, so that it takes 94 per cent of the range of $|\alpha|$ to remove one-half of the relative bias!

3.1.2.¹ By a similar procedure we can evaluate the bias for the case where the initial value is a fixed variate, here chosen as zero. Let

$$(3.13) \quad X_0 = \text{fixed},$$

$$(3.14) \quad X_t = \alpha X_{t-1} + u_t, \quad t = 1, 2, \dots, T,$$

where α and the u 's have the same properties as before. T now denotes the number of *stochastic* observations. The likelihood function becomes

$$(3.15) \quad (2\pi)^{-T/2} \exp\left\{-\frac{1}{2} \sum_{t=1}^T (X_t - \alpha X_{t-1})^2\right\}.$$

Then the *least-squares*² estimate $\tilde{\alpha}_T^*$ of α , which in this case [unlike that specified by (2.9)] is also a *maximum-likelihood* estimate, is given by

$$(3.16) \quad \tilde{\alpha}_T^* = \frac{\sum_{t=1}^T X_t X_{t-1}}{\sum_{t=1}^T X_{t-1}^2}.$$

Now consider the special case where

$$(3.17) \quad X_0 = 0.$$

Here we find that $\tilde{\alpha}_T^*$ as defined in (3.16) equals $\tilde{\alpha}_T$ as defined in (2.3).

¹T. W. Anderson has made helpful suggestions in connection with the problem treated in this section.

²The asterisk in $\tilde{\alpha}_T^*$ indicates that $\tilde{\alpha}_T^*$ is the least-squares estimate of α for the sample (X_0, X_1, \dots, X_T) while $\tilde{\alpha}_T$ is the least-squares estimate for the sample (X_1, \dots, X_T) .

Denote by $\mathcal{E}^*(\tilde{a}_T)$ the expectation of \tilde{a}_T evaluated on the basis of (3.15) with $X_0 = 0$; also, write

$$(3.18) \quad N_T^* \equiv N_T^*(\beta) \equiv \frac{\mathcal{E}^*(\tilde{a}_T)}{\alpha}.$$

Then we find that

$$(3.19) \quad N_3^*(\beta) = \frac{3 + \beta}{4 + \beta}.$$

Thus $N_3^*(0) = 0.75 = N_3(0)$; ¹ but while $N_3(1) = 1$, we have here $N_3^*(1) = 4/5$: there is a bias even at $|\alpha| = 1$!

The case where $X_0 \neq 0$ has not been treated, but one might conjecture that, for given T and β , $N_T^*(\beta) \rightarrow 1$ as X_0 becomes numerically large.

3.2.1. When $T = 4, 5$ the integrals to be evaluated are of the elliptic type. For $T \geq 6$ they are hyperelliptic. Only for the case of $T = 4$ has the bias been evaluated in closed form.

To perform the integration in the elliptic case it is necessary to factor the T -rowed determinant $C_0^{(T)}$ where ² [cf. (3.4) above]

$$(3.20) \quad C_0^{(T)} = (1 - \beta) \varphi^{-2}(t_1, t_2) = \begin{vmatrix} y & f & 0 & \dots & 0 & 0 & 0 \\ 1 & z & f & \dots & 0 & 0 & 0 \\ 0 & 1 & z & \dots & 0 & 0 & 0 \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & \dots & 1 & z & f \\ 0 & 0 & 0 & \dots & 0 & 1 & 1 \end{vmatrix}_{(T)}$$

in a sample of T observations and

$$(3.21) \quad C_0^{(T)} = C^{(T)} \Big|_{t_1=0}.$$

We may write

¹This is an example of a more general phenomenon: $N_T^*(0) = N_T(0)$ for all T , cf. (4.4).

²For the definitions of y, z, f , see (3.5).

$$(3.22) \quad C_0^{(T)} = \begin{vmatrix} z-\beta & \beta & 0 & \dots & 0 & 0 & 0 \\ 1 & z & \beta & \dots & 0 & 0 & 0 \\ 0 & 1 & z & \dots & 0 & 0 & 0 \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & \dots & 1 & z & \beta \\ 0 & 0 & 0 & \dots & 0 & 1 & 1 \end{vmatrix}_{(T)}$$

For $\mu = 2, 3, \dots$ we have

$$(3.23) \quad \begin{aligned} C_0^{(2\mu+1)} &= R_\mu^2 - \beta R_{\mu-1}^2 = (R_\mu - \alpha R_{\mu-1})(R_\mu + \alpha R_{\mu-1}), \\ C_0^{(2\mu)} &= R_{\mu-1}(z R_{\mu-1} - 2\beta R_{\mu-2}), \end{aligned}$$

where

$$(3.24) \quad \begin{aligned} R_\mu &= \begin{vmatrix} z & \beta & 0 & \dots & 0 & 0 & 0 \\ 1 & z & \beta & \dots & 0 & 0 & 0 \\ 0 & 1 & z & \dots & 0 & 0 & 0 \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & \dots & 1 & 1 & z & \beta \\ 0 & 0 & 0 & \dots & 0 & 1 & 1 \end{vmatrix}_{(\mu+1)} \\ &= \begin{vmatrix} z-\beta & \beta & 0 & \dots & 0 & 0 & 0 \\ 1 & z & \beta & \dots & 0 & 0 & 0 \\ 0 & 1 & z & \dots & 0 & 0 & 0 \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & \dots & 1 & z & \beta \\ 0 & 0 & 0 & \dots & 0 & 1 & z \end{vmatrix}_{(\mu)} \end{aligned}$$

with

$$(3.25) \quad R_0 = 1; \quad R_1 = \begin{vmatrix} z & \beta \\ 1 & 1 \end{vmatrix}.$$

Equations (3.23) make it possible to find by radicals the roots of the determinant $C_0^{(T)}$ for $T \leq 9$. For example,

$$(3.26) \quad C_0^{(4)} = (z - \beta)[z(z - \beta) - 2\beta],$$

$$C_0^{(5)} = \{[z(z - \beta) - \beta - \alpha(z - \beta)]\{[z(z - \beta) - \beta] + \alpha(z - \beta)\},$$

$$C_0^{(6)} = [z(z - \beta) - \beta]\{z[z(z - \beta) - \beta] - 2\beta(z - \beta)\},$$

etc.

3.2.2. Now in order to obtain $N_4(\beta)$ we follow a procedure similar to that for $T = 3$. We have

$$(3.27) \quad \mathcal{E}(\tilde{a}_4) = -\frac{1}{2} \sqrt{1 - \beta} \int_{-\infty}^0 [C_0^{(4)}]^{-3/2} \frac{\partial C^{(4)}}{\partial t_1} \Big|_{t_1=0} dt_2,$$

where¹

$$(3.28) \quad C^{(4)} = \begin{vmatrix} y & f & 0 & 0 \\ 1 & z & f & 0 \\ 0 & 1 & z & f \\ 0 & 0 & 1 & 1 \end{vmatrix},$$

so that

$$(3.29) \quad C_0^{(4)} = y(y^2 + \beta y - 2\beta),$$

$$(3.30) \quad \frac{\partial C^{(4)}}{\partial t_1} \Big|_0 = -2\alpha [(y^2 + \beta y - 2\beta) + (2y + \beta)].$$

Hence, substituting these values into (3.27) and dividing both

¹The symbols used here are defined in (3.5).

sides by α we have

$$(3.31) \quad N_4(\beta) = \sqrt{1-\beta} \int_{-\infty}^0 \frac{(y^2 + \beta y - 2\beta) + (2y + \beta)}{[y(y^2 + \beta y - 2\beta)]^{3/2}} dt_2 .$$

Splitting into partial fractions and making the substitutions

$$(3.32) \quad y = -2t_2 + 1 ,$$

$$(3.33) \quad w^2 = \frac{y - y_1}{y} , \quad y_1 = \frac{1}{2} (\sqrt{\beta^2 + 8\beta} - \beta) ,$$

and, finally,

$$(3.34) \quad u = \operatorname{sn}^{-1} w ,$$

we obtain

$$(3.35) \quad \frac{4\sqrt{\beta^2 + 8\beta}}{\sqrt{1-\beta}} N_4(\beta) = \frac{1}{y_1^2} \int cs^2 u \, du + \frac{1}{2y_1} \int cn^2 u \, du \\ + \frac{1}{y_3(y_1 - y_3)} \int cd^2 u \, du ,$$

$$(3.36) \quad y_1, y_3 = \frac{1}{2} (\pm \sqrt{\beta^2 + 8\beta} - \beta) ,$$

where Glaisher's notation is used ($csu \equiv cnu / \operatorname{sn} u$, $cd u \equiv cnu / \operatorname{dn} u$); the upper limit of integration is $\operatorname{sn}^{-1}(1)$ and the lower limit $\operatorname{sn}^{-1}[(1 - y_1)^{1/2}]$.

With the help of formulae given by [Whittaker and Watson, 22.72, Ex. 3], we perform the integration, thus obtaining, after simplification,

$$(3.37) \quad N_4(\beta) = \frac{1}{2\beta} \left[(1 + \beta) - \frac{\sqrt{1-\beta}}{4\sqrt{\beta^2 + 8\beta}} \Delta F \right] ,$$

where

$$(3.38) \quad \Delta F = F_{90^\circ}(\theta) - F_\varphi(\theta).$$

$F_\varphi(\theta)$ is the incomplete elliptic integral of the first kind with modular angle θ and amplitude φ ; $F_{90^\circ}(\theta)$ is the corresponding complete integral. The angles θ and φ are given by

$$(3.39) \quad \begin{aligned} \sin \varphi &= \sqrt{1 - y_1}, \\ \sin \theta &= k, \quad 2k^2 = 1 + \frac{\beta}{\sqrt{\beta^2 + 8\beta}}. \end{aligned}$$

3.2.3. It is easily seen that

$$(3.40) \quad N_4(1) \equiv \lim_{\beta \rightarrow 1} N_4(\beta) = 1.$$

From this and the behavior of $N_3(1)$, one might conjecture that $N_T(1) = 1$ for all T .

It will be shown later in (4.4) that $N_4(0) \equiv \lim_{\beta \rightarrow 0} N_4(\beta) = 11/15 < N_3(0)$. As can be seen in Table 1 and Figure 1, $N_4(\beta)$ is a monotonic function of β . Whether this is the property of $N_T(\beta)$ for all T is not known.

3.2.4. The numerical values¹ computed from (3.37) are given in Table 1 and plotted in Figure 1. It is of interest to note that $N_3(\beta) > N_4(\beta)$ for all β . It will be seen later [from (4.4)] that $N_T(0)$ has a minimum for $T = 4$ (if T is an integer)². The convergence of $N_4(\beta)$ to 1 as $\beta \rightarrow 1$ is very slow, as was also the case for $N_3(\beta)$. To reduce the relative bias to one-half of its maximal value (so that $N_4 = 0.866$) we must have $|\alpha| = 0.95$. Whether the situation is quite as serious for larger samples is not known.

4.1. By expanding the integrand of (3.27), with T replacing

¹Obtained with the aid of Miss Estelle Mass.

²For $T \geq 2$ and real, but not necessarily an integer, the minimum is at $T = 2 + \sqrt{3}$; $N_{2+\sqrt{3}}(0) = \sqrt{3} - 1 = 0.7321$.

TABLE 1.

Ratio of Estimate Expectation to True Parameter Value for Samples of 3 and 4 Observations

α	$\beta \equiv \alpha^2$	$N_3 \equiv \mathcal{E}(\tilde{a}_3) / \alpha$	$N_3^* \equiv \mathcal{E}^*(\tilde{a}_3) / \alpha$	$N_4 \equiv \mathcal{E}(\tilde{a}_4) / \alpha$
0	0	0.7500	0.7500	0.7333
0.1	0.01	0.7506	0.7506	0.7340
0.2	0.04	0.7526	0.7525	0.7359
0.3	0.09	0.7559	0.7555	0.7388
0.4	0.16	0.7609	0.7596	0.7434
0.5	0.25	0.7679	0.7647	0.7501
0.6	0.36	0.7778	0.7706	0.7595
0.7	0.49	0.7918	0.7773	0.7730
0.8	0.64	0.8125	0.7845	0.7938
0.9	0.81	0.8482	0.7921	0.8307
0.95	0.9025	0.8810	0.7960	0.8656
0.99	0.9801	0.9382	0.7992	0.9289
1.00	1.0000	1.0000	0.8000	1.0000

N and N^* refer to stochastic and fixed (zero) initial-value cases, respectively. The subscript indicates the number of (stochastic) observations in the sample.

TABLE 2.

The Limit (for Small Parameter Values) of the Ratio of Estimate Expectation to True Parameter Value: Stochastic Initial-Value Case

Sample Size T	$N_T(0) \equiv \lim_{\alpha \rightarrow 0} \mathcal{E}(\tilde{a}_T) / \alpha$	Sample Size T	$N_T(0) \equiv \lim_{ \alpha \rightarrow 0} \mathcal{E}(\tilde{a}_T) / \alpha$
22	1.0000	13	0.8690
3	0.7500	14	0.8769
4	0.7333	15	0.8839
5	0.7500	16	0.8902
6	0.7714	17	0.8958
7	0.7917	18	0.9009
8	0.8095	19	0.9056
9	0.8250	20	0.9098
10	0.8384	50	0.9616
11	0.8500	100	0.9804
12	0.8601	500	0.9960

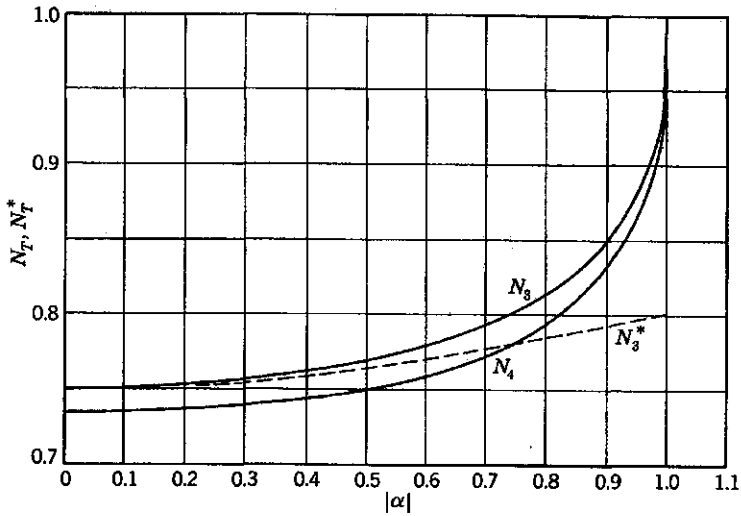


FIGURE 1. Ratio of Estimate Expectation to True Parameter Value for Samples of 3 and 4 Observations. (For explanation of symbols see Table 1.)

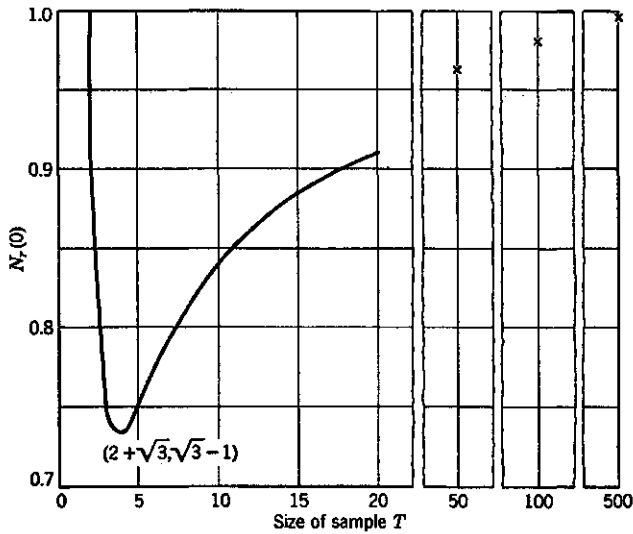


FIGURE 2. The Limit (for Small Parameter Values) of the Ratio of Estimate Expectation to True Parameter Value: Stochastic Initial-Value Case. (For explanation of symbols see Table 2.)

the affix 4, in a Maclaurin series and integrating termwise, it is possible to obtain an expansion¹ for $N_T(\beta)$ in powers of β .

It will now be shown how the first term of the expansion is obtained.² We first find

$$(4.1) \quad \lim_{\beta \rightarrow 0} C_0^{(T)} = y^{T-1},$$

$$(4.2) \quad \lim_{\beta \rightarrow 0} \frac{1}{\alpha} \frac{\partial C^{(T)}}{\partial t_1} \Big|_{t_1=0} = -2y^{T-3} [y + (T-2)]_0.$$

Hence

$$(4.3) \quad N_T(0) \equiv \lim_{\beta \rightarrow 0} \frac{\mathcal{E}(\tilde{\alpha}_T)}{\alpha} = \int_{-\infty}^0 \frac{y^{T-3} [y + (T-2)]}{(y^{T-1})^{3/2}} dt_2,$$

and, upon evaluation, this yields³

$$(4.4) \quad N_T(0) = \frac{T^2 - 2T + 3}{(T-1)(T+1)}.$$

Hence the bias exists for all finite-sized samples except $T = 2$. The values of $N_T(0)$ for some T are given in Table 2 and plotted in Figure 2.

It will be noted that the relative bias is 9 per cent for a sample of 20 observations and 2 per cent for a sample of 100 observations.

The second and third terms of the expansion of $N_T(\beta)$ in powers of β are obtained by similar methods, although the procedure becomes quite laborious. Writing

¹As before, $N_T(\beta) \equiv \mathcal{E}(\tilde{\alpha}_T)/\alpha$ and $\tilde{\alpha}_T = \sum_{t=2}^T X_t X_{t-1} / \sum_{t=2}^T X_{t-1}^2$ is the least-squares estimate of α .

²Valuable suggestions in connection with this problem were made by Professor John von Neumann, Institute for Advanced Study.

³The same formula holds for $N_T^*(0)$, i.e., for the fixed initial-value case with $X_0 = 0$; see section 3.1.2 above.

$$(4.5) \quad C_{ij}^{(T)} = \left[\frac{\partial^j}{\partial \beta^j} \left(\frac{\partial^i C^{(T)}}{\partial f^i} \right)_{f=\beta} \right]_{\beta=0} .$$

we have

$$(4.6) \quad \begin{aligned} C_{00}^{(T)} &= y^{T-1} , \\ C_{01}^{(T)} &= y^{T-3} [(T-3)y - (T-2)] , \\ C_{02}^{(T)} &= (T-3)(T-4)y^{T-5} (y^2 - 2y + 1) , \\ C_{10}^{(T)} &= -y^{T-3} [y + (T-2)] , \\ C_{11}^{(T)} &= -(T-3)y^{T-5} [y^2 + (T-5)y - (T-4)] , \\ C_{12}^{(T)} &= -(T-4)y^{T-7} [(T-3)y^3 + (T-4)(T-7)y^2 \\ &\quad - (T-5)(2T-11)y + (T-5)(T-6)] , \end{aligned}$$

where the expression for $C_{12}^{(T)}$ is valid only for $T > 3$. With the help of (4.6) we find that, for $T > 3$, the first three coefficients of the expansion

$$(4.7) \quad N_T(\beta) = N_T(0) + N_T'(0)\beta + \frac{1}{2} N_T''(0)\beta^2 + \dots$$

are given by

$$(4.8) \quad N_T(0) = \frac{T^2 - 2T + 3}{(T-1)(T+1)} ,$$

$$(4.9) \quad N_T'(0) = \frac{2(T^2 - 8T + 21)}{(T-1)(T+1)(T+3)(T+5)} ,$$

$$(4.10) \quad N_T''(0) = \frac{4(T^4 + 24T^3 + 98T^2 - 264T - 99)}{(T-1)(T+1)(T+3)(T+5)(T+7)(T+9)} .$$

Equations (4.8) and (4.9) are valid for $T \geq 2$ and $T \geq 3$, respectively, equation (4.10) for $T > 3$ only.

4.2. It is quite evident that the first three terms of the expansion for $N_T(\beta)$ do not give a very good approximation for values of β near 1. For example, $N_4(1) = 1$, while the first three terms of the expansion give only 0.823. For $|\alpha| = 0.5$, however, while¹ $N_4(0.25) = 0.7501$, the first three terms of the expansion give 0.7497 which is correct to the third decimal digit. If the same phenomenon should exist for $T > 4$, which remains to be proved, the first three terms of the expansion could not only be used successfully for values of $|\alpha|$ below 0.5 but also for values of $|\alpha|$ somewhat (though not too much) above 0.5.

Moreover, it would seem that these three expansion terms will always give values lower than the true $N_T(\beta)$. It would be desirable to examine the correctness of this conjecture and also to obtain an upper bound for $N_T(\beta)$ in terms of the expansion. It seems possible that such a bound is given by the expression

$$(4.11) \quad 1 - (1 - \beta)N_T'(0) - (1 - \beta^2)\frac{1}{2}N_T''(0).$$

However, even if proved correct, this would only be useful for values of β in the neighborhood of 1. For example, for $T = 4$ and $|\alpha| = 0.9$, the upper bound would be given by 0.978, while the correct value is only 0.831. It may be observed that the first three terms of the expansion give a much better approximation, namely 0.801, despite the high value of $|\alpha|$.

4.3. The following more general propositions concerning first-order stochastic linear difference equations still await proof:

$$(4.12) \quad N_T(\beta) \leq 1 \quad \text{for } \beta < 1 \text{ and all } T,$$

$$(4.13) \quad \frac{\partial N_T(\beta)}{\partial \beta} > 0 \quad \text{for } \beta < 1 \text{ and all } T,$$

and hence

$$(4.14) \quad \lim_{\beta \rightarrow 1} N_T(\beta) = 1 \quad \text{for all } T.$$

¹It will be remembered that the argument of $N_T(\cdot)$ is β , not $|\alpha|$.

It would also be desirable to investigate analogous problems first for higher-order difference equations and then for equation systems.¹ It would be of interest to see whether the following two propositions are generally true: (1) that a system becomes more strongly damped when estimate expectations are substituted for the structural parameters, and (2) that the relative bias is lower in systems with stronger damping and tends to zero as the characteristic roots of the system approach 1 in absolute value.

These investigations should be carried out for both the stochastic and fixed-variate initial-value case, without the fixed initial value being necessarily zero.

Finally, higher moments of the sampling distributions of the estimates should be investigated. It is probable that these problems will have to be studied with the help of more powerful tools, especially that of the approximate sampling distributions,² provided the upper bound of the approximation error can be determined.

1

At this point one would have to distinguish between the bias of the (maximum-likelihood) estimates of the structural coefficients and the bias of the (least-squares equivalent to maximum-likelihood for the fixed-variate initial-value case) estimates of the regression coefficients. The distributions of the conditional variances and of the disturbance covariance matrix also deserve attention in small samples.

²See [Koopmans, 1942], [Dixon], and [Leipnik].

XVI. MODELS INVOLVING A CONTINUOUS TIME VARIABLE

BY TJALLING C. KOOPMANS

1. In the opinion of this author, the simultaneous-equations method to which the first part of this volume is devoted constitutes an important advance over single-equation methods in the measurement of economic relations. Reasons for this opinion have been stated elsewhere [Koopmans, 1945]. In this note it is my intention to point out certain shortcomings of the new methods in their present stage.

Limitations to the usefulness of the new methods arise from the combination of the following two aspects in the specification of the distribution of the variables:

- (a) the treatment of time as a discrete variable,
- (b) the assumption that disturbances at different points in time are independent.

These aspects are also found in the single-equation methods that have been generally used, and some of the limitations to be mentioned below therefore apply equally to simultaneous-equations and single-equation methods. However, the development of the simultaneous-equations method has brought to light further disadvantages of the specifications (a) and (b) above, and thus has made their revision even more desirable than before.

2. It was found in another article¹ that in the simultaneous-equations method the procedure for estimating the coefficient α_{11} of the endogenous variable x_1 in the first structural equation, say, is different according to whether that variable is "predetermined" or is one of the "jointly dependent variables." A given variable may be in one category or the other depending solely on the timing of that variable in the equation concerned. If an endogenous variable x_1 appears in the first equation only with a lag of one unit behind the most recent timing of that variable in the

¹[II-3.1.3]. For the definition of "endogenous," "predetermined," "dependent," variables, see also [XVII].

equation system as a whole, then in the first equation it will be classified as predetermined. If the time lag is reduced to zero, the variable is in general to be treated as one of the jointly dependent variables. This situation can be reduced *ad absurdum* by making the time unit of measurement smaller and smaller and at the same time reducing the time lag to zero. Then suddenly at the moment the time lag reaches zero the status of the variable x_1 in the first equation is changed, and with it the "unbiased" estimate of the coefficient in question is changed. The solution of this paradox is, of course, that such a procedure is illegitimate. The assumption of independent disturbances in successive observations can be maintained only if the size of the time unit to which these observations refer is not made too small. Therefore, the independence assumption makes the distinction between predetermined and dependent variables appear as absolute instead of a matter of degree, which it would be in a more refined model.

The main source of disturbances is the erratic element in economic behavior. Some causes of erratic behavior (not already represented by measurable variables), like the weather affecting the amount and direction of consumers' expenditure, may be so variable as to reverse themselves in a few days. Other causes like fads and fashions affecting consumption, confidence or lack of confidence affecting investment, may lead to deviations in the same direction for a whole year or even longer. Therefore, as the time unit of observation is reduced in size, a situation in which serial correlation of the disturbances in a given equation can no longer be neglected is bound to arise at some stage.

Methods based on the independence assumption thus involve a lower bound on the permissible size of the time unit of observation. This precludes adequate treatment of a number of important statistical problems in the measurement of economic relations. The time lags occurring in economic behavior are not always integral multiples of one time unit of a size compatible with the independence assumption. They are almost always distributed lags, with the lower limit to the range of lags sometimes practically equal to zero. There is therefore a need for methods of estimating the parameters that characterize lag distributions.

Another problem that cannot be studied adequately under the independence assumption is that of the most economic time unit of observation. Whether it is best to use annual, quarterly, or monthly, data depends on a comparison of the cost of collection of such data (if not already available), the cost of calculating

the necessary estimates, and the information gained (in the sense of smaller sampling errors of estimated parameters) by a given reduction in the size of the unit period of observation. The latter gain in information is likely to be reduced by serial correlation in the disturbances. For time units below a certain size, this gain in information can therefore not be analyzed theoretically on the basis of the independence assumption.

3. An adequate model for the study of the foregoing problem is obtained by considering the disturbances (and therefore the economic variables) as generated by a stochastic process with a continuous time variable. Let us for simplicity assume that this process is normal and stationary, i.e., the values of the variables $u_g(t)$ at any set of time points t_1, t_2, \dots, t_S have a joint normal distribution depending only on the differences $t_2 - t_1, \dots, t_S - t_{S-1}$. In that case the process is entirely characterized [Doob, Theorem 4.3] by the elements

$$(1) \quad \mathcal{E} u_g(t) u_h(t) = \sigma_{gh}(0)$$

of the covariance matrix $\Sigma(0)$, and by a matrix Ξ , determining the lagged covariance matrix

$$(2) \quad \mathcal{E} u_g(t) u_h(t + \tau) = \sigma_{gh}(\tau)$$

through the formula

$$(3) \quad \Sigma(\tau) = \Sigma(0) e^{-\tau \Xi}, \quad \tau > 0.$$

If economic variables $x_n(t)$ are considered as determined by equations in which quantities $u_g(t)$ of this nature are the only random elements, it is necessary to define further the way in which these variables are observed. In practice, the method of observation is again a discrete procedure. One method of observation is to make readings

$$(4) \quad x_n(t_0), x_n(t_1), \dots, x_n(t_S),$$

$$t_S - t_{S-1} = h, \quad s = 1, \dots, S,$$

at equidistant points in time. Price variables are sometimes observed in this way. Quantities of goods and flows of money are

usually observed through averages

$$(5) \quad \bar{x}_n(t_s) = \int_{t_s-h}^{t_s} x_n(\tau) d\tau$$

over a period of observation of length h . For any specified method of obtaining a finite number of observations for each variable in the system, the set of all observations on all variables becomes subject to a joint probability distribution derivable from (3) (or from whatever other process of a continuous time variable t is specified).

Although the mathematical difficulties involved may be considerable, a model of this kind would provide a means of studying the estimation of lag distributions and the choice of the most economic time unit of observation.

4. Perhaps the most important advantage of a continuous treatment of time has not yet been mentioned. It is explained in another article [II-2.5.6] that, in the discrete case, as soon as the independence assumption for successive disturbances is dropped, the problem of identification of the structural equations is greatly complicated. For the removal of that assumption may open up a new group of transformations of the equations (involving shifts along the time axis) that preserve the probability distribution of the variables.

The introduction of a continuous time variable is perhaps the best way to study fully all aspects of the identification problem of relations between economic time series. Consider for instance a system of one linear equation containing only one variable x . In the discrete formulation this equation would be of the type

$$(6) \quad \mathcal{L}x(t) \equiv x(t) + \alpha_1 x(t-1) + \cdots + \alpha_\tau x(t-\tau) = u(t).$$

Under the independence assumption for $u(t)$,

$$(7) \quad \mathcal{E}u(t)u(t+\theta) = 0 \quad \text{if } \theta \neq 0,$$

there is no identification problem. For any linear combination of the type

$$(8) \quad \lambda_0 \mathcal{L}x(t) + \lambda_1 \mathcal{L}x(t-1) + \cdots + \lambda_K \mathcal{L}x(t-K) = \bar{v}(t) = \sum_{k=0}^K \lambda_k u(t-k)$$

introduces serial correlation into $v(t)$ unless all but one of the quantities λ_k vanish.

In the continuous formulation the equation is of the type

$$(9) \quad \mathcal{L}^*x(t) \equiv x(t) - \int_{-\infty}^t \varphi(t - \tau) x(\tau) d\tau = u(t)$$

where the disturbance process, if normal and stationary, is described by

$$(10) \quad \mathcal{E}u(t) u(t + \theta) = f(|\theta|).$$

Now there are infinitely many transformations in the space of the functions φ and f which preserve the form of (9). The simplest of these is obtained by substituting

$$(11) \quad \int_{-\infty}^t \varphi(t - \tau) x(\tau) d\tau + u(t)$$

for $x(t)$ under the integral sign in (9). Another possibility is first to write

$$(12) \quad \varphi(t - \tau) = \varphi_1(t - \tau) + \varphi_2(t - \tau),$$

and to substitute (11) only in one of the two integrals so obtained, etc. The only invariants of all these transformations are the autocovariance function

$$(13) \quad g(\theta) = \mathcal{E}x(t) x(t + \theta)$$

of the variable $x(t)$ and all its functions and functionals. These are therefore the only identifiable characteristics of the process (9), and only these are subject to estimation.

The identification difficulties just described are absent from the process

$$(14) \quad x(t) = \int_{-\infty}^t \varphi(t - \tau) x(\tau) d\tau + \alpha y(t) + u(t)$$

containing an observable exogenous variable $y(t)$. They reappear if $y(t)$ occurs in the form

$$(15) \quad x(t) = \int_{-\infty}^t \varphi(t - \tau) x(\tau) d\tau + \int_{-\infty}^t \psi(t - \tau) y(\tau) d\tau + u(t)$$

with unknown lag-distribution functions φ and ψ .

The process (9) becomes a system of equations if $x(t)$ and $u(t)$ are interpreted as column vectors, $\varphi(t - \tau)$ as a matrix. In this case, the identification of individual equations is aided if economic considerations require or permit the specification that all diagonal elements of $\varphi(t - \tau)$ shall vanish for all values of $t - \tau$. In two dimensions this leads to the system

$$(16) \quad \begin{aligned} x_1(t) &= \int_{-\infty}^t \varphi_{12}(t - \tau) x_2(\tau) d\tau + u_1(t), \\ x_2(t) &= \int_{-\infty}^t \varphi_{21}(t - \tau) x_1(\tau) d\tau + u_2(t). \end{aligned}$$

The equations (16) are completely identified, i.e., there is no transformation other than the identity, in the space of the functions φ_{12} , φ_{21} and f_1 , f_2 [defined as in (10)] which preserves the form of (16).

PART THREE

SPECIFICATION OF HYPOTHESES

XVII. WHEN IS AN EQUATION SYSTEM COMPLETE FOR STATISTICAL PURPOSES?

BY TJALLING C. KOOPMANS

	Page
1. Endogenous and Exogenous Variables	393
2. Exogenous Variables in Economic Theory	393
3. Statistical Definition of Exogenous Variables	394
4. Endogenous and Exogenous Variables in Systems Containing Time Lags	399
5. The Nature of Exogenous Variables	402
6. Predetermined Variables and Jointly Dependent Variables	402
7. Summary of the Classification of Variables	405
8. Insufficiency of the "Approximate" Causal Principle	405
9. Time Lags as a Criterion of Classification	407
10. Identification Problems in Systems Containing Exogenous or Other Predetermined Variables	407

1. *Endogenous and Exogenous Variables*

In static or dynamic economic theory, the criteria employed in determining whether or not a system of equations is complete are derived from the purpose for which such systems are constructed: the explanation of economic phenomena. In each case a distinction is drawn between the *endogenous* variables that the economist sets out to explain and the *exogenous* variables that he takes as given. The number of equations required for the explanation of the values or, in the dynamic case, of the movements of the endogenous variables, then equals the number of such variables.

2. *Exogenous Variables in Economic Theory*

In determining which variables are set aside as exogenous, two main principles are implicitly or explicitly applied in economic literature. They might be described as the departmental principle and the causal principle. The departmental principle treats as exogenous those variables which are wholly or partly outside the scope of economics, like weather and climate, earthquakes, popula-

tion, technological change, political events. The causal principle, which does not always lead to the same result, regards as exogenous those variables which influence the remaining (endogenous) variables but are not influenced thereby.

The causal principle is often used also if it applies only approximately, that is, if the influence of the endogenous variables on those treated as exogenous is presumed to be small. For instance, in explaining the level of employment in a country which has only a small share in world trade, the shifts in the schedules of foreign demand for its exports and of foreign supply of its imports are sometimes treated as exogenous in first approximation. Another example is found in the formation of quantity and price of a consumers' good that attracts only a small fraction of consumers' expenditure. In such cases, consumers' income is often taken as an exogenous variable, operating at the demand side, although of course consumers' income itself depends on the demand for *all* commodities. In order to distinguish between cases where the causal principle of classification of variables is strictly or only approximately applicable, we shall in what follows speak of the strict causal principle and the approximate causal principle.

There is no sharp line of demarcation between the application of the approximate causal principle and what deserves mention as a third principle or consideration: the purpose of exposition. At a certain stage of the analysis, variables are often treated as exogenous to facilitate understanding of the model studied, reserving for later elaboration their inclusion among the endogenous variables.

3. *Statistical Definition of Exogenous Variables*

One of the main purposes of the present volume is to study the statistical implications of the fact that economic data are governed by a system of simultaneous equations. It is therefore necessary to review the foregoing principles, and such other considerations as may present themselves, in relation to the purpose of statistical estimation of the equations of dynamic economics. The question which, if any, are exogenous variables must be raised afresh from the statistical point of view.

It will be clear that the departmental principle is of no value in this connection. Suppose that rainfall and temperature enter as determining factors in one or more equations of the system. Is it possible that a certain method of estimation, applied to the economic relations with disregard of purely physical relations between

these and other meteorological variables, for that very reason leads to inconsistent and asymptotically biased estimates¹ of the parameters of the economic relations? The answer to this question obviously does not depend on which particular department is in charge of studying such relations.

The causal definition of exogenous variables does permit an answer to this question, with specific reference to the maximum-likelihood method of estimation. The answer is based on a mathematical observation which we shall now formulate. Suppose that

$$(1) \quad \varphi_n(\alpha_{n1}, \dots, \alpha_{nQ_n}; x_1, \dots, x_N) = u_n, \quad n = 1, \dots, N,$$

represents the complete system of all structural equations (linear or otherwise) between *all* variables, economic or noneconomic, that enter directly or indirectly into the explanation of economic variables. Suppose, for the present, that all these variables enter without time lags, and let the random terms u_n have for each time t a joint continuous distribution

$$(2) \quad f(u_1, \dots, u_N) du_1 \cdots du_N,$$

these distributions being independent for successive values $t = 1, 2, \dots$ of t . Suppose finally that both the equations (1) and the variables x_1, \dots, x_N , can be separated into two sets, with numbering $n = 1, \dots, G$ and $n = G + 1, \dots, N$ for the first and second sets, respectively, such that the following three assumptions are satisfied.

ASSUMPTION 3.1. *The first set of "endogenous variables" does not occur in the second set of equations:*

$$(3) \quad \begin{cases} (3a) & \varphi_n(\alpha_{n1}, \dots, \alpha_{nQ_n}; x_1, \dots, x_G; x_{G+1}, \dots, x_N) = u_n, \\ & n = 1, \dots, G, \\ (3b) & \varphi_n(\alpha_{n1}, \dots, \alpha_{nQ_n}; x_{G+1}, \dots, x_N) = u_n, \\ & n = G + 1, \dots, N. \end{cases}$$

¹An estimate of a parameter α , derived from a sample of size T , is called consistent if, for any $\epsilon > 0$, $\lim_{T \rightarrow \infty} P(|a - \alpha| > \epsilon) = 0$ if $P(E)$ denotes the

ASSUMPTION 3.2. *The distribution function in (2) can be factorized as follows:*

$$(4) \quad f(u_1, \dots, u_N) = f_1(u_1, \dots, u_G) f_2(u_{G+1}, \dots, u_N).$$

ASSUMPTION 3.3. *In any point x_1, \dots, x_N permitted by the distribution function (2), the Jacobian $\partial(u_1, \dots, u_N)/\partial(x_1, \dots, x_N)$ of the "transformation" (3) differs from zero everywhere in the space of the parameters α_{nq} , $q = 1, \dots, Q_n$; $n = 1, \dots, N$.*

The economic meaning of these assumptions is that (Assumption 3.1) we isolate in (3b) the equations that connect only exogenous variables, that (Assumption 3.2) we assume that the random elements entering those equations – and hence the (exogenous) variables of the second set themselves – are distributed independently from the random disturbances in the equations (3a) explaining the endogenous variables, and that (Assumption 3.3) no situation can arise in which small changes in the disturbances u_1, \dots, u_N lead to very large changes in the variables x_1, \dots, x_N .

Since the equation system (3) together with (4) completely specifies the joint distribution of the variables x_1, \dots, x_N , maximum-likelihood estimation based on (3) and (4) may be expected to lead to consistent and asymptotically unbiased estimates of all identifiable¹ parameters.² Assuming this to be the case, we shall now prove that in consequence *the identifiable parameters of the equation (3a) can be estimated, by the method of maximum likelihood, i.e., consistently and without bias in large samples, from those equations only, treating the "exogenous" variables x_{G+1}, \dots, x_N as if they were fixed in repeated samples.* To prove this point, let us regard (3a) and (3b) as one system of equations (3) and write down the joint distribution function h of a complete set of observations for one value of t . This function is given by

probability of an event E . The estimate a is called asymptotically unbiased if $\lim_{T \rightarrow \infty} \mathbb{E}a = \alpha$, where \mathbb{E} denotes the mathematical expectation.

¹For the concept of identifiability, see [II-2], also [IV].

²This has been proved rigorously for linear systems; see [II-3.3].

$$\begin{aligned}
 & h(\alpha_{11}, \dots, \alpha_{N0_N}; x_1, \dots, x_N) dx_1 \cdots dx_N \\
 (5) \quad & = f(u_1, \dots, u_N) \frac{\partial(u_1, \dots, u_N)}{\partial(x_1, \dots, x_N)} dx_1 \cdots dx_N,
 \end{aligned}$$

in which the quantities u_1, \dots, u_N are regarded as functions (1) of x_1, \dots, x_N . Because of the assumed form of (3b), the Jacobian in (5) factorizes as follows:

$$\begin{aligned}
 (6) \quad \frac{\partial(u_1, \dots, u_N)}{\partial(x_1, \dots, x_N)} &= \begin{vmatrix} \frac{\partial u_1}{\partial x_1} & \cdots & \frac{\partial u_G}{\partial x_1} & 0 & \cdots & 0 \\ \cdot & \cdots & \cdot & \cdot & \cdots & \cdot \\ \frac{\partial u_1}{\partial x_G} & \cdots & \frac{\partial u_G}{\partial x_G} & 0 & \cdots & 0 \\ \frac{\partial u_1}{\partial x_{G+1}} & \cdots & \frac{\partial u_G}{\partial x_{G+1}} & \frac{\partial u_{G+1}}{\partial x_{G+1}} & \cdots & \frac{\partial u_N}{\partial x_{G+1}} \\ \cdot & \cdots & \cdot & \cdot & \cdots & \cdot \\ \frac{\partial u_1}{\partial x_N} & \cdots & \frac{\partial u_G}{\partial x_N} & \frac{\partial u_{G+1}}{\partial x_N} & \cdots & \frac{\partial u_N}{\partial x_N} \end{vmatrix} \\
 &= \frac{\partial(u_1, \dots, u_G)}{\partial(x_1, \dots, x_G)} \cdot \frac{\partial(u_{G+1}, \dots, u_N)}{\partial(x_{G+1}, \dots, x_N)}.
 \end{aligned}$$

Combining corresponding factors in (6) and (4), we obtain the following factorization for h :

$$(7) \quad h(\alpha_{11}, \dots, \alpha_{N Q_N}; x_1, \dots, x_N) \\ = h_1(\alpha_{11}, \dots, \alpha_{G Q_G}; x_1, \dots, x_N) h_2(\alpha_{G+1,1}, \dots, \alpha_{N Q_N}; x_{G+1}, \dots, x_N),$$

in which the parameters α_{nq} occurring in each factor of (6) have been determined from (3).

What happens if the likelihood function

$$(8) \quad F \equiv \prod_{t=1}^T h(\alpha_{11}, \dots, \alpha_{N Q_N}; x_1(t), \dots, x_N(t))$$

following from (7) is maximized with respect to the parameters α_{nq} ? To answer this question, it is necessary to remember that the factors f_1 and f_2 in (4) are nonnegative, and are positive somewhere within their range, in view of their nature as probability densities. Since the nonvanishing Jacobians in (6) enter in (5) and hence in (7) only with their absolute values, the factors h_1 and h_2 in (4) are likewise nonnegative and somewhere positive. Furthermore, the two factors $F_1 = \prod_t h_1$ and $F_2 = \prod_t h_2$ of (8) depend on entirely different sets of parameters, $\alpha_{11}, \dots, \alpha_{G Q_G}$ and $\alpha_{G+1,1}, \dots, \alpha_{N Q_N}$, respectively. It follows that, for the product $F = F_1 F_2$ to be as large as possible, it is necessary and sufficient that each of the factors F_1 and F_2 separately be as large as possible.

The problem of (asymptotically unbiased) estimation of the parameters by the method of maximum likelihood has thereby been split into two separate estimation problems for the two sets of parameters occurring in (3a) and (3b), respectively. In particular, it is possible to disregard entirely the estimation of the equations (3b) and to estimate only the parameters of (3a) by maximizing F_1 . The latter problem is identical with the estimation problem met with if the exogenous variables x_{G+1}, \dots, x_N in the equations (3a) are regarded not as determined by (3b), but as given in advance and fixed in repeated samples.

It might be observed that although under the present assumptions the maximum-likelihood estimates of the parameters occurring

in (3a) are the same functions of the observations whether the exogenous variables are regarded as subject to a probability distribution or as fixed in repeated samples, a difference arises again between these two cases in the distribution of the maximum-likelihood estimates at least in small samples. However, a procedure is available whereby this difference can be removed. In the case where the exogenous variables are subject to a probability distribution, it is both permissible and useful to consider, for the construction of confidence intervals, only that subclass of all possible samples in which the exogenous variables have values equal to those observed.¹

The foregoing proposition provides a justification for using the concept of exogenous variables defined according to the causal principle. In equation systems in which all variables enter simultaneously (the only case covered so far), a clear separation has now been obtained between the endogenous variables that must be explained by an equal number of relations (all of which are in principle relevant to the estimation of any one of them) and the exogenous variables that may be accepted without explanation. It is worth stressing again that for purposes of statistical estimation the concept of exogenous variables must be defined more strictly and narrowly than for some purposes of economic theory. Whereas the theorist may at a certain stage choose, for reasons of approximation or exposition, to disregard possible influences exerted by variables "inside" his system on variables regarded by him as "outside," the statistician must be convinced (on a priori grounds or as a result of statistical test) of the absence of such influence before he can declare the "outside" variable to be exogenous in the foregoing sense.

4. *Endogenous and Exogenous Variables in Systems Containing Time Lags*

We must now take into consideration the fact that the action of one variable on another is often subject to a time lag. Suppose therefore that the variables x_1, \dots, x_n occur in the equations (1) not only with timing t , but also with various time lags τ which are integral nonnegative multiples of the chosen unit of time:

¹For a discussion of this device, see [Hotelling, 1940].

$$\varphi_n\{\alpha_{n1}, \dots, \alpha_{nQ_n}; x_1(t), \dots, x_N(t); x_1(t-1), \dots, x_N(t-1); \dots\}$$

$$(9) \quad = u_n(t), \quad n = 1, \dots, N; \quad t = 1, \dots, T.$$

The joint distribution of all observations $x_n(t)$, $n = 1, \dots, N$; $t = 1, \dots, T$, again follows from that of the disturbances $u_n(t)$, $n = 1, \dots, N$; $t = 1, \dots, T$, provided we specify how any values $x_n(t-\tau)$ with $t-\tau \leq 0$, occurring in the equations (1) for $t \geq 1$, are distributed. We shall assume such values to be constant in repeated samples and equal to the values observed in the sample at hand.

The first question requiring an answer under the present assumption is whether maximum-likelihood estimation based on the joint distribution function of all observations $x_n(t)$, $n = 1, \dots, N$; $t = 1, \dots, T$, still leads to consistent and asymptotically unbiased estimates. This does not follow from the general theory of maximum-likelihood estimation because the observations $x_n(t)$ for successive values of t are no longer stochastically independent. However, Mann and Wald [1943] have proved that for linear systems of the form (9), maximum-likelihood estimates are consistent and asymptotically unbiased. Presuming that this result can be extended to nonlinear systems, we shall therefore continue the analysis on the basis of the maximum-likelihood method of estimation.

It is easily seen that the Jacobian

$$(10) \quad J_T \equiv \frac{\partial \{u_1(1), \dots, u_N(1); \dots; u_1(T), \dots, u_N(T)\}}{\partial \{x_1(1), \dots, x_N(1); \dots; x_1(T), \dots, x_N(T)\}}$$

of the transformation from all disturbances to all observations (that are not constants) factorizes as follows

$$J_T = \begin{vmatrix} \frac{\partial\{u_1(1), \dots, u_N(1)\}}{\partial\{x_1(1), \dots, x_N(1)\}} & \frac{\partial\{u_1(2), \dots, u_N(2)\}}{\partial\{x_1(1), \dots, x_N(1)\}} & \dots & \dots \\ 0 & \frac{\partial\{u_1(2), \dots, u_N(2)\}}{\partial\{x_1(2), \dots, x_N(2)\}} & \dots & \dots \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \frac{\partial\{u_1(T), \dots, u_N(T)\}}{\partial\{x_1(T), \dots, x_N(T)\}} \end{vmatrix} \quad (11)$$

$$= \prod_{t=1}^T \frac{\partial\{u_1(t), \dots, u_N(t)\}}{\partial\{x_1(t), \dots, x_N(t)\}} = \prod_{t=1}^T J(t),$$

say, all entries to the right of the main diagonal in J_T being immaterial. Consequently, the joint distribution function of the observations is again given by the expressions (8) and (5), the only difference being that the arguments u_1, \dots, u_N of f in (5) are now given by (9) instead of (1).

By going through the previous reasoning after this change in the structural equations, it will be seen that the previous result remains true: If (9) has the form (3), where u_n now stands for $u_n(t)$, and each of the symbols x_n stands for a sequence of values $x_n(t), x_n(t-1), \dots$, maximum-likelihood estimation of the parameters $\alpha_{11}, \dots, \alpha_{GQ_G}$ can be performed by maximizing only the function

$$(12) \quad F_1 = \prod_{t=1}^T f_1\{u_1(t), \dots, u_G(t)\} \left| \frac{\partial\{u_1(t), \dots, u_G(t)\}}{\partial\{x_1(t), \dots, x_G(t)\}} \right|,$$

with $u_g(t)$ given by (3a).

Thus the distinction between endogenous and exogenous variables

is maintained equally well in systems containing time lags. Further analysis needs to be concerned only with the function (12) which represents the conditional probability density in the space of all values of the endogenous variables x_1, \dots, x_G within the period $t = 1, \dots, T$, for given values of the exogenous variables x_{G+1}, \dots, x_N within that period, and of all variables x_1, \dots, x_N previous to that period ($t \leq 0$). Distribution theory of the estimates of the parameters α_{gg} , $g = 1, \dots, G$, can also be based entirely on the conditional distribution function (12). It is therefore unnecessary from now on to carry the "given" values of the exogenous variables x_{G+1}, \dots, x_N along in all formulae, and we shall rewrite (3a) as

$$\varphi_g\{\alpha_{g1}, \dots, \alpha_{gG}; x_1(t), \dots, x_G(t); x_1(t-1), \dots, x_G(t-1); \dots\} = u_g(t),$$

(13)

$$g = 1, \dots, G; \quad t = 1, \dots, T.$$

5. *The Nature of Exogenous Variables*

Before continuing the analysis on the basis of (12) and (13), the question is in order as to the nature of the variables so set aside as "exogenous." Which factors in man's physical and historical environment are not influenced by his economic activity? If the question is put in this way one can think of little else besides changes in weather, climate, geology, and geography that are brought about by natural causes. There remains a host of sociological, political, and psychological factors that are in continuous interaction with economic activity, and therefore cannot, on any grounds so far adduced, be accepted as they come without incorporating the explanation of their fluctuations in the system of equations.

6. *Predetermined Variables and Jointly Dependent Variables*

A further delimitation of the area that needs to be covered by an equation system, for purposes of unbiased estimation of parameters in large samples, can sometimes be obtained through closer study of the manner in which time lags occur in the structural equations (13).

The assumption that the variables x_1, \dots, x_G are endogenous implies that each variable occurs without time lag in at least one equation, because otherwise the equation system (13) would not describe the determination of each variable. The earlier values $x_i(t-1), x_i(t-2), \dots, i = 1, \dots, G$, of the endogenous variables may accordingly be called predetermined variables. At time t their values are already given by the equations (13) for earlier values of t . Any possibility of again splitting off from the system (13) a subset of equations for separate treatment depends, as before, on further factorizing properties of the Jacobian $J(t)$ referring to one single time point t as defined in the last line of (11). This leads immediately to the following extension of the observations made above. Suppose that both the equations (13) and the (endogenous) variables $x_i(t), i = 1, \dots, G$, can again be separated into two sets, (13a) and (13b), such that instead of the previous Assumptions 3.1, 3.2, 3.3 the following assumptions are satisfied:

ASSUMPTION 6.1. *The first set of variables $x_i(t), i = 1, \dots, G_1$, does not occur in the second set ($g = G_1 + 1, \dots, G$) of equations (13) except with a time lag $\tau \geq 1$:*

$$\begin{aligned} & \varphi_g \{ \alpha_{g1}, \dots, \alpha_{gQ_g}; x_1(t), \dots, x_{G_1}(t); \\ (13a) \quad & x_1(t-1), \dots; x_{G_1}(t-1); \dots \} = u_g(t), \\ & g = 1, \dots, G_1; \quad t = 1, \dots, T, \end{aligned}$$

(13)

$$\begin{aligned} & \varphi_g \{ \alpha_{g1}, \dots, \alpha_{gQ_g}; x_{G_1+1}(t), \dots, x_G(t); \\ (13b) \quad & x_1(t-1), \dots, x_{G_1}(t-1); \dots \} = u_g(t), \\ & g = G_1 + 1, \dots, G; \quad t = 1, \dots, T, \end{aligned}$$

ASSUMPTIONS 6.2 and 6.3. *The previous Assumptions 3.2 (factorizing of the distribution function of disturbances) and 3.3 (non-vanishing Jacobian), respectively, are again satisfied for the new subdivisions of variables and equations.*

Then it follows that *the parameters of the two sets of equations can again be estimated by the maximum-likelihood method in two separate steps.* In particular, the first set of equations can be estimated without any knowledge of the second set beyond that expressed by the Assumptions 6.1, 6.2, 6.3.

If the conditions of this proposition are satisfied, the endogenous variables of the second set $x_{G_1+1}(t), \dots, x_G(t)$ will be called the predetermined variables even where they occur without time lags. For even their values $x_g(t)$ without lags are determined by *earlier* values of the variables of both sets and by disturbances independent of the *simultaneous* values of the variables of the first set. These characteristics permit the predetermined variables to be treated in the first set of equations *as if* they were exogenous variables for the purpose of the separate determination of maximum-likelihood estimates of the parameters α_{gq} , $q = 1, \dots, Q_g$, $g = 1, \dots, G$, of the first set of equations.

With regard to the distribution of these estimates the position is slightly more complicated than in the case where only the exogenous variables and the equations explaining these variables are split off. In the present case the small-sample distribution of the maximum-likelihood estimates a_{gq} of the parameters α_{gq} , $g = 1, \dots, G_1$, of the first set of equations cannot be made independent of the parameters α_{gq} , $g = G_1 + 1, \dots, G$, of the second set of equations by any such construction or restriction as was applied in the case of exogenous variables. The asymptotic theory of maximum-likelihood estimation for stochastic difference equations, reviewed in [II-3.3], implies, however, that at least for linear systems with normally distributed disturbances, the limiting distribution of the a_{gq} , $g = 1, \dots, G_1$, for infinitely large samples is independent of the parameters α_{gq} , $g = G_1 + 1, \dots, G$, and is such as to make the estimates involved both consistent and asymptotically unbiased.

The distinction of predetermined variables thus makes possible a great further reduction in the size and scope of the equation systems needed for statistical purposes in cases where large-sample approximations are adequate.¹ In order to be "complete" in this "large-sample" statistical sense, the equation system then only needs to cover with reasonable completeness the formation of those

¹See, however, [XV] where L. Hurwicz discusses in a simple case the bias involved in the application of large-sample approximations to samples of moderate size.

economic (and possibly some noneconomic) variables that are connected through a process of *instantaneous* interactions. The variables between which these instantaneous interactions occur will be called the *jointly dependent variables*, or, briefly, the *dependent variables*. Any further determining variables affecting this process of interaction, which variables themselves depend only on earlier values of the endogenous variables and on chance variation independent of that process, need not be "explained" through additional equations, but can be treated as predetermined variables. In particular, slow changes in environment affecting consumers' preferences, technological changes, and other factors usually described as "trends" can probably be treated as predetermined variables in this sense with a reasonable degree of accuracy. Likewise, the interactions between economic variables and political developments is often subject to sizable time lags, and in such cases need not be taken into account for the purposes here considered.

7. Summary of the Classification of Variables

The foregoing distinctions are summarized in the following table¹ which also provides a comparison of the notation of the present article with that introduced in two different places in [II].

8. Insufficiency of the "Approximate" Causal Principle

Further comment is needed on some of the principles underlying this classification. Both the distinction between exogenous and endogenous variables and that between predetermined and dependent variables are based on a subdivision of the complete set of equations "explaining" the formation of all variables into two subsets of equations. In both cases it is necessary to stipulate (see Assumption 3.3) that the disturbances affecting equations thereby placed into different subsets should be independently distributed. The necessity of this assumption definitely limits the extent to which equation systems can be reduced in size without sacrificing "completeness" in the statistical sense. In particular, this assumption frequently prohibits the application of the "approximate" causal principle in cutting down the size of the system. This point is of great importance, for instance, to the statistical measurement of demand and supply curves for individual commodities. For this reason, and because of the need for further mathematical demonstration of the insufficiency of the "approximate" causal

¹See p. 406.

TABLE 7.1.

Notation introduced in II-1.3	CLASSIFICATION OF VARIABLES WITH NOTATION USED IN THIS ARTICLE		Notation introduced in II-2.1.2
	<u>Distinction independent of timing of variable</u>	Distinctions dependent on timing of variable in the equations	
<p>$y_i(t-\tau), \tau \geq 0$ $i = 1, \dots, G$</p> <p>I. Endogenous variables (affected by variables under I and II), $x_n, n = 1, \dots, G$</p> <p>$z_k(t-\tau), \tau \geq 0$ $k = 1, \dots, K$</p> <p>II. Exogenous variables (not affected by the variables under I), $x_n, n = G+1, \dots, N$</p>	<p>A. occurring without time lag, $x_i(t), i = 1, \dots, G$</p> <p>B. occurring with a time lag, $x_i(t-\tau), \tau \geq 1,$ $i = 1, \dots, G$</p>	<p>a. in instantaneous interaction with each other, $x_i(t), i = 1, \dots, G_1$</p> <p>b. affected only by variables II and IB and by random disturbances independent of variables IAa, $x_i(t), i = G_1+1, \dots, G$</p>	<p>Jointly dependent variables $y(t)$</p> <p>Predetermined variables $z(t)$</p>

principle, a separate discussion of this problem will be given elsewhere in due course.

9. *Time Lags as a Criterion of Classification*

The question must be raised whether the concept of a time lag is sharp enough to serve as the basis for a classification of variables that determines their treatment in the estimation procedure. After all, time lags may occur in all sizes, varying continuously down to zero.

It is possible to give a provisional answer to this objection which safeguards the statistical procedures developed in this volume as logically defensible and consistent. It has been assumed in those procedures that disturbances in successive observations (that is, in successive time points or time intervals) are independent. That assumption cannot continue to be valid if the unit of observation is made smaller and smaller. Therefore, the term "time lag" in the previous discussion means a time lag not smaller than the smallest time unit of observation for which the independence assumption is still tenable as a fair approximation to reality.

Nevertheless, such a discrete treatment of what is a continuous time variable should not be accepted as the final word of statistical theory applied to economic time series. In this writer's opinion, an extension of the "discrete" methods of this volume to permit continuously distributed lags, ranging down to zero wherever appropriate, is the next important refinement needed in the development of statistical methods adapted to the analysis of relations between economic time series. In [XVI] the need for and the nature of such a generalization are discussed further.

10. *Identification Problems in Systems Containing Exogenous or Other Predetermined Variables*

In the foregoing discussion we have considered two successive reductions of the size of the original equation system (3), where x_n now, as in section 4, stands for a sequence of values $x_n(t)$, $x_n(t-1)$, In the first reduction the equations (3b) connecting only exogenous variables are omitted, and the equations (3a) or (13) are retained. In the second reduction the subsystem (13b) "explaining" the predetermined variables is omitted, and the remain-

ing equations (13a) are retained. In order to provide a link with the discussion of identification problems in II-2, the following question should be raised: Could any equation of (13a), or any parameter of such an equation, which appears to be identifiable on the basis of a study of transformations in the space of the parameters occurring in (13a) only, be found nonidentifiable if possible transformations in the space of all parameters entering in the complete system (3), or in its first subsystem (13), were taken into consideration? The answer is that, under the various independence assumptions stated above, this cannot occur as long as the distribution of the disturbances $u_g(t)$ is nonsingular, in the sense that no exact relationship of the type

$$(14) \quad \psi\{u_{g_1}(t), u_{g_2}(t), \dots\} = 0$$

is satisfied by the disturbances of either of the subsystems (3b) or (13b). For, if the distribution of the disturbances is nonsingular, any attempt to combine equations from (3b) or (13b) with those of (13a) to produce new equations of the form (13a) will violate the independence assumptions as between the disturbances in (3a) and (3b) or in (13a) and (13b), respectively. Under the present assumptions, therefore, nonsingularity of the distribution of disturbances is the only condition that needs to be met by the complete system (3) in order that questions of identifiability can be settled by study of the appropriate subsystem (3a) or (13a). In order to visualize the meaning of this condition, let us assume that a relation (14) holds for the disturbances of the equations (3b) "explaining" the exogenous variables. By inserting the left hand members of (3b) for $u_g(t)$ in (14), we find that then an exact relation

$$(15) \quad \chi\{x_{g+1}(t), \dots, x_n(t)\} = 0$$

connects the exogenous variables. This case was excluded from the discussion of identification problems in II-2 by the wording of Definition 2.1.2. Whether the presence of a relation (15) does in fact destroy identifiability of otherwise identifiable parameters of (3a) depends on the functional form of (15) in connection with any a priori information about the functional form of the equations (3a). For instance, if (15) is linear, the identifiability of any linear equation in (3a), containing the same variables as (15) with-

out restrictions on their coefficients, is destroyed. The occurrence of such a case in applications should reveal itself, theoretically through infinite, practically through very high calculated, standard errors of the estimates of parameters whose identification is destroyed.

XVIII. SYSTEMS WITH NONADDITIVE DISTURBANCES¹

BY LEONID HURWICZ

0.1. In this note the need is shown for generalization of the manner in which disturbances are introduced in stochastic equation systems. While at present in the systems treated each equation contains only *one* disturbance, in more realistic systems there might be equations each containing *several* disturbances. Moreover, some of these disturbances might enter in a nonadditive fashion as coefficients of the observed variables.

A model with nonadditive disturbances, while often appropriate on grounds of a priori (economic) knowledge, is much more difficult to handle mathematically. In fact the difficulties arise at an extremely early stage: it has not yet been possible, in general, to derive formulae on which the computation of estimates could be based. The maximum-likelihood estimates have been obtained for a lagless single-equation system, i.e., where all variables but one are fixed². Even in this simplest case the equation for the unknown estimates is, in general, nonlinear.

But when a simultaneous system is considered, even the likelihood function cannot, in general, be obtained except in the form of an integral. Hence, unless more powerful analytical methods can be found, the only practical solution (if the maximum-likelihood principle is to be used) would seem to be to resort to numerical methods³ in evaluating the likelihood function and locating its maximum.

The same would appear to hold for the method of moments, although it may be that this problem deserves a more careful investigation. One may hope that other criteria for estimation can be found, but if not, the numerical methods applied to the maximum-likelihood criterion of estimation would seem to be the best way out.

¹Part of the work done on this paper was done in 1945-46 during the author's tenure of the Guggenheim Memorial Fellowship.

²See [XIX].

³Perhaps with the aid of the recently designed electronic computing machines.

The next stage would be that of investigating the sampling properties of the estimates used. This would also be a work of considerable complexity.

0.2. This note is divided into three sections. Section 1 is devoted to the specification of the nonadditive disturbance model. Section 2 provides some examples, mostly from the field of economics, of nonadditive disturbance models. In section 3 the likelihood function is derived for a special and relatively simple case of a two-equation model where one of the equations contains a nonadditive disturbance.

1.1. The simultaneous linear systems treated by methods of statistical analysis as a rule belong to the category¹

$$(1.1) \quad \varphi_g(y', z') \equiv \varphi_g(x') = u_g, \quad g = 1, 2, \dots, G,$$

where $y = (y_1, \dots, y_{K_y})$, $z = (z_1, \dots, z_{K_z})$, $G = K_y$; y is stochastic and z a fixed variate. We observe $x \equiv (y, z)$ but not the disturbances u_g . The function (vector) $\Phi = (\varphi_1, \dots, \varphi_G)$ and the covariance matrix Σ^u of the disturbances are partly specified by a priori knowledge (the identifying restrictions) and are partly to be estimated from the observations on x . In the linear case (1.1) becomes

$$(1.2) \quad \alpha_g x' = u_g, \quad g = 1, 2, \dots, G,$$

where α_g is the g th row of the structural coefficient matrix A .

Now (1.1) is but a special case of

$$(1.3) \quad \psi_g(x', u') = 0, \quad g = 1, 2, \dots, G = K_y,$$

where $u = (u_1, \dots, u_{K_u})$. (1.1) can be obtained from (1.3) by assuming

$$(1.4) \quad K_u = K_y,$$

$$(1.5) \quad \frac{\partial \psi_g}{\partial u_p} = \delta_{gp},$$

¹ y' is the transpose of the row vector y , etc.

where δ_{gi} is the Kronecker symbol.

If (1.5) were abandoned but (1.4) retained, it would still be possible to solve for the u 's and rewrite (1.3) as

$$(1.6) \quad \psi_g^*(x') = u_g, \quad g = 1, 2, \dots, G,$$

although equations (1.6) would no longer be the original behavior equations. This type of situation, because of its nonlinear nature, would be of considerable mathematical complexity. The Jacobian of the transformation (from u to y) becomes a function of y , rather than a constant as in (1.2). Thus, while it is not difficult to obtain the likelihood function, it will, in general, be quite difficult to find the maximizing values of the unknown parameters.

But when K_u is permitted to exceed K_y ($=G$) an additional difficulty arises since the distribution of y is then essentially a marginal one; an integration must be performed in order to obtain the likelihood function. The difficulty of this integration is due to the fact that the integrand contains the absolute value of a Jacobian and the latter is a function of the variables of integration (see below, section 3.2).

1.2. Let (1.3) hold with $K_u > K_y$ so that the equations (1.4) and (1.5) are not satisfied. Nevertheless, the disturbances might still enter (1.3) in a purely additive manner. For instance, we might have the system

$$(1.7) \quad \begin{aligned} y_1 + \alpha_{12} y_2 + u_{11} + \varepsilon u_{12} &= 0, \\ \alpha_{21} y_1 + y_2 + u_{21} &= 0. \end{aligned}$$

This, however, is a trivial case and it will henceforth be assumed that any number of additive disturbances is treated as one single disturbance; thus in the above case we would define $u_1^* \equiv u_{11} + \varepsilon u_{12}$ and say that $K_y = K_u = 2$.

1.3. The simplest general case of nonadditive disturbances is given by the equation system

$$(1.8) \quad (\alpha_g^y + u_g^y) x' \equiv (\alpha_g^y + u_g^y) y' + (\alpha_g^z + u_g^z) z' = 0, \\ g = 1, 2, \dots, G = K_y,$$

where u_g is the g th row of the disturbance matrix $U \equiv [u_{gs}]$, $g = 1, 2, \dots, G$, $1 \leq s \leq K_x$; it is assumed, without loss of generality, that $E(U) = 0$; α_g is, as before, the g th row of the structural coefficient matrix A .

The linear additive disturbance case (1.2) can be obtained from (1.8) by making all the elements of U except one column in u^z corresponding to, say, $z_{K_x} \equiv 1$, vanish.

2.1. Why a nonadditive disturbance model should be the appropriate one cannot be seen without examining the manner of introducing the stochastic element into those models.

Suppose the *undisturbed* model is given by

$$(2.1) \quad \varphi_g(x'; \theta_g) = 0, \quad g = 1, 2, \dots, G.$$

To transform this into a *stochastic (disturbed)* model we must consider some component $\theta_g^{(1)}$ of θ_g as stochastic. Just which component of θ_g should become stochastic cannot be decided arbitrarily: it depends on the composition of the universe of objects from which the observation is assumed to be drawn. Thus suppose the theory, as embodied in (2.1), refers to an individual firm, while the observations (say in a cross-section study) describe a group of firms (the universe being some larger aggregate of existing and, possibly, non-existent firms). Let $\theta_g^{(1)}$ be the component of θ_g which characterizes an individual firm. Then $\theta_g^{(1)}$ may depend 1) on x , or 2) on nonobservable fixed variates, or 3) it may be regarded as stochastic. 1) creates no new problem, 2) may imply loss of identification, 3) converts $\theta_g^{(1)}$ into a disturbance u_g .

Suppose the third possibility holds so that $\theta_g^{(1)}$ may be regarded as stochastic. Then, if $\theta_g^{(1)}$ has more than one component, we are dealing with the case of nonadditive disturbances except for the trivial case mentioned in section 1.2. Such situations are extremely frequent in economic models, although, unfortunately, the second of the above three possibilities is not uncommon.

Some economic models of this type are described in article [I], others in section 2.2 below.

2.2. An example of (1.3) can be found in the estimation of a firm's production function in a cross-section study, as discussed on p. 415. Other examples inserted here in small print may help

clarify the nature of the problem and show the great range of possibilities.

(A) *A game of chance (fixed variates)*. Let X_{it} , $i = 1, 2, \dots, k$, be k numbers arbitrarily chosen by the player at the time t and α_{it} the number given by throwing the i th set of dice. (The α 's are independently distributed in time.) The player wins each time the amount

$$(1) \quad Z_t = \sum_i \alpha_{it} X_{it}.$$

We observe the X 's and the Z 's but not the α 's, and we wish to estimate the distribution of the latter, say its first two moments.

(B) *A game of chance (time series; discrete)*. Let α_{it} , $i = 0, 1$, be numbers given by throwing the i th set of dice at the time t . (The α 's are again independently distributed in time.) X_0 is an arbitrary initial constant, and the player's winnings at time t , denoted by X_t , are determined by

$$(2) \quad X_t = \alpha_{it} X_{t-1} + \alpha_{0t}, \quad t = 1, 2, \dots$$

Here again the problem would consist in estimating the unknown distribution of the α 's (which are not observable) given the observations on X_t .

(C) *An economic system (time series; discrete)*. Let there be a system of equations

$$(3) \quad \begin{aligned} \alpha_{11t} x_{1t} + \alpha_{12t} x_{2,t-1} + \alpha_{10t} &= u_{1t}, \\ \alpha_{21t} x_{1t} + \alpha_{22t} x_{2t} + \alpha_{20t} &= u_{2t}, \end{aligned}$$

which may be interpreted as supply and demand equations, respectively, the x 's representing quantity and price, with a lag in the first equation as in the case of the "hog cycle." The α 's have the subscript t since, for instance, they may vary seasonally or exhibit a trend. But it seems quite plausible to consider them as also having variations of a random nature. In that case the problem of estimating the distribution of the α 's would again arise.

(D) *Continuous systems*. The above system may also be presented in continuous form. However, in order to simplify exposition, the continuous example given here will correspond to the one-variable discrete case, say to

$$(4) \quad x_t = \sum_{\tau=1}^{\tau_m} \alpha_{\tau t} x_{t-\tau} + \alpha_{0t} + u_t.$$

The corresponding (slightly simplified) continuous situation would be

$$(5) \quad x(t) = \int_0^{\infty} \varphi(\tau) x(t-\tau) d\tau + u(t),$$

where $u(t)$ is a continuous normal process with a specified autocorrelation pattern.

The difficulties of mathematical manipulation are not to be underestimated. But in principle it would seem possible to assign stochastic properties to φ as well as (or instead of) u . This would require making φ a function of t as well as of τ ; for instance,

$$(6) \quad x(t) = \int_0^{\infty} \varphi(t, \tau) x(t-\tau) d\tau + u(t),$$

where φ (as a function of t), or u , or both, are specified continuous stochastic processes.

When the production function is of the Cobb-Douglas type it may be written as

$$(2.2) \quad X_0^{(j)} = \sum_{i=1}^K \alpha_i^{(j)} X_i^{(j)} + \alpha_0^{(j)}.$$

Here X_0 is the logarithm of the product and the other X 's are logarithms of the factors. The superscript refers to the j th firm.

Now the $\alpha^{(j)}$ (the "productivities") may differ from firm to firm and the firm in existence may be considered as a sample from the universe of all possible firms with $\alpha^{(j)}$ which are stochastic variates.

The current practice does this, but only for $\alpha_0^{(j)}$. It is difficult to see any reason, except that of expediency, for regarding other α 's as fixed.

On the other hand, if the α 's are random variables, it is reasonable to try to estimate their joint distribution. This, however, seems to be a rather difficult task. In the following section we show what problems are involved in such an estimation procedure. The example treated is highly oversimplified and should be regarded only as illustrating the mathematics involved.

3.1. Let the nonadditive disturbance model be

$$(3.1) \quad x_1 + u_3 x_2 = u_1, \quad \alpha x_1 + x_2 = u_2,$$

where the x 's are observed, the u 's are the disturbances (with u_3 entering in a nonadditive fashion), and α is a (known or unknown) constant.

The problem in general is that of estimating α and the distribution of the u 's. This can be accomplished by the maximum-likelihood method if the likelihood function is known. This function will now be obtained for the case of jointly normally distributed disturbances; the sample is assumed to be random so that equation (3.4) holds. [In (3.4), $h(x_1, x_2)$ is the joint probability density function of the observed variates.]

It should be noted that the problem thus defined would in general be indeterminate; without additional *a priori* knowledge only certain relations between the parameters can be estimated. It will be assumed, however, that after the likelihood function is formed, these restrictions will be formulated and taken into account in the procedure of maximizing the likelihood function.

It should also be noted that it may not be very realistic to endow the disturbances with a normal distribution. It might, for instance, be more appropriate to make the marginal distribution of u_3 one of χ^2 type, etc.

3.2. Thus let the joint probability density function of the disturbances be

$$(3.2) \quad f(u_1, u_2, u_3) = (2\pi)^{-3/2} (\det [\sigma^{ij}])^{1/2} \exp(-\frac{1}{2} Q)$$

with

$$(3.3) \quad Q = \sum_{i,j=1}^3 \sigma^{ij} (u_i - \alpha_i)(u_j - \alpha_j).$$

We shall now proceed to obtain the joint distribution $h(x_1, x_2, u_3)$ and then integrate out u_3 . The result of the integration, say $g(x_1, x_2)$ will yield the logarithm of the likelihood function

$$(3.4) \quad \sum_{t=1}^T \log g(x_{1t}, x_{2t}).$$

The computation of $h(x_1, x_2, u_3)$ is straightforward.

For convenience of notation we may add to (3.1) a third equation

$$(3.5) \quad x_3 = u_3$$

involving the auxiliary variable x_3 and consider (3.1) with (3.5) as a transformation from (u_1, u_2, u_3) to (x_1, x_2, x_3) . In the result u_3 may again be written instead of x_3 .

The Jacobian of the transformation is

$$(3.6) \quad |J| = |1 - \alpha x_3|,$$

where $| \quad |$ is the absolute-value symbol.

Making the appropriate substitutions we find

$$(3.7) \quad h(x_1, x_2, x_3) = (2\pi)^{-3/2} (\det [\sigma^{ij}])^{1/2} |1 - \alpha u_3| \\ \times \exp\left\{-\frac{1}{2} \frac{(u_3 + \psi)^2}{\varphi^2} - \frac{1}{2} \lambda\right\},$$

where φ , ψ , and λ do not contain u_3 .¹

¹The values of φ , ψ , and λ are as follows:

$$(f.1) \quad \begin{aligned} \varphi &= 1/\sqrt{B_2}, \\ \psi &= B_1/B_2, \\ \lambda &= B_0 - B_1^2/B_2, \end{aligned}$$

where the B 's are given by

$$(f.2.1) \quad \begin{aligned} B_2 &= \sigma^{11} x_2^2 + 2\sigma^{13} x_2 + \sigma^{33}, \\ B_1 &= \sigma^{11} w_1 x_2 + \sigma^{12} x_2 w_2 + \sigma^{13} (w_1 - x_2 \alpha_3) \\ &\quad + \sigma^{23} w_2 - \sigma^{33} \alpha_3, \\ B_0 &= \sigma^{11} w_1^2 + 2\sigma^{12} w_1 w_2 + \sigma^{22} w_2^2 - 2\sigma^{13} w_1 \alpha_3 \\ &\quad - 2\sigma^{23} w_2 \alpha_3 + \sigma^{33} \alpha_3^2. \end{aligned}$$

The w 's in (f.2.1) have the following meaning:

$$(f.2.2) \quad \begin{aligned} w_1 &= x_1 - \alpha_1, \\ w_2 &= \alpha x_1 + x_2 - \alpha_2. \end{aligned}$$

Hence

$$(3.8) \quad \begin{aligned} & g(x_1, x_2) \\ &= (2\pi)^{-\frac{3}{2}} (\det [\sigma^{ij}])^{\frac{1}{2}} e^{-\frac{1}{2}\lambda} \int_{-\infty}^{\infty} |1 - \alpha u_3| e^{-\frac{1}{2} \frac{(u_3 + \psi)^2}{\varphi^2}} du_3. \end{aligned}$$

It will be observed that for $\alpha = 0$, this integral may be evaluated in terms of elementary functions. In this case we have¹

$$(3.9) \quad g_0(x_1, x_2) = (2\pi)^{-1} (\det [\sigma^{ij}])^{\frac{1}{2}} e^{-\frac{1}{2}\lambda} \varphi.$$

¹When the u 's are independent with zero means and unit variances, $g_0(x_1, x_2)$ becomes

$$(f.3) \quad g_{00}(x_1, x_2) = (2\pi)^{-1} (1 + x_2^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(x_1^2 + x_2^2 - \frac{x_1^2 x_2^2}{1 + x_2^2})\right\}.$$

XIX. NOTE ON RANDOM COEFFICIENTS

BY HERMAN RUBIN

Let us consider an equation of the form

$$(1) \quad y_t = \sum_{k=1}^K a_{kt} x_{kt}, \quad t = 1, \dots, T,$$

which defines the dependent variable y_t in terms of the fixed variables x_{kt} and the coefficients a_{kt} which are assumed to be normally and independently distributed with mean α_k and variance $\sigma_k^2 \varphi_{kt}$, where the inverse weights φ_{kt} are known functions of the fixed variates.¹ We could estimate the parameters α_k by least squares, i.e., by the formula

$$\bar{\alpha}_k = \sum_{l=1}^L m_{yx_l} m^{lk},$$

where

$$m_{yx_l} = \sum_{t=1}^T y_t x_{lt},$$

and m^{lk} is the element in the l th row and k th column of

$$(m_{kl})^{-1} = \left(\sum_{t=1}^T x_{kt} x_{lt} \right)^{-1}.$$

Although this method can be shown, under certain conditions, to be consistent, it would not be efficient.

We shall derive equations defining maximum-likelihood estimates of the α_i and σ_j^2 . Let us consider the distribution of the y_t . We see that the quantities

$$(2) \quad z_t = y_t - \sum_{k=1}^K \alpha_k x_{kt}, \quad t = 1, \dots, T,$$

¹See also [XVIII].

are normally and independently distributed with mean 0 and variance

$$(3) \quad \sum x_{kt}^2 \sigma_k^2 \varphi_{kt}.$$

Then the likelihood function of the observation becomes

$$(4) \quad L = (2\pi)^{-\frac{T}{2}} \prod_{t=1}^T \left(\sum x_{kt}^2 \sigma_k^2 \varphi_{kt} \right)^{-\frac{1}{2}} \\ \times \exp \left\{ -\frac{1}{2} \frac{\sum_{t=1}^T (y_t - \sum_{k=1}^K \alpha_k x_{kt})^2}{\sum_{k=1}^K \sum x_{kt}^2 \sigma_k^2 \varphi_{kt}} \right\}.$$

Let us maximize $\log L$ with respect to α_i and σ_j^2 , subject to the restriction $\sigma_j^2 \geq 0$. We obtain

$$(5) \quad \sum_{t=1}^T \frac{x_{it} (y_t - \sum_{k=1}^K \alpha_k x_{kt})}{\sum_{k=1}^K x_{kt}^2 \sigma_k^2 \varphi_{kt}} = 0,$$

$$(6) \quad \sum_{t=1}^T \frac{x_{jt}^2 \varphi_{jt}}{\sum_{k=1}^K x_{kt}^2 \sigma_k^2 \varphi_{kt}} = \sum_{t=1}^T \frac{x_{jt}^2 \varphi_{jt} (y_t - \sum_{k=1}^K \alpha_k x_{kt})^2}{(\sum_{k=1}^K x_{kt}^2 \sigma_k^2 \varphi_{kt})^2} - \lambda_j,$$

$$(7) \quad \lambda_j \sigma_j^2 = 0,$$

$$(8) \quad \sigma_j^2 \geq 0.$$

It is necessary to introduce the Lagrange multipliers λ_j because the solutions of (5) and (6) with $\lambda_j = 0$ might give negative σ_j^2 . We take

that solution of (5), (6), (7), and (8) which gives to (4) its highest value. It should be noted that it is unnecessary to consider solutions having λ_j 's $\neq 0$ if there is a solution with λ_j 's = 0, and, in general, if there is a solution with only $\lambda_{j_1}, \dots, \lambda_{j_\beta} \neq 0$, it is not necessary to consider solutions with those λ_j 's $\neq 0$ and other λ_j 's $\neq 0$.

REFERENCES

- ANDERSON, R. L., "Use of Variance Components in the Analysis of Hog Prices in Two Markets," *Journal of the American Statistical Association*, Vol. 42, December, 1947, pp. 612-634.
- ANDERSON, T. W., and LEONID HURWICZ, "Errors and Shocks in Economic Relationships," Paper presented at the meeting of The Econometric Society, September 6-18, 1947, in Washington, D.C.; abstract in *Econometrica*, Vol. 16, January, 1948, pp. 36-37.
- ANDERSON, T. W., and HERMAN RUBIN, "Estimation of the Parameters of a Single Stochastic Difference Equation in a Complete System," *Annals of Mathematical Statistics*, Vol. 20, March, 1949, pp. 46-63 (and to be included in *Cowles Commission Paper, New Series, No. 36*).
- _____, "Asymptotic Properties of Estimates of the Parameters of a Single Equation in a Complete System of Stochastic Equations," to be published (and to be included in *Cowles Commission Paper, New Series, No. 36*).
- ANDREWS, WILLIAM H., JR., see JACOB MARSCHAK and WILLIAM H. ANDREWS, JR.
- BENTZEL, H., and H. WOLD, "On Statistical Demand Analysis from the Viewpoint of Simultaneous Equations," *Skandinavisk Aktuarietidskrift*, Vol. 29, Nos. 1-2, 1946, pp. 95-114.
- BIRKHOFF, GARRETT, and SAUNDERS MACLANE, *A Survey of Modern Algebra*, New York: The Macmillan Co., 1941, 450 pp.
- BROOKNER, RALPH J., "Choice of One Among Several Statistical Hypotheses," *Annals of Mathematical Statistics*, Vol. 16, September, 1945, pp. 221-242.
- CHERNOFF, HERMAN, "Gradient Methods of Maximization in Estimating Economic Parameters," Paper presented at the Madison meeting of The Econometric Society, September 7-10, 1948; abstract in *Econometrica*, Vol. 17, January, 1949, pp. 75-76.
- CRAMÉR, HARALD, *Random Variables and Probability Distributions*, London: Cambridge, England: The University Press, 1937, 120 pp.
- DAVID, F. N., and J. NEYMAN, "Extension of the Markoff Theorem on Least Squares," *Statistical Research Memoirs*, Vol. II, London: Department of Statistics, University of London, University College, 1938, pp. 105-116.
- DIXON, WILFRID J., "Further Contributions to the Problem of Serial Correlation," *Annals of Mathematical Statistics*, Vol. 15, June, 1944, pp. 119-144.
- DOOB, J. L., "The Elementary Gaussian Processes," *Annals of Mathematical Statistics*, Vol. 15, September, 1944, pp. 229-282.
- FRISCH, RAGNAR, 1929, "Correlation and Scatter in Statistical Variables," *Nordic Statistical Journal*, Vol. 1, 1929, pp. 36-102.

- _____, 1933, *Pitfalls in the Construction of Statistical Demand Curves*, Veröffentlichungen der Frankfurter Gesellschaft für Konjunkturforschung, Neue Folge, Heft 5, Leipzig: Hans Buske Verlag, 1933, 39 pp.
- _____, 1934, *Statistical Confluence Analysis by Means of Complete Regression Systems*, Oslo: Universitetets Økonomiske Institutt, 1934, 192 pp.
- _____, 1938, "Statistical Versus Theoretical Relations in Economic Macrodynamics," Memorandum prepared for a conference in Cambridge, England, July 18-20, 1938, to discuss drafts of Tinbergen's League of Nations Publications; mimeographed.
- FRISCH, RAGNAR, and BRUCE D. MUDGETT, "Statistical Correlation and the Theory of Cluster Types," *Journal of the American Statistical Association*, Vol. 26, December, 1931, pp. 375-392.
- GEARY, R. C., 1942, "Inherent Relations between Random Variables," *Proceedings of the Royal Irish Academy*, Vol. 47, Section A, March, 1942, pp. 63-196.
- _____, 1943, "Relations between Statistics: The General and the Sampling Problem when the Samples Are Large," *Proceedings of the Royal Irish Academy*, Vol. 49, Section A, December, 1943, pp. 177-196.
- GIRSHICK, M. A., and TRYGVE HAAVELMO, "Statistical Analysis of the Demand for Food: Examples of Simultaneous Estimation of Structural Equations," *Econometrica*, Vol. 15, April, 1947, pp. 79-110 (and reprinted as *Cowles Commission Paper, New Series, No. 24*).
- HAAVELMO, TRYGVE, 1943, "The Statistical Implications of a System of Simultaneous Equations," *Econometrica*, Vol. 11, January, 1943, pp. 1-12.
- _____, 1944, "The Probability Approach in Econometrics," *Econometrica*, Vol. 12, Supplement, July, 1944, 118 pp. (and reprinted as *Cowles Commission Paper, New Series, No. 4*).
- _____, 1947-A, "Methods of Measuring the Marginal Propensity to Consume," *Journal of the American Statistical Association*, Vol. 42, March, 1947, pp. 105-122 (and reprinted as *Cowles Commission Paper, New Series, No. 22*).
- _____, 1947-B, "Quantitative Research in Agricultural Economics: The Interdependence Between Agriculture and the National Economy," *Journal of Farm Economics*, Vol. 24, 1947, pp. 910-924 (and included in *Cowles Commission Paper, New Series, No. 27*).
- _____, see also M. A. GIRSHICK and TRYGVE HAAVELMO.
- HARDY, G. H., J. E. LITTLEWOOD, and G. POLYA, *Inequalities*, Cambridge, England: The University Press, 1934, 314 pp.
- HOTELLING, HAROLD, 1933, "Analysis of a Complex of Statistical Variables in Principal Components," *Journal of Educational Psychology*, Vol. 24, September and October, 1933, pp. 417-444 and 498-520.
- _____, 1936-A, "Simplified Calculation of Principal Components," *Psychometrika*, Vol. 1, March, 1936, pp. 27-35.
- _____, 1936-B, "Relations Between Two Sets of Variables," *Biometrika*, Vol. 28, December, 1936, pp. 321-377.
- _____, 1940, "The Selection of Variates for Use in Prediction with Some Comments on the Problem of Nuisance Parameters," *Annals of Mathematical Statistics*, Vol. 11, September, 1940, pp. 271-283.

- _____, 1943-A, "Some New Methods in Matrix Calculation," *Annals of Mathematical Statistics*, Vol. 14, March, 1943, pp. 1-34.
- _____; 1943-B, "Further Points on Matrix Calculation and Simultaneous Equations," *Annals of Mathematical Statistics*, Vol. 14, December, 1943, pp. 440-441.
- _____, 1949, "Practical Problems of Matrix Calculation," *Proceedings of the Berkeley Symposium on Probability and Statistics*, Jerzy Neyman, ed., Berkeley and Los Angeles: University of California Press, 1949, pp. 275-293.
- HURWICZ, LEONID, 1944, "Stochastic Models of Economic Fluctuations," *Econometrica*, Vol. 12, April, 1944, pp. 114-124 (and reprinted as *Cowles Commission Paper, New Series, No. 3*).
- _____, 1946, "Theory of the Firm and of Investment," *Econometrica*, Vol. 14, April, 1946, pp. 109-136 (and reprinted as *Cowles Commission Paper, New Series, No. 16*).
- _____, 1947, "Some Problems Arising in Estimating Economic Relations," *Econometrica*, Vol. 15, July, 1947, pp. 236-240.
- _____, see also T. W. ANDERSON and LEONID HURWICZ.
- _____, see also JACOB MARSCHAK; LEONID HURWICZ; et al.
- JOHNSON, EVAN, JR., "Estimates of Parameters by Means of Least Squares," *Annals of Mathematical Statistics*, Vol. 11, December, 1940, pp. 453-456.
- KENDALL, M. G., "A New Measure of Rank Correlation," *Biometrika*, Vol. 30, June, 1938, pp. 81-93.
- KLEIN, LAWRENCE R., 1946-A, "Macroeconomics and the Theory of Rational Behavior," *Econometrica*, Vol. 14, April, 1946, pp. 93-108 (and reprinted as *Cowles Commission Paper, New Series, No. 14*).
- _____, 1946-B, "A Post-Mortem on Transition Predictions of National Product," *Journal of Political Economy*, Vol. 54, August, 1946, pp. 289-308 (and reprinted as *Cowles Commission Paper, New Series, No. 18*).
- _____, 1947, "The Use of Econometric Models as a Guide to Economic Policy," *Econometrica*, Vol. 15, April, 1947, pp. 111-151 (and reprinted as *Cowles Commission Paper, New Series, No. 23*).
- _____, 1950, *Economic Fluctuations in the United States, 1921-1941*, Cowles Commission Monograph No. 11, New York: John Wiley & Sons, 1950, about 170 pp.
- KOOPMANS, TJALLING C., 1937, *Linear Regression Analysis of Economic Time Series*, Haarlem: De Erven F. Bohn N. V., 1937, 150 pp.
- _____, 1942, "Serial Correlation and Quadratic Forms in Normal Variables," *Annals of Mathematical Statistics*, Vol. 13, March, 1942, pp. 14-33.
- _____, 1945, "Statistical Estimation of Simultaneous Economic Relations," *Journal of the American Statistical Association*, Vol. 40, December, 1945, pp. 448-466 (and reprinted as *Cowles Commission Paper, New Series, No. 11*).
- _____, 1949, "Identification Problems in Economic Model Construction," *Econometrica*, Vol. 17, April, 1949, pp. 125-144 (and reprinted as *Cowles Commission Paper, New Series, No. 31*).
- _____, see also JACOB MARSCHAK; LEONID HURWICZ; et al.

- LEIPNIK, ROY B., "Distribution of the Serial Correlation Coefficient in a Circularly Correlated Universe," *Annals of Mathematical Statistics*, Vol. 18, March, 1947, pp. 80-87 (and included in *Cowles Commission Paper, New Series, No. 21*).
- _____, see also JACOB MARSCHAK; LEONID HURWICZ; et al.
- LITTLEWOOD, J. E., see G. H. HARDY, J. E. LITTLEWOOD, and G. PÓLYA.
- MACDUFFEE, CYRUS COLTON, *Vectors and Matrices*, Menasha, Wisconsin: Mathematical Association of America, 1943, 192 pp.
- MACLANE, SAUNDERS, see GARRETT BIRKHOFF and SAUNDERS MACLANE.
- MADOW, WILLIAM G., "Note on the Distribution of the Serial Correlation Coefficient," *Annals of Mathematical Statistics*, Vol. 16, September, 1945, pp. 308-310.
- MANN, HENRY B., "Nonparametric Tests against Trend," *Econometrica*, Vol. 13, July, 1945, pp. 245-259.
- MANN, H. B., and A. WALD, "On the Statistical Treatment of Linear Stochastic Difference Equations," *Econometrica*, Vol. 11, July-October, 1943, pp. 173-220.
- MARSCHAK, JACOB, 1947-A, "Economic Structure, Path, Policy, and Prediction," *American Economic Review, Papers and Proceedings of the 59th Annual Meeting of the American Economic Association*, Vol. 37, May, 1947, pp. 81-84.
- _____, 1947-B, "Statistical Inference from Nonexperimental Observations: An Economic Example," Paper presented at the meeting of the International Statistical Institute and The Econometric Society, September 6-18, 1947, in Washington, D.C.; to be published in 1950 in *Proceedings of the International Statistical Conference* (and to be reprinted as *Cowles Commission Paper, New Series, No. 32*); abstract in *Econometrica*, Vol. 16, January, 1948, pp. 53-55.
- MARSCHAK, JACOB, and WILLIAM H. ANDREWS, JR., "Random Simultaneous Equations and the Theory of Production," *Econometrica*, Vol. 12, July-October, 1944, pp. 143-205 (and reprinted as *Cowles Commission Paper, New Series, No. 5*).
- MARSCHAK, JACOB; LEONID HURWICZ; TJALLING C. KOOPMANS and ROY B. LEIPNIK; "Estimating Relations from Nonexperimental Observations," Papers presented at the Cleveland meeting of The Econometric Society, January 24-27, 1946; abstracts in *Econometrica*, Vol. 14, April, 1946, pp. 165-172 (and included in *Cowles Commission Paper, New Series, No. 17*)
- MARSHALL, ALFRED, *Principles of Economics*, Eighth Edition, London: Macmillan & Co., 1920, 754 pp.
- MUDGEY, BRUCE D., see RAGNAR FRISCH and BRUCE D. MUDGEY.
- NEYMAN, J., "Outline of a Theory of Statistical Estimation Based on the Classical Theory of Probability," *Philosophical Transactions of the Royal Society of London, Series A*, Vol. 236, August 30, 1937, pp. 330-380.
- _____, see also F. N. DAVID and J. NEYMAN.
- PIERCE, B. O., *A Short Table of Integrals*, New York: Ginn and Co., 1929, 156 pp.
- PÓLYA, G., see G. H. HARDY, J. E. LITTLEWOOD, and G. PÓLYA.
- REIERSØL, OLAV, "Residual Variables in Regression and Confluence Analysis,"

- Skandinavisk Aktuarietidskrift*, Vol. 28, Nos. 3-4, 1945, pp. 201-217.
- RUBIN, HERMAN, 1945, "On the Distribution of the Serial Correlation Coefficient," *Annals of Mathematical Statistics*, Vol. 16, June, 1945, pp. 211-215 (and reprinted as *Cowles Commission Paper, New Series, No. 10*).
- _____, 1946, "Asymptotic Distribution of Moments from a System of Linear Stochastic Difference Equations," Abstract No. 332 in *Bulletin of the American Mathematical Society*, Vol. 52, September, 1946, pp. 827-828.
- _____, 1948, "Some Results on the Asymptotic Distribution of Maximum- and Quasi-Maximum-Likelihood Estimates," Abstract No. 529 in *Bulletin of the American Mathematical Society*, Vol. 54, November, 1948, p. 1080, and in *Annals of Mathematical Statistics*, Vol. 19, December, 1948, p. 598.
- _____, 1949, "Properties of Maximum- and Quasi-Maximum-Likelihood Estimates of Parameters of a System of Linear Stochastic Difference Equations with Serially Correlated Disturbances," abstract in *Annals of Mathematical Statistics*, Vol. 20, March, 1949, p. 137.
- _____, see also T. W. ANDERSON and HERMAN RUBIN.
- SNEDECOR, GEORGE W., *Calculation and Interpretation of Analysis of Variance and Covariance*, Ames, Iowa: Collegiate Press, Inc., 1934, 96 pp.
- TINBERGEN, J., *Statistical Testing of Business-Cycle Theories: Vol. I, A Method and Its Application to Investment Activity; Vol. II, Business Cycles in the United States of America, 1919-1932*; Geneva: League of Nations, 1939, 164 and 244 pp.
- TINTNER, GERHARD, 1942, "A Contribution to the Nonstatic Theory of Production," in *Studies in Mathematical Economics and Econometrics, in Memory of Henry Schultz*, O. Lange, F. McIntyre, T. O. Yntema, eds., Chicago: The University of Chicago Press, 1942, pp. 92-109.
- _____, 1946, "Multiple Regression for Systems of Equations," *Econometrica*, Vol. 14, January, 1946, pp. 5-36.
- ULLMAN, JOSEPH, "The Probability of Convergence of an Iterative Process of Inverting a Matrix," *Annals of Mathematical Statistics*, Vol. 15, June, 1944, pp. 205-213.
- WALD, ABRAHAM, 1940, "The Fitting of Straight Lines If Both Variables Are Subject to Error," *Annals of Mathematical Statistics*, Vol. 11, September, 1940, pp. 284-300.
- _____, 1942, *On the Principles of Statistical Inference*, Notre Dame Mathematical Lectures, No. 1, South Bend, Indiana: University of Notre Dame, 1942, 47 pp.
- _____, 1944, "Note on a Lemma," *Annals of Mathematical Statistics*, Vol. 15, September, 1944, pp. 330-333.
- _____, see also H. B. MANN and A. WALD.
- WALD, A., and J. WOLFOWITZ, "Confidence Limits for Continuous Distribution Functions," *Annals of Mathematical Statistics*, Vol. 10, June, 1939, pp. 105-118.
- WATSON, G. N., see E. T. WHITTAKER and G. N. WATSON.
- WAUGH, F. V., "A Note Concerning Hotelling's Method of Inverting a Partitioned Matrix," *Annals of Mathematical Statistics*, Vol. 16, June, 1945, pp. 216-217.

- WHITTAKER, E. T., and G. N. WATSON, *A Course of Modern Analysis*, American Edition, Cambridge, England: The University Press; New York: The Macmillan Co., 1945, 608 pp.
- WILKS, S. S., *Mathematical Statistics*, Princeton: Princeton University Press, 1943, 284 pp.
- WILLIAMS, J. D., "Moments of the Ratio of the Mean Square Successive Differences to the Mean Square Difference in Samples from a Normal Universe," *Annals of Mathematical Statistics*, Vol. 12, June, 1941, pp. 239-241.
- WINTNER, A., *Lectures by Aurel Wintner on Asymptotic Distributions and Infinite Convolutions*, Ann Arbor, Michigan: Edwards Brothers, 1938, 54 pp.
- WOLD, HERMAN, *A Study in the Analysis of Stationary Time Series*, Uppsala: Almqvist & Wiksells, 1938, 214 pp.
- _____, see also H. BENTZEL and H. WOLD.
- WOLFOWITZ, J., see A. WALD and J. WOLFOWITZ.

INDEX

- A priori information, *see* Information
A priori restrictions, *see* Restrictions
Absolutely unbiased estimate, 281, 367
Aggregation, optimum, 7
Aggregative model, 7
Albert, A. Adrian, 5
Anderson, R. L., 2
Anderson, T. W., 2, 4, 5, 58, 111, 312, 338, 372
Andrews, William H., Jr., 5, 32, 97, 105
Approximations, successive, in the computation of maximum-likelihood estimates, 153
Autocorrelation coefficient, 48
Autocorrelation of disturbances, *see* Disturbances
Autonomous relation, 263
Autoregressive time series, *see* Time series
- Basic matrices, *see* Matrix
Behavior, 7, 13, 19, 21, 44
 of buyers, 36, 37
 of consumers, 31
 economic, 4
 erratic element in, 385
 law of, 63
 of lenders, 31
 of producers, 36, 37
Behavior equation, 54
Bentzel, H., 35
Best policy, *see* Policy
Best unbiased estimate, *see* Estimate
Bias, 265, 365, 370
 magnitude of in small-sample estimation, 366
 relative, 365, 380
 single-equation, 277
 see also Estimate
Bilinear restriction, *see* Restrictions
Brookner, Ralph F., 44
Bunch map analysis, 46, 261, 265
- Calculus of variations, 284
Cauchy-Schwarz inequality, 357
Cayley numbers, 356
Central limit theorem, 358
Change of structure, *see* Structure
Characteristic function, 370
Chernoff, Herman, 47
Chi-square distribution, 308, 318, 320
Class frequencies, disproportionate, 324
Cobb, C. W., 415
Cobb-Douglas production function, 415
Coefficient,
 of regression, *see* Regression coefficient
 structural, 35, 36
 biased estimate of, 365
 vs. predictive estimation, 273
 see also Parameters
Complement, orthogonal, 89
Complementary equations, *see* Equations
Computation, 5, 7
 cost of, 230
Confidence coefficient, 307
Confidence interval, 399
Confidence region, 30, 43, 307, 308
 diameter of, 309
 for small samples, 317
Confidence set, 307
Confluence analysis, 258

- Confluent relation, 70
 Consistent estimate, 41, 306
 Continuous process, 49
 Correlation, serial, 310
 see also Disturbances and Residuals
 Covariance matrix, *see Disturbances and Lag-covariance matrix*
 Cramér, Harald, 357, 358, 362
 Critical region, 307, 346
 in K-test, 349
 Cyclical fluctuations, 34
- David, F. N., 259, 285
 Demand, 14, 15, 29, 36
 shift of, 36, 50
 Demand curve, 21
 Demand equation, 36, 41, 50
 Demand and supply analysis, pitfalls in, 4
 Dependent variable, *see Variables*
 Design of social experiments, 355
 see also Experiments
 Design, randomized block, 355
 Determination, 6
 predictive, 12, 16, 17, 25
 structural, 13, 15
 Difference equations, 56, 330
 with autocorrelated disturbances, 343
 linear, 41, 55
 stochastic, 272, 298, 332, 354, 368, 382
 noncircular, 365
 Discrete model, 33, 49
 Distribution,
 complete, 275
 derived, 275
 of disturbances, *see Disturbances*
 normal, 28, 57, 272, 301, 359
 of observables, 20
 past, 30
 see also Probability distribution
 Distribution function, cumulative, 274, 280, 305
 Disturbances, 4, 18, 56, 262, 276, 312, 338, 386
 additive linear, 413
 autocorrelation of, 337
 autoregressive, 337
 moving-average, 337, 339
 correlated, 211, 231
 covariance matrix of, 247, 275, 276, 313, 320, 411
 diagonality restriction on, 183, 191
 nonsingularity of, 58
 distribution of,
 normal, 246, 314, 416
 nonsingular, 408
 in equations, 56, 338
 see also Shocks
 independence of in successive observations, 314, 384, 385, 387, 407
 nonadditive, 50, 247, 410, 412, 415
 nonautocorrelated, 330, 336
 nonobservable, 19
 serial correlation of, 385
 uncorrelated, 159, 166, 231
 in variables, 57, 338
 see also Errors
 Diurnal fluctuations, 329, 332
 Dixon, Wilfrid J., 48, 370, 383
 Doob, J. L., 386
 Doolittle method, 323
 Douglas, P. H., 415
 Dwyer, Paul S., 323
- Economic theory, 2, 9, 13, 47
 Endogenous variable, *see Variables*
 Equations,
 complementary, 43
 definitional, 22
 disturbances in, *see Disturbances and Shocks*
 final (Tinbergen), 34
 forecast, 21
 institutional, 54
 normal, 47
 stochastic, 56, 305, 410
 structural, 8, 27, 32, 63, 314, 384, 402
 complete system of, 395
 completed subset of, 98, 100

- identifiability of, 78, 176
 - linearity of, 108
- subsidiary, 8, 22
- system of,
 - a priori information on, 7, 19
 - complete, 393, 395
 - in reduced form, *see* Reduced form
 - simultaneous, 2, 4, 55, 394, 411
 - stable, 133
 - technical, 54
- see also* Difference equations and Jacobian determinant
- Equilibrium, multiple, 9
- Errors,
 - matrix of, in computations, 324
 - of measurement or observation, 2, 18, 32, 50, 57, 262
 - additive, 20
 - nonadditive, 21
- see also* Disturbances in variables
- Estimate,
 - absolutely unbiased, 281, 367
 - asymptotic normality of, 139, 317
 - asymptotically unbiased, 395, 396
 - best absolutely unbiased linear, 283
 - best unbiased linear, 273, 279
 - conditionally unbiased, 281, 367
 - consistent, 41, 306
 - efficient, 41
 - least-squares, *see* Least-squares estimate
 - limited-information, *see* Limited-information method of estimation
 - maximum-likelihood, *see* Maximum-likelihood estimate
 - Markoff, *see* Markoff estimate
 - optimal properties of, 25
 - reduced-form maximum-likelihood, *see* Limited-information method of estimation
 - unbiased, 38, 259, 281, 367, 385
- Estimation, 6, 70
 - in incomplete systems of equations, 305
 - for large samples, 4
 - nonparametric, 45, 46
 - of parameters, 62
 - predictive, 25, 34, 38, 266, 270, 277, 366
 - under changed structure, 269
 - under unchanged structure, 278
 - structural, 3, 25, 248, 269, 271, 366
 - partial or incomplete, 41, 42, 46
- Exogenous variable, *see* Variables
- Experiment, 2, 3, 6, 17, 32, 45, 355
- Explosive system, 48, 356
- F -distribution, 310, 319
- Fluctuations,
 - cyclical, 34
 - diurnal, 329, 332
 - economic, 54
 - seasonal, 32, 47, 329, 332, 334, 343
- Forecast, 245, 268
- Forecast equation, 21
- Fourier series, 366
- Frisch, Ragnar, 4, 20, 32, 69, 258, 260
- Gain (welfare), 16, 26, 38
 - maximization of, 11
 - function, 11, 16
 - functional, 26, 31
- Geary, R. C., 20
- Girshick, M. A., 5, 43, 111
- Haavelmo, T., 2, 4, 5, 36, 40, 41, 43, 55, 63, 70, 262, 354
- Hardy, G. H., 291
- Homoscedasticity, 284
- Hotelling, H., 2, 47, 154, 175, 201, 210, 230, 233, 323, 399
- Hurwicz, L., 2, 4, 5, 26, 58, 336, 338
- Hypothesis, 2
 - considered, 49
 - maintained, 49
 - null, 345
 - statistical, 4

- Identifiable parameter, 85, 239
396
- Identifiability, 159
almost everywhere in the parameter space, 82
in spite of autocorrelation, 338, 340
conditions for, 243
criteria for,
 based on counting, 101, 102, 104
 in linear dynamic systems, 46
incomplete, 15, 96, 205
of a linear form, 239
local, 239
 condition for, 243
of structural equations, 77, 176
 subset of, 96
unique, 15
- Identification, 6, 10, 14, 20, 30, 62, 69, 78, 122, 238, 320, 387, 407
complete, 96, 180, 188
under linear restrictions, 78
under linear and bilinear restrictions, 93
multiple, 96
partial, 15, 31
unique, 96, 154
- Identification power, 245, 269, 278
definition of, 248
complete, 248
incomplete, 246, 248
multiple, 248, 255, 256
partially unique, 246
unique, 245, 248
various forms of (summary), 255
- Identifying model, *see* Model
- Income, 26
national, 14, 21, 29, 36
- Incomplete system, 49
- Incomplete (partial) model, 8, 16
- Independence of disturbances, *see* Disturbances
- Index numbers, 7
- Inference, 20
statistical, 3, 5, 70
- Information,
 additional, 31
 a priori, 7, 19, 22, 37, 41
 disregard of, 111
 limited, *see* Limited-information method of estimation
- Information-preserving maximum-likelihood estimation, *see* Maximum-likelihood estimation
- Initial value, 298, 366
- Initial-value case, 369
- Integral,
 elliptic, 373, 377
 hyperelliptic, 373
 Stieltjes, 281
- Intercorrelation of independent variables, 258
- Invariance of a function of the parameters, 242
- Invariants, fundamental set of, 244
- Iterative methods for maximizing the likelihood function, 153, 155, 323, 324
Newton method, 190, 196, 200, 203, 210, 222, 227
 computational procedure of, 206
 numerical experiment with, 227
 numerical illustration of, 209
 speed of convergence in, 205, 209
- \mathcal{P}_1 , \mathcal{P}_h , and \mathcal{P}_{h_n} methods (processes)
168, 169, 170, 172, 190, 218
asymptotic convergence properties of, 172, 183, 189, 215, 221
choice of h (in \mathcal{P}_h), 182
computations for, 192, 196, 218, 222
modification of by normalization, 220
numerical illustration of, 222
speed of convergence in, 189
see also Revision, iterative
- Jacobian determinant, 72, 244, 339, 396, 398, 400, 412, 417
factorization of, 397, 400, 403
nonsingularity of, 60
partitioning of, 234
- Johnson, Evan, Jr., 288

- K*-test, 348
 critical region in, 349
 Kronecker symbol, 163, 285, 412
 Kendall, M. G., 346
 Klein, L. R., 2, 5, 40, 43
 Koopmans, T. C., 2, 4, 5, 20, 43,
 48, 196, 245, 260, 280, 294,
 298, 369, 383, 384
- Lags, *see* Time lags and Variables
 Lag-correlation coefficient, 332
 Lag-covariance matrix, 342, 386
 Lagrange multiplier, 287, 332, 420
 Least-squares estimate, 40, 42,
 236, 258, 304, 315, 324, 372,
 419
 bias of, 365
 consistency of, 272
 efficiency of, 272
 and maximum likelihood, equiva-
 lence of, 189, 293, 301, 334,
 335, 366
 optimal properties of, 39
 of regression coefficients, 39,
 46, 272, 335
 single-equation, 117, 192
 small-sample properties of, 272
 variance of, 290
 Leipnik, R. B., 2, 5, 383
 Likelihood function, 110, 120, 154,
 303, 312, 321, 331, 333, 352,
 354, 369, 372, 398, 412, 416
 factorization of, 159, 397
 logarithmic, 235
 maximization of, 155, 304
 see also Iterative methods
 maximum of,
 absolute, 177
 depressed by the restrictions,
 177, 216
 first-order conditions for, 47,
 156, 166
 highest, 235
 restricted, 167, 169
 restricted, 120
 unrestricted, 110
 see also Taylor expansion of the
 likelihood function
 Likelihood-ratio criterion, 352
 Likelihood-ratio test, 320
 Limited-information method of esti-
 mation, 5, 111, 311, 313
 advantages and disadvantages of,
 321
 Linear model, 9
 Linear transformation in parameter
 space, 76, 247, 303
 Linearity, 39, 55
 Littlewood, J. E., 291
 Local identifiability, *see* Identi-
 fiability
- MacDuffee, Cyrus Colton, 163
 Maclaurin expansion, 369, 380
 Madow, W. G., 2, 5, 48
 Mann, H. B., 2, 4, 32, 41, 42, 55,
 56, 67, 115, 133, 135, 148,
 345, 356, 360, 400
 Markoff, 274, 279
 Markoff estimate, 284, 288
 absolutely unbiased, 285
 conditionally unbiased, 285
 variance of, 290
 Markoff theorem, 258, 259, 285,
 366
 Marschak, Jacob, 2, 5, 11, 32, 97,
 105
 Marshall, Alfred, 9
 Matrix,
 basic, 160, 164, 191
 canonical form of, 197, 200
 orthogonality of, 206, 222
 orthogonalization of, 162,
 172, 196, 200, 218
 inverse, 323
 norm of, 324
 partitioning of, 323
 rank of, 79
 restriction, 160, 164
 Maximal set of structures, 252
 Maximum-likelihood estimate, 41,
 191, 279, 296, 312, 314, 329,
 333, 353, 356, 372, 410, 416
 asymptotically unbiased, 400
 bias of, 365
 computation of, 153
 see also Iterative methods for
 maximizing the likelihood

- function
 - consistency of, 356, 400
 - in difference-equation systems, 366
 - first-order conditions for, 213
 - large-sample distribution of, 404
 - and least-squares, equivalence of, 189, 293, 301, 334, 335, 366
 - of regression coefficients, 39
 - quasi-, 134
 - sampling variance and covariance of, 42, 153, 209, 228
 - in small samples, 366, 404
- Maximum-likelihood estimation, 110, 189, 395
 - information-preserving, 42
 - limited-information, or reduced-form, *see* Limited-information method of estimation
- Measurement error, *see* Errors
- Missing observations, 354
- Model, 4, 5, 7, 14, 245, 248, 276
 - aggregative, 7
 - choice of, 44
 - complete, 7, 16, 23, 311
 - discrete, 33, 34, 49
 - dynamic, 32
 - identifying, 245, 278
 - partially, 256
 - totally, 256
 - uniquely, 249, 253
 - incomplete (partial), 8, 16
 - linear, 9
 - multiequational, unitemporal, 39
 - multitemporal, 32, 39, 40, 42
 - nonlinear, 9, 47
 - nonstochastic, 5, 19
 - partitionable, 35
 - sectional, 7, 22
 - self-contained, 7, 22
 - shock, 21
 - shock-and-error, 20
 - simultaneous-equations, 4
 - stochastic, 18, 19, 265, 413
 - continuous, 34
 - structure-identifying, 245, 269
 - subsidiary, 23
 - uniequational, 35
 - complete, 28, 32, 38
 - multitemporal, 40, 42
 - unitemporal, 32, 39
- Moment, 112
- Moment matrix, 112, 154, 276, 315
- Moving-average disturbance, 337, 339
- Mudgett, Bruce D., 20
- Multicollinearity, 46, 262
- Multiple regression, *see* Regression
- Neumann, John von, 5, 47, 154, 172, 380
- Newton method, *see* Iterative methods for maximizing the likelihood function
- Neyman, J., 259, 285, 307
- Nonlinear model, 9, 47
- Nonstochastic model, 5, 19
- Normal distribution, *see* Distribution
- Normal equations, 47
- Normality, asymptotic,
 - of moments, 136
 - of estimates, 139, 317
- Normalization, 68, 154, 158, 162, 189
- Notation, xiii, xiv, 70, 112, 405
 - conformity in, 5
 - of operators, 165
 - of vectors, 81
- Null hypothesis, 345
- Observable variable, *see* Variables
- Observations, 9, 14, 267, 317
 - errors of, *see* Errors
 - missing, 354
 - number of, 10, 15
 - passive, 64, 70
 - successive, 345
 - dependence between, 345
 - independence of, 280
 - see also* Disturbances
- Observation method, 386
- Observational reduced form, *see* Reduced form
- Observational structure, *see* Structure
- Orthogonal complement, 89

- Parameters, 6
 basic, 160
 estimation of, 62
 identifiable, 85, 239, 396
 structural, 8, 27, 276
 unrestricted set of, 161
 complete, 160
 see also Coefficient
- Parameter space, 76, 240, 244
 observationally equivalent points
 in, 77
 restricted, 77
 projection on, 163, 165
- Path, 3
- Pierce, B. O., 371
- Policy, 2, 3, 6, 41, 49, 271
 best, 12, 16, 27, 29, 31
 national, 30
 nonstructural, 11
 optimal, 38
 private, 30
 structural, 11, 13, 14, 26, 271
- Policy-maker, 11, 26
- Pólya, G., 291
- Polynomial series, 366
- Predetermined variable, *see* Variables
- Predictand, 27, 38, 266, 277
 conditional distribution of, 272
 probability distribution of, 267
- Prediction, 41, 263, 266, 267, 277
 generalized, 271
 under unchanged structure, 257
see also Determination and Estimation
- Predictor, 27, 29, 40, 267, 277, 323
- Predictor-predictand relationship, 268
- Predictor set, 39
- Price, 14, 26, 36
 prospective, 21
- Price control, 11
- Probability distribution, 2, 3, 18, 19, 55, 56, 387
 of income, 26
see also Likelihood function
- Probability measure, 307
- Process,
 continuous, 49, 50
 random (stochastic), 4
 cumulative, 48
 explosive, 44
 random (stochastic), 4, 47, 369, 386
 discrete, 343
 stationary, 41, 48
- Production conditions (technology), 7
- Production function, Cobb-Douglas, 415
- Profits, prospective, 21
- Propensity to spend, marginal, 37
- Property, structural, 22
- Quaternions, 356
- Random elements, 260
- Random series, 43
- Random variable, *see* Variables
- Randomized block design, 355
- Recursion, 346
- Reduced form, 9, 11, 14, 20, 24, 27, 31, 34, 40, 189
 observational, 10
- Reduced-form method of estimation, *see* Limited-information method of estimation
- Regression, 117, 271
 multiple, 47
- Regression coefficient, 27, 32, 35, 38, 279, 293
 biased estimate of, 365
 normally distributed, 318
 significance of, 346
 vs. structural coefficient, 273
see also Least-squares estimate and Maximum-likelihood estimate
- Regression equation, 29, 35
- Regression function, 27, 272
 estimation of, 272
 linear, 28
 homoscedastic, 279, 281
- Regular point, 244
- Reiersøl, Olav, 20
- Relation,
 anonymous, 14

- autonomous, 263
- confluent, 70
- economic, 6
- functional, 263
- nonlinear, 50
- stochastic, 6, 263
 - single, estimation of, 3
 - structural, linearity of, 246
- Representation,
 - of the distribution of the variables, 62
 - observationally equivalent, 63
 - parametric, 247
 - structural, 63
- Residuals, 43, 117
 - serial correlation of, 354
- Resolved form, 34, 40
- Restrictions, a priori, 7, 44, 58, 64, 110, 154, 164, 238
 - on \tilde{A} , 178
 - generalization of, 230, 234
 - bilinear, 66, 93, 95
 - counting of, 97, 98
 - on the distribution of disturbances, 66, 183, 191
 - dummy, 85, 159
 - identifying, 176, 248, 411
 - inequalities as, 67
 - just adequate in number and variety, 132
 - normalizing, 68
 - number and variety of, 98
 - single-parameter, 65, 154, 190, 192
 - statistical testing of, 66, 320
- Restriction matrix, 160, 164
- Reverse arrangement, 346
- Revision, iterative,
 - principle of, 157
 - complete, of B_0 , 157
 - simultaneous, of the rows of B_0 , 158
 - successive, of the rows of B_0 , 157
 - see also* Iterative methods for maximizing the likelihood function
- Rubin, H., 2, 5, 111, 136, 245, 312
- Sample,
 - finite, 35, 246
 - infinite, 38, 245
 - large, 39
 - repeated, 24
 - small, 39, 42
 - bias of estimates in, 365
 - confidence region, 317
- Sample space, 265
- Sampling fluctuations, 190
- Sampling variance and covariance
 - of maximum-likelihood estimates, 42, 153, 209, 228
- Scale factors, 158
 - oscillation of, 182
- Stochastic function, 272
- Science,
 - experimental, 2
 - nonexperimental, 17, 32
- Seasonality, *see* Fluctuations
- Sectional model, 7, 22
- Separated form, 34, 40
- Sequence of random variables, 345
- Serial correlation, 310
- Series,
 - circular, 48
 - Fourier, 366
 - polynomial, 366
 - random, 43
 - see also* Time series
- Shifts,
 - of demand, 36
 - in time, 109
- Shocks, 2, 3, 19, 58
 - distribution of, 31, 32, 36
 - independence of, 34, 35, 36, 41
 - nonadditivity of, 21
 - noncorrelation of, 38
 - successive, 33, 39
 - see also* Disturbances
- Shock model, 21
- Shock-and-error model, 20
- Shock structure, 22
- Significance level, 30
- Significance tests in time series, 352
- Simultaneous equations, *see* Equations
- Single-equation approach, 262, 276

- standard error of estimated parameter in, 277
 Single-parameter restrictions, *see* Restrictions, a priori
 Snedecor, G. W., 319, 352
 Snedecor's distribution, 352
 Stable system, *see* System
 Stationary process, 41, 48
 Stieltjes integral, 281
 Stochastic equations, *see* Equations and Difference equations
 Stochastic process, *see* Process
 Stochastic variables, *see* Variables
 Structural equation, *see* Equations
 Structural estimation, *see* Estimation
 Structural parameter, *see* Parameters
 Structural policy, *see* Policy
 Structure, 2, 8, 14, 19, 73, 245, 247, 267
 change of, 3, 11, 16, 25, 38, 267, 275
 controlled, 11, 26
 uncontrolled, 11, 13, 26, 27
 definition of, 8
 economic, 2, 3
 dynamic character of, 3
 equivalence, observational, of, 72, 73, 76
 future, 17
 identifiable, 15, 44
 incomplete, 42
 and models, 246
 observational, 10, 13, 17, 27, 35
 original, 17
 overidentifiable, 44
 shock-, 22
 uniquational, 17, 18
 linear, 16
 nonlinear, 16
 Structure-identifying model, *see* Model
 Supply, 14, 36
 System,
 explosive, without disturbances, 356
 incomplete, 49
 lagless, 273
 nonlinear, 317
 stable, without disturbances, 356, 358
 Systematic part, 261, 262, 264
 T-distribution, 309
 difference equation for the moments of, 346
 ultimate normality of, 346
 T-test (Kendall), 346
 unbiasedness of, 348
 Taste, 19, 50
 Taylor expansion of the likelihood function, 166, 172, 174, 203, 214
 characteristic values of, 175, 187, 217
 multiplicities of, 179
 characteristic vectors of, 179, 186, 217
 orthogonal set of, 176, 178
 Tchebycheff inequality, 359
 Technical equation, 54
 Technology, 14
 Test,
 likelihood ratio, 320
 nonparametric, 47, 345
 based on ranks, 346
 based on small-sample theory, 320
 of significance in time series, 352
 against trend, 345
 of the assumption of zero coefficients, 320
 see also T-test
 Time as a variable, 47, 329
 continuous, 384, 387
 Time lags, 30, 33, 34, 49, 55, 59, 64, 399, 402, 407
 distributed, 385, 387
 see also Lag-correlation coefficients and Lag-covariance matrix
 Time series, 3, 365
 autoregressive, 272
 noncircular, 366
 short, 40, 42, 48
 significance tests in, 352
 Time unit of observation, most

- economic, 385
- Tinbergen, J., 34, 55
- Tintner, G., 2, 4, 20, 26
- Transformation,
 - admissible, 239
 - identical, 11
 - linear, 76, 247, 303
 - nonsingular, 240
- Trend, 32, 47, 48, 329, 330, 331, 333, 347
 - downward, 345
 - test against, 345
 - upward, 345
- "True" variable, *see* Variables

- Ullman, Joseph, 323
- Unbiased estimate, *see* Estimate
- Unequational model, *see* Model
- Unique identification, *see* Identification
- Unitemporal model, 32, 39

- Variables (variates),
 - actual, 21
 - dependent, 32, 35, 55, 405
 - see also* Variables, jointly dependent
 - disturbances in, *see* Disturbances and Errors
 - economic, 49, 262
 - exogenous, 29, 30, 56, 70, 393, 395, 399, 406
 - lagged, 33, 247, 272
 - nonlagged, 32
 - observable, 8, 23
 - exogenous, 29, 30, 33, 41, 49, 56, 70, 263, 276, 393, 399, 402, 406
 - controllable, 11
 - in economic theory, 393
 - observable, 8, 23, 388
 - statistical definition of, 394
 - uncontrollable, 11, 13, 27
 - expected value of, 260
 - fixed, 247, 274, 298, 419
 - intercorrelation of, 258
 - jointly dependent, 33, 59, 70, 111, 312, 314, 384, 402, 405, 406
 - lagged, 33
 - noneconomic, 49
 - nonobservable, 2, 4
 - see also* Disturbances, Errors, and Shocks
 - observable, 3, 6, 8, 19, 23, 32, 268, 271
 - predetermined, 33, 36, 38, 49, 59, 70, 111, 312, 314, 333, 384, 402, 403, 404, 406
 - prospective, 21
 - random (stochastic), 2, 3, 6, 259, 356, 367
 - independent, 259
 - normally distributed, 308
 - sequence of, 345
 - representation of the distribution of, 62
 - theoretical, economic meaning of, 262
 - timing of, 62
 - "true", 262
 - see also* Time as a variable
- Variance, analysis of, 324
- Variates, *see* Variables
- Vector, 81
 - elementary, 175
- Vries, B. A. de, 53

- Wald, A., 2, 4, 20, 32, 41, 42, 44, 45, 55, 56, 67, 103, 115, 133, 135, 148, 356, 366, 400
- Watson, G. N., 372
- Waugh, F. V., 323
- Welfare, *see* Gain
- Whittaker, E. T., 376
- Wilks, S., 352
- Williams, J. D., 370
- Wintner, A., 362, 363
- Wold, Herman, 35, 337
- Wolfowitz, J., 45