

# Linear Inverse Problems in Structural Econometrics Estimation based on spectral decomposition and regularization\*

Marine Carrasco  
University of Rochester

Jean-Pierre Florens  
Université de Toulouse (GREMAQ and IDEI)

Eric Renault  
University of North Carolina, Chapel Hill

---

\*We thank Richard Blundell, Xiaohong Chen, Serge Darolles, James Heckman, Jan Johannes, Oliver Linton, Jean-Michel Loubes, Enno Mammen, Costas Meghir, Whitney Newey, Jean-Francois Richard, Anne Vanhems, and Ed Vytlačil for helpful discussions. Carrasco gratefully acknowledges financial support from the National Science Foundation, grant # SES-0211418.

## Abstract

Inverse problems can be described as functional equations where the value of the function is known or easily estimable but the argument is unknown. Many problems in econometrics can be stated in the form of inverse problems where the argument itself is a function. For example, consider a nonlinear regression where the functional form is the object of interest. One can readily estimate the conditional expectation of the dependent variable given a vector of instruments. From this estimate, one would like to recover the unknown functional form.

This chapter provides an introduction to the estimation of the solution to inverse problems. This chapter focuses mainly on integral equations of the first kind. Solving these equations is particularly challenging as the solution does not necessarily exist, may not be unique, and is not continuous. As a result, a regularized (or smoothed) solution needs to be implemented. We review different regularization methods and study the properties of the estimator. Integral equations of the first kind appear, for example, in the generalized method of moments when the number of moment conditions is infinite, and in the nonparametric estimation of instrumental variable regressions. In the last section of this chapter, we investigate integral equations of the second kind, whose solutions may not be unique but are continuous. Such equations arise when additive models and measurement error models are estimated nonparametrically.

Keywords: Additive models, Integral equation, Generalized Method of Moments, Instrumental variables, Many regressors, Nonparametric estimation, Tikhonov and Landweber-Fridman regularizations.

JEL: C13, C14, C20.

# 1. Introduction

## 1.1. Structural models and functional estimation

The objective of this chapter is to analyze functional estimation in structural econometric models. Different approaches exist to structural inference in econometrics and our presentation may be viewed as a nonparametric extension of the basic example of structural models, namely the static linear simultaneous equations model (SEM). Let us consider  $Y$  a vector of random endogenous variables and  $Z$  a vector of exogenous random variables. A SEM is characterized by a system

$$B_\theta Y + C_\theta Z = U \tag{1.1}$$

where  $B_\theta$  and  $C_\theta$  are matrices that are functions of an unknown “structural” parameter  $\theta$  and  $E[U|Z] = 0$ . The reduced form is a multivariate regression model

$$Y = \Pi Z + V \tag{1.2}$$

where  $\Pi$  is the matrix of ordinary regression coefficients. The relation between reduced and structural form is, in the absence of higher moments restrictions, characterized by:

$$B_\theta \Pi + C_\theta = 0. \tag{1.3}$$

The two essential issues of structural modeling, the identification and the overidentification problems, follow from the consideration of Equation (1.3). The uniqueness of the solution in  $\theta$  for given  $\Pi$  defines the identification problem. The existence of a solution (or restrictions imposed on  $\Pi$  to guarantee the existence) defines the overidentification question. The reduced form parameter  $\Pi$  can be estimated by OLS and if a unique solution in  $\theta$  exists for any  $\Pi$ , it provides the Indirect Least Square estimate of  $\theta$ . If the solution does not exist for any  $\Pi$ ,  $\theta$  can be estimated by a suitable minimization of  $B_\theta \hat{\Pi} + C_\theta$  where  $\hat{\Pi}$  is an estimator of  $\Pi$ .

In this chapter, we address the issue of functional extension of this construction. The data generating process (DGP) is described by a stationary ergodic stochastic process which generates a sequence of observed realizations of a random vector  $X$ .

The structural econometric models considered in this chapter are about the stationary distribution of  $X$ . This distribution is characterized by its cumulative distribution function (c.d.f.)  $F$ , while the functional parameter of interest is an element  $\varphi$  of some infinite dimensional Hilbert space. Following the notation of Florens (2003), the structural econometric model defines the connection between  $\varphi$  and  $F$  under the form of a functional equation:

$$A(\varphi, F) = 0. \tag{1.4}$$

This equation extends Equation (1.3) and the definitions of identification (uniqueness of this solution) and of overidentification (constraints on  $F$  such that a solution exists) are analogous to the SEM case. The estimation is also performed along the same line:  $F$

can be estimated by the empirical distribution of the sample or by a more sophisticated estimator (like kernel smoothing) belonging to the domain of  $A$ .  $\varphi$  is estimated by solving (1.4) or, in the presence of overidentification, by a minimization of a suitable norm of  $A(\varphi, F)$  after plugging in the estimator of  $F$ .

This framework may be clarified by some remarks.

1. All the variables are treated as random in our model and this construction seems to differ from the basic econometric models which are based on a distinction between exogenous or conditioning variables and endogenous variables. Actually this distinction may be used in our framework. Let  $X$  be decomposed into  $Y$  and  $Z$  and  $F$  into  $F_Y(\cdot|Z = z)$  the conditional c.d.f. of  $Y$  given  $Z = z$ , and  $F_Z$  the marginal c.d.f. of  $Z$ . Then, the exogeneity of  $Z$  is tantamount to the conjunction of two conditions. Firstly, the solution  $\varphi$  of (1.4) only depends on  $F_Y(\cdot|Z = z)$  and  $\varphi$  is identified by the conditional model only. Secondly if  $F_Y(\cdot|Z = z)$  and  $F_Z$  are “variations free” in a given statistical model defined by a family of sampling distributions (intuitively no restrictions link  $F_Y(\cdot|Z = z)$  and  $F_Z$ ), no information on  $F_Y(\cdot|Z = z)$  (and then on  $\varphi$ ) is lost by neglecting the estimation of  $F_Z$ . This definition fully encompasses the usual definition of exogeneity in terms of cuts (see Engle, Hendry and Richard (1983), Florens and Mouchart (1985)). Extension of that approach to sequential models and then to sequential or weak exogeneity is straightforward.
2. Our construction does not explicitly involve residuals or other unobservable variables. As will be illustrated in the examples below, most of the structural econometric models are formalized by a relationship between observable and unobservable random elements. A first step in the analysis of these models is to express the relationship between the functional parameters of interest and the DGP, or, in our terminology, to specify the relation  $A(\varphi, F) = 0$ . We start our presentation at the second step of this approach and our analysis is devoted to the study of this equation and to its use for estimation.
3. The overidentification is handled by extending the definition of the parameter in order to estimate overidentified models. Even if  $A(\varphi, F) = 0$  does not have a solution for a given  $F$ , the parameter  $\varphi$  is still defined as the minimum of a norm of  $A(\varphi, F)$ . Then  $\varphi$  can be estimated from an estimation of  $F$ , which does not satisfy the overidentification constraints. This approach extends the original Generalized Method of Moments (GMM) treatment of overidentification. Another way to take into account overidentification constraints consists in estimating  $F$  under these constraints (the estimator of  $F$  is the nearest distribution to the empirical distribution for which there exists a solution,  $\varphi$ , of  $A(\varphi, F) = 0$ ). This method extends the new approach to GMM called the empirical likelihood analysis (see Owen (2001) and references therein). In this chapter, we remain true to the first approach: if the equation  $A(\varphi, F) = 0$  has no solution it will be replaced by the first order condition of the minimization of a norm of  $A(\varphi, F)$ . In that case, this first order condition defines a functional equation usually still denoted  $A(\varphi, F) = 0$ .

## 1.2. Notation

In this chapter,  $X$  is a random element of a finite or infinite dimensional space  $\mathcal{X}$ . In most of the examples,  $\mathcal{X}$  is a finite dimensional euclidean space ( $\mathcal{X} \subset \mathbb{R}^m$ ) and the distribution on  $X$ , denoted  $F$  is assumed to belong to a set  $\mathcal{F}$ . If  $F$  is absolutely continuous, its density is denoted by  $f$ . Usually,  $X$  is decomposed into several components,  $X = (Y, Z, W) \in \mathbb{R}^p \times \mathbb{R}^q \times \mathbb{R}^r$  ( $p+q+r = m$ ) and the marginal c.d.f. or probability density function (p.d.f.) are denoted by  $F_Y, F_Z, F_W$  and  $f_Y, f_X, f_W$  respectively. Conditional c.d.f. are denoted by  $F_Y(\cdot|Z = z)$  or  $F_Y(\cdot|z)$  and conditional density by  $f_Y(\cdot|Z = z)$  or  $f_Y(\cdot|z)$ . The sample may be an i.i.d. sample of  $X$  (denoted in that case  $(x_i)_{i=1,\dots,n}$ ) or weakly dependent time series sample denoted  $(x_t)_{t=1,\dots,T}$  in the dynamic case.

The paper focuses on the estimation of an infinite dimensional parameter denoted by  $\varphi$ , which is an element of a Hilbert space  $\mathcal{H}$  (mathematical concepts are recalled in Section 2). In some particular cases, finite dimensional parameters are considered and this feature is underlined by the notation  $\theta \in \Theta \subset \mathbb{R}^d$ .

The structural model is expressed by an operator  $A$  from  $\mathcal{H} \times \mathcal{F}$  into an Hilbert space  $\mathcal{E}$  and defines the equation  $A(\varphi, F) = 0$ . The (possibly local) solution of this equation is denoted by:

$$\varphi = \Psi(F). \tag{1.5}$$

For statistical discussions, a specific notation for the true value is helpful and  $F_0$  will denote the true c.d.f. (associated with the density  $f_0$  and with the true parameter  $\varphi_0$  (or  $\theta_0$ )). The estimators of the c.d.f. will be denoted by  $F_n$  in an i.i.d. setting or  $F_T$  in a dynamic environment.

The operator  $A$  may take various forms. Particular cases are linear operators with respect to  $F$  or to  $\varphi$ . The first case will be illustrated in the GMM example but most of the paper will be devoted to the study of linear operator relatively to  $\varphi$ . In that case, equation  $A(\varphi, F) = 0$  can be rewritten :

$$A(\varphi, F) = K\varphi - r = 0 \tag{1.6}$$

where  $K$  is a linear operator from  $\mathcal{H}$  to  $\mathcal{E}$  depending on  $F$  and  $r$  is an element of  $\mathcal{E}$  and is also a function of  $F$ . The properties of  $K$  are essential and we will present different examples of integral or differential operators. More generally,  $A$  may be nonlinear either with respect to  $F$  or to  $\varphi$ , but as usual in functional analysis, most of the analysis of nonlinear operators may be done locally (around the true value typically) and reduces to the linear case. Game theoretic models or surplus estimation give examples of nonlinear models.

The problem of solving Equation (1.4) enters in the class of inverse problems. An inverse problem consists of the resolution of an equation where the elements of the equations are imperfectly known. In the linear case, the equation is  $K\varphi = r$  and  $F$  is not exactly known but only estimated. Thus,  $r$  is also imperfectly known. The econometric situation is more complex than most of the inverse problems studied in the statistical literature

because  $K$  is also only imperfectly known. According to the classification proposed by Vapnik (1998), the stochastic inverse problems of interest in this chapter are more often than not characterized by equations where both the operator and the right-hand side term need to be estimated. Inverse problems are said to be well-posed if a unique solution exists and depends continuously on the imperfectly known elements of the equation. In our notation, this means that  $\Psi$  in (1.5) exists as a function of  $F$  and is continuous. Then if  $F$  is replaced by  $F_n$ , the solution  $\varphi_n$  of  $A(\varphi_n, F_n) = 0$  exists and the convergence of  $F_n$  to  $F_0$  implies the convergence of  $\varphi_n$  to  $\varphi_0$  by continuity. Unfortunately a large class of inverse problems relevant to econometric applications are not well-posed (they are then said to be ill-posed in the Hadamard sense, see e.g. Kress (1999), Vapnik (1998)). In this case, a regularization method needs to be implemented to stabilize the solution. Our treatment of ill-posed problems is closed to that of Van Rooij and Ryumgaart (1999).

### 1.3. Examples

This section presents various examples of inverse problems motivated by structural econometric models. We will start with the GMM example, which is the most familiar to econometricians. Subsequently, we present several examples of linear (w.r.t.  $\varphi$ ) inverse problems. The last three examples are devoted to nonlinear inverse problems.

#### 1.3.1. Generalized Method of Moments (GMM)

Let us assume that  $X$  is  $m$  dimensional and the parameter of interest  $\theta$  is also finite dimensional ( $\theta \in \Theta \subset \mathbb{R}^d$ ). We consider a function

$$h : \mathbb{R}^m \times \Theta \rightarrow \mathcal{E} \tag{1.7}$$

and the equation connecting  $\theta$  and  $F$  is defined by:

$$A(\theta, F) = E^F(h(X, \theta)) = 0 \tag{1.8}$$

A particular case is given by  $h(X, \theta) = \mu(X) - \theta$  where  $\theta$  is exactly the expectation of a transformation  $\mu$  of the data. More generally,  $\theta$  may be replaced by an infinite dimensional parameter  $\varphi$  but we do not consider this extension here.

The GMM method was introduced by Hansen (1982) and has received numerous extensions (see Ai and Chen (2003) for the case of an infinite dimensional parameter). GMM consists in estimating  $\theta$  by solving an inverse problem linear in  $F$  but nonlinear in  $\theta$ . It is usually assumed that  $\theta$  is identified i.e. that  $\theta$  is uniquely characterized by Equation (1.8). Econometric specifications are generally overidentified and a solution to (1.8) only exists for some particular  $F$ , including the true DGP  $F_0$ , under the hypothesis of correct specification of the model. The c.d.f  $F$  is estimated by the empirical distribution and the equation (1.8) becomes:

$$\frac{1}{n} \sum_{i=1}^n h(x_i, \theta) = 0, \tag{1.9}$$

which has no solution in general. Overidentification is treated by an extension of the definition of  $\theta$  as follows:

$$\theta = \arg \min_{\theta} \|BE^F(h)\|^2 \quad (1.10)$$

where  $B$  is a linear operator in  $\mathcal{E}$  and  $\|\cdot\|$  denotes the norm in  $\mathcal{E}$ . This definition coincides with (1.8) if  $F$  satisfies the overidentification constraints. Following Equation (1.10), the estimator is:

$$\hat{\theta}_n = \arg \min_{\theta} \left\| B_n \left( \frac{1}{n} \sum_{i=1}^n h(x_i, \theta) \right) \right\|^2 \quad (1.11)$$

where  $B_n$  is a sequence of operators converging to  $B$ . If the number of moment conditions is finite,  $B_n$  and  $B$  are square matrices.

As  $\theta$  is finite dimensional, the inverse problem generated by the first order conditions of (1.10) or (1.11) is well-posed and consistency of the estimators follows from standard regularity conditions. As it will be illustrated in Section 6, an ill-posed inverse problem arises if the number of moment conditions is infinite and if optimal GMM is used. In finite dimensions, optimal GMM is obtained using a specific weighting matrix,  $B = \Sigma^{-\frac{1}{2}}$ , where  $\Sigma$  is the asymptotic variance of  $\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n h(x_i, \theta) \right)$  ( $\Sigma = Var(h)$  in i.i.d. sampling). In the general case, optimal GMM requires the minimization of  $\|g\|^2$  where

$$\Sigma^{\frac{1}{2}} g = E^F(h) \quad (1.12)$$

The function  $g$  is then the solution of a linear inverse problem. If the dimension of  $h$  is not finite, Equation (1.12) defines an ill-posed inverse problem, which requires a regularization scheme (see Section 3).

### 1.3.2. Instrumental variables

Instrumental regression is a possible strategy to perform nonparametric estimation when explanatory variables are endogenous. Let us decompose  $X$  into  $(Y, X, W)$  where  $Y \in \mathbb{R}$ ,  $Z \in \mathbb{R}^q$ ,  $W \in \mathbb{R}^r$ . The subvectors  $Z$  and  $W$  may have elements in common. The econometrician starts with a relation

$$Y = \varphi(Z) + U \quad (1.13)$$

where  $U$  is a random term which does not satisfy  $E(U|Z) = 0$ . This assumption is replaced by the more general hypothesis

$$E(U|W) = 0 \quad (1.14)$$

and  $W$  is called the set of instrumental variables. Condition (1.14) defines  $\varphi$  as the solution of an integral equation. In terms of density, (1.14) means that

$$A(\varphi, F) = \int \varphi(z) f_Z(z|W = w) dz - \int y f_Y(y|W = w) dy = 0 \quad (1.15)$$

Using previous notation, the first part of (1.15) is denoted  $K\varphi$  and the second part is equal to  $r$ .

This expression is linear in  $\varphi$  and can be made linear in  $F$  by eliminating the denominator through a multiplication by  $f_W(w)$ . However, as will be seen later, this problem is essentially nonlinear in  $F$  because the treatment of overidentification and of regularization will necessarily reintroduce the denominator in (1.15).

Instrumental regression introduced in (1.15) can be generalized to local instrumental regression and to generalized local instrumental regression. These extensions are relevant in more complex models than (1.13), where in particular the error term may enter the equation in non additive ways (see for such a treatment, Florens, Heckman, Meghir, and Vytlacil (2002)). For example, consider the equation

$$Y = \varphi(Z) + Z\varepsilon + U \tag{1.16}$$

where  $Z$  is scalar and  $\varepsilon$  is a random unobservable heterogeneity component. It can be proved that, under a set of identification assumptions,  $\varphi$  satisfies the equations :

$$A_j(\varphi, F) = E^F \left( \frac{\partial \varphi(Z)}{\partial Z} | W = w \right) - \frac{\frac{\partial}{\partial W_j} E(Y|W = w)}{\frac{\partial}{\partial W_j} E(Z|W = w)} = 0 \tag{1.17}$$

for any  $j = 1, \dots, r$ . This equation, linear with respect to  $\varphi$ , combines integral and differential operators.

Instrumental variable estimation and its local extension define ill-posed inverse problems as will be seen in Section 5.

### 1.3.3. Deconvolution

Another classical example of ill-posed inverse problem is given by the deconvolution problem. Let us assume that  $X, Y, Z$  be three scalar random elements such that

$$Y = X + Z \tag{1.18}$$

Only  $Y$  is observable. The two components  $X$  and  $Z$  are independent. The density of  $Z$  (the error term) is known and denoted  $g$ . The parameter of interest is the density  $\varphi$  of  $X$ . Then  $\varphi$  is solution of:

$$\begin{aligned} A(\varphi, F) &= \int \varphi(y)g(x - y)dy - f(x) = 0 \\ &= K\varphi - r \end{aligned} \tag{1.19}$$

This example is comparable to the instrumental variables case but only the r.h.s.  $r = f$  is unknown whereas the operator  $K$  is given.



### 1.3.4. Regression with many regressors

This example also constitutes a case of linear ill-posed inverse problems. Let us consider a regression model where the regressors are indexed by  $\tau$  belonging to an infinite index set provided with a measure  $\Pi$ . The model says:

$$Y = \int Z(\tau)\varphi(\tau)\Pi(d\tau) + U \quad (1.20)$$

where  $E(U|(Z(\tau))_\tau) = 0$  and  $\varphi$  is the parameter of interest and is infinite dimensional. Examples of regression with many regressors are now common in macroeconomics (see Stock and Watson (2002) or Forni and Reichlin (1998) for two presentations of this topic).

Let us assume that  $Y$  and  $(Z(\tau))_\tau$  are observable. Various treatments of (1.20) can be done and we just consider the following analysis. The conditional moment equation  $E(U|(Z(\tau))_\tau) = 0$  implies an infinite number of conditions indexed by  $\tau$ :

$$E(Z(\tau)U) = 0, \quad \forall \tau$$

or equivalently

$$\int E^F(Z(\tau)Z(\rho))\varphi(\rho)\Pi(d\rho) - E^F(YZ(\tau)) = 0, \quad \forall \tau \quad (1.21)$$

This equation generalizes the usual normal equations of the linear regression to an infinite number of regressors. The inverse problem defined in (1.21) is linear in both  $F$  and  $\varphi$  but it is ill posed. An intuitive argument to illustrate this issue is to consider the estimation using a finite number of observations of the second moment operator  $E^F(Z(\tau)Z(\rho))$  which is infinite dimensional. The resulting multicollinearity problem is solved by a ridge regression. The “infinite matrix”  $E^F(Z(\cdot)Z(\cdot))$  is replaced by  $\alpha I + E^F(Z(\cdot)Z(\cdot))$  where  $I$  is the identity and  $\alpha$  a positive number, or by a reduction of the set of regressors to the first principal components. These two solutions are particular examples of regularization methods (namely the Tikhonov and the spectral cut-off regularizations), which will be introduced in Section 3.

### 1.3.5. Additive models

The properties of the integral equations generated by this example and by the next one are very different from that of the three previous examples. We consider an additive regression model:

$$Y = \varphi(Z) + \psi(W) + U \quad (1.22)$$

where  $E(U|Z, W) = 0$  and  $X = (Y, Z, W)$  is the observable element. The parameters of interest are the two functions  $\varphi$  and  $\psi$ . The approach we propose here is related to the backfitting approach (see Hastie and Tibshirani (1990)). Other treatments of additive

models have been considered in the literature (see Pagan and Ullah (1999)). Equation (1.22) implies

$$\begin{cases} E^F(Y|Z = z) = \varphi(z) + E^F(\psi(W)|Z = z) \\ E^F(Y|W = w) = E^F(\varphi(Z)|W = w) + \psi(w) \end{cases} \quad (1.23)$$

and by substitution

$$\begin{aligned} \varphi(z) - E^F(E^F(\varphi(Z)|W)|Z = z) \\ = E^F(Y|Z = z) - E^F(E^F(Y|W)|Z = z) \end{aligned} \quad (1.24)$$

or, in our notations:

$$(I - K)\varphi = r$$

where  $K = E^F(E^F(\cdot |W)|Z)$ . Backfitting refers to the iterative method to solve Equation (1.23).

An analogous equation characterizes  $\psi$ . Actually even if (1.22) is not well specified, these equations provide the best approximation of the regression of  $Y$  given  $Z$  and  $W$  by an additive form. Equation (1.24) is a linear integral equation and even if this inverse problem is ill-posed because  $K$  is not one-to-one ( $\varphi$  is only determined up to a constant term), the solution is still continuous and therefore the difficulty is not as important as that of the previous examples.

### 1.3.6. Measurement-error models or nonparametric analysis of panel data

We denote  $\eta$  to be an unobservable random variable for which two measurements  $Y_1$  and  $Y_2$  are available. These measurements are affected by a bias dependent on observable variables  $Z_1$  and  $Z_2$ . More formally:

$$\begin{cases} Y_1 = \eta + \varphi(Z_1) + U_1 & E(U_1|\eta, Z_1, Z_2) = 0 \\ Y_2 = \eta + \varphi(Z_2) + U_2 & E(U_2|\eta, Z_1, Z_2) = 0 \end{cases} \quad (1.25)$$

An i.i.d. sample  $(y_{1i}, y_{2i}, \eta_i, z_{1i}, z_{2i})$  is drawn but the  $\eta_i$  are unobservable. Equivalently this model may be seen as a two period panel data with individual effects  $\eta_i$ .

The parameter of interest is the “bias function”  $\varphi$ , identical for the two observations. In the measurement context, it is natural to assume that the joint distribution of the observables is independent of the order of the observations, or equivalently  $(Y_1, Z_1, Y_2, Z_2)$  are distributed as  $(Y_2, Z_2, Y_1, Z_1)$ . This assumption is not relevant in a dynamic context.

The model is transformed in order to eliminate the unobservable variable by difference:

$$Y = \varphi(Z_2) - \varphi(Z_1) + U \quad (1.26)$$

where  $Y = Y_2 - Y_1$ ,  $U = U_2 - U_1$ , and  $E(U|Z_1, Z_2) = 0$ .

This model is similar to an additive model except for the symmetry between the variables, and the fact that with the notation of (1.22),  $\varphi$  and  $\psi$  are identical. An application of this model may be found in Gaspar and Florens (1998) where  $y_{1i}$  and  $y_{2i}$  are two measurements of the level of the ocean in location  $i$  by a satellite radar altimeter,  $\eta_i$  is the true level and  $\varphi$  is the “sea state bias” depending on the waves’ height and the wind speed ( $Z_{1i}$  and  $Z_{2i}$  are both two dimensional).

The model is treated through the relation:

$$E(Y|Z_2 = z_2) = \varphi(z_2) - E(\varphi(Z_1)|Z_2 = z_2) \quad (1.27)$$

which defines an integral equation  $K\varphi = r$ . The exchangeable property between the variables implies that conditioning on  $Z_1$  gives the same equation (where  $Z_1$  and  $Z_2$  are exchanged).

### 1.3.7. Game theoretic model

This example and the next ones present economic models formalized by nonlinear inverse problems. As the focus of this chapter is on linear equations, these examples are given for illustration and will not be treated outside of this section. The analysis of nonlinear functional equations raises numerous questions: uniqueness and existence of the solution, asymptotic properties of the estimator, implementation of the estimation procedure and numerical computation of the solution.

Most of these questions are usually solved locally by a linear approximation of the nonlinear problem deduced from a suitable concept of derivative. A strong concept of derivation (typically Frechet derivative) is needed to deal with the implicit form of the model, which requires the use of the Implicit Function theorem.

The first example of nonlinear inverse problems follows from the strategic behavior of the players in a game. Let us assume that for each game, each player receives a random signal or type denoted by  $\xi$  and plays an action  $X$ . The signal is generated by a probability described by its c.d.f.  $\varphi$ , and the players all adopt a strategy  $\sigma$  dependent on  $\varphi$  which associates  $X$  with  $\xi$ , i.e.

$$X = \sigma_\varphi(\xi).$$

The strategy  $\sigma_\varphi$  is determined as an equilibrium of the game (e.g. Nash equilibrium) or by an approximation of the equilibrium (bounded rationality behavior). The signal  $\xi$  is private knowledge for the player but is unobserved by the econometrician, and the c.d.f.  $\varphi$  is common knowledge for the players but is unknown for the statistician. The strategy  $\sigma_\varphi$  is determined from the rules of the game and by the assumptions on the behavior of the players. The essential feature of the game theoretic model from a statistical viewpoint is that the relation between the unobservable and the observable variables depends on the distribution of the unobservable component. The parameter of interest is the c.d.f.  $\varphi$  of the signals.

Let us restrict our attention to cases where  $\xi$  and  $X$  are scalar and where  $\sigma_\varphi$  is strictly increasing. Then the c.d.f.  $F$  of the observable  $X$  is connected with  $\varphi$  by:

$$A(\varphi, F) = F \circ \sigma_\varphi - \varphi = 0 \quad (1.28)$$

If the signals are i.i.d. across the different players and different games,  $F$  can be estimated by a smooth transformation of the empirical distribution and Equation (1.28) is solved in  $\varphi$ . The complexity of this relation can be illustrated by the auction model. In the private value first price auction model,  $\xi$  is the value of the object and  $X$  the bid. If the number of bidders is  $N + 1$  the strategy function is equal to:

$$X = \xi - \frac{\int_{\xi}^{\xi} \varphi^N(u) du}{\varphi^N(\xi)} \quad (1.29)$$

where  $[\underline{\xi}, \bar{\xi}]$  is the support of  $\xi$  and  $\varphi^N(u) = [\varphi(u)]^N$  is the c.d.f. of the maximum private value among  $N$  players.

Model (1.28) may be extended to a non iid setting (depending on exogenous variables) or to the case where  $\sigma_\varphi$  is partially unknown. The analysis of this model has been done by Guerre, Perrigne and Vuong (2000) in a nonparametric context. The framework of inverse problem is used by Florens, Protopopescu and Richard (1997).

### 1.3.8. Solution of a differential equation

In several models like the analysis of the consumer surplus, the function of interest is the solution of a differential equation depending on the data generating process.

Consider for example a class of problems where  $X = (Y, Z, W) \in \mathbb{R}^3$  is i.i.d.,  $F$  is the c.d.f. of  $X$  and the parameter  $\varphi$  verifies:

$$\frac{d}{dz}\varphi(z) = m_F(z, \varphi(z)) \quad (1.30)$$

when  $m_F$  is a regular function depending on  $F$ . A first example is

$$m_F(z, w) = E^F(Y|Z = z, W = w) \quad (1.31)$$

but more complex examples may be constructed in order to take into account the endogeneity of one or two variables. For example,  $Z$  may be endogenous and  $m_F$  may be defined by:

$$E(Y|W_1 = w_1, W_2 = w_2) = E(m_F(Z, W_1)|W_1 = w_1, W_2 = w_2) \quad (1.32)$$

Economic applications can be found in Hausman (1981, 1985) and Hausman and Newey (1995) and a theoretical treatment of these two problems is given by Vanhems (2000) and Loubes and Vanhems (2001).

### 1.3.9. Instrumental variables in a nonseparable model

Another example of a nonlinear inverse problem is provided by the following model:

$$Y = \varphi(Z, U) \tag{1.33}$$

where  $Z$  is an endogenous variable. The function  $\varphi$  is the parameter of interest. Denote  $\varphi_z(u) = \varphi(z, u)$ . Assume that  $\varphi_z(u)$  is an increasing function of  $u$  for each  $z$ . Moreover, the distribution,  $F_U$  of  $U$  is assumed to be known for identification purposes. Model (1.33) may arise in a duration model where  $Y$  is the duration (see Equation (2.2) of Horowitz 1999). One difference with Horowitz (1999) is the presence of an endogenous variable here. There is a vector of instruments  $W$ , which are independent of  $U$ . Because  $U$  and  $W$  are independent, we have

$$P(U \leq u | W = w) = P(U \leq u) = F_U(u). \tag{1.34}$$

Denote  $f$  the density of  $(Y, Z)$  and

$$F(y, z | w) = \int_{-\infty}^y f(t, z | w) dt.$$

$F$  can be estimated using the observations  $(y_i, z_i, w_i)$ ,  $i = 1, 2, \dots, n$ . By a slight abuse of notation, we use the notation  $P(Y \leq y, Z = z | W = w)$  for  $F(y, z | w)$ . We have

$$\begin{aligned} P(U \leq u, Z = z | W = w) &= P(\varphi_z(Y)^{-1} \leq u, Z = z | W = w) \\ &= P(Y \leq \varphi_z(u), Z = z | W = w) \\ &= F(\varphi_z(u), z | w). \end{aligned} \tag{1.35}$$

Combining Equations (1.34) and (1.35), we obtain

$$\int F(\varphi_z(u), z | w) dz = F_U(u). \tag{1.36}$$

Equation (1.36) belongs to the class of Urysohn equations of Type I (Polyanin and Manzhirov, 1998). The estimation of the solution of Equation (1.36) is discussed in Florens (2005).

## 1.4. Organization of the chapter

Section 2 reviews the basic definitions and properties of operators in Hilbert spaces. The focus is on compact operators because they have the advantage of having a discrete spectrum. We recall some laws of large numbers and central limit theorems for Hilbert valued random elements. Finally, we discuss how to estimate the spectrum of a compact operator and how to estimate the operators themselves.

Section 3 is devoted to solving integral equations of the first kind. As these equations are ill-posed, the solution needs to be regularized (or smoothed). We investigate the properties of the regularized solutions for different types of regularizations.

In Section 4, we show under suitable assumptions the consistency and asymptotic normality of regularized solutions.

Section 5 detail five examples: the ridge regression, the factor model, the infinite number of regressors, the deconvolution, and the instrumental variables estimation.

Section 6 has two parts. First, it recalls the main results relative to reproducing kernels. Reproducing kernel theory is closely related to that of the integral equations of the first kind. Second, we explain the extension of GMM to a continuum of moment conditions and show how the GMM objective function reduces to the norm of the moment functions in a specific reproducing kernel Hilbert space. Several examples are provided.

Section 7 tackles the problem of solving integral equations of the second kind. A typical example of such a problem is the additive model introduced earlier.

## 2. Spaces and Operators

The purpose of this section is to introduce terminology and to state the main properties of operators in Hilbert spaces that are used in our econometric applications. Most of these results can be found in Debnath and Mikusinsky (1999) and Kress (1999). Ait-Sahalia, Hansen, and Scheinkman (2004) provide an excellent survey of operator methods for the purpose of financial econometrics.

### 2.1. Hilbert spaces

We start by recalling some of the basic concepts of analysis. In the sequel,  $\mathbf{C}$  denotes the set of complex numbers. A vector space equipped by a norm is called a normed space. A sequence  $(\varphi_n)$  of elements in a normed space is called a Cauchy sequence if for every  $\varepsilon > 0$  there exists an integer  $N(\varepsilon)$  such that

$$\|\varphi_n - \varphi_m\| < \varepsilon$$

for all  $n, m \geq N(\varepsilon)$ , i.e, if  $\lim_{n,m \rightarrow \infty} \|\varphi_n - \varphi_m\| = 0$ . A space  $S$  is complete if every Cauchy sequence converges to an element in  $S$ . A complete normed vector space is called a Banach space.

Let  $(E, \mathcal{E}, \Pi)$  be a probability space and

$$L_C^p(E, \mathcal{E}, \Pi) = \left\{ f : E \rightarrow \mathbf{C} \text{ measurable s.t. } \|f\| \equiv \left( \int |f|^p d\Pi \right)^{1/p} < \infty \right\}, p \geq 1.$$

Then,  $L_C^p(E, \mathcal{E}, \Pi)$  is a Banach space. If we only consider functions valued in  $\mathbb{R}$  this space is still a Banach space and is denoted in that case by  $L^p$  (we drop the subscript  $C$ ). In the sequel, we also use the following notation. If  $E$  is a subset of  $\mathbf{R}^p$ , then the  $\sigma$ -field  $\mathcal{E}$  will always be the Borel  $\sigma$ -field and will be omitted in the notation  $L^p(\mathbf{R}^p, \Pi)$ . If  $\Pi$  has a density  $\pi$  with respect to Lebesgue measure,  $\Pi$  will be replaced by  $\pi$ . If  $\pi$  is uniform, it will be omitted in the notation.

**Definition 2.1 (Inner product).** Let  $H$  be a complex vector space. A mapping  $\langle, \rangle : H \times H \rightarrow \mathbf{C}$  is called an inner product in  $H$  if for any  $\varphi, \psi, \xi \in H$  and  $\alpha, \beta \in \mathbf{C}$  the following conditions are satisfied:

- (a)  $\langle \varphi, \psi \rangle = \overline{\langle \psi, \varphi \rangle}$  (the bar denotes the complex conjugate),
- (b)  $\langle \alpha\varphi + \beta\psi, \xi \rangle = \alpha \langle \varphi, \xi \rangle + \beta \langle \psi, \xi \rangle$ ,
- (c)  $\langle \varphi, \varphi \rangle \geq 0$  and  $\langle \varphi, \varphi \rangle = 0 \iff \varphi = 0$ .

A vector space equipped by an inner product is called an inner product space.

**Example.** The space  $\mathbf{C}^N$  of ordered  $N$ -tuples  $x = (x_1, \dots, x_N)$  of complex numbers, with the inner product defined by

$$\langle x, y \rangle = \sum_{l=1}^N x_l \overline{y_l}$$

is an inner product space.

**Example.** The space  $l^2$  of all sequences  $(x_1, x_2, \dots)$  of complex numbers such that  $\sum_{j=1}^{\infty} |x_j|^2 < \infty$  with the inner product defined by  $\langle x, y \rangle = \sum_{j=1}^{\infty} x_j \overline{y_j}$  for  $x = (x_1, x_2, \dots)$  and  $y = (y_1, y_2, \dots)$  is an infinite dimensional inner product space.

**Example.** The space  $L_C^2(E, \mathcal{E}, \Pi)$  associated with the inner product defined by

$$\langle \varphi, \psi \rangle = \int \varphi \overline{\psi} d\Pi$$

is an inner product space. On the other hand,  $L_C^p(E, \mathcal{E}, \Pi)$  is not a inner product space if  $p \neq 2$ .

An inner product satisfies the Cauchy-Schwartz inequality, that is,

$$|\langle \varphi, \psi \rangle|^2 \leq \langle \varphi, \varphi \rangle \langle \psi, \psi \rangle$$

for all  $\varphi, \psi \in H$ . Remark that  $\langle \varphi, \varphi \rangle$  is real because  $\langle \varphi, \varphi \rangle = \overline{\langle \varphi, \varphi \rangle}$ . It actually defines a norm  $\|\varphi\| = \langle \varphi, \varphi \rangle^{1/2}$  (this is the norm induced by the inner product  $\langle, \rangle$ ).

**Definition 2.2 (Hilbert space).** If an inner product space is complete in the induced norm, it is called a Hilbert space.

A standard theorem in functional analysis guarantees that every inner product space  $H$  can be completed to form a Hilbert space  $\mathcal{H}$ . Such a Hilbert space is said to be the completion of  $H$ .

**Example.**  $\mathbf{C}^N$ ,  $l^2$  and  $L^2(\mathbb{R}, \Pi)$  are Hilbert spaces.

**Example.** (Sobolev space) Let  $\Omega = [a, b]$  be an interval of  $\mathbb{R}$ . Denote by  $\tilde{H}^m(\Omega)$ ,  $m = 1, 2, \dots$ , the space of all complex-valued functions  $\varphi \in \mathcal{C}^m$  such that for all  $|l| \leq m$ ,  $\varphi^{(l)} = \partial^l \varphi(\tau) / \partial \tau^l \in L^2(\Omega)$ . The inner product on  $\tilde{H}^m(\Omega)$  is

$$\langle \varphi, \psi \rangle = \int_a^b \sum_{l=0}^m \varphi^{(l)}(\tau) \overline{\psi^{(l)}(\tau)} d\tau.$$

$\tilde{H}^m(\Omega)$  is an inner product space but it is not a Hilbert space because it is not complete. The completion of  $\tilde{H}^m(\Omega)$ , denoted  $H^m(\Omega)$ , is a Hilbert space.

**Definition 2.3 (Convergence).** A sequence  $(\varphi_n)$  of vectors in an inner product space  $H$ , is called strongly convergent to a vector  $\varphi \in H$  if  $\|\varphi_n - \varphi\| \rightarrow 0$  as  $n \rightarrow \infty$ .

Remark that if  $(\varphi_n)$  converges strongly to  $\varphi$  in  $H$  then  $\langle \varphi_n, \psi \rangle \rightarrow \langle \varphi, \psi \rangle$  as  $n \rightarrow \infty$ , for every  $\psi \in H$ . The converse is false.

**Definition 2.4.** Let  $H$  be an inner product space. A sequence  $(\varphi_n)$  of nonzero vectors in  $H$  is called an orthogonal sequence if  $\langle \varphi_m, \varphi_n \rangle = 0$  for  $n \neq m$ . If in addition  $\|\varphi_n\| = 1$  for all  $n$ , it is called an orthonormal sequence.

**Example.** Let  $\pi(x)$  be the pdf of a normal with mean  $\mu$  and variance  $\sigma^2$ . Denote by  $\phi_j$  the Hermite polynomials of degree  $j$ :

$$\phi_j(x) = (-1)^j \frac{d^j \pi}{dx^j}. \quad (2.1)$$

The functions  $\phi_j(x)$  form an orthogonal system in  $L^2(\mathbb{R}, \pi)$ .

Any sequence of vectors  $(\psi_j)$  in an inner product space that is linearly independent, i.e.,

$$\sum_{j=1}^{\infty} \alpha_j \psi_j = 0 \Rightarrow \alpha_j = 0 \quad \forall j = 1, 2, \dots$$

can be transformed into an orthonormal sequence by the method called Gram-Schmidt orthonormalization process. This process consists of the following steps. Given  $(\psi_j)$ , define a sequence  $(\varphi_j)$  inductively as

$$\begin{aligned} \varphi_1 &= \frac{\psi_1}{\|\psi_1\|}, \\ \varphi_2 &= \frac{\psi_2 - \langle \psi_2, \varphi_1 \rangle \varphi_1}{\|\psi_2 - \langle \psi_2, \varphi_1 \rangle \varphi_1\|} \\ &\vdots \\ \varphi_n &= \frac{\psi_n - \sum_{l=1}^{n-1} \langle \psi_n, \varphi_l \rangle \varphi_l}{\|\psi_n - \sum_{l=1}^{n-1} \langle \psi_n, \varphi_l \rangle \varphi_l\|}. \end{aligned}$$

As a result,  $(\varphi_j)$  is orthonormal and any linear combinations of vectors  $\varphi_1, \dots, \varphi_n$  is also a linear combinations of  $\psi_1, \dots, \psi_n$  and vice versa.

**Theorem 2.5 (Pythagorean formula).** If  $\varphi_1, \dots, \varphi_n$  are orthogonal vectors in an inner product space, then

$$\left\| \sum_{j=1}^n \varphi_j \right\|^2 = \sum_{j=1}^n \|\varphi_j\|^2.$$



From the Pythagorean formula, it can be seen that the  $\alpha_l$  that minimize

$$\left\| \varphi - \sum_{j=1}^n \alpha_j \varphi_j \right\|$$

are such that  $\alpha_j = \langle \varphi, \varphi_j \rangle$ . Moreover

$$\sum_{j=1}^n |\langle \varphi, \varphi_j \rangle|^2 \leq \|\varphi\|^2. \quad (2.2)$$

Hence the series  $\sum_{j=1}^{\infty} |\langle \varphi, \varphi_j \rangle|^2$  converges for every  $\varphi \in H$ . The expansion

$$\varphi = \sum_{j=1}^{\infty} \langle \varphi, \varphi_j \rangle \varphi_j \quad (2.3)$$

is called a generalized Fourier series of  $\varphi$ . In general, we do not know whether the series in (2.3) is convergent. Below we give a sufficient condition for convergence.

**Definition 2.6 (Complete orthonormal sequence).** An orthonormal sequence  $(\varphi_j)$  in an inner product space  $H$  is said to be complete if for every  $\varphi \in H$  we have

$$\varphi = \sum_{j=1}^{\infty} \langle \varphi, \varphi_j \rangle \varphi_j$$

where the equality means

$$\lim_{n \rightarrow \infty} \left\| \varphi - \sum_{j=1}^n \langle \varphi, \varphi_j \rangle \varphi_j \right\| = 0$$

where  $\|\cdot\|$  is the norm in  $H$ .

A complete orthonormal sequence  $(\varphi_j)$  in an inner product space  $H$  is an orthonormal basis in  $H$ , that is every  $\varphi \in H$  has a unique representation  $\varphi = \sum_{j=1}^{\infty} \alpha_j \varphi_j$  where  $\alpha_l \in \mathbf{C}$ . If  $(\varphi_j)$  is a complete orthonormal sequence in an inner product space  $H$  then the set

$$\text{span} \{\varphi_1, \varphi_2, \dots\} = \left\{ \sum_{j=1}^n \alpha_j \varphi_j : \forall n \in \mathbf{N}, \forall \alpha_1, \dots, \alpha_n \in \mathbf{C} \right\}$$

is dense in  $H$ .

**Theorem 2.7.** An orthonormal sequence  $(\varphi_j)$  in a Hilbert space  $\mathcal{H}$  is complete if and only if  $\langle \varphi, \varphi_j \rangle = 0$  for all  $j = 1, 2, \dots$  implies  $\varphi = 0$ .

**Theorem 2.8 (Parseval's formula).** *An orthonormal sequence  $(\varphi_j)$  in a Hilbert space  $\mathcal{H}$  is complete if and only if*

$$\|\varphi\|^2 = \sum_{j=1}^{\infty} |\langle \varphi, \varphi_j \rangle|^2 \quad (2.4)$$

for every  $\varphi \in \mathcal{H}$ .

**Definition 2.9 (Separable space).** *A Hilbert space is called separable if it contains a complete orthonormal sequence.*

**Example.** A complete orthonormal sequence in  $L^2([-\pi, \pi])$  is given by

$$\phi_j(x) = \frac{e^{ijx}}{\sqrt{2\pi}}, \quad j = \dots, -1, 0, 1, \dots$$

Hence, the space  $L^2([-\pi, \pi])$  is separable.

**Theorem 2.10.** *Every separable Hilbert space contains a countably dense subset.*

## 2.2. Definitions and basic properties of operators

In the sequel, we denote  $K : \mathcal{H} \rightarrow \mathcal{E}$  the operator that maps a Hilbert space  $\mathcal{H}$  (with norm  $\|\cdot\|_{\mathcal{H}}$ ) into a Hilbert space  $\mathcal{E}$  (with norm  $\|\cdot\|_{\mathcal{E}}$ ).

**Definition 2.11.** *An operator  $K : \mathcal{H} \rightarrow \mathcal{E}$  is called linear if*

$$K(\alpha\varphi + \beta\psi) = \alpha K\varphi + \beta K\psi$$

for all  $\varphi, \psi \in \mathcal{H}$  and all  $\alpha, \beta \in \mathbf{C}$ .

**Definition 2.12.** (i) *The null space of  $K : \mathcal{H} \rightarrow \mathcal{E}$  is the set  $\mathcal{N}(K) = \{\varphi \in \mathcal{H} : K\varphi = 0\}$ .*

(ii) *The range of  $K : \mathcal{H} \rightarrow \mathcal{E}$  is the set  $\mathcal{R}(K) = \{\psi \in \mathcal{E} : \psi = K\varphi \text{ for some } \varphi \in \mathcal{H}\}$ .*

(iii) *The domain of  $K : \mathcal{H} \rightarrow \mathcal{E}$  is the subset of  $\mathcal{H}$  denoted  $\mathcal{D}(K)$  on which  $K$  is defined.*

(iv) *An operator is called finite dimensional if its range is of finite dimension.*

**Theorem 2.13.** *A linear operator is continuous if it is continuous at one element.*

**Definition 2.14.** *A linear operator  $K : \mathcal{H} \rightarrow \mathcal{E}$  is called bounded if there exists a positive number  $C$  such that*

$$\|K\varphi\|_{\mathcal{E}} \leq C \|\varphi\|_{\mathcal{H}}$$

for all  $\varphi \in \mathcal{H}$ .

**Definition 2.15.** The norm of a bounded operator  $K$  is defined as

$$\|K\| \equiv \sup_{\|\varphi\| \leq 1} \|K\varphi\|_{\mathcal{E}}$$

**Theorem 2.16.** A linear operator is continuous if and only if it is bounded.

**Example.** The identity operator defined by  $\mathcal{I}\varphi = \varphi$  for all  $\varphi \in \mathcal{H}$  is bounded with  $\|\mathcal{I}\| = 1$ .

**Example.** Consider the differential operator:

$$(D\varphi)(x) = \frac{d\varphi(\tau)}{d\tau} = \varphi'(\tau)$$

defined on the space  $E_1 = \{\varphi \in L^2([-\pi, \pi]) : \varphi' \in L^2([-\pi, \pi])\}$  with norm  $\|\varphi\| = \sqrt{\int_{-\pi}^{\pi} |f(\tau)|^2 d\tau}$ . For  $\varphi_j(\tau) = \sin j\tau$ ,  $j = 1, 2, \dots$ , we have  $\|\varphi_j\| = \sqrt{\int_{-\pi}^{\pi} |\sin(j\tau)|^2 d\tau} = \sqrt{\pi}$  and  $\|D\varphi_j\| = \sqrt{\int_{-\pi}^{\pi} |j \cos(j\tau)|^2 d\tau} = j\sqrt{\pi}$ . Therefore  $\|D\varphi_j\| = j\|\varphi_j\|$  proving that the differential operator is not bounded.

**Theorem 2.17.** Each linear operator  $K$  from a finite dimensional normed space  $\mathcal{H}$  into a normed space  $\mathcal{E}$  is bounded.

An important class of linear operators are valued in  $\mathbf{C}$  and they are characterized by Riesz theorem. From (2.2), we know that for any fixed vector  $g$  in an inner product space  $H$ , the formula  $G(\varphi) = \langle \varphi, g \rangle$  defines a bounded linear functional on  $H$ . It turns out that if  $H$  is a Hilbert space, then every bounded linear functional is of this form.

**Theorem 2.18 (Riesz).** Let  $\mathcal{H}$  be a Hilbert space. Then for each bounded linear function  $G : \mathcal{H} \rightarrow \mathbf{C}$  there exists a unique element  $g \in \mathcal{H}$  such that

$$G(\varphi) = \langle \varphi, g \rangle$$

for all  $\varphi \in \mathcal{H}$ . The norms of the element  $g$  and the linear function  $F$  coincide

$$\|g\|_{\mathcal{H}} = \|G\|$$

where  $\|\cdot\|_{\mathcal{H}}$  is the norm in  $\mathcal{H}$  and  $\|\cdot\|$  is the operator norm.

**Definition 2.19 (Hilbert space isomorphism).** A Hilbert space  $\mathcal{H}_1$  is said to be isometrically isomorphic (congruent) to a Hilbert space  $\mathcal{H}_2$  if there exists a one-to-one linear mapping  $J$  from  $\mathcal{H}_1$  to  $\mathcal{H}_2$  such that

$$\langle J(\varphi), J(\psi) \rangle_{\mathcal{H}_2} = \langle \varphi, \psi \rangle_{\mathcal{H}_1}$$

for all  $\varphi, \psi \in \mathcal{H}_1$ . Such a mapping  $J$  is called a Hilbert space isomorphism (or congruence) from  $\mathcal{H}_1$  to  $\mathcal{H}_2$ .

The terminology “congruence” is used by Parzen (1959, 1970).

**Theorem 2.20.** *Let  $\mathcal{H}$  be a separable Hilbert space.*

- (a) *If  $\mathcal{H}$  is infinite dimensional, then it is isometrically isomorphic to  $l^2$ .*
- (b) *If  $\mathcal{H}$  has a dimension  $N$ , then it is isometrically isomorphic to  $\mathbf{C}^N$ .*

A consequence of Theorem 2.20 is that two separable Hilbert spaces of the same dimension (finite or infinite) are isometrically isomorphic.

**Theorem 2.21.** *Let  $\mathcal{H}$  and  $\mathcal{E}$  be Hilbert spaces and let  $K : \mathcal{H} \rightarrow \mathcal{E}$  be a bounded operator. Then there exists a uniquely determined linear operator  $K^* : \mathcal{E} \rightarrow \mathcal{H}$  with the property*

$$\langle K\varphi, \psi \rangle_{\mathcal{E}} = \langle \varphi, K^*\psi \rangle_{\mathcal{H}}$$

for all  $\varphi \in \mathcal{H}$  and  $\psi \in \mathcal{E}$ . Moreover, the operator  $K^*$  is bounded and  $\|K\| = \|K^*\|$ .  $K^*$  is called the adjoint operator of  $K$ .

Riesz Theorem 2.18 implies that, in Hilbert spaces, the adjoint of a bounded operator always exists.

**Example 2.1. (discrete case)** Let  $\pi$  and  $\rho$  be two discrete probability density functions on  $\mathbb{N}$ . Let  $\mathcal{H} = L^2(\mathbb{N}, \pi) = \{\varphi : \mathbb{N} \rightarrow \mathbb{R}, \varphi = (\varphi_l)_{l \in \mathbb{N}} \text{ such that } \sum_{l \in \mathbb{N}} \varphi_l^2 \pi(l) < \infty\}$  and  $\mathcal{E} = L^2(\mathbb{N}, \rho)$ . The operator  $K$  that associates to elements  $(\varphi_l)_{l \in \mathbb{N}}$  of  $\mathcal{H}$  elements  $(\psi_p)_{p \in \mathbb{N}}$  of  $\mathcal{E}$  such that

$$(K\varphi)_p = \psi_p = \sum_{l \in \mathbb{N}} k(p, l) \varphi_l \pi(l)$$

is an infinite dimensional matrix. If  $\mathcal{H}$  and  $\mathcal{E}$  are finite dimensional, then  $K$  is simply a matrix.

**Example 2.2. (integral operator)** An important kind of operator is the integral operator. Let  $\mathcal{H} = L^2_C(\mathbb{R}^q, \pi)$  and  $\mathcal{E} = L^2_C(\mathbb{R}^r, \rho)$  where  $\pi$  and  $\rho$  are pdf. The integral operator  $K : \mathcal{H} \rightarrow \mathcal{E}$  is defined as

$$K\varphi(\tau) = \int k(\tau, s) \varphi(s) \pi(s) ds. \tag{2.5}$$

The function  $k$  is called the kernel of the operator. If  $k$  satisfies

$$\int \int |k(\tau, s)|^2 \pi(s) \rho(\tau) ds d\tau < \infty \tag{2.6}$$

( $k$  is said to be a  $L^2$ -kernel) then  $K$  is a bounded operator and

$$\|K\| \leq \sqrt{\int \int |k(\tau, s)|^2 \pi(s) \rho(\tau) ds d\tau}.$$

Indeed for any  $\varphi \in \mathcal{H}$ , we have

$$\begin{aligned}\|K\varphi\|_{\mathcal{E}}^2 &= \int \left| \int k(\tau, s) \varphi(s) \pi(s) ds \right|^2 \rho(\tau) d\tau \\ &= \int |\langle k(\tau, \cdot), \varphi(\cdot) \rangle_{\mathcal{H}}|^2 \rho(\tau) d\tau \\ &\leq \int \|k(\tau, \cdot)\|_{\mathcal{H}}^2 \|\varphi\|_{\mathcal{H}}^2 \rho(\tau) d\tau\end{aligned}$$

by Cauchy-Schwarz inequality. Hence we have

$$\begin{aligned}\|K\varphi\|_{\mathcal{E}}^2 &\leq \|\varphi\|_{\mathcal{H}}^2 \int \|k(\tau, \cdot)\|_{\mathcal{H}}^2 \rho(\tau) d\tau \\ &= \|\varphi\|_{\mathcal{H}}^2 \int \int |k(\tau, s)|^2 \pi(s) \rho(\tau) ds d\tau.\end{aligned}$$

The upperbound for  $\|K\|$  follows.

The adjoint  $K^*$  of the operator  $K$  is also an integral operator

$$K^*\psi(s) = \int k^*(s, \tau) \psi(\tau) \rho(\tau) d\tau \quad (2.7)$$

with  $k^*(s, \tau) = \overline{k(\tau, s)}$ . Indeed, we have

$$\begin{aligned}\langle K\varphi, \psi \rangle_{\mathcal{E}} &= \int (K\varphi)(\tau) \overline{\psi(\tau)} \rho(\tau) d\tau \\ &= \int \left( \int k(\tau, s) \varphi(s) \pi(s) ds \right) \overline{\psi(\tau)} \rho(\tau) d\tau \\ &= \int \varphi(s) \left( \int k(\tau, s) \overline{\psi(\tau)} \rho(\tau) \right) \pi(s) ds \\ &= \int \varphi(s) \overline{\left( \int k^*(s, \tau) \psi(\tau) \rho(\tau) \right)} \pi(s) ds \\ &= \langle \varphi, K^*\psi \rangle_{\mathcal{H}}.\end{aligned}$$

There are two types of integral operators we are interested in, the covariance operator and the conditional expectation operator.

**Example 2.3. (conditional expectation operator)** When  $K$  is a conditional expectation operator, it is natural to define the spaces of reference as a function of unknown pdfs. Let  $(Z, W) \in \mathbb{R}^q \times \mathbb{R}^r$  be a r.v. with distribution  $F_{Z,W}$ , let  $F_Z$ , and  $F_W$  be the marginal distributions of  $Z$  and  $W$  respectively. The corresponding pdfs are denoted  $f_{Z,W}$ ,  $f_Z$ , and  $f_W$ . Define

$$\begin{aligned}\mathcal{H} &= L^2(\mathbb{R}^q, f_Z) \equiv L_Z^2, \\ \mathcal{E} &= L^2(\mathbb{R}^r, f_W) \equiv L_W^2.\end{aligned}$$

Let  $K$  be the conditional expectation operator:

$$\begin{aligned} K & : L_Z^2 \rightarrow L_W^2 \\ \varphi & \rightarrow E[\varphi(Z)|W]. \end{aligned} \tag{2.8}$$

$K$  is an integral operator with kernel

$$k(w, z) = \frac{f_{Z,W}(z, w)}{f_Z(z) f_W(w)}$$

By Equation (2.7), its adjoint  $K^*$  has kernel  $k^*(z, w) = k(w, z)$  and is also a conditional expectation operator:

$$\begin{aligned} K^* & : L_W^2 \rightarrow L_Z^2 \\ \psi & \rightarrow E[\psi(W)|Z]. \end{aligned}$$

**Example 2.4. (Restriction of an operator on a subset of  $\mathcal{H}$ )** Let  $K : \mathcal{H} \rightarrow \mathcal{E}$  and consider the restriction denoted  $K_0$  of  $K$  on a subspace  $\mathcal{H}_0$  of  $\mathcal{H}$ .  $K_0 : \mathcal{H}_0 \rightarrow \mathcal{E}$  is such that  $K_0$  and  $K$  coincide on  $\mathcal{H}_0$ . It can be shown that the adjoint  $K_0^*$  of  $K_0$  is the operator mapping  $\mathcal{E}$  into  $\mathcal{H}_0$  such that

$$K_0^* = PK^* \tag{2.9}$$

where  $P$  is the projection on  $\mathcal{H}_0$ . The expression of  $K_0^*$  will reflect the extra information contained in  $\mathcal{H}_0$ .

To prove (2.9), we use the definition of  $K^*$  :

$$\begin{aligned} \langle K\varphi, \psi \rangle_{\mathcal{E}} & = \langle \varphi, K^*\psi \rangle_{\mathcal{H}} \text{ for all } \varphi \in \mathcal{H}_0 \\ & = \langle \varphi, K_0^*\psi \rangle_{\mathcal{H}_0} \text{ for all } \varphi \in \mathcal{H}_0 \\ \Leftrightarrow & \langle \varphi, K^*\psi - K_0^*\psi \rangle_{\mathcal{H}} = 0 \text{ for all } \varphi \in \mathcal{H}_0 \\ \Leftrightarrow & K^*\psi - K_0^*\psi \in \mathcal{H}_0^\perp \\ \Leftrightarrow & K_0^*\psi = PK^*\psi. \end{aligned}$$

A potential application of this result to the conditional expectation in Example 2.3 is the case where  $\varphi$  is known to be additive. Let  $Z = (Z_1, Z_2)$ . Then

$$\mathcal{H}_0 = \{\varphi(Z) = \varphi_1(Z_1) + \varphi_2(Z_2) : \varphi_1 \in L_{Z_1}^2, \varphi_2 \in L_{Z_2}^2\}.$$

Assume that  $E[\varphi_1(Z_1)] = E[\varphi_2(Z_2)] = 0$ . We have  $P\varphi = (\varphi_1, \varphi_2)$  with

$$\begin{aligned} \varphi_1 & = (I - P_1P_2)^{-1}(P_1 - P_1P_2)\varphi, \\ \varphi_2 & = (I - P_1P_2)^{-1}(P_2 - P_1P_2)\varphi, \end{aligned}$$

where  $P_1$  and  $P_2$  are the projection operators on  $L_{Z_1}^2$  and  $L_{Z_2}^2$  respectively. If the two spaces  $L_{Z_1}^2$  and  $L_{Z_2}^2$  are orthogonal, then  $\varphi_1 = P_1\varphi$  and  $\varphi_2 = P_2\varphi$ .

**Definition 2.22 (Self-adjoint).** If  $K = K^*$  then  $K$  is called self-adjoint (or Hermitian).

Remark that if  $K$  is a self-adjoint integral operator, then  $k(s, \tau) = \overline{k(\tau, s)}$ .

**Theorem 2.23.** Let  $K : \mathcal{H} \rightarrow \mathcal{H}$  be a self-adjoint operator then

$$\|K\| = \sup_{\|\varphi\|=1} |\langle K\varphi, \varphi \rangle_{\mathcal{H}}|.$$

**Definition 2.24 (Positive operator).** An operator  $K : \mathcal{H} \rightarrow \mathcal{H}$  is called positive if it is self-adjoint and  $\langle K\varphi, \varphi \rangle_{\mathcal{H}} \geq 0$ .

**Definition 2.25.** A sequence  $(K_n)$  of operators  $K_n : \mathcal{H} \rightarrow \mathcal{E}$  is called pointwise convergent if for every  $\varphi \in \mathcal{H}$ , the sequence  $K_n\varphi$  converges in  $\mathcal{E}$ . A sequence  $(K_n)$  of bounded operators converges in norm to a bounded operator  $K$  if  $\|K_n - K\| \rightarrow 0$  as  $n \rightarrow \infty$ .

**Definition 2.26 (Compact operator).** A linear operator  $K : \mathcal{H} \rightarrow \mathcal{E}$  is called a compact operator if for every bounded sequence  $(\varphi_n)$  in  $\mathcal{H}$ , the sequence  $(K\varphi_n)$  contains a convergent subsequence in  $\mathcal{E}$ .

**Theorem 2.27.** Compact linear operators are bounded.

Not every bounded operator is compact. An example is given by the identity operator on an infinite dimensional space  $\mathcal{H}$ . Consider an orthonormal sequence  $(e_n)$  in  $\mathcal{H}$ . Then the sequence  $\mathcal{I}e_n = e_n$  does not contain a convergent subsequence.

**Theorem 2.28.** Finite dimensional operators are compact.

**Theorem 2.29.** If the sequence  $K_n : \mathcal{H} \rightarrow \mathcal{E}$  of compact linear operators are norm convergent to a linear operator  $K : \mathcal{H} \rightarrow \mathcal{E}$ , i.e.,  $\|K_n - K\| \rightarrow 0$  as  $n \rightarrow \infty$ , then  $K$  is compact. Moreover, every compact operator is the limit of a sequence of operators with finite dimensional range.

Hilbert Schmidt operators are discussed in Dunford and Schwartz (1988, p. 1009), Dautray and Lyons (1988, Vol 5, p.41, chapter VIII).

**Definition 2.30 (Hilbert-Schmidt operator).** Let  $\{\varphi_j, j = 1, 2, \dots\}$  be a complete orthonormal set in a Hilbert space  $\mathcal{H}$ . An operator  $K : \mathcal{H} \rightarrow \mathcal{E}$  is said to be a Hilbert-Schmidt operator if the quantity  $\|\cdot\|_{HS}$  defined by

$$\|K\|_{HS} = \left\{ \sum_{j=1}^{\infty} \|K\varphi_j\|_{\mathcal{E}}^2 \right\}^{1/2}$$

is finite. The number  $\|K\|_{HS}$  is called the Hilbert-Schmidt norm of  $K$ . Moreover

$$\|K\| \leq \|K\|_{HS} \tag{2.10}$$

and hence  $K$  is bounded.

From (2.10), it follows that HS norm convergence implies (operator) norm convergence.

**Theorem 2.31.** *The Hilbert-Schmidt norm is independent of the orthonormal basis used in its definition.*

**Theorem 2.32.** *Every Hilbert-Schmidt operator is compact.*

**Theorem 2.33.** *The adjoint of a Hilbert-Schmidt operator is itself a Hilbert-Schmidt operator and  $\|K\|_{HS} = \|K^*\|_{HS}$ .*

Theorem 2.32 implies that Hilbert-Schmidt (HS) operators can be approached by a sequence of finite dimensional operators, which is an attractive feature when it comes to estimating  $K$ . Remark that the integral operator  $K$  defined by (2.5) and (2.6) is a Hilbert-Schmidt (HS) operator and its adjoint is also a HS operator. Actually, all Hilbert-Schmidt operators of  $L^2(\mathbb{R}^q, \pi)$  in  $L^2(\mathbb{R}^r, \rho)$  are integral operators. The following theorem is proved in Dautray and Lions (Vol. 5, p. 45).

**Theorem 2.34.** *An operator of  $L^2(\mathbb{R}^q, \pi)$  in  $L^2(\mathbb{R}^r, \rho)$  is Hilbert-Schmidt if and only if it admits a kernel representation (2.5) conformable to (2.6). In this case, the kernel  $k$  is unique.*

**Example 2.1 (continued).** Let  $K$  from  $L^2(\mathbb{N}, \pi)$  in  $L^2(\mathbb{N}, \rho)$  with kernel  $k(l, p)$ .  $K$  is a Hilbert-Schmidt operator if  $\sum \sum k(l, p)^2 \pi(l) \rho(p) < \infty$ . In particular, the operator defined by  $(K\varphi)_1 = \varphi_1$  and  $(K\varphi)_p = \varphi_p - \varphi_{p-1}$ ,  $p = 2, 3, \dots$  is not a Hilbert-Schmidt operator; it is not even compact.

**Example 2.3 (continued).** By Theorem 2.34, a sufficient condition for  $K$  and  $K^*$  to be compact is

$$\int \int \left[ \frac{f_{Z,W}(z, w)}{f_Z(z) f_W(w)} \right]^2 f_Z(z) f_W(w) dz dw < \infty.$$

**Example 2.5 (Conditional expectation with common elements).** Consider a conditional expectation operator from  $L^2(X, Z)$  into  $L^2(X, W)$  defined by

$$(K\varphi)(x, w) = E[\varphi(X, Z) | X = x, W = w].$$

Because there are common elements between the conditioning variable and the argument of the function  $\varphi$ , the operator  $K$  is not compact. Indeed, let  $\varphi(X)$  be such that  $E(\varphi^2) = 1$ , we have  $K\varphi = \varphi$ . It follows that the image of the unit circle in  $L^2(X, Z)$  contains the unit circle of  $L^2(X)$  and hence is not compact. Therefore,  $K$  is not compact.

**Example 2.6 (Restriction).** For illustration, we consider the effect of restricting  $K$  on a subset of  $L_C^2(\mathbb{R}^q, \pi)$ . Consider  $\tilde{K}$  the operator defined by

$$\begin{aligned} \tilde{K} & : L_C^2(\mathbb{R}^q, \tilde{\pi}) \rightarrow L_C^2(\mathbb{R}^r, \tilde{\rho}) \\ \tilde{K}\varphi & = K\varphi \end{aligned}$$



for every  $\varphi \in L_C^2(\mathbb{R}^q, \tilde{\pi})$ , where  $L_C^2(\mathbb{R}^q, \tilde{\pi}) \subset L_C^2(\mathbb{R}^q, \pi)$  and  $L_C^2(\mathbb{R}^r, \tilde{\rho}) \supset L_C^2(\mathbb{R}^r, \rho)$ . Assume that  $K$  is a HS operator defined by (2.5). Under which conditions is  $\tilde{K}$  an HS operator? Let

$$\begin{aligned}\tilde{K}\varphi(s) &= \int k(\tau, s) \varphi(s) \pi(s) ds \\ &= \int k(\tau, s) \frac{\pi(s)}{\tilde{\pi}(s)} \varphi(s) \tilde{\pi}(s) ds \\ &\equiv \int \tilde{k}(\tau, s) \varphi(s) \tilde{\pi}(s) ds.\end{aligned}$$

Assume that  $\tilde{\pi}(s) = 0$  implies  $\pi(s) = 0$  and  $\rho(\tau) = 0$  implies  $\tilde{\rho}(\tau) = 0$ . Note that

$$\begin{aligned}& \int \left| \tilde{k}(\tau, s) \right|^2 \tilde{\pi}(s) \tilde{\rho}(\tau) ds d\tau \\ &= \int |k(\tau, s)|^2 \frac{\pi(s)}{\tilde{\pi}(s)} \frac{\tilde{\rho}(\tau)}{\rho(\tau)} \pi(s) \rho(\tau) ds d\tau \\ &< \sup_s \left| \frac{\pi(s)}{\tilde{\pi}(s)} \right| \sup_\tau \left| \frac{\tilde{\rho}(\tau)}{\rho(\tau)} \right| \int |k(\tau, s)|^2 \pi(s) \rho(\tau) ds d\tau.\end{aligned}$$

Hence the HS property is preserved if (a) there is a constant  $c > 0$  such that  $\pi(s) \leq c\tilde{\pi}(s)$  for all  $s \in \mathbb{R}^q$  and (b) there is a constant  $d$  such that  $\tilde{\rho}(\tau) \leq d\rho(\tau)$  for all  $\tau \in \mathbb{R}^r$ .

### 2.3. Spectral decomposition of compact operators

For compact operators, spectral analysis reduces to the analysis of eigenvalues and eigenfunctions. Let  $K : \mathcal{H} \rightarrow \mathcal{H}$  be a compact linear operator.

**Definition 2.35.**  $\lambda$  is an eigenvalue of  $K$  if there is a nonzero vector  $\phi \in \mathcal{H}$  such that  $K\phi = \lambda\phi$ .  $\phi$  is called the eigenfunction of  $K$  corresponding to  $\lambda$ .

**Theorem 2.36.** All eigenvalues of a self-adjoint operator are real and eigenfunctions corresponding to different eigenvalues are orthogonal.

**Theorem 2.37.** All eigenvalues of a positive operator are nonnegative.

**Theorem 2.38.** For every eigenvalue  $\lambda$  of a bounded operator  $K$ , we have  $|\lambda| \leq \|K\|$ .

**Theorem 2.39.** Let  $K$  be a self-adjoint compact operator, the set of its eigenvalues  $(\lambda_j)$  is countable and its eigenvectors  $(\phi_j)$  can be orthonormalized. Its largest eigenvalue (in absolute value) satisfies  $|\lambda_1| = \|K\|$ . If  $K$  has infinitely many eigenvalues  $|\lambda_1| \geq |\lambda_2| \geq \dots$ , then  $\lim_{j \rightarrow \infty} \lambda_j = 0$ .

Let  $K : \mathcal{H} \rightarrow \mathcal{E}$ ,  $K^*K$  and  $KK^*$  are self-adjoint positive operators on  $\mathcal{H}$  and  $\mathcal{E}$  respectively. Hence their eigenvalues are nonnegative by Theorem 2.37.

**Definition 2.40.** Let  $\mathcal{H}$  and  $\mathcal{E}$  be Hilbert spaces,  $K : \mathcal{H} \rightarrow \mathcal{E}$  be a compact linear operator and  $K^* : \mathcal{E} \rightarrow \mathcal{H}$  be its adjoint. The square roots of the eigenvalues of the nonnegative self-adjoint compact operator  $K^*K : \mathcal{H} \rightarrow \mathcal{H}$  are called the singular values of  $K$ .

The following results (Kress, 1999, Theorem 15.16) apply to operators that are not necessarily self-adjoint.

**Theorem 2.41.** Let  $(\lambda_j)$  denote the sequence of the nonzero singular values of the compact linear operator  $K$  repeated according to their multiplicity. Then there exist orthonormal sequences  $\phi_j$  of  $\mathcal{H}$  and  $\psi_j$  of  $\mathcal{E}$  such that

$$K\phi_j = \lambda_j\psi_j, \quad K^*\psi_j = \lambda_j\phi_j \quad (2.11)$$

for all  $j \in N$ . For each  $\varphi \in \mathcal{H}$  we have the singular value decomposition

$$\varphi = \sum_{j=1}^{\infty} \langle \varphi, \phi_j \rangle \phi_j + Q\varphi \quad (2.12)$$

with the orthogonal projection operator  $Q : \mathcal{H} \rightarrow \mathcal{N}(K)$  and

$$K\varphi = \sum_{j=1}^{\infty} \lambda_j \langle \varphi, \phi_j \rangle \psi_j. \quad (2.13)$$

$\{\lambda_j, \phi_j, \psi_j\}$  is called the singular system of  $K$ . Note that  $\lambda_j^2$  are the nonzero eigenvalues of  $KK^*$  and  $K^*K$  associated with the eigenfunctions  $\psi_j$  and  $\phi_j$  respectively.

**Theorem 2.42.** Let  $K$  be the integral operator defined by (2.5) and assume Condition (2.6) holds. Let  $\{\lambda_j, \phi_j, \psi_j\}$  be as in (2.11). Then:

(i) The Hilbert Schmidt norm of  $K$  can be written as

$$\|K\|_{HS} = \left\{ \sum_{j \in N} |\lambda_j|^2 \right\}^{1/2} = \left\{ \int \int |k(\tau, s)|^2 \pi(s) \rho(\tau) ds d\tau \right\}^{1/2}$$

where each  $\lambda_j$  is repeated according to its multiplicity.

(ii) (Mercer's formula)  $k(\tau, s) = \sum_{j=1}^{\infty} \lambda_j \psi_j(\tau) \overline{\phi_j(s)}$ .

**Example (degenerate operator).** Consider an integral operator defined on  $L^2([a, b])$  with a Pincherle-Goursat kernel

$$Kf(\tau) = \int_a^b k(\tau, s) f(s) ds,$$

$$k(\tau, s) = \sum_{l=1}^n a_l(\tau) b_l(s).$$

Assume that  $a_l$  and  $b_l$  belong to  $L^2([a, b])$  for all  $l$ . By (2.6), it follows that  $K$  is bounded. Moreover, as  $K$  is finite dimensional, we have  $K$  compact by Theorem 2.28. Assume that the set of functions  $(a_l)$  is linearly independent. The equality  $K\phi = \lambda\phi$  yields

$$\sum_{l=1}^n a_l(\tau) \int b_l(s) \phi(s) ds = \lambda\phi(\tau),$$

hence  $\phi(\tau)$  is necessarily of the form  $\sum_{l=1}^n c_l a_l(\tau)$ . The dimension of the range of  $K$  is therefore  $n$ , there are at most  $n$  nonzero eigenvalues.

**Example.** Let  $\mathcal{H} = L^2([0, 1])$  and the integral operator  $Kf(\tau) = \int_0^1 (\tau \wedge s) f(s) ds$  where  $\tau \wedge s = \min(\tau, s)$ . It is possible to explicitly compute the eigenvalues and eigenfunctions of  $K$  by solving  $K\phi = \lambda\phi \iff \int_0^\tau s\phi(s) ds + \tau \int_\tau^1 \phi(s) ds = \lambda\phi(\tau)$ . Using two successive differentiations with respect to  $\tau$ , we obtain a differential equation  $\phi(\tau) = -\lambda\phi''(\tau)$  with boundary conditions  $\phi(0) = 0$  and  $\phi'(1) = 0$ . Hence the set of orthonormal eigenvectors is  $\phi_j(\tau) = \sqrt{2} \sin((\pi j \tau)/2)$  associated with the eigenvalues  $\lambda_j = 4/(\pi^2 j^2)$ ,  $j = 1, 3, 5, \dots$ . We can see that the eigenvalues converge to zero at an arithmetic rate.

**Example.** Let  $\pi$  be the pdf of the standard normal distribution and  $\mathcal{H} = L^2(\mathbb{R}, \pi)$ . Define  $K$  as the integral operator with kernel

$$k(\tau, s) = \frac{l(\tau, s)}{\pi(\tau)\pi(s)}$$

where  $l(\tau, s)$  is the joint pdf of the bivariate normal  $\mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$ . Then  $K$  is a self-adjoint operator with eigenvalues  $\lambda_j = \rho^j$  and has eigenfunctions that take the Hermite polynomial form  $\phi_j$ ,  $j = 1, 2, \dots$  defined in (2.1). This is an example where the eigenvalues decay exponentially fast.

## 2.4. Random element in Hilbert spaces

### 2.4.1. Definitions

Let  $\mathcal{H}$  be a real separable Hilbert space with norm  $\|\cdot\|$  induced by the inner product  $\langle \cdot, \cdot \rangle$ . Let  $(\Omega, \mathcal{F}, P)$  be a complete probability space. Let  $X : \Omega \rightarrow \mathcal{H}$  be a Hilbert space-valued random element (an  $\mathcal{H}$ -r.e.).  $X$  is integrable or has finite expectation  $E(X)$  if  $E(\|X\|) = \int_\Omega \|X\| dP < \infty$ , in that case  $E(X)$  satisfies  $E(X) \in \mathcal{H}$  and  $E[\langle X, \varphi \rangle] = \langle E(X), \varphi \rangle$  for all  $\varphi \in \mathcal{H}$ . An  $\mathcal{H}$ -r.e.  $X$  is weakly second order if  $E[\langle X, \varphi \rangle^2] < \infty$  for all  $\varphi \in \mathcal{H}$ . For a weakly second order  $\mathcal{H}$ -r.e.  $X$  with expectation  $E(X)$ , we define the covariance operator  $K$  as

$$\begin{aligned} K & : \mathcal{H} \rightarrow \mathcal{H} \\ K\varphi & = E[\langle X - E(X), \varphi \rangle (X - E(X))] \end{aligned}$$

for all  $\varphi \in \mathcal{H}$ . Note that  $\text{var} \langle X, \varphi \rangle = \langle K\varphi, \varphi \rangle$ .

**Example.** Let  $\mathcal{H} = L^2([0, 1])$  with  $\|g\| = \left[ \int_0^1 g(\tau)^2 d\tau \right]^{1/2}$  and  $X = h(\tau, Y)$  where  $Y$  is a random variable and  $h(\cdot, Y) \in L^2([0, 1])$  with probability one. Assume  $E(h(\tau, Y)) = 0$ , then the covariance operator takes the form:

$$\begin{aligned} K\varphi(\tau) &= E[\langle h(\cdot, Y), \varphi \rangle h(\tau, Y)] \\ &= E\left[\left(\int h(s, Y)\varphi(s) ds\right) h(\tau, Y)\right] \\ &= \int E[h(\tau, Y)h(s, Y)]\varphi(s) ds \\ &\equiv \int k(\tau, s)\varphi(s) ds. \end{aligned}$$

Moreover, if  $h(\tau, Y) = I\{Y \leq \tau\} - F(\tau)$  then  $k(\tau, s) = F(\tau \wedge s) - F(\tau)F(s)$ .

**Definition 2.43.** An  $\mathcal{H}$ -r.e.  $Y$  has a Gaussian distribution on  $\mathcal{H}$  if for all  $\varphi \in \mathcal{H}$  the real-valued r.v.  $\langle \varphi, Y \rangle$  has a Gaussian distribution on  $\mathbb{R}$ .

**Definition 2.44 (strong mixing).** Let  $\{X_{i,n}, i = \dots, -1, 0, 1, \dots; n \geq 1\}$  be an array of  $\mathcal{H}$ -r.e., defined on the probability space  $(\Omega, \mathcal{F}, P)$  and define  $\mathcal{A}_{n,a}^{n,b} = \sigma(X_{i,n}, a \leq i \leq b)$  for all  $-\infty \leq a \leq b \leq +\infty$ , and  $n \geq 1$ . The array  $\{X_{i,n}\}$  is called a strong or  $\alpha$ -mixing array of  $\mathcal{H}$ -r.e. if  $\lim_{j \rightarrow \infty} \alpha(j) = 0$  where

$$\alpha(j) = \sup_{n \geq 1} \sup_l \sup_{A, B} \left[ |P(A \cap B) - P(A)P(B)| : A \in \mathcal{A}_{n,-\infty}^{n,l}, B \in \mathcal{A}_{n,l+j}^{n,+\infty} \right].$$

#### 2.4.2. Central limit theorem for mixing processes

We want to study the asymptotic properties of  $Z_n = n^{-1/2} \sum_{i=1}^n X_{i,n}$  where  $\{X_{i,n} : 1 \leq i \leq n\}$  is an array of  $\mathcal{H}$ -r.e.. Weak and strong laws of large numbers for near epoch dependent (NED) processes can be found in Chen and White (1996). Here we provide sufficient conditions for the weak convergence of processes to be denoted  $\Rightarrow$  (see Davidson, 1994, for a definition). Weak convergence is stronger than the standard central limit theorem (CLT) as illustrated by a simple example. Let  $(X_i)$  be an iid sequence of zero mean weakly second order elements of  $\mathcal{H}$ . Then for any  $Z$  in  $\mathcal{H}$ ,  $\langle X_i, Z \rangle$  is an iid zero mean sequence of  $\mathbf{C}$  with finite variance  $\langle KZ, Z \rangle$ . Then standard CLT implies the asymptotic normality of  $\frac{1}{\sqrt{n}} \sum_{i=1}^n \langle X_i, Z \rangle$ . The weak convergence of  $\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i$  to a Gaussian process  $\mathcal{N}(0, K)$  in  $\mathcal{H}$  requires an extra assumption, namely  $E\|X_1\|^2 < \infty$ . Weak convergence theorems for NED processes that might have trending mean (hence are not covariance stationary) are provided by Chen and White (1998). Here, we report results for mixing processes proved by Politis and Romano (1994). See also van der Vaart and Wellner (1996) for iid sequences.

**Theorem 2.45.** Let  $\{X_{i,n} : 1 \leq i \leq n\}$  be a double array of stationary mixing  $\mathcal{H}$ -r.e. with zero mean, such that for all  $n$ ,  $\|X_{i,n}\| < B$  with probability one, and  $\sum_{j=1}^m j^2 \alpha(j) < \infty$

$Km^r$  for all  $1 \leq m \leq n$  and  $n$ , and some  $r < 3/2$ . Assume, for any integer  $l \geq 1$ , that  $(X_{1,n}, \dots, X_{l,n})$ , regarded as a r.e. of  $\mathcal{H}^l$ , converges in distribution to say,  $(X_1, \dots, X_l)$ . Moreover, assume  $E[\langle X_{1,n}, X_{l,n} \rangle] \rightarrow E[\langle X_1, X_l \rangle]$  as  $n \rightarrow \infty$  and

$$\lim_{n \rightarrow \infty} \sum_{l=1}^n E[\langle X_{1,n}, X_{l,n} \rangle] = \sum_{l=1}^{\infty} E[\langle X_1, X_l \rangle] < \infty.$$

Let  $Z_n = n^{-1/2} \sum_{i=1}^n X_{i,n}$ . For any  $\varphi \in \mathcal{H}$ , let  $\sigma_{\varphi,n}^2$  denote the variance of  $\langle Z_n, \varphi \rangle$ . Assume

$$\sigma_{\varphi,n}^2 \xrightarrow{n \rightarrow \infty} \sigma_{\varphi}^2 \equiv \text{Var}(\langle X_1, \varphi \rangle) + 2 \sum_{i=1}^{\infty} \text{cov}(\langle X_1, \varphi \rangle, \langle X_{1+i}, \varphi \rangle). \quad (2.14)$$

Then  $Z_n$  converges weakly to a Gaussian process  $\mathcal{N}(0, K)$  in  $\mathcal{H}$ , with zero mean and covariance operator  $K$  satisfying  $\langle K\varphi, \varphi \rangle = \sigma_{\varphi}^2$  for each  $\varphi \in \mathcal{H}$ .

In the special case when the  $X_{i,n} = X_i$  form a stationary sequence, the conditions simplify considerably:

**Theorem 2.46.** Assume  $X_1, X_2, \dots$  is a stationary sequence of  $\mathcal{H}$ -r.e. with mean  $\mu$  and mixing coefficient  $\alpha$ . Let  $Z_n = n^{-1/2} \sum_{i=1}^n (X_i - \mu)$ .

(i) If  $E(\|X_1\|^{2+\delta}) < \infty$  for some  $\delta > 0$ , and  $\sum_j [\alpha(j)]^{\delta/(2+\delta)} < \infty$

(ii) or if  $X_1, X_2, \dots$  is iid and  $E\|X_1\|^2 < \infty$

Then  $Z_n$  converges weakly to a Gaussian process  $G \sim \mathcal{N}(0, K)$  in  $\mathcal{H}$ . The distribution of  $G$  is determined by the distribution of its marginals  $\langle G, \varphi \rangle$  which are  $\mathcal{N}(0, \sigma_{\varphi}^2)$  distributed for every  $\varphi \in \mathcal{H}$  where  $\sigma_{\varphi}^2$  is defined in (2.14).

Let  $\{e_l\}$  be a complete orthonormal basis of  $\mathcal{H}$ . Then  $\|X_1\|^2 = \sum_{l=1}^{\infty} \langle X_1, e_l \rangle^2$  and hence in the iid case, it suffices to check that  $E\|X_1\|^2 = \sum_{l=1}^{\infty} E[\langle X_1, e_l \rangle^2] < \infty$ .

The following theorem is stated in more general terms in Chen and White (1992).

**Theorem 2.47.** Let  $A_n$  be a random bounded linear operator from  $\mathcal{H}$  to  $\mathcal{H}$  and  $A \neq 0$  be a nonrandom bounded linear operator from  $\mathcal{H}$  to  $\mathcal{H}$ . If  $\|A_n - A\| \rightarrow 0$  in probability as  $n \rightarrow \infty$  and  $Y_n \Rightarrow Y \sim \mathcal{N}(0, K)$  in  $\mathcal{H}$ . Then  $A_n Y_n \Rightarrow AY \sim \mathcal{N}(0, AK A^*)$ .

In Theorem 2.47, the boundedness of  $A$  is crucial. In most of our applications,  $A$  will not be bounded and we will not be able to apply Theorem 2.47. Instead we will have to check the Liapunov condition (Davidson 1994) ‘‘by hand’’.

**Theorem 2.48.** Let the array  $\{X_{i,n}\}$  be independent with zero mean and variance sequence  $\{\sigma_{i,n}^2\}$  satisfying  $\sum_{i=1}^n \sigma_{i,n}^2 = 1$ . Then  $\sum_{i=1}^n X_{i,n} \xrightarrow{d} \mathcal{N}(0, 1)$  if

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n E[|X_{i,n}|^{2+\delta}] = 0 \quad (\text{Liapunov condition})$$

for some  $\delta > 0$ .

## 2.5. Estimation of an operator and its adjoint

### 2.5.1. Estimation of an operator

In many cases of interest, an estimator of the compact operator,  $K$ , is given by a degenerate operator of the form

$$\hat{K}_n \varphi = \sum_{l=1}^{L_n} a_l(\varphi) \varepsilon_l \quad (2.15)$$

where  $\varepsilon_l \in \mathcal{E}$ ,  $a_l(\varphi)$  is linear in  $\varphi$ .

Examples:

1 - Covariance operator

$$K\varphi(\tau_1) = \int E[h(\tau_1, X)h(\tau_2, X)]\varphi(\tau_2) d\tau_2.$$

Replacing the expectation by the sample mean, one obtains an estimator of  $K$  :

$$\begin{aligned} \hat{K}_n \varphi(\tau_1) &= \int \left( \frac{1}{n} \sum_{i=1}^n h(\tau_1, x_i) h(\tau_2, x_i) \right) \varphi(\tau_2) d\tau_2 \\ &= \sum_{i=1}^n a_i(\varphi) \varepsilon_i \end{aligned}$$

with

$$a_i(\varphi) = \frac{1}{n} \int h(\tau_2, x_i) \varphi(\tau_2) d\tau_2 \text{ and } \varepsilon_i = h(\tau_1, x_i).$$

Note that here  $K$  is self-adjoint and the rate of convergence of  $\hat{K}_n$  is parametric.

2 - Conditional expectation operator

$$K\varphi(w) = E[\varphi(Z) | W = w].$$

The kernel estimator of  $K$  with kernel  $\omega$  and bandwidth  $c_n$  is given by

$$\begin{aligned} \hat{K}_n \varphi(w) &= \frac{\sum_{i=1}^n \varphi(z_i) \omega\left(\frac{w-w_i}{c_n}\right)}{\sum_{i=1}^n \omega\left(\frac{w-w_i}{c_n}\right)} \\ &= \sum_{i=1}^n a_i(\varphi) \varepsilon_i \end{aligned}$$

where

$$a_i(\varphi) = \varphi(z_i) \text{ and } \varepsilon_i = \left[ \frac{\omega\left(\frac{w-w_i}{c_n}\right)}{\sum_{i=1}^n \omega\left(\frac{w-w_i}{c_n}\right)} \right].$$

In this case, the rate of convergence of  $\hat{K}_n$  is nonparametric, see Subsection 4.1.

### 2.5.2. Estimation of the adjoint of a conditional expectation operator

Consider a conditional expectation operator as described in Example 2.3. Let  $K : L_Z^2 \rightarrow L_W^2$  be such that  $(K\varphi)(w) = E[\varphi(Z) | W = w]$  and its adjoint is  $K^* : L_W^2 \rightarrow L_Z^2$  with  $(K^*\psi)(z) = E[\psi(W) | Z = z]$ . Let  $\hat{f}_{Z,W}$ ,  $\hat{f}_Z(z)$ , and  $\hat{f}_W(w)$  be nonparametric estimators of  $f_{Z,W}$ ,  $f_Z(z)$ , and  $f_W(w)$  obtained either by kernel or sieves estimators. Assume that  $K$  and  $K^*$  are estimated by replacing the unknown pdfs by their estimators, that is:

$$\begin{aligned}\hat{K}_n\varphi(w) &= \int \frac{\hat{f}_{Z,W}(z,w)}{\hat{f}_Z(z)}\varphi(z) dz, \\ (\widehat{K^*})_n\psi(z) &= \int \frac{\hat{f}_{Z,W}(z,w)}{\hat{f}_W(w)}\psi(w) dw.\end{aligned}$$

Then we have  $(\widehat{K^*})_n \neq (\hat{K}_n)^*$  for  $\mathcal{H} = L_Z^2$  and  $\mathcal{E} = L_W^2$ . Indeed, we do not have

$$\left\langle \hat{K}_n\varphi, \psi \right\rangle_{\mathcal{E}} = \left\langle \varphi, (\widehat{K^*})_n\psi \right\rangle_{\mathcal{H}}. \quad (2.16)$$

There are two solutions to this problem. First, we choose as space of references  $\mathcal{H}_n = L^2(\mathbb{R}^q, \hat{f}_Z)$  and  $\mathcal{E}_n = L^2(\mathbb{R}^r, \hat{f}_W)$ . In which case,  $(\widehat{K^*})_n = (\hat{K}_n)^*$  for  $\mathcal{H}_n$  and  $\mathcal{E}_n$  because

$$\left\langle \hat{K}_n\varphi, \psi \right\rangle_{\mathcal{E}_n} = \left\langle \varphi, (\widehat{K^*})_n\psi \right\rangle_{\mathcal{H}_n}. \quad (2.17)$$

The new spaces  $\mathcal{H}_n$  and  $\mathcal{E}_n$  depend on the sample size and on the estimation procedure. Another approach consists in defining  $\mathcal{H} = L^2(\mathbb{R}^q, \pi)$  and  $\mathcal{E} = L^2(\mathbb{R}^r, \rho)$  where  $\pi$  and  $\rho$  are known and satisfy: There exist  $c, c' > 0$  such that  $f_Z(z) \leq c\pi(z)$  and  $f_W(w) \geq c'\rho(w)$ . Then

$$\begin{aligned}K^*\psi(z) &= \int \frac{f_{Z,W}(z,w)}{f_W(w)} \frac{\rho(w)}{\pi(z)} \psi(w) dw \\ &\neq E[\psi(W) | Z = z].\end{aligned}$$

In that case,  $(\widehat{K^*})_n = (\hat{K}_n)^*$  for  $\mathcal{H}$  and  $\mathcal{E}$  but the choice of  $\pi$  and  $\rho$  require some knowledge on the support and the tails of the distributions of  $W$  and  $Z$ .

An alternative solution to estimating  $K$  and  $K^*$  by kernel is to estimate the spectrum of  $K$  and to apply Mercer's formula. Let  $\mathcal{H} = L_Z^2$  and  $\mathcal{E} = L_W^2$ . The singular system  $\{\lambda_j, \phi_j, \psi_j\}$  of  $K$  satisfies

$$\lambda_j = \sup_{\phi_j, \psi_j} E[\phi_j(Z)\psi_j(W)], \quad j = 1, 2, \dots \quad (2.18)$$

subject to  $\|\phi_j\|_{\mathcal{H}} = 1, \langle \phi_j, \phi_l \rangle_{\mathcal{H}} = 0, l = 1, 2, \dots, j-1, \|\psi_j\|_{\mathcal{E}} = 1, \langle \psi_j, \psi_l \rangle_{\mathcal{E}} = 0, l = 1, 2, \dots, j-1$ . Assume the econometrician observes a sample  $\{w_i, z_i : i = 1, \dots, n\}$ . To

estimate  $\{\lambda_j, \phi_j, \psi_j\}$ , one can either estimate (2.18) by replacing the expectation by the sample mean or by replacing the joint pdf by a nonparametric estimator.

The first approach was adopted by Darolles, Florens, and Renault (1998). Let

$$\begin{aligned}\mathcal{H}_n &= \left\{ \varphi : \mathbb{R}^q \rightarrow \mathbb{R}, \int \varphi(z)^2 d\widehat{F}_Z(z) < \infty \right\}, \\ \mathcal{E}_n &= \left\{ \psi : \mathbb{R}^r \rightarrow \mathbb{R}, \int \psi(w)^2 d\widehat{F}_W(w) < \infty \right\}\end{aligned}$$

where  $\widehat{F}_Z$  and  $\widehat{F}_W$  are the empirical distributions of  $Z$  and  $W$ . That is  $\|\varphi\|_{\mathcal{H}_n}^2 = \frac{1}{n} \sum_{i=1}^n \varphi(z_i)^2$  and  $\|\psi\|_{\mathcal{E}_n}^2 = \frac{1}{n} \sum_{i=1}^n \psi(w_i)^2$ . Darolles, Florens, and Renault (1998) propose to estimate  $\{\lambda_j, \phi_j, \psi_j\}$  by solving

$$\hat{\lambda}_j = \sup_{\hat{\phi}_j, \hat{\psi}_j} \frac{1}{n} \sum_{i=1}^n \left[ \hat{\phi}_j(z_i) \hat{\psi}_j(w_i) \right], \quad j = 1, 2, \dots \quad (2.19)$$

subject to  $\|\hat{\phi}_j\|_{\mathcal{H}_n} = 1, \langle \hat{\phi}_j, \hat{\phi}_l \rangle_{\mathcal{H}_n} = 0, l = 1, 2, \dots, j-1, \|\hat{\psi}_j\|_{\mathcal{E}_n} = 1, \langle \hat{\psi}_j, \hat{\psi}_l \rangle_{\mathcal{E}_n} = 0, l = 1, 2, \dots, j-1$  where  $\hat{\phi}_j$  and  $\hat{\psi}_j$  are elements of increasing dimensional spaces

$$\begin{aligned}\hat{\phi}_j(z) &= \sum_{j=1}^J \alpha_j a_j(z), \\ \hat{\psi}_j(w) &= \sum_{j=1}^J \beta_j b_j(w)\end{aligned}$$

for some basis  $\{a_j\}$  and  $\{b_j\}$ . By Mercer's formula (2.13),  $K$  can be estimated by

$$\begin{aligned}\widehat{K}_n \varphi(w) &= \sum \hat{\lambda}_j \left( \int \hat{\phi}_j(z) \varphi(z) d\widehat{F}_Z \right) \hat{\psi}_j(w) \\ (\widehat{K^*})_n \psi(z) &= \sum \hat{\lambda}_j \left( \int \hat{\psi}_j(w) \psi(w) d\widehat{F}_W \right) \hat{\phi}_j(z).\end{aligned}$$

Hence  $(\widehat{K^*})_n = (\widehat{K}_n)^*$  for  $\mathcal{H}_n$  and  $\mathcal{E}_n$ .

The second approach consists in replacing  $f_{Z,W}$  by a nonparametric estimator  $\hat{f}_{Z,W}$ . Darolles, Florens, and Gourieroux (2004) use a kernel estimator, whereas Chen, Hansen and Scheinkman (1998) use B-spline wavelets. Let  $\mathcal{H}_n = L^2(\mathbb{R}^q, \hat{f}_Z)$  and  $\mathcal{E}_n = L^2(\mathbb{R}^r, \hat{f}_W)$  where  $\hat{f}_Z$  and  $\hat{f}_W$  are the marginals of  $\hat{f}_{Z,W}$ . (2.18) can be replaced

$$\hat{\lambda}_j = \sup_{\phi_j, \psi_j} \int \phi_j(z) \psi_j(w) \hat{f}_{Z,W}(z, w) dz dw, \quad j = 1, 2, \dots \quad (2.20)$$



subject to  $\|\phi_j\|_{\mathcal{H}_n} = 1, \langle \phi_j, \phi_l \rangle_{\mathcal{H}_n} = 0, l = 1, 2, \dots, j-1, \|\psi_j\|_{\mathcal{E}_n} = 1, \langle \psi_j, \psi_l \rangle_{\mathcal{E}_n} = 0, l = 1, 2, \dots, j-1$ . Denote  $\{\hat{\lambda}_j, \hat{\phi}_j, \hat{\psi}_j\}$  the resulting estimators of  $\{\lambda_j, \phi_j, \psi_j\}$ . By Mercer's formula,  $K$  can be approached by

$$\begin{aligned}\hat{K}_n \varphi(w) &= \sum \hat{\lambda}_j \left( \int \hat{\phi}_j(z) \varphi(z) \hat{f}_Z(z) dz \right) \hat{\psi}_j(w) \\ (\widehat{K^*})_n \psi(z) &= \sum \hat{\lambda}_j \left( \int \hat{\psi}_j(w) \psi(w) \hat{f}_W(w) dw \right) \hat{\phi}_j(z).\end{aligned}$$

Hence  $(\widehat{K^*})_n = (\hat{K}_n)^*$  for  $\mathcal{H}_n$  and  $\mathcal{E}_n$ . Note that in the three articles mentioned above,  $Z = X_{t+1}$  and  $W = X_t$  where  $\{X_t\}$  is a Markov process. These papers are mainly concerned with estimation. When the data are the discrete observations of a diffusion process, the nonparametric estimations of a single eigenvalue-eigenfunction pair and of the marginal distribution are enough to recover a nonparametric estimate of the diffusion coefficient. The techniques described here can also be used for testing the reversibility of the process  $\{X_t\}$ , see Darolles, Florens, and Gouriéroux (2004).

### 2.5.3. Computation of eigenvalues and eigenfunctions of finite dimensional operators

Here, we assume that we have some estimators of  $K$  and  $K^*$ , denoted  $\hat{K}_n$  and  $\hat{K}_n^*$  such that  $\hat{K}_n$  and  $\hat{K}_n^*$  have finite range and satisfy

$$\hat{K}_n \varphi = \sum_{l=1}^{L_n} a_l(\varphi) \varepsilon_l \quad (2.21)$$

$$\hat{K}_n^* \psi = \sum_{l=1}^{L_n} b_l(\psi) \eta_l \quad (2.22)$$

where  $\varepsilon_l \in \mathcal{E}, \eta_l \in \mathcal{H}$ ,  $a_l(\varphi)$  is linear in  $\varphi$  and  $b_l(\psi)$  is linear in  $\psi$ . Moreover the  $\{\varepsilon_l\}$  and  $\{\eta_l\}$  are assumed to be linearly independent. It follows that

$$\begin{aligned}\hat{K}_n^* \hat{K}_n \varphi &= \sum_{l=1}^{L_n} b_l \left( \sum_{l'=1}^{L_n} a_{l'}(\varphi) \varepsilon_{l'} \right) \eta_l \\ &= \sum_{l, l'=1}^{L_n} a_{l'}(\varphi) b_l(\varepsilon_{l'}) \eta_l.\end{aligned} \quad (2.23)$$

We calculate the eigenvalues and eigenfunctions of  $\hat{K}_n^* \hat{K}_n$  by solving

$$\hat{K}_n^* \hat{K}_n \phi = \lambda^2 \phi.$$

Hence  $\phi$  is necessarily of the form:  $\phi = \sum_l \beta_l \eta_l$ . Replacing in (2.23), we have

$$\lambda^2 \beta_l = \sum_{\nu, j=1}^{L_n} \beta_j a_\nu (\eta_j) b_l (\varepsilon_\nu). \quad (2.24)$$

Denote  $\underline{\hat{\beta}} = [\beta_1, \dots, \beta_{L_n}]$  the solution of (2.24). Solving (2.24) is equivalent to finding the  $L_n$  nonzero eigenvalues  $\hat{\lambda}_1^2, \dots, \hat{\lambda}_{L_n}^2$  and eigenvectors  $\underline{\hat{\beta}}^1, \dots, \underline{\hat{\beta}}^{L_n}$  of an  $L_n \times L_n$ -matrix  $C$  with principle element

$$c_{l,j} = \sum_{\nu=1}^{L_n} a_\nu (\eta_j) b_l (\varepsilon_\nu).$$

The eigenfunctions of  $\hat{K}_n^* \hat{K}_n$  are

$$\hat{\phi}_j = \sum_{l=1}^{L_n} \hat{\beta}_l^j \eta_l, \quad j = 1, \dots, L_n$$

associated with  $\hat{\lambda}_1^2, \dots, \hat{\lambda}_{L_n}^2$ .  $\{\hat{\phi}_j : j = 1, \dots, L_n\}$  need to be orthonormalized. The estimators of the singular values are  $\hat{\lambda}_j = \sqrt{\hat{\lambda}_j^2}$ .

#### 2.5.4. Estimation of noncompact operators

This chapter mainly focuses on compact operators, because compact operators can be approached by a sequence of finite dimensional operators and therefore can be easily estimated. However, it is possible to estimate a noncompact operator by an estimator, which is infinitely dimensional. A simple example is provided by the conditional expectation operator with common elements.

**Example 2.5 (continued).** This example is discussed in Hall and Horowitz (2004). Assume that the dimension of  $Z$  is  $p$ . The conditional expectation operator  $K$  can be estimated by a kernel estimator with kernel  $\omega$  and bandwidth  $c_n$

$$\left(\hat{K}\varphi\right)(x, w) = \frac{\sum_{i=1}^n \left[ \int \frac{1}{c_n^p} \varphi(x, z) \omega\left(\frac{z-z_i}{c_n}\right) dz \right] \omega\left(\frac{x-x_i}{c_n}\right) \omega\left(\frac{w-w_i}{c_n}\right)}{\sum_{i=1}^n \omega\left(\frac{x-x_i}{c_n}\right) \omega\left(\frac{w-w_i}{c_n}\right)}.$$

We can see that  $\hat{K}$  is an infinite dimensional operator because all functions  $\varphi(x)$  that depend only on  $x$  are in the range of  $\hat{K}$ .

### 3. Regularized solutions of integral equations of the first kind

Let  $\mathcal{H}$  and  $\mathcal{E}$  be two Hilbert spaces considered only over the real scalars for the sake of notational simplicity. Let  $K$  be a linear operator on  $\mathcal{D}(K) \subset \mathcal{H}$  into  $\mathcal{E}$ . This section

discusses the properties of operator equations (also called Fredholm equations) of the first kind

$$K\varphi = r \tag{3.1}$$

where  $K$  is typically an integral compact operator. Such an equation in  $\varphi$  is in general an ill-posed problem by opposition to a well-posed problem. Equation (3.1) is said to be well-posed if (i) (*existence*) a solution exists, (ii) (*uniqueness*) the solution is unique, and (iii) (*stability*) the solution is continuous in  $r$ , that is  $\varphi$  is stable with respect to small changes in  $r$ . Whenever one of these conditions is not satisfied, the problem is said to be ill-posed. The lack of stability is particularly problematic and needs to be addressed by a regularization scheme. Following Wahba (1973) and Nashed and Wahba (1974), we introduce generalized inverses of operators in Reproducing Kernel Hilbert Spaces (RKHS). Properties of RKHS will be studied more extensively in Section 6.

### 3.1. Ill-posed and well-posed problems

This introductory subsection gives an overview of the problems encountered when solving an equation  $K\varphi = r$  where  $K$  is a linear operator, not necessarily compact. A more detailed encounter can be found in Groetsch (1993). We start with a formal definition of a well-posed problem.

**Definition 3.1.** *Let  $K : \mathcal{H} \rightarrow \mathcal{E}$ . The equation*

$$K\varphi = r \tag{3.2}$$

*is called well-posed if  $K$  is bijective and the inverse operator  $K^{-1} : \mathcal{E} \rightarrow \mathcal{H}$  is continuous. Otherwise, the equation is called ill-posed.*

Note that  $K$  is injective means  $\mathcal{N}(K) = \{0\}$ , and  $K$  is surjective means  $\mathcal{R}(K) = \mathcal{E}$ . In this section, we will restrict ourselves to the case where  $K$  is a bounded (and therefore continuous) linear operator. By Banach theorem (Kress, 1999, page 266), if  $K : \mathcal{H} \rightarrow \mathcal{E}$  is a bounded linear operator,  $K$  bijective implies that  $K^{-1} : \mathcal{E} \rightarrow \mathcal{H}$  is bounded and therefore continuous. In this case,  $K\varphi = r$  is well-posed.

An example of a well-posed problem is given by

$$(I - C)\varphi = r$$

where  $C : \mathcal{H} \rightarrow \mathcal{H}$  is a compact operator and 1 is not an eigenvalue of  $C$ . This is an example of integral equation of the second kind that will be studied in Section 7.

We now turn our attention to ill-posed problems.

**Problem of uniqueness:**

If  $\mathcal{N}(K) \neq \{0\}$ , then to any solution of  $\varphi$  of (3.2), one can add an element  $\varphi_1$  of  $\mathcal{N}(K)$ , so that  $\varphi + \varphi_1$  is also a solution. A way to achieve uniqueness is to look for the solution with minimal norm.

**Problem of existence:**

A solution to (3.2) exists if and only if

$$r \in \mathcal{R}(K).$$

Since  $K$  is linear,  $\mathcal{R}(K)$  is a subspace of  $\mathcal{E}$ , however it generally does not exhaust  $\mathcal{E}$ . Therefore, a traditional solution of (3.2) exists only for a restricted class of functions  $r$ . If we are willing to broaden our notion of solution, we may enlarge the class of functions  $r$  for which a type of generalized solution exists to a dense subspace of functions of  $\mathcal{E}$ .

**Definition 3.2.** An element  $\tilde{\varphi} \in \mathcal{H}$  is said to be a least squares solution of (3.2) if:

$$\|K\tilde{\varphi} - r\| \leq \|Kf - r\|, \text{ for any } f \in \mathcal{H} \quad (3.3)$$

If the set  $S_r$  of all least squares solutions of (3.2) for a given  $r \in \mathcal{E}$  is not empty and admits an element  $\varphi$  of minimal norm, then  $\varphi$  is called a pseudosolution of (3.2).

The pseudosolution, when it exists, is denoted  $\varphi = K^\dagger r$  where  $K^\dagger$  is by definition the Moore Penrose generalized inverse of  $K$ . However, the pseudosolution does not necessarily exist. The pseudosolution exists if and only if  $Pr \in \mathcal{R}(K)$  where  $P$  is the projection operator on  $\overline{\mathcal{R}(K)}$ , the closure of the range of  $K$ . Note that  $Pr \in \mathcal{R}(K)$  if and only if

$$r = Pr + (1 - P)r \in \mathcal{R}(K) + \mathcal{R}(K)^\perp. \quad (3.4)$$

Therefore, a pseudosolution exists if and only if  $r$  lies in the dense subspace  $\mathcal{R}(K) + \mathcal{R}(K)^\perp$  of  $\mathcal{E}$ .

We distinguish two cases:

1.  $\mathcal{R}(K)$  is closed.

For any  $r \in \mathcal{E}$ ,  $\varphi = K^\dagger r$  exists and is continuous in  $r$ .

**Example.**  $(I - C)\varphi = r$  where  $K$  is compact and 1 is an eigenvalue of  $K$ . The problem is ill-posed because the solution is not unique but it is not severally ill-posed because the pseudosolution exists and is continuous.

2.  $\mathcal{R}(K)$  is not closed.

The pseudosolution exists if and only if  $r \in \mathcal{R}(K) + \mathcal{R}(K)^\perp$ . But here,  $\varphi = K^\dagger r$  is not continuous in  $r$ .

**Example.**  $K$  is a compact infinitely dimensional operator.

For the purpose of econometric applications, condition (3.4) will be easy to maintain since:

Either  $(K, r)$  denotes the true unknown population value, and then the assumption  $r \in \mathcal{R}(K)$  means that the structural econometric model is well-specified. Inverse problems with specification errors are beyond the scope of this chapter.

Or  $(K, r)$  denotes some estimators computed from a finite sample of size  $n$ . Then, insofar as the chosen estimation procedure is such that  $\mathcal{R}(K)$  is closed (for instance because it is finite dimensional as in Subsection 2.5.1), we have  $\mathcal{R}(K) + \mathcal{R}(K)^\perp = \mathcal{E}$ .

The continuity assumption of  $K$  will come in general with the compactness assumption for population values and, for sample counterparts, with the finite dimensional property. Moreover, the true unknown value  $K_0$  of  $K$  will be endowed with the *identification assumption*:

$$\mathcal{N}(K_0) = \{0\} \tag{3.5}$$

and the *well-specification assumption*:

$$r_0 \in \mathcal{R}(K_0). \tag{3.6}$$

(3.5) and (3.6) ensure the existence of a unique true unknown value  $\varphi_0$  of  $\varphi$  defined as the (pseudo) solution of the operator equation  $K_0\varphi_0 = r_0$ . Moreover, this solution is not going to depend on the choice of topologies on the two spaces  $\mathcal{H}$  and  $\mathcal{E}$ .

It turns out that a compact operator  $K$  with infinite-dimensional range is a prototype of an operator for which  $\mathcal{R}(K)$  is not closed. Therefore, as soon as one tries to generalize structural econometric estimation from a parametric setting ( $K$  finite dimensional) to a non-parametric one, which can be seen as a limit of finite dimensional problems ( $K$  compact), one is faced with an ill-posed inverse problem. This is a serious issue for the purpose of consistent estimation, since in general one does not know the true value  $r_0$  of  $r$  but only a consistent estimator  $\hat{r}_n$ . Therefore, there is no hope to get a consistent estimator  $\hat{\varphi}_n$  of  $\varphi$  by solving  $K\hat{\varphi}_n = \hat{r}_n$  that is  $\hat{\varphi}_n = K^\dagger\hat{r}_n$ , when  $K^\dagger$  is not continuous. In general, the issue to address will be even more involved since  $K^\dagger$  and  $K$  must also be estimated.

Let us finally recall a useful characterization of the Moore-Penrose generalized inverse of  $K$ .

**Proposition 3.3.** *Under (3.4),  $K^\dagger r$  is the unique solution of minimal norm of the equation  $K^*K\varphi = K^*r$ .*

In other words, the pseudosolution  $\varphi$  of (3.2) can be written in two ways:

$$\varphi = K^\dagger r = (K^*K)^\dagger K^*r$$

For  $r \in \mathcal{R}(K)$  (well-specification assumption in the case of true unknown values),  $K^*r \in \mathcal{R}(K^*K)$  and then  $(K^*K)^{-1}K^*r$  is well defined. The pseudosolution can then be represented from the singular value decomposition of  $K$  as

$$\varphi = K^\dagger r = (K^*K)^{-1}K^*r = \sum_{j=1}^{\infty} \frac{\langle r, \psi_j \rangle}{\lambda_j} \phi_j \tag{3.7}$$

It is worth noticing that the spectral decomposition (3.7) is also valid for any  $r \in \mathcal{R}(K) + \mathcal{R}(K)^\perp$  to represent the pseudosolution  $\varphi = K^\dagger r = (K^*K)^\dagger K^*r$  since  $r \in \mathcal{R}(K)^\perp$  is equivalent to  $K^\dagger r = 0$ .

Formula (3.7) clearly demonstrates the ill-posed nature of the equation  $K\varphi = r$ . If we perturb the right-hand side  $r$  by  $r^\delta = r + \delta\psi_j$ , we obtain the solution  $\varphi^\delta = \varphi + \delta\phi_j/\lambda_j$ . Hence, the ratio  $\|\varphi^\delta - \varphi\| / \|r^\delta - r\| = 1/\lambda_j$  can be made arbitrarily large due to the fact that the singular values tend to zero. Since the influence of estimation errors in  $r$  is controlled by the rate of this convergence, Kress (1999, p. 280) says that the equation is “mildly ill-posed” if the singular values decay slowly to zero and that it is “severely ill-posed” if they decay rapidly. Actually, the critical property is the relative decay rate of the sequence  $\langle r, \psi_j \rangle$  with respect to the decay of the sequence  $\lambda_j$ . To see this, note that the solution  $\varphi$  has to be determined from its Fourier coefficients by solving the equations

$$\lambda_j \langle \varphi, \phi_j \rangle = \langle r, \psi_j \rangle, \text{ for all } j.$$

Then, we may expect high instability of the solution  $\varphi$  if  $\lambda_j$  goes to zero faster than  $\langle \varphi, \phi_j \rangle$ . The properties of regularity spaces introduced below precisely document this intuition.

### 3.2. Regularity spaces

As stressed by Nashed and Wahba (1974), an ill-posed problem relative to  $\mathcal{H}$  and  $\mathcal{E}$  may be recast as a well-posed problem relative to new spaces  $\mathcal{H}' \subset \mathcal{H}$  and  $\mathcal{E}' \subset \mathcal{E}$ , with topologies on  $\mathcal{H}'$  and  $\mathcal{E}'$ , which are different from the topologies on  $\mathcal{H}$  and  $\mathcal{E}$  respectively. While Nashed and Wahba (1974) generally build these Hilbert spaces  $\mathcal{H}'$  and  $\mathcal{E}'$  as RKHS associated with an arbitrary self-adjoint Hilbert-Schmidt operator, we focus here on the RKHS associated with  $(K^*K)^\beta$ , for some positive  $\beta$ . More precisely, assuming that  $K$  is Hilbert-Schmidt and denoting  $(\lambda_j, \phi_j, \psi_j)$  its singular system, we define the self-adjoint operator  $(K^*K)^\beta$  by:

$$(K^*K)^\beta \varphi = \sum_{j=1}^{\infty} \lambda_j^{2\beta} \langle \varphi, \phi_j \rangle \phi_j.$$

**Definition 3.4.** *The  $\beta$ -regularity space of the compact operator  $K$  is defined for all  $\beta > 0$ , as the RKHS associated with  $(K^*K)^\beta$ . That is, the space:*

$$\Phi_\beta = \left\{ \varphi \in \mathcal{N}(K)^\perp \text{ such that } \sum_{j=1}^{\infty} \frac{\langle \varphi, \phi_j \rangle^2}{\lambda_j^{2\beta}} < \infty \right\} \quad (3.8)$$

where a Hilbert structure is being defined through the inner product

$$\langle f, g \rangle_\beta = \sum_{j=1}^{\infty} \frac{\langle f, \phi_j \rangle \langle g, \phi_j \rangle}{\lambda_j^{2\beta}}$$

for  $f$  and  $g \in \Phi_\beta$ .

Note however that the construction of RKHS considered here is slightly more general than the one put forward in Nashed and Wahba (1974) since we start from elements of a general Hilbert space, not limited to be a  $L^2$  space of functions defined on some interval of the real line. This latter example will be made explicit in Section 6. Moreover, the focus of our interest here will only be the regularity spaces associated with the true unknown value  $K_0$  of the operator  $K$ . Then, the identification assumption will ensure that all the regularity spaces are dense in  $\mathcal{H}$  :

**Proposition 3.5.** *Under the identification assumption  $\mathcal{N}(K) = \{0\}$ , the sequence of eigenfunctions  $\{\phi_j\}$  associated with the non-zero singular values  $\lambda_j$  defines a Hilbert basis of  $\mathcal{H}$ . In particular, all the regularity spaces  $\Phi_\beta$ ,  $\beta > 0$ , contain the vectorial space spanned by the  $\{\phi_j\}$  and, as such, are dense in  $\mathcal{H}$ .*

Proposition 3.5. is a direct consequence of the singular value decomposition (2.9). Generally speaking, when  $\beta$  increases,  $\Phi_\beta$ ,  $\beta > 0$ , is a decreasing family of subspaces of  $\mathcal{H}$ . Hence,  $\beta$  may actually be interpreted as the regularity level of the functions  $\varphi$ , as illustrated by the following result.

**Proposition 3.6.** *Under the identification assumption ( $\mathcal{N}(K) = \{0\}$ ), for any  $\beta > 0$ ,*

$$\Phi_\beta = \mathcal{R} \left[ (K^*K)^{\frac{\beta}{2}} \right].$$

In particular,  $\Phi_1 = \mathcal{R}(K^*)$ .

**Proof.** By definition, the elements of the range of  $(K^*K)^{\frac{\beta}{2}}$  can be written  $f = \sum_{j=1}^{\infty} \lambda_j^\beta \langle \varphi, \phi_j \rangle \phi_j$  for some  $\varphi \in \mathcal{H}$ . Note that this decomposition also describes the range of  $K^*$  for  $\beta = 1$ . Then:

$$\|f\|_\beta^2 = \sum_{j=1}^{\infty} \frac{\langle \varphi, \phi_j \rangle^2}{\lambda_j^{2\beta}} \lambda_j^{2\beta} = \sum_{j=1}^{\infty} \langle \varphi, \phi_j \rangle^2 = \|\varphi\|^2 < \infty.$$

Hence  $\mathcal{R} \left[ (K^*K)^{\frac{\beta}{2}} \right] \subset \Phi_\beta$ .

Conversely, for any  $\varphi \in \Phi_\beta$ , one can define:

$$f = \sum_{j=1}^{\infty} \frac{\langle \varphi, \phi_j \rangle}{\lambda_j^\beta} \phi_j$$

and then  $(K^*K)^{\frac{\beta}{2}} f = \sum_{j=1}^{\infty} \langle \varphi, \phi_j \rangle \phi_j = \varphi$  since  $\mathcal{N}(K) = \{0\}$ . Hence,  $\Phi_\beta \subset \mathcal{R} \left[ (K^*K)^{\frac{\beta}{2}} \right]$ .

Since we mainly consider operators,  $K$ , which are integral operators with continuous kernels, applying the operator  $(K^*K)^{\frac{\beta}{2}}$  has a smoothing effect, which is stronger for larger values of  $\beta$ . This is the reason why the condition  $\varphi \in \Phi_\beta$  qualifies the level,  $\beta$

of regularity or smoothness of  $\varphi$ . The associated smoothness properties are studied in further details in Loubes and Vanhems (2003). The space  $\Phi_1$  of functions is also put forward in Schaumburg (2004) when  $K$  denotes the conditional expectation operator for a continuous time Markov process  $X_t$  with Levy type generator sampled in discrete time. He shows that whenever a transformation  $\varphi(X_t)$  of the diffusion process is considered with  $\varphi \in \Phi_1$ , the conditional expectation operator  $E[\varphi(X_{t+h}) | X_t]$  admits a convergent power series expansion as the exponential of the infinitesimal generator.

The regularity spaces  $\Phi_\beta$  are of interest here as Hilbert spaces (included in  $\mathcal{H}$  but endowed with another scalar product) where our operator equation (3.2) is going to become well-posed. More precisely, let us also consider the family of regularity spaces  $\Psi_\beta$  associated with the compact operator  $K^*$ :

$$\Psi_\beta = \left\{ \psi \in \mathcal{N}(K^*)^\perp \text{ such that } \sum_{j=1}^{\infty} \frac{\langle \psi, \psi_j \rangle^2}{\lambda_j^{2\beta}} < \infty \right\}$$

$\Psi_\beta$  is a Hilbert space endowed with the inner product:

**Definition 3.7.**  $\langle F, G \rangle_\beta = \sum_{j=1}^{\infty} \frac{\langle F, \psi_j \rangle \langle G, \psi_j \rangle}{\lambda_j^{2\beta}}$  for  $F$  and  $G \in \Psi_\beta$ .

Note that the spaces  $\Psi_\beta$  are not in general dense in  $\mathcal{E}$  since  $\mathcal{N}(K^*) \neq \{0\}$ . But they describe well the range of  $K$  when  $K$  is restricted to some regularity space:

**Proposition 3.8.** *Under the identification assumption  $\mathcal{N}(K) = \{0\}$ ,  $K(\Phi_\beta) = \Psi_{\beta+1}$  for all positive  $\beta$ . In particular,  $\Psi_1 = \mathcal{R}(K)$ .*

**Proof.** We know from Proposition 3.6 that when  $\varphi \in \Phi_\beta$ , it can be written:

$\varphi = \sum_{j=1}^{\infty} \lambda_j^\beta \langle f, \phi_j \rangle \phi_j$  for some  $f \in \mathcal{H}$ . Then,  $K\varphi = \sum_{j=1}^{\infty} \lambda_j^{\beta+1} \langle f, \phi_j \rangle \psi_j \in \Psi_{\beta+1}$ . Hence  $K(\Phi_\beta) \subset \Psi_{\beta+1}$ .

Conversely, since according to a singular value decomposition like (2.9), the sequence  $\{\psi_j\}$  defines a basis of  $\mathcal{N}(K^*)^\perp$ , any element of  $\Psi_{\beta+1}$  can be written as

$$\psi = \sum_{j=1}^{\infty} \langle \psi, \psi_j \rangle \psi_j \text{ with } \sum_{j=1}^{\infty} \frac{\langle \psi, \psi_j \rangle^2}{\lambda_j^{2\beta+2}} < \infty.$$

Let us then define  $\varphi = \sum_{j=1}^{\infty} (1/\lambda_j) \langle \psi, \psi_j \rangle \phi_j$ . We have

$$\sum_{j=1}^{\infty} \frac{\langle \varphi, \phi_j \rangle^2}{\lambda_j^{2\beta}} = \sum_{j=1}^{\infty} \frac{\langle \psi, \psi_j \rangle^2}{\lambda_j^{2\beta+2}} < \infty$$

and thus  $\varphi \in \Phi_\beta$ . Moreover,  $K\varphi = \sum_{j=1}^{\infty} \langle \psi, \psi_j \rangle \psi_j = \psi$ . This proves that  $\Psi_{\beta+1} \subset K(\Phi_\beta)$ . ■



Therefore, when viewed as an operator from  $\Phi_\beta$  into  $\Psi_{\beta+1}$ ,  $K$  has a closed range defined by the space  $\Psi_{\beta+1}$  itself. It follows that the ill-posed problem

$$\begin{aligned} K &: \mathcal{H} \rightarrow \mathcal{E} \\ K\varphi &= r \end{aligned}$$

may be viewed as well-posed relative to the subspaces  $\Phi_\beta$  into  $\Psi_{\beta+1}$  and their associated norms. This means that

(i) First, we think about the pseudosolution  $\varphi = K^\dagger r$  as a function of  $r$  evolving in  $\Psi_{\beta+1}$ , for some positive  $\beta$ .

(ii) Second, continuity of  $\varphi = K^\dagger r$  with respect to  $r$  must be understood with respect to the norms  $\|r\|_{\beta+1} = \langle r, r \rangle_{\beta+1}^{1/2}$  and  $\|\varphi\|_\beta = \langle \varphi, \varphi \rangle_\beta^{1/2}$

To get the intuition of this result, it is worth noticing that these new topologies define another adjoint operator  $K_\beta^*$  of  $K$  characterized by:

$$\langle K\varphi, \psi \rangle_{\beta+1} = \langle \varphi, K_\beta^* \psi \rangle_\beta,$$

and thus:

$$K_\beta^* \psi = \sum_{j=1}^{\infty} (1/\lambda_j) \langle \psi, \psi_j \rangle \phi_j.$$

In particular,  $K_\beta^* \psi_j = (1/\lambda_j) \phi_j$ . In other words, all the eigenvalues of  $K_\beta^* K$  and  $KK_\beta^*$  are now equal to one and the pseudosolution is defined as:

$$\varphi = K_\beta^\dagger r = K_\beta^* r = \sum_{j=1}^{\infty} \frac{\langle r, \psi_j \rangle}{\lambda_j} \phi_j.$$

The pseudosolution depends continuously on  $r$  because  $K_\beta^\dagger = K_\beta^*$  is a bounded operator for the chosen topologies; it actually has a unit norm.

For the purpose of econometric estimation, we may be ready to assume that the true unknown value  $\varphi_0$  belongs to some regularity space  $\Phi_\beta$ . This just amounts to an additional smoothness condition about our structural functional parameter of interest. Then, we are going to take advantage of this regularity assumption through the rate of convergence of some regularization bias as characterized in the next subsection.

Note finally that assuming  $\varphi_0 \in \Phi_\beta$ , that is  $r_0 \in \Psi_{\beta+1}$  for some positive  $\beta$ , is nothing but a small reinforcement of the common criterion of existence of a solution, known as the Picard's theorem (see e.g. Kress, 1999, page 279), which states that  $r_0 \in \Psi_1 = \mathcal{R}(K)$ . The spaces  $\Phi_\beta$  and  $\Psi_\beta$  are strongly related to the concept of Hilbert scales, see Natterer (1984), Engl, Hanke, and Neubauer (1996), and Tautenhahn (1996).

### 3.3. Regularization schemes

As pointed out in Subsection 3.1, the ill-posedness of an equation of the first kind with a compact operator stems from the behavior of the sequence of singular values, which

converge to zero. This suggests trying to regularize the equation by damping the explosive asymptotic effect of the inversion of singular values. This may be done in at least two ways:

A first estimation strategy consists in taking advantage of the well-posedness of the problem when reconsidered within regularity spaces. Typically, a sieve kind of approach may be designed, under the maintained assumption that the true unknown value  $r_0 \in \Psi_{\beta+1}$  for some positive  $\beta$ , in such a way that the estimator  $\hat{r}_n$  evolves when  $n$  increases, in an increasing sequence of finite dimensional subspaces of  $\Psi_{\beta+1}$ . Note however that when the operator  $K$  is unknown, the constraint  $\hat{r}_n \in \mathcal{N}(K^*)^\perp$  may be difficult to implement. Hence, we will not pursue this route any further.

The approach adopted in this chapter follows the general regularization framework of Kress (1999, Theorem 15.21). It consists in replacing a sequence  $\{1/\mu_j\}$  of explosive inverse singular values by a sequence  $\{q(\alpha, \mu_j)/\mu_j\}$  where the *damping function*  $q(\alpha, \mu)$  is chosen such that:

- (i)  $\{q(\alpha, \mu)/\mu\}$  remains bounded when  $\mu$  goes to zero (damping effect),
- (ii) for any given  $\mu$  :  $\lim_{\alpha \rightarrow 0} q(\alpha, \mu) = 1$  (asymptotic unbiasedness).

Since our inverse problem of interest can be addressed in two different ways:

$$\varphi = K^\dagger r = (K^*K)^\dagger K^*r,$$

the regularization scheme can be applied either to  $K^\dagger$  ( $\mu_j = \lambda_j$ ) or to  $(K^*K)^\dagger$  ( $\mu_j = \lambda_j^2$ ). The latter approach is better suited for our purpose since estimation errors will be considered below at the level of  $(K^*K)$  and  $K^*r$  respectively. We maintain in this subsection the identification assumption  $\mathcal{N}(K) = \{0\}$ . We then define:

**Definition 3.9.** A regularized version  $\varphi_\alpha = A_\alpha K^*r$  of the pseudosolution  $\varphi = (K^*K)^\dagger K^*r$  is defined as:

$$\begin{aligned} \varphi_\alpha &= \sum_{j=1}^{\infty} \frac{1}{\lambda_j^2} q(\alpha, \lambda_j^2) \langle K^*r, \phi_j \rangle \phi_j = \sum_{j=1}^{\infty} \frac{1}{\lambda_j} q(\alpha, \lambda_j^2) \langle r, \psi_j \rangle \phi_j \\ &= \sum_{j=1}^{\infty} q(\alpha, \lambda_j^2) \langle \varphi, \phi_j \rangle \phi_j \end{aligned} \quad (3.9)$$

where the real-valued function,  $q$ , is such that

$$\begin{aligned} |q(\alpha, \mu)| &\leq d(\alpha) \mu \\ \lim_{\alpha \rightarrow 0} q(\alpha, \mu) &= 1. \end{aligned} \quad (3.10)$$

Note that (3.9) leaves unconstrained the values of the operator  $A_\alpha$  on the space  $\mathcal{R}(K^*)^\perp = \mathcal{N}(K)$ . However, since  $\mathcal{N}(K) = \{0\}$ ,  $A_\alpha$  is uniquely defined as

$$A_\alpha \varphi = \sum_{j=1}^{\infty} \frac{1}{\lambda_j^2} q(\alpha, \lambda_j^2) \langle \varphi, \phi_j \rangle \phi_j \quad (3.11)$$

for all  $\varphi \in \mathcal{H}$ . Note that as  $q$  is real,  $A_\alpha$  is self-adjoint. Then by (3.10),  $A_\alpha$  is a bounded operator from  $\mathcal{H}$  into  $\mathcal{H}$  with

$$\|A_\alpha\| \leq d(\alpha). \quad (3.12)$$

In the following, we will always normalize the exponent of the regularization parameter  $\alpha$  such that  $\alpha d(\alpha)$  has a positive finite limit  $c$  when  $\alpha$  goes to zero. By construction,  $A_\alpha K^* K \varphi \rightarrow \varphi$  as  $\alpha$  goes to zero. When a genuine solution exists ( $r = K\varphi$ ), the regularization induces a bias:

$$\varphi - \varphi_\alpha = \sum_{j=1}^{\infty} [1 - q(\alpha, \lambda_j^2)] \langle r, \psi_j \rangle (\phi_j / \lambda_j) = \sum_{j=1}^{\infty} [1 - q(\alpha, \lambda_j^2)] \langle \varphi, \phi_j \rangle \phi_j \quad (3.13)$$

The squared regularization bias is

$$\|\varphi - \varphi_\alpha\|^2 = \sum_{j=1}^{\infty} b^2(\alpha, \lambda_j^2) \langle \varphi, \phi_j \rangle^2, \quad (3.14)$$

where  $b(\alpha, \lambda_j^2) = 1 - q(\alpha, \lambda_j^2)$  is the *bias function* characterizing the weight of the Fourier coefficient  $\langle \varphi, \phi_j \rangle$ . Below, we show that the most common regularization schemes fulfill the above conditions. We characterize these schemes through the definitions of the *damping weights*  $q(\alpha, \mu)$  or equivalently, of the *bias function*  $b(\alpha, \mu)$ .

**Example (Spectral cut-off).**

The spectral cut-off regularized solution is

$$\varphi_\alpha = \sum_{\lambda_j^2 \geq \alpha/c} \frac{1}{\lambda_j} \langle r, \psi_j \rangle \phi_j.$$

The explosive influence of the factor  $(1/\mu)$  is filtered out by imposing  $q(\alpha, \mu) = 0$  for small  $\mu$ , that is  $|\mu| < \alpha/c$ .  $\alpha$  is a positive regularization parameter such that no bias is introduced when  $|\mu|$  exceeds the threshold  $\alpha/c$ :

$$q(\alpha, \mu) = I\{|\mu| \geq \alpha/c\} = \begin{cases} 1 & \text{if } |\mu| \geq \alpha/c, \\ 0 & \text{otherwise.} \end{cases}$$

For any given scaling factor  $c$ , the two conditions of Definition 3.9. are then satisfied (with  $d(\alpha) = c/\alpha$ ) and we get a bias function  $b(\alpha, \lambda^2)$  which is maximized (equal to 1) when  $\lambda^2 < \alpha/c$  and minimized (equal to 0) when  $\lambda^2 \geq \alpha/c$ .

**Example (Landweber-Fridman).**

Landweber-Fridman regularization is characterized by

$$\begin{aligned} A_\alpha &= c \sum_{l=0}^{1/\alpha-1} (I - cK^*K)^l K^*, \\ \varphi_\alpha &= c \sum_{l=0}^{1/\alpha-1} (I - cK^*K)^l K^*r. \end{aligned}$$

The basic idea is similar to spectral cut-off but with a smooth bias function. Of course, one way to make the bias function continuous while meeting the conditions  $b(\alpha, 0) = 1$  and  $b(\alpha, \lambda^2) = 0$  for  $\lambda^2 > \alpha/c$  would be to consider a piecewise linear bias function with  $b(\alpha, \lambda^2) = 1 - (c/\alpha)\lambda^2$  for  $\lambda^2 \leq \alpha/c$ . Landweber-Fridman regularization makes it smooth, while keeping the same level and the same slope at  $\lambda^2 = 0$  and zero bias for large  $\lambda^2$ ,  $b(\alpha, \lambda^2) = (1 - c\lambda^2)^{1/\alpha}$  for  $\lambda^2 \leq 1/c$  and zero otherwise, that is

$$q(\alpha, \mu) = \begin{cases} 1 & \text{if } |\mu| > 1/c, \\ 1 - (1 - c\mu)^{1/\alpha} & \text{for } |\mu| \leq 1/c. \end{cases}$$

For any given scaling factor  $c$ , the two conditions of Definition 3.9 are then satisfied with again  $d(\alpha) = c/\alpha$ .

**Example (Tikhonov regularization).**

Here, we have

$$\begin{aligned} A_\alpha &= \left( \frac{\alpha}{c} I + K^* K \right)^{-1}, \\ \varphi_\alpha &= \left( \frac{\alpha}{c} I + K^* K \right)^{-1} K^* r \\ &= \sum_{j=1}^{\infty} \frac{\lambda_j}{\lambda_j^2 + \alpha/c} \langle r, \psi_j \rangle \phi_j \end{aligned}$$

where  $c$  is some scaling factor. In contrast to the two previous examples, the bias function is never zero but decreases toward zero at a hyperbolic rate (when  $\lambda^2$  becomes infinitely large), while still starting from 1 for  $\lambda^2 = 0$ :

$$b(\alpha, \lambda^2) = \frac{(\alpha/c)}{(\alpha/c) + \lambda^2}.$$

that is:

$$q(\alpha, \lambda^2) = \frac{\lambda^2}{(\alpha/c) + \lambda^2}$$

For any given scaling factor  $c$ , the two conditions of Definition 3.9 are again satisfied with  $d(\alpha) = c/\alpha$ .

We are going to show now that the regularity spaces  $\Phi_\beta$  introduced in the previous subsection are well-suited for controlling the regularization bias. The basic idea is a straightforward consequence of (3.15):

$$\|\varphi - \varphi_\alpha\|^2 \leq [\sup_j b^2(\alpha, \lambda_j^2) \lambda_j^{2\beta}] \|\varphi\|_\beta^2 \tag{3.15}$$

Therefore, the rate of convergence (when the regularization parameter  $\alpha$  goes to zero) of the regularization bias will be controlled, for  $\varphi \in \Phi_\beta$ , by the rate of convergence of

$$M_\beta(\alpha) = \sup_j b^2(\alpha, \lambda_j^2) \lambda_j^{2\beta}$$

The following definition is useful to characterize the regularization schemes.

**Definition 3.10 (Geometrically unbiased regularization).** *A regularization scheme characterized by a bias function  $b(\alpha, \lambda^2)$  is said to be geometrically unbiased at order  $\beta > 0$  if:*

$$\sup_{\lambda^2} b^2(\alpha, \lambda^2) \lambda^{2\beta} = O(\alpha^\beta).$$

**Proposition 3.11.** *The spectral cut-off and the Landweber-Fridman regularization schemes are geometrically unbiased at any positive order  $\beta$ , whereas Tikhonov regularization scheme is geometrically unbiased at all order  $\beta \in (0, 2]$ . For Tikhonov regularization and  $\beta \geq 2$ , we have*

$$M_\beta(\alpha) = O(\alpha^2).$$

**Proof.** In the spectral cut-off case, there is no bias for  $\lambda_j^2 > \alpha/c$  while the bias is maximum, equal to one, for smaller  $\lambda_j^2$ . Therefore:

$$M_\beta(\alpha) \leq (\alpha/c)^\beta.$$

In the Landweber-Fridman case, there is no bias for  $\lambda_j^2 > 1/c$  but a decreasing bias  $(1 - c\lambda_j^2)^{1/\alpha}$  for  $\lambda_j^2$  increasing from zero to  $(1/c)$ . Therefore,  $M_\beta(\alpha) \leq [\text{Sup}_{\lambda^2 \leq (1/c)} (1 - c\lambda^2)^{2/\alpha} \lambda^{2\beta}]$ . The supremum is reached for  $\lambda^2 = (\beta/c)[\beta + (2/\alpha)]^{-1}$  and gives:

$$M_\beta(\alpha) \leq (\beta/c)^\beta [\beta + (2/\alpha)]^{-\beta} \leq (\beta/2)^\beta (\alpha/c)^\beta.$$

In the Tikhonov case, the bias decreases hyperbolically and then  $M_\beta(\alpha) \leq \sup_{\lambda^2} [\frac{(\alpha/c)}{(\alpha/c) + \lambda^2}]^2 \lambda^{2\beta}$ . For  $\beta < 2$ , the supremum is reached for  $\lambda^2 = (\beta\alpha/c)[2 - \beta]^{-1}$  and thus

$$M_\beta(\alpha) \leq \lambda^{2\beta} \leq [\beta/(2 - \beta)]^\beta (\alpha/c)^\beta.$$

As  $K$  is bounded, its largest eigenvalue is bounded. Therefore, for  $\beta \geq 2$ , we have

$$M_\beta(\alpha) \leq (\alpha/c)^2 \sup_j \lambda_j^{2(\beta-2)}.$$

■

**Proposition 3.12.** *Let  $K : \mathcal{H} \rightarrow \mathcal{E}$  be an injective compact operator. Let us assume that the solution  $\varphi$  of  $K\varphi = r$  lies in the  $\beta$ -regularity space  $\Phi_\beta$  of operator  $K$ , for some positive  $\beta$ . Then, if  $\varphi_\alpha$  is defined by a regularization scheme geometrically unbiased at order  $\beta$ , we have*

$$\|\varphi_\alpha - \varphi\|^2 = O(\alpha^\beta).$$

Therefore, the smoother the function  $\varphi$  of interest ( $\varphi \in \Phi_\beta$  for larger  $\beta$ ) is, the faster the rate of convergence to zero of the regularization bias will be. However, a degree of smoothness larger than or equal to 2 (corresponding to the case  $\varphi \in \mathcal{R}[(K^*K)]$ ) may be useless in the Tikhonov case. Indeed, for Tikhonov, we have  $\|\varphi_\alpha - \varphi\|^2 = O(\alpha^{\min(\beta, 2)})$ . This is basically the price to pay for a regularization procedure, which is simple to implement and rather intuitive (see Subsection 3.4 below) but introduces a regularization bias which never vanishes completely.

Both the operator interpretation and the practical implementation of smooth regularization schemes (Tikhonov and Landweber-Fridman) are discussed below.

### 3.4. Operator interpretation and implementation

In contrast to spectral cut-off, the advantage of Tikhonov and Landweber-Fridman regularization schemes is that they can be interpreted in terms of operators. Their algebraic expressions only depend on the global value of  $(K^*K)$  and  $(K^*r)$ , and not of the singular value decomposition. An attractive feature is that it implies that they can be implemented from the computation of sample counterparts  $(\hat{K}_n \hat{K}_n^*)$  and  $(\hat{K}_n^* \hat{r}_n)$  without resorting to an estimation of eigenvalues and eigenfunctions.

**The Tikhonov regularization** is based on

$$\begin{aligned} (\alpha_n I + K^*K) \varphi_{\alpha_n} &= K^*r \Leftrightarrow \\ \varphi_{\alpha_n} &= \sum_{j=1}^{\infty} \frac{\lambda_j}{\lambda_j^2 + \alpha_n} \langle r, \psi_j \rangle \phi_j \end{aligned}$$

for a penalization term  $\alpha_n$  and  $\lambda_j = \sqrt{\lambda_j^2}$ , while, for notational simplicity, the scaling factor  $c$  has been chosen equal to 1.

The interpretation of  $\alpha_n$  as a penalization term comes from the fact that  $\varphi_\alpha$  can be seen as the solution of

$$\varphi_\alpha = \arg \min_{\varphi} \|K\varphi - r\|^2 + \alpha \|\varphi\|^2 = \langle \varphi, K^*K\varphi + \alpha\varphi - 2K^*r \rangle + \|r\|^2.$$

To see this, just compute the Frechet derivative of the above expression and note that it is zero only for  $K^*K\varphi + \alpha\varphi = K^*r$ .

This interpretation of Tikhonov regularization in terms of penalization may suggest to look for quasi-solutions (see Kress, 1999, section 16-3), that is solutions of the minimization of  $\|K\varphi - r\|$  subject to the constraint that the norm is bounded by  $\|\varphi\| \leq \rho$  for

given  $\rho$ . For the purpose of econometric estimation, the quasi-solution may actually be the genuine solution if the specification of the structural econometric model entails that the function of interest  $\varphi$  lies in some compact set (Newey and Powell, 2003).

If one wants to solve directly the first order conditions of the above minimization, it is worth mentioning that the inversion of the operator  $(\alpha I + K^*K)$  is not directly well-suited for iterative approaches since, typically for small  $\alpha$ , the series expansion of  $[I + (1/\alpha)K^*K]^{-1}$  does not converge. However, a convenient choice of the estimators  $\hat{K}_n$  and  $\hat{K}_n^*$  may allow us to replace the inversion of infinite dimensional operators by the inversion of finite dimensional matrices.

More precisely, when  $\hat{K}_n$  and  $\hat{K}_n^*$  can be written as in (2.21) and (2.22), one can directly write the finite sample problem as:

$$\begin{aligned} (\alpha_n I + \hat{K}_n^* \hat{K}_n) \varphi &= \hat{K}_n^* r \Leftrightarrow \\ \alpha_n \varphi + \sum_{l,l'=1}^{L_n} a_{l'}(\varphi) b_l(\varepsilon_{l'}) \eta_l &= \sum_{l=1}^{L_n} b_l(r) \eta_l \end{aligned} \quad (3.16)$$

1) First we compute  $a_l(\varphi)$ :

Apply  $a_j$  to (3.16):

$$\alpha_n a_j(\varphi) + \sum_{l,l'=1}^{L_n} a_{l'}(\varphi) b_l(\varepsilon_{l'}) a_j(\eta_l) = \sum_{l=1}^{L_n} b_l(r) a_j(\eta_l) \quad (3.17)$$

(3.17) can be rewritten as

$$(\alpha_n I + A) \underline{a} = \underline{b}$$

where  $\underline{a} = [a_1(\varphi) \ a_2(\varphi) \ \cdots \ a_{L_n}(\varphi)]'$ ,  $A$  is the  $L_n \times L_n$ -matrix with principal element

$$A_{j,l'} = \sum_{l=1}^{L_n} b_l(\varepsilon_{l'}) a_j(\eta_l)$$

and

$$\underline{b} = \begin{bmatrix} \sum_l b_l(r) a_1(\eta_l) \\ \vdots \\ \sum_l b_l(r) a_{L_n}(\eta_l) \end{bmatrix}.$$

2) From (3.16), we have

$$\hat{\varphi}_n = \frac{1}{\alpha_n} \left[ \sum_{l=1}^{L_n} b_l(r) \eta_l - \sum_{l,l'=1}^{L_n} a_{l'}(\varphi) b_l(\varepsilon_{l'}) \eta_l \right].$$

### Landweber-Fridman regularization

The great advantage of this regularization scheme is not only that it can be written directly in terms of quantities  $(K^*K)$  and  $(K^*r)$ , but also the resulting operator problem can be solved by a simple iterative procedure, with a finite number of steps. To get this, one has to first choose a sequence of regularization parameters,  $\alpha_n$ , such that  $(1/\alpha_n)$  is an integer and second the scaling factor  $c$  so that  $0 < c < 1/\|K\|^2$ . This latter condition may be problematic to implement since the norm of the operator  $K$  may be unknown. The refinements of an asymptotic theory, that enables us to accommodate a first step estimation of  $\|K\|$  before the selection of an appropriate  $c$ , is beyond the scope of this chapter. Note however, that in several cases of interest,  $\|K\|$  is known a priori even though the operator  $K$  itself is unknown. For example, if  $K$  is the conditional expectation operator,  $\|K\| = 1$ .

The advantage of the condition  $c < 1/\|K\|^2$  is to guarantee a unique expression for the bias function  $b(\alpha, \lambda^2) = (1 - c\lambda^2)^{1/\alpha}$  since, for all eigenvalues,  $\lambda^2 \leq 1/c$ . Thus, when  $(1/\alpha)$  is an integer :

$$\varphi_\alpha = \sum_{j=1}^{\infty} \frac{1}{\lambda_j} q(\alpha, \lambda_j^2) \langle r, \psi_j \rangle \phi_j$$

with

$$\begin{aligned} q(\alpha, \lambda_j^2) &= 1 - (1 - c\lambda_j^2)^{1/\alpha} \\ &= c\lambda_j^2 \sum_{l=0}^{1/\alpha-1} (1 - c\lambda_j^2)^l. \end{aligned}$$

Thus,

$$\begin{aligned} \varphi_\alpha &= c \sum_{l=0}^{1/\alpha-1} \sum_{j=1}^{\infty} \lambda_j (1 - c\lambda_j^2)^l \langle r, \psi_j \rangle \phi_j \\ &= c \sum_{l=0}^{1/\alpha-1} \sum_{j=1}^{\infty} \lambda_j^2 (1 - c\lambda_j^2)^l \langle \varphi, \phi_j \rangle \phi_j \\ &= c \sum_{l=0}^{1/\alpha-1} (I - cK^*K)^l K^*K\varphi. \end{aligned}$$

Therefore, the estimation procedure will only resort to estimators of  $K^*K$  and of  $K^*r$ , without need for either the singular value decomposition nor any inversion of operators. For a given  $c$  and regularization parameter  $\alpha_n$ , the estimator of  $\varphi$  is

$$\hat{\varphi}_n = c \sum_{l=0}^{1/\alpha_n-1} \left( I - c\hat{K}_n^*\hat{K}_n \right)^l \hat{K}_n^*\hat{r}_n.$$



$\hat{\varphi}_n$  can be computed recursively by

$$\hat{\varphi}_{l,n} = \left( I - c\hat{K}_n^*\hat{K}_n \right) \hat{\varphi}_{l-1,n} + c\hat{K}_n^*\hat{r}_n, \quad l = 1, 2, \dots, 1/\alpha_n - 1.$$

starting with  $\hat{\varphi}_{0,n} = c\hat{K}_n^*\hat{r}_n$ . This scheme is known as the Landweber-Fridman iteration (see Kress, 1999, p. 287).

### 3.5. Estimation bias

Regularization schemes have precisely been introduced because the right hand side  $r$  of the inverse problem  $K\varphi = r$  is generally unknown and replaced by an estimator. Let us denote by  $\hat{r}_n$  an estimator computed from an observed sample of size  $n$ . As announced in the introduction, a number of relevant inverse problems in econometrics are even more complicated since the operator  $K$  itself is unknown. Actually, in order to apply a regularization scheme, we may not need only an estimator of  $K$  but also of its adjoint  $K^*$  and of its singular system  $\{\lambda_j, \phi_j, \psi_j : j = 1, 2, \dots\}$ . In this subsection, we consider such estimators  $\hat{K}_n$ ,  $\hat{K}_n^*$ , and  $\{\hat{\lambda}_j, \hat{\phi}_j, \hat{\psi}_j : j = 1, \dots, L_n\}$  as given. We also maintain the identification assumption, so that the equation  $K\varphi = r$  defines without ambiguity a true unknown value  $\varphi_0$ .

If  $\varphi_\alpha = A_\alpha K^* r$  is the chosen regularized solution, the proposed estimator  $\hat{\varphi}_n$  of  $\varphi_0$  is defined by

$$\hat{\varphi}_n = \hat{A}_{\alpha_n} \hat{K}_n^* \hat{r}_n. \quad (3.18)$$

Note that the definition of this estimator involves two decisions. First, we need to select a sequence  $(\alpha_n)$  of regularization parameters so that  $\lim_{n \rightarrow \infty} \alpha_n = 0$  (possibly in a stochastic sense in the case of a data-driven regularization) in order to get a consistent estimator of  $\varphi_0$ . Second, for a given  $\alpha_n$ , we estimate the second order regularization scheme  $A_{\alpha_n} K^*$  by  $\hat{A}_{\alpha_n} \hat{K}_n^*$ . Generally speaking,  $\hat{A}_{\alpha_n}$  is defined from (3.9) where the singular values are replaced by their estimators and the inner products  $\langle \varphi, \phi_j \rangle$  are replaced by their empirical counterparts (see Subsection 2.5.3). Yet, we have seen above that in some cases, the estimation of the regularized solution does not involve the estimators  $\hat{\lambda}_j$  but only the estimators  $\hat{K}_n$  and  $\hat{K}_n^*$ .

In any case, the resulting estimator bias  $\hat{\varphi}_n - \varphi_0$  has two components:

$$\hat{\varphi}_n - \varphi_0 = \hat{\varphi}_n - \varphi_{\alpha_n} + \varphi_{\alpha_n} - \varphi_0. \quad (3.19)$$

While the second component  $\varphi_{\alpha_n} - \varphi_0$  defines the regularization bias characterized in Subsection 3.3, the first component  $\hat{\varphi}_n - \varphi_{\alpha_n}$  is the bias corresponding to the estimation of the regularized solution of  $\varphi_{\alpha_n}$ . The goal of this subsection is to point out a set of statistical assumptions about the estimators  $\hat{K}_n$ ,  $\hat{K}_n^*$ , and  $\hat{r}_n$  that gives an (asymptotic) upper bound to the specific estimation bias magnitude,  $\|\hat{\varphi}_n - \varphi_{\alpha_n}\|$  when the regularization bias  $\|\varphi_{\alpha_n} - \varphi_0\|$  is controlled.

**Definition 3.13 (Smooth regularization).** A regularization scheme is said to be smooth when

$$\left\| \left( \hat{A}_{\alpha_n} \hat{K}_n^* \hat{K}_n - A_{\alpha_n} K^* K \right) \varphi_0 \right\| \leq d(\alpha_n) \left\| \hat{K}_n^* \hat{K}_n - K^* K \right\| \|\varphi_{\alpha_n} - \varphi_0\| (1 + \varepsilon_n) \quad (3.20)$$

with  $\varepsilon_n = O\left(\left\|\hat{K}_n^* \hat{K}_n - K^* K\right\|\right)$ .

**Proposition 3.14 (Estimation bias).** If  $\varphi_\alpha = A_\alpha K^* r$  is the regularized solution conformable to Definition 3.9 and  $\hat{\varphi}_n = \hat{A}_{\alpha_n} \hat{K}_n^* \hat{r}_n$ , then

$$\begin{aligned} & \left\| \hat{\varphi}_n - \varphi_{\alpha_n} \right\| \\ & \leq d(\alpha_n) \left\| \hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi_0 \right\| + \left\| \left( \hat{A}_{\alpha_n} \hat{K}_n^* \hat{K}_n - A_{\alpha_n} K^* K \right) \varphi_0 \right\| \end{aligned} \quad (3.21)$$

In addition, both the Tikhonov and Landweber-Fridman regularization schemes are smooth. In the Tikhonov case,  $\varepsilon_n = 0$  identically.

**Proof.**

$$\begin{aligned} \hat{\varphi}_n - \varphi_{\alpha_n} &= \hat{A}_{\alpha_n} \hat{K}_n^* \hat{r}_n - A_{\alpha_n} K^* r \\ &= \hat{A}_{\alpha_n} \hat{K}_n^* \left( \hat{r}_n - \hat{K}_n \varphi_0 \right) + \hat{A}_{\alpha_n} \hat{K}_n^* \hat{K}_n \varphi_0 - A_{\alpha_n} K^* K \varphi_0 \end{aligned}$$

Thus,

$$\left\| \hat{\varphi}_n - \varphi_{\alpha_n} \right\| \leq d(\alpha_n) \left\| \hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi_0 \right\| + \left\| \hat{A}_{\alpha_n} \hat{K}_n^* \hat{K}_n \varphi_0 - A_{\alpha_n} K^* K \varphi_0 \right\|.$$

1) Case of Tikhonov regularization:

$$\begin{aligned} & \hat{A}_{\alpha_n} \hat{K}_n^* \hat{K}_n \varphi_0 - A_{\alpha_n} K^* K \varphi_0 \\ &= \hat{A}_{\alpha_n} \left( \hat{K}_n^* \hat{K}_n - K^* K \right) \varphi_0 + \left( \hat{A}_{\alpha_n} - A_{\alpha_n} \right) K^* K \varphi_0. \end{aligned} \quad (3.22)$$

Since in this case,

$$A_\alpha = (\alpha I + K^* K)^{-1},$$

the identity

$$B^{-1} - C^{-1} = B^{-1}(C - B)C^{-1}$$

gives

$$\hat{A}_{\alpha_n} - A_{\alpha_n} = \hat{A}_{\alpha_n} \left( K^* K - \hat{K}_n^* \hat{K}_n \right) A_{\alpha_n}$$

and thus,

$$\begin{aligned} (\hat{A}_{\alpha_n} - A_{\alpha_n}) K^* K \varphi_0 &= \hat{A}_{\alpha_n} (K^* K - \hat{K}_n^* \hat{K}_n) A_{\alpha_n} K^* K \varphi_0 \\ &= \hat{A}_{\alpha_n} (K^* K - \hat{K}_n^* \hat{K}_n) \varphi_{\alpha_n}. \end{aligned} \quad (3.23)$$

(3.22) and (3.23) together give

$$\begin{aligned} &\hat{A}_{\alpha_n} \hat{K}_n^* \hat{K}_n \varphi_0 - A_{\alpha_n} K^* K \varphi_0 \\ &= \hat{A}_{\alpha_n} (\hat{K}_n^* \hat{K}_n - K^* K) (\varphi_0 - \varphi_{\alpha_n}), \end{aligned}$$

which shows that Tikhonov regularization is smooth with  $\varepsilon_n = 0$ .

2) Case of Landweber-Fridman regularization:

In this case,

$$\begin{aligned} \varphi_\alpha &= \sum_{j=1}^{\infty} \left[ 1 - (1 - c\lambda_j^2)^{1/\alpha} \right] \langle \varphi_0, \phi_j \rangle \phi_j \\ &= \left[ I - (I - cK^*K)^{1/\alpha} \right] \varphi_0. \end{aligned}$$

Thus,

$$\begin{aligned} &\hat{A}_{\alpha_n} \hat{K}_n^* \hat{K}_n \varphi_0 - A_{\alpha_n} K^* K \varphi_0 \\ &= \left[ (I - cK^*K)^{1/\alpha_n} - (I - c\hat{K}_n^* \hat{K}_n)^{1/\alpha_n} \right] \varphi_0 \\ &\quad + \left[ I - (I - c\hat{K}_n^* \hat{K}_n)^{1/\alpha_n} (I - cK^*K)^{-1/\alpha_n} \right] (I - cK^*K)^{1/\alpha_n} \varphi_0 \\ &\quad + \left[ I - (I - c\hat{K}_n^* \hat{K}_n)^{1/\alpha_n} (I - cK^*K)^{-1/\alpha_n} \right] (\varphi_0 - \varphi_{\alpha_n}). \end{aligned}$$

Then, a Taylor expansion gives:

$$\begin{aligned} &\left\| I - (I - c\hat{K}_n^* \hat{K}_n)^{1/\alpha_n} (I - cK^*K)^{-1/\alpha_n} \right\| \\ &= \left\| \frac{c}{\alpha_n} (\hat{K}_n^* \hat{K}_n - K^*K) \right\| (1 + \varepsilon_n) \end{aligned}$$

with  $\varepsilon_n = O\left(\left\| \hat{K}_n^* \hat{K}_n - K^*K \right\|\right)$ .

The result follows with  $d(\alpha) = c/\alpha$ . ■

Note that we are not able to establish (3.20) for the spectral cut-off regularization. In that case, the threshold introduces a lack of smoothness, which precludes a similar Taylor expansion based argument as above.

The result of Proposition 3.14 jointly with (3.19) shows that two ingredients matter in controlling the estimation bias  $\|\hat{\varphi}_n - \varphi_0\|$ . First, the choice of a sequence of regularization

parameters,  $\alpha_n$ , will govern the speed of convergence to zero of the regularization bias  $\|\varphi_{\alpha_n} - \varphi_0\|$  (for  $\varphi_0$  in a given  $\Phi_\beta$ ) and the speed of convergence to infinity of  $d(\alpha_n)$ . Second, nonparametric estimation of  $K$  and  $r$  will determine the rates of convergence of  $\|\hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi_0\|$  and  $\|\hat{K}_n^* \hat{K}_n - K^* K\|$ .

## 4. Asymptotic properties of solutions of integral equations of the first kind

### 4.1. Consistency

Let  $\varphi_0$  be the solution of  $K\varphi = r$ . By abuse of notation, we denote  $X_n = O(c_n)$  for positive sequences  $\{X_n\}$  and  $\{c_n\}$ , if the sequence  $X_n/c_n$  is upper bounded.

We maintain the following assumptions:

**A1.**  $\hat{K}_n, \hat{r}_n$  are consistent estimators of  $K$  and  $r$ .

**A2.**  $\|\hat{K}_n^* \hat{K}_n - K^* K\| = O\left(\frac{1}{a_n}\right)$

**A3.**  $\|\hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi_0\| = O\left(\frac{1}{b_n}\right)$

As before  $\varphi_\alpha = A_\alpha K^* r$  is the regularized solution where  $A_\alpha$  is a second order regularization scheme and  $\hat{\varphi}_n = \hat{A}_{\alpha_n} \hat{K}_n^* \hat{r}_n$ . Proposition 4.1 below follows directly from Proposition 3.9 and Definition 3.6 (with the associated normalization rule  $\alpha d(\alpha) = O(1)$ ):

**Proposition 4.1.** *When applying a smooth regularization scheme, we get:*

$$\begin{aligned} & \|\hat{\varphi}_n - \varphi_0\| \\ &= O\left(\frac{1}{\alpha_n b_n} + \left(\frac{1}{\alpha_n a_n} + 1\right) \|\varphi_{\alpha_n} - \varphi_0\|\right). \end{aligned}$$

#### Discussion on the rate of convergence:

The general idea is that the fastest possible rate of convergence in probability of  $\|\hat{\varphi}_n - \varphi_0\|$  to zero should be the rate of convergence of the regularization bias  $\|\varphi_{\alpha_n} - \varphi_0\|$ . Proposition 4.1 shows that these two rates of convergence will precisely coincide when the rate of convergence to zero of the regularization parameter,  $\alpha_n$ , is chosen sufficiently slow with respect to both the rate of convergence  $a_n$  of the sequence of approximations of the true operator, and the rate of convergence  $b_n$  of the estimator of the right-hand side of the operator equation. This is actually a common strategy when both the operator and the right-hand side of the inverse problem have to be estimated (see e.g. Vapnik (1998), corollary p. 299).

To get this, it is first obvious that  $\alpha_n b_n$  must go to infinity at least as fast as  $\|\varphi_{\alpha_n} - \varphi_0\|^{-1}$ . For  $\varphi_0 \in \Phi_\beta$ ,  $\beta > 0$ , and a geometrically unbiased regularization scheme, this means that:

$$\alpha_n^2 b_n^2 \geq \alpha_n^{-\beta}$$

that is  $\alpha_n \geq b_n^{-\frac{2}{\beta+2}}$ . To get the fastest possible rate of convergence under this constraint, we will choose:

$$\alpha_n = b_n^{-\frac{2}{\beta+2}}.$$

Then, the rate of convergence of  $\|\hat{\varphi}_n - \varphi_0\|$  and  $\|\varphi_{\alpha_n} - \varphi_0\|$  will coincide if and only if  $a_n b_n^{-\frac{2}{\beta+2}}$  is bounded away from zero. Thus, we have proved:

**Proposition 4.2.** *Consider a smooth regularization scheme, which is geometrically unbiased of order  $\beta > 0$  with estimators of  $K$  and  $r$  conformable to Assumptions A1, A2, A3, and  $a_n b_n^{-\frac{2}{\beta+2}}$  bounded away from zero. For  $\varphi_0 \in \Phi_\beta$ , the optimal choice of the regularization parameter is  $\alpha_n = b_n^{-\frac{2}{\beta+2}}$ , and then*

$$\|\hat{\varphi}_n - \varphi_0\| = O\left(b_n^{-\frac{\beta}{\beta+2}}\right).$$

For Tikhonov regularization, when  $\varphi_0 \in \Phi_\beta$ ,  $\beta > 0$ , provided  $a_n b_n^{-\min(\frac{2}{\beta+2}, \frac{1}{2})}$  is bounded away from zero and  $\alpha_n = b_n^{-\min(\frac{2}{\beta+2}, \frac{1}{2})}$ , we have

$$\|\hat{\varphi}_n - \varphi_0\| = O\left(b_n^{-\min(\frac{\beta}{\beta+2}, \frac{1}{2})}\right).$$

Note that the only condition about the estimator of the operator  $K^*K$  is that its rate of convergence,  $a_n$ , is sufficiently fast to be greater than  $b_n^{\frac{2}{\beta+2}}$ . Under this condition, the rate of convergence of  $\hat{\varphi}_n$  does not depend upon the accuracy of the estimator of  $K^*K$ . Of course, the more regular the unknown function  $\varphi_0$  is, the larger  $\beta$  is and the easier it will be to meet the required condition. Generally speaking, the condition will involve the relative bandwidth sizes in the nonparametric estimation of  $K^*K$  and  $K^*r$ . Note that if, as it is generally the case for a convenient bandwidth choice (see e.g. subsection 5.4),  $b_n$  is the parametric rate ( $b_n = \sqrt{n}$ ),  $a_n$  must be at least  $n^{1/(\beta+2)}$ . For  $\beta$  not too small, this condition will be fulfilled by optimal nonparametric rates. For instance, the optimal unidimensional nonparametric rate,  $n^{2/5}$ , will work as soon as  $\beta \geq 1/2$ .

The larger  $\beta$  is, the faster the rate of convergence of  $\hat{\varphi}_n$  is. In the case where  $\varphi_0$  is a finite linear combination of  $\{\phi_j\}$  (case where  $\beta$  is infinite), and  $b_n = \sqrt{n}$ , an estimator based on a geometrically unbiased regularization scheme (such as Landweber-Fridman) achieves the parametric rate of convergence. We are not able to obtain such a fast rate for Tikhonov, therefore it seems that if the function  $\varphi_0$  is suspected to be very regular, Landweber-Fridman is preferable to Tikhonov. However, it should be noted that the rates of convergence in Proposition 4.2 are upperbounds and could possibly be improved upon.

## 4.2. Asymptotic normality

Asymptotic normality of

$$\begin{aligned}\hat{\varphi}_n - \varphi_0 &= \hat{\varphi}_n - \varphi_{\alpha_n} + \varphi_{\alpha_n} - \varphi_0 \\ &= \hat{A}_{\alpha_n} \hat{K}_n^* \hat{r}_n - A_{\alpha_n} K^* K \varphi_0 + \varphi_{\alpha_n} - \varphi_0\end{aligned}$$

can be deduced from a functional central limit theorem applied to  $\hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi_0$ . Therefore, we must reinforce Assumption A3 by assuming a weak convergence in  $\mathcal{H}$ :

**Assumption WC:**

$$b_n \left( \hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi_0 \right) \Rightarrow \mathcal{N}(0, \Sigma) \text{ in } \mathcal{H}.$$

According to (3.21), (3.22), and (3.23), we have in the case of Tikhonov regularization:

$$b_n (\hat{\varphi}_n - \varphi_0) = b_n \hat{A}_{\alpha_n} \left[ \hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi_0 \right] \quad (4.1)$$

$$+ b_n \hat{A}_{\alpha_n} \left[ \hat{K}_n^* \hat{K}_n - K^* K \right] (\varphi_0 - \varphi_{\alpha_n}) \quad (4.2)$$

while an additional term corresponding  $\varepsilon_n$  in (3.20) should be added for general regularization schemes. The term (4.1) can be rewritten as

$$\hat{A}_{\alpha_n} \xi + \hat{A}_{\alpha_n} (\xi_n - \xi)$$

where  $\xi$  denotes the random variable  $\mathcal{N}(0, \Sigma)$  in  $\mathcal{H}$  and

$$\xi_n = b_n \left( \hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi_0 \right).$$

By definition:

$$\frac{\langle \hat{A}_{\alpha_n} \xi, g \rangle}{\left\| \Sigma^{1/2} \hat{A}_{\alpha_n} g \right\|} \xrightarrow{d} \mathcal{N}(0, 1)$$

for all  $g \in \mathcal{H}$ . Then, we may hope to get a standardized normal asymptotic probability distribution for

$$\frac{\langle b_n (\hat{\varphi}_n - \varphi_0), g \rangle}{\left\| \Sigma^{1/2} \hat{A}_{\alpha_n} g \right\|}$$

for vectors  $g$  conformable to the following assumption:

**Assumption G**

$$\frac{\left\| \hat{A}_{\alpha_n} g \right\|}{\left\| \Sigma^{1/2} \hat{A}_{\alpha_n} g \right\|} = O(1).$$

Indeed, we have in this case:

$$\frac{\left| \left\langle \hat{A}_{\alpha_n} (\xi_n - \xi), g \right\rangle \right|}{\left\| \Sigma^{1/2} \hat{A}_{\alpha_n} g \right\|} \leq \frac{\|\xi_n - \xi\| \left\| \hat{A}_{\alpha_n} g \right\|}{\left\| \Sigma^{1/2} \hat{A}_{\alpha_n} g \right\|},$$

which converges to zero in probability because  $\|\xi_n - \xi\| \xrightarrow{P} 0$  by WC. We are then able to show:

**Proposition 4.3.** *Consider a Tikhonov regularization. Suppose Assumptions A1, A2, A3, and WC hold and  $\varphi_0 \in \Phi_\beta$ ,  $\beta > 0$ , with  $b_n \alpha_n^{\min(\beta/2, 1)} \xrightarrow[n \rightarrow \infty]{} 0$ , we have for any  $g$  conformable to  $G$ :*

$$\frac{\langle b_n (\hat{\varphi}_n - \varphi_0), g \rangle}{\left\| \Sigma^{1/2} \hat{A}_{\alpha_n} g \right\|} \xrightarrow{d} \mathcal{N}(0, 1).$$

**Proof.** From the proof of Proposition 3.9, we have:

$$\begin{aligned} & \langle b_n (\hat{\varphi}_n - \varphi_{\alpha_n}), g \rangle \\ = & \left\langle \hat{A}_{\alpha_n} \xi, g \right\rangle \\ & + \left\langle \hat{A}_{\alpha_n} (\xi_n - \xi), g \right\rangle \\ & + \left\langle b_n \hat{A}_{\alpha_n} \left[ \hat{K}_n^* \hat{K}_n - K^* K \right] (\varphi_0 - \varphi_{\alpha_n}), g \right\rangle \end{aligned} \quad (4.3)$$

in the case of Tikhonov regularization. We already took care of the terms in  $\xi$  and  $\xi_n$ , it remains to deal with the bias term corresponding to (4.3):

$$\begin{aligned} & \frac{b_n \left\langle \hat{A}_{\alpha_n} \left( \hat{K}_n^* \hat{K}_n - K^* K \right) (\varphi_0 - \varphi_{\alpha_n}), g \right\rangle}{\left\| \Sigma^{1/2} \hat{A}_{\alpha_n} g \right\|} \\ \leq & \frac{b_n \left\langle \left( \hat{K}_n^* \hat{K}_n - K^* K \right) (\varphi_0 - \varphi_{\alpha_n}), \hat{A}_{\alpha_n} g \right\rangle}{\left\| \Sigma^{1/2} \hat{A}_{\alpha_n} g \right\|} \\ \leq & b_n \left\| \hat{K}_n^* \hat{K}_n - K^* K \right\| \|\varphi_0 - \varphi_{\alpha_n}\| \frac{\left\| \hat{A}_{\alpha_n} g \right\|}{\left\| \Sigma^{1/2} \hat{A}_{\alpha_n} g \right\|} \\ = & O \left( \frac{b_n \alpha_n^{\min(\beta/2, 1)}}{a_n} \right). \end{aligned}$$

■

**Discussion of Proposition 4.3.**

(i) It is worth noticing that Proposition 4.3 does not in general deliver a weak convergence result for  $b_n(\hat{\varphi}_n - \varphi_0)$  because it does not hold for all  $g \in \mathcal{H}$ . However, the condition G is not so restrictive. It just amounts to assume that the multiplication by  $\Sigma^{1/2}$  does not modify the rate of convergence of  $\hat{A}_{\alpha_n}g$ .

(ii) We remark that for  $g = K^*Kh$ ,  $\hat{A}_{\alpha_n}g$  and  $\Sigma^{1/2}\hat{A}_{\alpha_n}g$  converge respectively to  $h$  and  $\Sigma^{1/2}h$ . Moreover, if  $g \neq 0$ ,  $\Sigma^{1/2}h = \Sigma^{1/2}(K^*K)^{-1}g \neq 0$ . Therefore, in this case, not only is the condition G fulfilled but we get asymptotic normality with rate of convergence  $b_n$ , that is typically root  $n$ . This result is conformable to the theory of asymptotic efficiency of inverse estimators as recently developed by Van Rooij, Ruymgaart and Van Zwet (2000). They show that there is a dense linear submanifold of functionals for which the estimators are asymptotically normal at the root  $n$  rate with optimal variance (in the sense of minimum variance in the class of the moment estimators). We do get optimal variance in Proposition 4.3 since in this case (using heuristic notations as if we were in finite dimension) the asymptotic variance is:

$$\begin{aligned} & \lim_{n \rightarrow \infty} g' A_{\alpha_n} \Sigma A_{\alpha_n} \\ &= g' (K^*K)^{-1} \Sigma (K^*K)^{-1} g. \end{aligned}$$

Moreover, we get this result in particular for any nonzero  $g$  in  $\mathcal{R}(K^*K)$  while we know that  $\mathcal{R}(K^*)$  is dense in  $\mathcal{H}$  (identification condition). Generally speaking, Van Rooij, Ruymgaart and Van Zwet (2000) stress that the inner products do not converge weakly for every vector  $g$  in  $\mathcal{H}$  at the same rate, if they converge at all.

(iii) The condition  $b_n \alpha_n^{\min(\beta/2, 1)} \rightarrow 0$  imposes a convergence to zero of the regularization coefficient  $\alpha_n$  faster than the rate  $\alpha_n = b_n^{-\min(\frac{2}{\beta+2}, \frac{1}{2})}$  required for the consistency. This stronger condition is needed to show that the regularization bias multiplied by  $b_n$  converges to zero. A fortiori, the estimation bias term vanishes asymptotically.

The results of Proposition 4.3 are established under strong assumptions: convergence in  $\mathcal{H}$  and restriction on  $g$ . An alternative method consists in establishing the normality of  $\hat{\varphi}_n$  by the Liapunov condition (Davidson, 1994), see the example on deconvolution in Section 5 below.

## 5. Applications

A well-known example is that of the kernel estimator of the density. Indeed, the estimation of the pdf  $f$  of a random variable  $X$  can be seen as solving an integral equation of the first kind

$$Kf(x) = \int_{-\infty}^{+\infty} I(u \leq x) f(u) du = F(x) \tag{5.1}$$

where  $F$  is the cdf of  $X$ . Applying the Tikhonov regularization to (5.1), one obtains a kernel estimator of  $f$ . This example is detailed in Hardle and Linton (1994) and in Vapnik (1998, pages 308-311) and will not be discussed further.



This section reviews the standard examples of the ridge regression and factor models and less standard examples such as the regression with an infinite number of regressors, the deconvolution and the instrumental variable estimation.

### 5.1. Ridge regression

The Tikhonov regularization discussed in Section 3 can be seen as an extension of the well-known ridge regression. The ridge regression was introduced by Hoerl and Kennard (1970). It was initially motivated by the fact that in the presence of near multicollinearity of the regressors, the least squares estimator may vary dramatically as the result of a small perturbation in the data. The ridge estimator is more stable and may have a lower risk than the conventional least squares estimator. For a review of this method, see Judge, Griffiths, Hill, Lutkepohl, and Lee (1980) and for a discussion in the context of inverse problems, see Ruymgaart (2001).

Consider the linear model (the notation of this paragraph is specific and corresponds to general notations of linear models):

$$y = X\theta + \varepsilon \tag{5.2}$$

where  $y$  and  $\varepsilon$  are  $n \times 1$ -random vectors,  $X$  is a  $n \times q$  matrix of regressors of full rank, and  $\theta$  is an unknown  $q \times 1$ -vector of parameters. The number of explanatory variables,  $q$ , is assumed to be constant and  $q < n$ . Assume that  $X$  is exogenous and all the expectations are taken conditionally on  $X$ . The classical least-squares estimator of  $\theta$  is

$$\hat{\theta} = (X'X)^{-1} X'y.$$

There exists an orthogonal transformation such that  $X'X/n = P'DP$  with

$$D = \begin{pmatrix} \mu_1 & & 0 \\ & \ddots & \\ 0 & & \mu_q \end{pmatrix},$$

$\mu_j > 0$ , and  $P'P = I_q$ . Using the mean square error as measure of the risk, we get

$$\begin{aligned}
E \left\| \hat{\theta} - \theta \right\|^2 &= E \left\| (X'X)^{-1} (X' (X\theta + \varepsilon) - \theta) \right\|^2 \\
&= E \left\| (X'X)^{-1} X' \varepsilon \right\|^2 \\
&= E \left( \varepsilon' X (X'X)^{-2} X' \varepsilon \right) \\
&= \sigma^2 \text{trace} \left( X (X'X)^{-2} X' \right) \\
&= \frac{\sigma^2}{n} \text{trace} \left( \left( \frac{X'X}{n} \right)^{-1} \right) \\
&= \frac{\sigma^2}{n} \text{trace} (P'D^{-1}P) \\
&= \frac{\sigma^2}{n} \sum_{j=1}^q \frac{1}{\mu_j}.
\end{aligned}$$

If some of the columns of  $X$  are closely collinear, the eigenvalues may be very small and the risk very large. Moreover, when the number of regressors is infinite, the risk is no longer bounded.

A solution is to use the ridge regression estimator:

$$\begin{aligned}
\hat{\theta}_a &= \arg \min_{\theta} \|y - X\theta\|^2 + a \|\theta\|^2 \\
\Rightarrow \hat{\theta}_a &= (aI + X'X)^{-1} X'y
\end{aligned}$$

for  $a > 0$ . We prefer to introduce  $\alpha = a/n$  and define

$$\hat{\theta}_\alpha = \left( \alpha I + \frac{X'X}{n} \right)^{-1} \frac{X'y}{n}. \tag{5.3}$$

This way, the positive constant  $\alpha$  corresponds to the regularization parameter introduced in earlier sections.

The estimator  $\hat{\theta}_\alpha$  is no longer unbiased. We have

$$\theta_\alpha = E \left( \hat{\theta}_\alpha \right) = \left( \alpha I + \frac{X'X}{n} \right)^{-1} \frac{X'X}{n} \theta.$$

Using the fact that  $A^{-1} - B^{-1} = A^{-1} [B - A] B^{-1}$ . The bias can be rewritten as

$$\begin{aligned}
\theta_\alpha - \theta &= \left( \alpha I + \frac{X'X}{n} \right)^{-1} \frac{X'X}{n} \theta - \left( \frac{X'X}{n} \right)^{-1} \frac{X'X}{n} \theta \\
&= \alpha \left( \alpha I + \frac{X'X}{n} \right)^{-1} \theta.
\end{aligned}$$

The risk becomes

$$\begin{aligned}
E \left\| \hat{\theta}_\alpha - \theta \right\|^2 &= E \left\| \hat{\theta}_\alpha - \theta_\alpha \right\|^2 + \|\theta_\alpha - \theta\|^2 \\
&= E \left\| \left( \alpha I + \frac{X'X}{n} \right)^{-1} \frac{X'\varepsilon}{n} \right\|^2 + \alpha^2 \left\| \left( \alpha I + \frac{X'X}{n} \right)^{-1} \theta \right\|^2 \\
&= E \left( \frac{\varepsilon'X}{n} \left( \alpha I + \frac{X'X}{n} \right)^{-2} \frac{X'\varepsilon}{n} \right) + \alpha^2 \left\| \left( \alpha I + \frac{X'X}{n} \right)^{-1} \theta \right\|^2 \\
&= \frac{\sigma^2}{n} \text{trace} \left( \left( \alpha I + \frac{X'X}{n} \right)^{-2} \frac{X'X}{n} \right) + \alpha^2 \left\| \left( \alpha I + \frac{X'X}{n} \right)^{-1} \theta \right\|^2 \\
&= \frac{\sigma^2}{n} \sum_{j=1}^q \frac{\mu_j}{(\alpha + \mu_j)^2} + \alpha^2 \sum_{j=1}^q \frac{((P\theta)_j)^2}{(\alpha + \mu_j)^2}.
\end{aligned}$$

There is the usual trade-off between the variance (decreasing in  $\alpha$ ) and the bias (increasing in  $\alpha$ ). For each  $\theta$  and  $\sigma^2$ , there is a value of  $\alpha$  for which the risk of  $\hat{\theta}_\alpha$  is smaller than that of  $\hat{\theta}$ . As  $q$  is finite, we have  $E \left\| \hat{\theta}_\alpha - \theta_\alpha \right\|^2 \sim 1/n$  and  $\|\theta_\alpha - \theta\|^2 \sim \alpha^2$ . Hence, the MSE is minimized for  $\alpha_n \sim 1/\sqrt{n}$ . Let us compare this rate with that necessary to the asymptotic normality of  $\hat{\theta}_\alpha$ . We have

$$\hat{\theta}_\alpha - \theta = -\alpha \left( \alpha I + \frac{X'X}{n} \right)^{-1} \theta + \left( \alpha I + \frac{X'X}{n} \right)^{-1} \frac{X'\varepsilon}{n}.$$

Therefore, if  $X$  and  $\varepsilon$  satisfy standard assumptions of stationarity and mixing,  $\hat{\theta}_\alpha$  is consistent as soon as  $\alpha_n$  goes to zero and  $\sqrt{n}(\hat{\theta}_\alpha - \theta)$  is asymptotically centered normal provided  $\alpha_n = o(1/\sqrt{n})$ , which is a faster rate than that obtained in the minimization of the MSE. Data-dependent methods for selecting the value of  $\alpha$  are discussed in Judge et al. (1980).

Note that the ridge estimator (5.3) is the regularized inverse of the equation

$$y = X\theta, \tag{5.4}$$

where obviously  $\theta$  is overidentified as there are  $n$  equations for  $q$  unknowns. Let  $\mathcal{H}$  be  $\mathbb{R}^q$  endowed with the euclidean norm and  $\mathcal{E}$  be  $\mathbb{R}^n$  endowed with the norm,  $\|v\|^2 = v'v/n$ . Define  $K : \mathcal{H} \rightarrow \mathcal{E}$  such that  $Ku = Xu$  for any  $u \in \mathbb{R}^q$ . Solving  $\langle Ku, v \rangle = \langle u, K^*v \rangle$ , we find the adjoint of  $K$ ,  $K^* : \mathcal{E} \rightarrow \mathcal{H}$  where  $K^*v = X'v/n$  for any  $v \in \mathbb{R}^n$ . The Tikhonov regularized solution of (5.4) is given by

$$\hat{\theta}_\alpha = (\alpha I + K^*K)^{-1} K^*y,$$

which corresponds to (5.3). It is also interesting to look at the spectral cut-off regularization. Let  $\{P_1, P_2, \dots, P_q\}$  be the orthonormal eigenvectors of the  $q \times q$  matrix

$K^*K = X'X/n$  and  $\{Q_1, Q_2, \dots, Q_n\}$  be the orthonormal eigenvectors of the  $n \times n$  matrix  $KK^* = XX'/n$ . Let  $\lambda_j = \sqrt{\mu_j}$ . Then the spectral cut-off regularized estimator is

$$\hat{\theta}_\alpha = \sum_{\lambda_j \geq \alpha} \frac{1}{\lambda_j} \langle y, Q_j \rangle P_j = \sum_{\lambda_j \geq \alpha} \frac{1}{\lambda_j} \frac{y'Q_j}{n} P_j.$$

A variation on the spectral cut-off consists in keeping the  $l$  largest eigenvalues to obtain

$$\hat{\theta}_l = \sum_{j=1}^l \frac{1}{\lambda_j} \frac{y'Q_j}{n} P_j.$$

We will refer to this method as truncation. A forecast of  $y$  is given by

$$\hat{y} = K\hat{\theta}_l = \sum_{j=1}^l \frac{y'Q_j}{n} Q_j. \quad (5.5)$$

Equation (5.5) is particularly interesting for its connection with forecasting using factors described in the next subsection.

## 5.2. Principal components and factor models

Let  $X_{it}$  be the observed data for the  $i^{\text{th}}$  cross-section unit at time  $t$ , with  $i = 1, 2, \dots, q$  and  $t = 1, 2, \dots, T$ . Consider the following dynamic factor model

$$X_{it} = \delta_i' F_t + e_{it} \quad (5.6)$$

where  $F_t$  is an  $l \times 1$  vector of unobserved common factors and  $\delta_i$  is the vector of factor loadings associated with  $F_t$ . The factor model is used in finance, where  $X_{it}$  represents the return of asset  $i$  at time  $t$ , see Ross (1976). Here we focus on the application of (5.6) to forecasting a single time series using a large number of predictors as in Stock and Watson (1998, 2002), Forni and Reichlin (1998), and Forni, Hallin, Lippi, and Reichlin (2000). Stock and Watson (1998, 2002) consider the forecasting equation

$$y_{t+1} = \beta' F_t + \epsilon_{t+1}$$

where  $y_t$  is either the inflation or the industrial production and  $X_t$  in (5.6) comprises 224 macroeconomic time-series. If the number of factors  $l$  is known, then  $\Lambda = (\delta_1, \delta_2, \dots, \delta_q)$  and  $F = (F_1, F_2, \dots, F_T)'$  can be estimated from

$$\min_{\Lambda, F} \frac{1}{qT} \sum_{i=1}^q \sum_{t=1}^T (X_{it} - \delta_i' F_t)^2 \quad (5.7)$$

under the restriction  $F'F/T = I$ . The  $F$  solution of (5.7) are the eigenvectors of  $XX'/T$  associated with the  $l$  largest eigenvalues. Hence  $F = [Q_1 | \dots | Q_l]$  where  $Q_j$  is  $j$ th eigenvector of  $XX'/T$ . Using the compact notation  $y = (y_2, \dots, y_{T+1})'$ , a forecast of  $y$  is given by

$$\begin{aligned}\hat{y} &= F\hat{\beta} \\ &= F(F'F)^{-1}F'y \\ &= F\frac{F'y}{T} \\ &= \sum_{j=1}^l \frac{Q_j'y}{T}Q_j.\end{aligned}$$

We recognize (5.5). It means that forecasting using a factor model (5.6) is equivalent to forecasting  $Y$  from (5.4) using a regularized solution based on the truncation. The only difference is that in the factor literature, it is assumed that there exists a fixed number of common factors, whereas in the truncation approach (5.5), the number of factors grows with the sample size. This last assumption may seem more natural when the number of explanatory variables,  $q$  goes to infinity.

An important issue in factor analysis is the estimation of the number of factors. Stock and Watson (1998) propose to minimize the MSE of the forecast. Bai and Ng (2002) propose various BIC and AIC criteria that permit us to consistently estimate the number of factors, even when  $T$  and  $q$  go to infinity.

### 5.3. Regression with an infinity of regressors

Consider the following model

$$Y = \int Z(\tau)\varphi(\tau)\Pi(d\tau) + U \tag{5.8}$$

where  $Z$  is uncorrelated with  $U$  and may include lags of  $Y$ ,  $\Pi$  is a known measure (possibly with finite or discrete support). One observes  $(y_i, z_i(\tau))_{i=1, \dots, n}$ .

#### First approach: Ridge regression

(5.8) can be rewritten as

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \int z_1(\tau)\varphi(\tau)\Pi(d\tau) \\ \vdots \\ \int z_n(\tau)\varphi(\tau)\Pi(d\tau) \end{pmatrix} + \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix}$$

or equivalently

$$y = K\varphi + u$$

where the operator  $K$  is defined in the following manner

$$K : L^2(\Pi) \rightarrow R^n$$

$$K\varphi = \begin{pmatrix} \int z_1(\tau) \varphi(\tau) \Pi(d\tau) \\ \vdots \\ \int z_n(\tau) \varphi(\tau) \Pi(d\tau) \end{pmatrix}.$$

As is usual in the regression, the error term  $u$  is omitted and we solve

$$K\varphi = y$$

using a regularized inverse

$$\varphi^\alpha = (K^*K + \alpha I)^{-1} K^*y. \quad (5.9)$$

As an exercise, we compute  $K^*$  and  $K^*K$ . To compute  $K^*$ , we solve

$$\langle K\varphi, \psi \rangle = \langle \varphi, K^*\psi \rangle$$

for  $\psi = (\psi_1, \dots, \psi_n)$  and we obtain

$$(K^*y)(\tau) = \frac{1}{n} \sum_{i=1}^n y_i z_i(\tau),$$

$$K^*K\varphi(\tau) = \int \frac{1}{n} \sum_{i=1}^n z_i(\tau) z_i(s) \varphi(s) \Pi(ds).$$

The properties of the estimator (5.9) are further discussed in Van Rooij, Ruymgaart and Van Zwet (2000).

### Second approach: moment conditions

Alternatively, (5.8) can be rewritten as

$$E[Y - \langle Z, \varphi \rangle | Z(\tau)] = 0 \text{ for all } \tau \text{ in the support of } \Pi$$

Replacing the conditional moments by unconditional moments, we have

$$E[YZ(\tau) - \langle Z, \varphi \rangle Z(\tau)] = 0 \iff$$

$$\int E[Z(\tau)Z(s)] \varphi(s) \Pi(ds) = E[YZ(\tau)] \iff$$

$$T\varphi = r. \quad (5.10)$$

The operator  $T$  can be estimated by  $\hat{T}_n$ , the operator with kernel  $\frac{1}{n} \sum_{i=1}^n z_i(\tau) z_i(s)$  and  $r_F$  can be estimated by  $\hat{r}_n(\tau) = \frac{1}{n} \sum_{i=1}^n y_i z_i(\tau)$ . Hence (5.10) becomes

$$\hat{T}_n \varphi = \hat{r}_n, \quad (5.11)$$

which is equal to

$$K^* K \varphi = K^* y.$$

If one preconditions (5.11) by applying the operator  $\hat{T}_n^*$ , one gets the solution

$$\hat{\varphi}_n = \left( \alpha I + \hat{T}_n^* \hat{T}_n \right) \hat{T}_n^* \hat{r}_n \quad (5.12)$$

which differs from the solution (5.9). When  $\alpha$  goes to zero at an appropriate rate of convergence (different in both cases), the solutions of (5.9) and (5.12) will be asymptotically equivalent. Actually, the preconditioning by an operator in the Tikhonov regularization has the purpose of constructing an operator which is positive self-adjoint. Because  $\hat{T}_n = K^* K$  is already positive self-adjoint, there is no reason to precondition here. Sometimes preconditioning more than necessary is aimed at facilitating the calculations (see Ruymgaart, 2001).

Using the results of Section 4, we can establish the asymptotic normality of  $\hat{\varphi}_n$  defined in (5.12).

Assuming that

A1 -  $u_i$  has mean zero and variance  $\sigma^2$  and is uncorrelated with  $z_i(\tau)$  for all  $\tau$

A2 -  $u_i z_i(\cdot)$  is an iid process of  $L^2(\Pi)$ .

A3 -  $E \|u_i z_i(\cdot)\|^2 < \infty$ .

we have

(i)  $\left\| \hat{T}_n^2 - T^2 \right\| = O\left(\frac{1}{\sqrt{n}}\right)$

(ii)  $\sqrt{n} \left( \hat{T}_n \hat{r}_n - \hat{T}_n^2 \varphi_0 \right) \Rightarrow \mathcal{N}(0, \Sigma)$  in  $L^2(\Pi)$ .

(i) is straightforward. (ii) follows from

$$\begin{aligned} \hat{r}_n - \hat{T}_n \varphi_0 &= \frac{1}{n} \sum_{i=1}^n y_i z_i(\tau) - \int \frac{1}{n} \sum_{i=1}^n z_i(\tau) z_i(s) \varphi_0(s) \Pi(ds) \\ &= \frac{1}{n} \sum_{i=1}^n u_i z_i(\tau). \end{aligned}$$

Here,  $a_n = \sqrt{n}$  and  $b_n = \sqrt{n}$ . Under Assumptions A1 to A3, we have

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n u_i z_i(\tau) \Rightarrow \mathcal{N}(0, \sigma^2 T)$$

in  $L^2(\Pi)$  by Theorem 2.46. Hence

$$\sqrt{n} \left( \hat{T}_n \hat{r}_n - \hat{T}_n^2 \varphi_0 \right) \Rightarrow \mathcal{N}(0, \sigma^2 T^3).$$

Let us rewrite Condition G in terms of the eigenvalues  $\lambda_j$  and eigenfunctions  $\phi_j$  of  $T$ :

$$\frac{\left\| (T^2 + \alpha_n I)^{-1} g \right\|^2}{\left\| T^{3/2} (T^2 + \alpha_n I)^{-1} g \right\|^2} = O(1)$$

$$\Leftrightarrow \frac{\sum_{j=1}^{\infty} \frac{\langle g, \phi_j \rangle^2}{(\lambda_j^2 + \alpha)^2}}{\sum_{j=1}^{\infty} \frac{\lambda_j^3 \langle g, \phi_j \rangle^2}{(\lambda_j^2 + \alpha)^2}} = O(1).$$

Obviously condition G will not be satisfied for all  $g$  in  $L^2(\Pi)$ .

By Proposition 4.3, assuming that  $\varphi_0 \in \Phi_\beta$ ,  $0 < \beta < 2$  and  $\sqrt{n}\alpha_n^{\beta/2} \rightarrow 0$ , we have for  $g$  conformable with Condition G,

$$\frac{\langle \sqrt{n}(\hat{\varphi}_n - \varphi_0), g \rangle}{\|T^{3/2}(T^2 + \alpha_n I)^{-1}g\|} \xrightarrow{d} \mathcal{N}(0, 1).$$

The asymptotic variance is given by

$$\|T^{-1/2}g\|^2 = \sum_{j=1}^{\infty} \frac{\langle g, \phi_j \rangle^2}{\lambda_j}.$$

Whenever it is finite, that is whenever  $g \in \mathcal{R}(T^{-1/2})$ ,  $\langle (\hat{\varphi}_n - \varphi_0), g \rangle$  converges at the parametric rate.

A related but different model from (5.8) is the Hilbertian autoregressive model:

$$X_t = \rho(X_{t-1}) + \varepsilon_t \tag{5.13}$$

where  $X_t$  and  $\varepsilon_t$  are random elements in a Hilbert space and  $\rho$  is a compact linear operator. The difference between (5.13) and (5.8) is that in (5.8),  $Y$  is a random variable and not an element of a Hilbert space. Bosq (2000) proposes an estimator of  $\rho$  and studies its properties.

Kargin and Onatski (2004) are interested in the best prediction of the interest rate curve. They model the forward interest rate  $X_t(\tau)$  at maturity  $\tau$  by (5.13) where  $\rho$  is a Hilbert-Schmidt integral operator:

$$(\rho f)(\tau) = \int_0^\infty \rho(\tau, s) f(s) ds. \tag{5.14}$$

The operator  $\rho$  is identified from the covariance and cross-covariance of the process  $X_t$ . Let  $\Gamma_{11}$  be the covariance operator of random curve  $X_t$  and  $\Gamma_{12}$  the cross-covariance operator of  $X_t$  and  $X_{t+1}$ . For convenience, the kernels of  $\Gamma_{11}$  and  $\Gamma_{12}$  are denoted using the same notation. Equations (5.13) and (5.14) yield

$$\begin{aligned} \Gamma_{12}(\tau_1, \tau_2) &= E[X_{t+1}(\tau_1) X_t(\tau_2)] \\ &= E\left[\int \rho(\tau_1, s) X_t(s) X_t(\tau_2) ds\right] \\ &= \int \rho(\tau_1, s) \Gamma_{11}(s, \tau_2) ds. \end{aligned}$$



Hence,

$$\Gamma_{12} = \rho\Gamma_{11}.$$

Solving this equation requires a regularization because  $\Gamma_{11}$  is compact. Interestingly, Kargin and Onatski (2004) show that the best prediction of the interest rate curve in finite sample is not necessarily provided by the eigenfunctions of  $\Gamma_{11}$  associated with the largest eigenvalues. It means that the spectral cut-off does not provide satisfactory predictions in small samples. They propose a better predictor of the interest rate curve.

Continuous-time models have been extensively studied by Le Breton (1994) and Bergstrom (1988).

#### 5.4. Deconvolution

Assume we observe iid realizations  $y_1, \dots, y_n$  of a random variable  $Y$  with unknown pdf  $h$ , where  $Y$  satisfies

$$Y = X + Z$$

where  $X$  and  $Z$  are independent random variables with pdf  $\varphi$  and  $g$  respectively. The aim is to get an estimator of  $\varphi$  assuming  $g$  is known. This problem consists in solving in  $\varphi$  the equation:

$$h(y) = \int g(y-x)\varphi(x)dx. \quad (5.15)$$

(5.15) is an integral equation of the first kind where the operator  $K$  defined by  $(K\varphi)(y) = \int g(y-x)\varphi(x)dx$  has a known kernel and need not be estimated. Recall that the compactness property depends on the space of reference. If we define as space of reference,  $L^2$  with respect to Lebesgue measure, then  $K$  is not a compact operator and hence has a continuous spectrum. However, for a suitable choice of the reference spaces,  $K$  becomes compact. The most common approach to solving (5.15) is to use a deconvolution kernel estimator, this method was pioneered by Carroll and Hall (1988) and Stefanski and Carroll (1990). It is essentially equivalent to inverting Equation (5.15) by means of the continuous spectrum of  $K$ , see Carroll, Van Rooij, and Ruymgaart (1991) and Subsection 5.4.2 below. In a related paper, Van Rooij and Ruymgaart (1991) propose a regularized inverse to a convolution problem of the type (5.15) where  $g$  has the circle for support. They invert the operator  $K$  using its continuous spectrum.

##### 5.4.1. A new estimator based on Tikhonov regularization

The approach of Carrasco and Florens (2002) consists in defining two spaces of reference,  $L^2_{\pi_X}(\mathbb{R})$  and  $L^2_{\pi_Y}(\mathbb{R})$  as

$$\begin{aligned} L^2_{\pi_X}(\mathbb{R}) &= \left\{ \phi(x) \text{ such that } \int \phi(x)^2 \pi_X(x) dx < \infty \right\}, \\ L^2_{\pi_Y}(\mathbb{R}) &= \left\{ \psi(y) \text{ such that } \int \psi(y)^2 \pi_Y(y) dy < \infty \right\}, \end{aligned}$$

so that  $K$  is a Hilbert-Schmidt operator from  $L^2_{\pi_X}(\mathbb{R})$  to  $L^2_{\pi_Y}(\mathbb{R})$ , that is the following condition is satisfied

$$\int \int \left( \frac{\pi_Y(y) g(y-x)}{\pi_Y(y) \pi_X(x)} \right)^2 \pi_Y(y) \pi_X(x) dx dy < \infty.$$

As a result  $K$  has a discrete spectrum for these spaces of reference. Let  $\{\lambda_j, \phi_j, \psi_j\}$  denote its singular value decomposition. Equation (5.15) can be approximated by a well-posed problem using Tikhonov regularization

$$(\alpha_n I + K^* K) \varphi_{\alpha_n} = K^* h.$$

Hence we have

$$\begin{aligned} \varphi_{\alpha_n}(x) &= \sum_{j=1}^{\infty} \frac{1}{\alpha_n + \lambda_j^2} \langle K^* h, \phi_j \rangle \phi_j(x) \\ &= \sum_{j=1}^{\infty} \frac{1}{\alpha_n + \lambda_j^2} \langle h, K \phi_j \rangle \phi_j(x) \\ &= \sum_{j=1}^{\infty} \frac{\lambda_j}{\alpha_n + \lambda_j^2} \langle h, \psi_j \rangle \phi_j(x) \\ &= \sum_{j=1}^{\infty} \frac{\lambda_j}{\alpha_n + \lambda_j^2} E[\psi_j(Y_i) \pi_Y(Y_i)] \phi_j(x). \end{aligned}$$

The estimator of  $\varphi$  is obtained by replacing the expectation by a sample mean:

$$\hat{\varphi}_n = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{\infty} \frac{\lambda_j}{\alpha_n + \lambda_j^2} \psi_j(y_i) \pi_Y(y_i) \phi_j(x).$$

Note that we avoided estimating  $h$  by a kernel estimator. In some cases,  $\psi_j$  and  $\phi_j$  are known. For instance, if  $Z \sim \mathcal{N}(0, \sigma^2)$ ,  $\pi_Y(y) = \phi(y/\tau)$  and  $\pi_X(x) = \phi(x/\sqrt{\tau^2 + \sigma^2})$  then  $\psi_j$  and  $\phi_j$  are Hermite polynomials associated with  $\lambda_j = \rho^j$ . When  $\psi_j$  and  $\phi_j$  are unknown, they can be estimated via simulations. Since one can do as many simulations as one wishes, the error due to the estimation of  $\psi_j$  and  $\phi_j$  can be considered negligible.

Using the results of Section 3, one can establish the rate of convergence of  $\|\hat{\varphi}_n - \varphi_0\|$ . Assume that  $\varphi_0 \in \Phi_\beta$ ,  $0 < \beta < 2$ , that is

$$\sum_{j=1}^{\infty} \frac{\langle \varphi, \phi_j \rangle^2}{\lambda_j^{2\beta}} < \infty.$$

We have  $\|\varphi_{\alpha_n} - \varphi_0\| = O(\alpha_n^{\beta/2})$  and  $\|\hat{\varphi}_n - \varphi_{\alpha_n}\| = O(1/(\alpha_n \sqrt{n}))$  as here  $b_n = \sqrt{n}$ . For an optimal choice of  $\alpha_n = Cn^{-1/(\beta+2)}$ ,  $\|\hat{\varphi}_n - \varphi_0\|^2$  is  $O(n^{-\beta/(\beta+2)})$ . The mean integrated

square error (MISE) defined as  $E \|\hat{\varphi}_n - \varphi_0\|^2$  has the same rate of convergence. Fan (1993) provides the optimal rate of convergence for a minimax criterion on a Lipschitz class of functions. The optimal rate of the MISE when the error term is normally distributed is only  $(\ln n)^{-2}$  if  $\varphi$  is twice differentiable. On the contrary, here we get an arithmetic rate of convergence. The condition  $\varphi_0 \in \Phi_\beta$  has the effect of reducing the class of admissible functions and hence improves the rate of convergence. Which type of restriction does  $\varphi_0 \in \Phi_\beta$  impose? In Carrasco and Florens (2002), it is shown that  $\varphi_0 \in \Phi_1$  is satisfied if

$$\int \left| \frac{\psi_{\varphi_0}(t)}{\psi_g(t)} \right| dt < \infty \quad (5.16)$$

where  $\psi_{\varphi_0}$  and  $\psi_g$  denote the characteristic functions of  $\varphi_0$  and  $g$  respectively. This condition can be interpreted as the noise is “smaller” than the signal. Consider for example the case where  $\varphi_0$  and  $g$  are normal. Condition (5.16) is equivalent to the fact that the variance of  $g$  is smaller than that of  $\varphi_0$ . Note that the condition  $\varphi_0 \in \Phi_1$  relates  $\varphi_0$  and  $g$  while one usually imposes restrictions on  $\varphi_0$  independently of those on  $g$ .

#### 5.4.2. Comparison with the deconvolution kernel estimator

Let  $L^2_\lambda(\mathbb{R})$  be the space of square-integrable functions with respect to Lebesgue measure on  $\mathbb{R}$ . Let  $F$  denote the Fourier transform operator from  $L^2_\lambda(\mathbb{R})$  into  $L^2_\lambda(\mathbb{R})$  defined by

$$(Fq)(s) = \frac{1}{\sqrt{2\pi}} \int e^{isx} q(x) dx.$$

$F$  satisfies that  $F^* = F^{-1}$ . We see that

$$F(g * f) = \phi_g Ff$$

so that  $K$  admits the following spectral decomposition (see Carroll, van Rooij and Ruymgaart, 1991, Theorem 3.1.):

$$K = F^{-1} M_{\phi_g} F$$

where  $M_\rho$  is the multiplication operator  $M_\rho \varphi = \rho \varphi$ .

$$K^* K = F^{-1} M_{|\phi_g|^2} F.$$

We want to solve in  $f$  the equation:

$$K^* K f = K^* h.$$

Let us denote

$$q(x) = (K^* h)(x) = \int g(y-x) h(y) dy.$$

Then,

$$\hat{q}(x) = \frac{1}{n} \sum_{i=1}^n g(y_i - x)$$

is a  $\sqrt{n}$ -consistent estimator of  $q$ .

Using the spectral cut-off regularized inverse of  $K^*K$ , we get

$$\hat{f} = F^{-1} M_{\frac{1}{|\phi_g|^2} \{|\phi_g| > \alpha\}} F \hat{q}$$

Using the change of variables  $u = y_i - x$ , we have

$$\begin{aligned} (F\hat{q})(t) &= \frac{1}{n} \sum_{i=1}^n \int e^{itx} g(y_i - x) dx \\ &= \frac{1}{n} \sum_{i=1}^n \int e^{it(y_i - u)} g(u) du \\ &= \frac{1}{n} \sum_{i=1}^n \overline{\phi_g(t)} e^{ity_i}. \end{aligned}$$

$$\begin{aligned} \hat{f}(x) &= \frac{1}{2\pi} \int e^{-itx} I \{|\phi_g(t)| > \alpha\} \frac{1}{|\phi_g(t)|^2} (F\hat{q})(t) dt \\ &= \frac{1}{2\pi} \frac{1}{n} \sum_{i=1}^n \int e^{-it(y_i - x)} I \{|\phi_g(t)| > \alpha\} \frac{1}{\phi_g(t)} dt. \end{aligned}$$

Assuming that  $\phi_g > 0$  and strictly decreasing as  $|t|$  goes to infinity, we have  $I \{|\phi_g(t)| > \alpha\} = I \{-A \leq t \leq A\}$  for some  $A > 0$  so that

$$\hat{f}(x) = \frac{1}{2\pi} \frac{1}{n} \sum_{i=1}^n \int_{-A}^A \frac{e^{-it(y_i - x)}}{\phi_g(t)} dt.$$

Now compare this expression with the kernel estimator (see e.g. Stefanski and Carroll, 1990). For a smoothing parameter  $c$  and a kernel  $\omega$ , the kernel estimator is given by

$$\hat{f}_k(x) = \frac{1}{nc} \sum_{i=1}^n \frac{1}{2\pi} \int \frac{\phi_\omega(u)}{\phi_g(u/c)} e^{iu(y_i - x)/c} du. \quad (5.17)$$

Hence  $\hat{f}$  coincides with the kernel estimator when  $\phi_\omega(u) = I_{[-1,1]}(u)$ . This is the sinc kernel corresponding to  $\omega(x) = \sin c(x) = \sin(x)/x$ . This suggests that the kernel estimator is obtained by inverting an operator that has a continuous spectrum. Because this spectrum is given by the characteristic function of  $g$ , the speed of convergence of the

estimator depends on the behavior of  $\phi_g$  in the tails. For a formal exposition, see Carroll et al (1991, Example 3.1.). They assume in particular that the function to estimate is  $p$  differentiable and they obtain a rate of convergence (as a function of  $p$ ) that is of the same order as the rate of the kernel estimator.

By using the Tikhonov regularization instead of the spectral cut-off, we obtain

$$\hat{f}(y) = \frac{1}{n} \sum_{i=1}^n \int \frac{\overline{\phi_g(t)}}{|\phi_g(t)|^2 + \alpha} e^{-itx_i} e^{ity} dt.$$

We apply a change of variable  $u = -t$ ,

$$\hat{f}(y) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2\pi} \int \frac{\phi_g(u)}{|\phi_g(u)|^2 + \alpha} e^{iu(x_i - y)} du. \quad (5.18)$$

The formulas (5.18) and (5.17) differ only by the way the smoothing is applied.

## 5.5. Instrumental variables

This example is mainly based on Darolles, Florens and Renault (2002).

An economic relationship between a response variable  $Y$  and a vector  $Z$  of explanatory variables is often represented by an equation:

$$Y = \varphi(Z) + U, \quad (5.19)$$

where the function  $\varphi(\cdot)$  defines the parameter of interest while  $U$  is an error term. The relationship (5.19) does not characterize the function  $\varphi$  if the residual term is not constrained. This difficulty is solved if it is assumed that  $E[U | Z] = 0$ , or if equivalently  $\varphi(Z) = E[Y | Z]$ . However in numerous structural econometric models, the conditional expectation function is not the parameter of interest. The structural parameter is a relation between  $Y$  and  $Z$  where some of the  $Z$  components are endogenous. This is the case in various situations: simultaneous equations, error-in-variables models, and treatment models with endogenous selection etc.

The first issue is to add assumptions to Equation (5.19) in order to characterize  $\varphi$ . Two general strategies exist in the literature, at least for linear models. The first one consists in introducing some hypotheses on the joint distribution of  $U$  and  $Z$  (for example on the variance matrix). The second one consists in increasing the vector of observables from  $(Y, Z)$  to  $(Y, Z, W)$ , where  $W$  is a vector of instrumental variables. The first approach was essentially followed in the error-in-variables models and some similarities exist with the instrumental variables model (see e.g. Malinvaud (1970, ch. 9), Florens, Mouchart, Richard (1974) or Florens, Mouchart, Richard (1987) for the linear case). Instrumental variable analysis as a solution to an endogeneity problem was proposed by Reiersol (1941, 1945), and extended by Theil (1953), Basmann (1957), and Sargan (1958).

However, even in the instrumental variables framework, the definition of the functional parameter of interest remains ambiguous in the general nonlinear case. Three possible

definitions of  $\varphi$  have been proposed (see Florens, Heckman, Meghir and Vytlačil (2002) for a general comparison between these three concepts and their extensions to more general treatment models).

*i)* The first one replaces  $E[U | Z] = 0$  by  $E[U | W] = 0$ , or equivalently it defines  $\varphi$  as the solution of

$$E[Y - \varphi(Z) | W] = 0. \quad (5.20)$$

This definition was the foundation of the analysis of simultaneity in linear models or parametric nonlinear models (see Amemiya (1974)), but its extension to the nonparametric case raises new difficulties. The focus of this subsection is to show how to address this issue in the framework of ill-posed inverse problems. A first attempt was undertaken by Newey and Powell (2003), who prove consistency of a series estimator of  $\varphi$  in Equation (5.20). Florens (2003) and Blundell and Powell (2003) consider various nonparametric methods for estimating a nonlinear regression with endogenous regressors. Darolles, Florens, and Renault (2002) prove both the consistency and the asymptotic distribution of a kernel estimator of  $\varphi$ . Hall and Horowitz (2004) give the optimal rate of convergence of the kernel estimator under conditions which differ from those of Darolles, Florens, and Renault (2002). Finally, Blundell, Chen, and Kristensen (2003) propose a sieves estimator of the Engel curve.

*ii)* A second approach is now called *control function approach* and was systematized by Newey, Powell, and Vella (1999). This technique was previously developed in specific models (e.g. Mills ratio correction in some selection models for example). The starting point is to compute  $E[Y | Z, W]$  which satisfies:

$$E[Y | Z, W] = \varphi(Z) + h(Z, W), \quad (5.21)$$

where  $h(Z, W) = E[U | Z, W]$ . Equation (5.21) does not characterize  $\varphi$ . However we can assume that there exists a function  $V$  (the *control function*) of  $(Z, W)$  (typically  $Z - E[Z | W]$ ), which captures all the endogeneity of  $Z$  in the sense that  $E[U | W, V] = E[U | V] = \tilde{h}(V)$ . This implies that (5.21) may be rewritten as

$$E[Y | Z, W] = \varphi(Z) + \tilde{h}(V), \quad (5.22)$$

and under some conditions,  $\varphi$  may be identified from (5.22) up to an additive constant term. This model is an additive model where the  $V$  are not observed but are estimated.

*iii)* A third definition follows from the literature on treatment models (see e.g. Imbens, Angrist (1994), Heckman, Ichimura, Smith, Todd (1998) and Heckman, Vytlačil (2000)). We extremely simplify this analysis by considering  $Z$  and  $W$  as scalars. *Local instrument* is defined by  $\frac{\partial E[Y|W]}{\partial W} / \frac{\partial E[Z|W]}{\partial W}$ , and the function of interest  $\varphi$  is assumed to be characterized by the relation:

$$\frac{\frac{\partial E[Y|W]}{\partial W}}{\frac{\partial E[Z|W]}{\partial W}} = E \left[ \frac{\partial \varphi}{\partial Z} \mid W \right]. \quad (5.23)$$

Let us summarize the arguments, which justify Equation (5.23). Equation (5.19) is extended to a non separable model

$$Y = \varphi(Z) + Z\varepsilon + U \quad (5.24)$$

where  $\varepsilon$  and  $U$  are two random noises.

First, we assume that

$$E(U|W) = E(\varepsilon|W) = 0$$

This assumption extends the instrumental variable assumption but is not sufficient to identify the parameter of interest  $\varphi$ . From (5.24) we get:

$$E(Y|W = w) = \int [\varphi(z) + zr(z, w)] f_Z(z|w) dz$$

where  $f_Z(\cdot|\cdot)$  denote the conditional density of  $Z$  given  $W$  and  $r(z, w) = E(\varepsilon|Z = z, W = w)$ . Then

$$\begin{aligned} \frac{\partial}{\partial w} E(Y|W = w) &= \int \varphi(z) \frac{\partial}{\partial w} f_Z(z|w) dz + \int z \frac{\partial}{\partial w} r(z, w) f_Z(z|w) dz \\ &+ \int zr(z, w) \frac{\partial}{\partial w} f_Z(z|w) dz. \end{aligned}$$

We assume that the order of integration and derivative may commute (in particular the boundary of the distribution of  $Z$  given  $W = w$  does not depends on  $w$ ).

Second, we introduce the assumption that  $V = Z - E(Z|W)$  is independent of  $W$ . In terms of density, this assumption implies that  $f_Z(z|w) = \tilde{f}(z - m(w))$  where  $m(w) = E(Z|W = w)$  and  $\tilde{f}$  is the density of  $v$ . Then:

$$\begin{aligned} \frac{\partial}{\partial w} E(Y|W = w) &= -\frac{\partial m(w)}{\partial w} \int \varphi(z) \frac{\partial}{\partial z} f_Z(z|w) dz \\ &+ \int z \frac{\partial}{\partial w} r(z, w) f_Z(z|w) dz \\ &- \frac{\partial m(w)}{\partial w} \int zr(z, w) \frac{\partial}{\partial z} f_Z(z|w) dz \end{aligned}$$

An integration by parts of the first and the third integrals gives

$$\begin{aligned} \frac{\partial}{\partial w} E(Y|W = w) &= \frac{\partial m(w)}{\partial w} \int \frac{\partial}{\partial z} \varphi(z) f_Z(z|w) dz \\ &+ \int z \left( \frac{\partial r}{\partial w} + \frac{\partial m}{\partial w} \frac{\partial r}{\partial z} \right) f_Z(z|w) dz \\ &+ \frac{\partial m(w)}{\partial w} \int r(z, w) f_Z(z|w) dz \end{aligned}$$

The last integral is zero under  $E(\varepsilon|w) = 0$ . Finally, we need to assume that the second integral is zero. This is true in particular if there exists  $\tilde{r}$  such that  $r(z, w) = \tilde{r}(z - m(w))$ .

Hence, Equation (5.23) is verified.

These three concepts are identical in the linear normal case but differ in general. We concentrate our presentation in this chapter on the pure instrumental variable cases defined by equation (5.20).

For a general approach of Equation (5.20) in terms of inverse problems, we introduce the following notation:

$$K : L_F^2(Z) \rightarrow L_F^2(W) \quad \varphi \rightarrow K\varphi = E[\varphi(Z) | W],$$

$$K^* : L_F^2(W) \rightarrow L_F^2(Z) \quad \psi \rightarrow K^*\psi = E[\psi(W) | Z].$$

All these spaces are defined relatively to the true (unknown) DGP. The two linear operators  $K$  and  $K^*$  satisfy:

$$\langle \varphi(Z), \psi(W) \rangle = E[\varphi(Z) \psi(W)] = \langle K\varphi(W), \psi(W) \rangle_{L_F^2(W)} = \langle \varphi(Z), K^*\psi(Z) \rangle_{L_F^2(Z)}.$$

Therefore,  $K^*$  is the adjoint operator of  $K$ , and reciprocally. Using these notations, the unknown instrumental regression  $\varphi$  corresponds to any solution of the functional equation:

$$A(\varphi, F) = K\varphi - r = 0, \tag{5.25}$$

where  $r(W) = E[Y | W]$ .

In order to illustrate this construction and the central role played by the adjoint operator  $K^*$ , we first consider the example where  $Z$  is discrete, namely  $Z$  is binary. This model is considered by Das (2005) and Florens and Malavolti (2002). In that case, a function  $\varphi(Z)$  is characterized by two numbers  $\varphi(0)$  and  $\varphi(1)$  and  $L_Z^2$  is isomorphic to  $\mathbb{R}^2$ . Equation (5.20) becomes

$$\varphi(0) \text{Prob}(Z = 0 | W = w) + \varphi(1) \text{Prob}(Z = 1 | W = w) = E(Y | W = w).$$

The instruments  $W$  need to take at least two values in order to identify  $\varphi(0)$  and  $\varphi(1)$  from this equation. In general,  $\varphi$  is overidentified and overidentification is solved by replacing (5.25) by

$$K^*K\varphi = K^*r$$

or, in the binary case, by

$$\varphi(0) E(\text{Prob}(Z = 0 | W) | Z) + \varphi(1) E(\text{Prob}(Z = 1 | W) | Z) = E(E(Y | W) | Z).$$

In the latter case, we get two equations which in general have a unique solution.

This model can be extended by considering  $Z = (Z_1, Z_2)$  where  $Z_1$  is discrete ( $Z_1 \in \{0, 1\}$ ) and  $Z_2$  is exogenous (i.e.  $W = (W_1, Z_2)$ ). In this extended binary model,  $\varphi$  is characterized by two functions  $\varphi(0, z_2)$  and  $\varphi(1, z_2)$ , the solutions of

$$\begin{aligned} \varphi(0, z_2) E(\text{Prob}(Z_1 = 0 | W) | Z_1 = z_1, Z_2 = z_2) + \varphi(1, z_2) E(\text{Prob}(Z_1 = 1 | W) | Z_1 = z_1, Z_2 = z_2) \\ = E(E(Y | W) | Z_1 = z_1, Z_2 = z_2), \quad \text{for } z_1 = 0, 1. \end{aligned}$$



The properties of the estimator based on the previous equation are considered in Florens and Malavolti (2002). In this case, no regularization is needed because  $K^*K$  has a continuous inverse (since the dimension is finite in the pure binary case and  $K^*K$  is not compact in the extended binary model).

We can also illustrate our approach in the case when the Hilbert spaces are not necessarily  $L^2$  spaces. Consider the following semiparametric case. The function  $\varphi$  is constrained to be an element of

$$\mathcal{X} = \left\{ \varphi \text{ such that } \varphi = \sum_{l=1}^L \beta_l \varepsilon_l \right\}$$

where  $(\varepsilon_l)_{l=1, \dots, L}$  is a vector of fixed functions in  $L_F^2(Z)$ . Then  $\mathcal{X}$  is a finite dimensional Hilbert space. However, we keep the space  $\mathcal{E}$  equal to  $L_F^2(W)$ . The model is then partially parametric but the relation between  $Z$  and  $W$  is treated nonparametrically. In this case, it can easily be shown that  $K^*$  transforms any function  $\psi$  of  $L_F^2(W)$  into a function of  $\mathcal{X}$ , which is its best approximation in the  $L^2$  sense (see Example 2.4. in Section 2). Indeed:

$$\text{If } \psi \in L_F^2(W), \forall j \in \{1, \dots, L\}$$

$$E(\varepsilon_j \psi) = \langle K \varepsilon_j, \psi \rangle = \langle \varepsilon_j, K^* \psi \rangle.$$

Moreover,  $K^* \psi \in \mathcal{X} \implies K^* \psi = \sum_{l=1}^L \alpha_l \varepsilon_l$ , therefore

$$\begin{aligned} \left\langle \varepsilon_j, \sum_{l=1}^L \alpha_l \varepsilon_l \right\rangle &= E(\psi \varepsilon_j) \\ \Leftrightarrow \sum_{l=1}^L \alpha_l E(\varepsilon_j \varepsilon_l) &= E(\psi \varepsilon_j). \end{aligned}$$

The function  $\varphi$  defined as the solution of  $K\varphi = r$  is in general overidentified but the equation  $K^*K\varphi = K^*r$  always has a unique solution. The finite dimension of  $\mathcal{X}$  implies that  $(K^*K)^{-1}$  is a finite dimensional linear operator and is then continuous. No regularization is required.

Now we introduce an assumption which is only a regularity condition when  $Z$  and  $W$  have no element in common. However, this assumption cannot be satisfied if there are some elements in common between  $Z$  and  $W$ . Extensions to this latter case are discussed in Darolles, Florens and Renault (2002), see also Example 2.5. in Section 2.

**Assumption A.1:** *The joint distribution of  $(Z, W)$  is dominated by the product of its marginal distributions, and its density is square integrable w.r.t. the product of margins.*

Assumption A.1 ensures that  $K$  and  $K^*$  are Hilbert Schmidt operators, and is a sufficient condition for the compactness of  $K$ ,  $K^*$ ,  $KK^*$  and  $K^*K$  (see Lancaster (1968), Darolles, Florens, Renault (2002)) and Theorem 2.34.

Under Assumption A1, the instrumental regression  $\varphi$  is identifiable if and only if 0 is not an eigenvalue of  $K^*K$ . Then, for the sake of expositional simplicity, we focus on the i.i.d. context:

**Assumption A.2:** *The data  $(y_i, z_i, w_i)$   $i = 1, \dots, n$ , are i.i.d samples of  $(Y, Z, W)$ .*

We estimate the joint distribution  $F$  of  $(Y, Z, W)$  using a kernel smoothing of the empirical distribution. In the applications, the bandwidths differ, but they all have the same speed represented by the notation  $c_n$ .

For economic applications, one may be interested either by the unknown function  $\varphi(Z)$  itself, or only by its moments, including covariances with some known functions. These moments may for instance be useful for testing economic statements about scale economies, elasticities of substitutions, and so on.

For such tests, one will only need the empirical counterparts of these moments and their asymptotic probability distribution. An important advantage of the instrumental variable approach is that it permits us to estimate the covariance between  $\varphi(Z)$  and  $g(Z)$  for a large class of functions. Actually, the identification assumption amounts to ensure that the range  $\mathcal{R}(K^*)$  is dense in  $L_F^2(Z)$  and for any  $g$  in this range:

$$\exists \psi \in L_F^2(W), g(Z) = E[\psi(W) | Z],$$

and then  $Cov[\varphi(Z), g(Z)] = Cov[\varphi(Z), E[\psi(W) | Z]] = Cov[\varphi(Z), \psi(W)] = Cov[E[\varphi(Z) | W], \psi(W)] = Cov[Y, \psi(W)]$ , can be estimated with standard parametric techniques. For instance, if  $E[g(Z)] = 0$ , the empirical counterpart of  $Cov[Y, \psi(W)]$ , i.e.:

$$\frac{1}{n} \sum_{i=1}^n Y_i \psi(W_i),$$

is a root- $n$  consistent estimator of  $Cov[\varphi(Z), g(Z)]$ , and:

$$\sqrt{n} \left[ \frac{1}{n} \sum_{i=1}^n Y_i \psi(W_i) - Cov[\varphi(Z), g(Z)] \right] \xrightarrow{d} \mathcal{N}(0, Var[Y\psi(W)]),$$

where  $Var[Y\psi(W)]$  will also be estimated by its sample counterpart. However in practice this analysis has very limited interest because even if  $g$  is given,  $\psi$  is not known and must be estimated by solving the integral equation  $g(Z) = E[\psi(W) | Z]$ , where the conditional distribution of  $W$  given  $Z$  is also estimated.

Therefore, the real problem of interest is to estimate  $Cov[\varphi(Z), g(Z)]$ , or  $\langle \varphi, g \rangle$  by replacing  $\varphi$  by an estimator. This estimator will be constructed by solving a regularized version of the empirical counterpart of (5.25) where  $K$  and  $r$  are replaced by their estimators. In the case of kernel smoothing, the necessity of regularization appears obviously. Using the notation of 2.5.3, the equation

$$\hat{K}_n \varphi = \hat{r}_n$$

becomes

$$\frac{\sum_{i=1}^n \varphi(z_i) \omega\left(\frac{w-w_i}{c_n}\right)}{\sum_{i=1}^n \omega\left(\frac{w-w_i}{c_n}\right)} = \frac{\sum_{i=1}^n y_i \omega\left(\frac{w-w_i}{c_n}\right)}{\sum_{i=1}^n \omega\left(\frac{w-w_i}{c_n}\right)}.$$

The function  $\varphi$  can not be obtained from this equation except for the values  $\varphi(z_i)$  equal to  $y_i$ . This solution does not constitute a consistent estimate. The regularized Tikhonov solution is the solution of:

$$\alpha_n \varphi(z) + \frac{\sum_{j=1}^n \omega\left(\frac{z-z_j}{c_n}\right) \frac{\sum_{i=1}^n \varphi(z_i) \omega\left(\frac{w_j-w_i}{c_n}\right)}{\sum_{i=1}^n \omega\left(\frac{w_j-w_i}{c_n}\right)}}{\sum_{j=1}^n \omega\left(\frac{z-z_j}{c_n}\right)} = \frac{\sum_{j=1}^n \omega\left(\frac{z-z_j}{c_n}\right) \frac{\sum_{i=1}^n y_i \omega\left(\frac{w_j-w_i}{c_n}\right)}{\sum_{i=1}^n \omega\left(\frac{w_j-w_i}{c_n}\right)}}{\sum_{j=1}^n \omega\left(\frac{z-z_j}{c_n}\right)}.$$

This functional equation may be solved in two steps. First, the  $z$  variable is fixed to the values  $z_i$  and the system becomes an  $n \times n$  linear system, which can be solved in order to obtain the  $\varphi(z_i)$ . Second, the previous expression gives a value of  $\varphi(z)$  for any value of  $z$ .

If  $n$  is very large, this inversion method may be difficult to apply and may be replaced by a Landweber Fridman resolution (see Section 3). A first expression of  $\varphi(z)$  may be for instance the estimated conditional expectation  $E(E(Y|W)|Z)$  and this estimator will be modified a finite number of times by the formula

$$\hat{\varphi}_{l,n} = \left(I - c\hat{K}_n^* \hat{K}_n\right) \hat{\varphi}_{l-1,n} + c\hat{K}_n^* \hat{r}_n.$$

To simplify our analysis, we impose a relatively strong assumption:

**Assumption A.3:** The error term is homoskedastic, that is:

$$\text{Var}(U|W) = \sigma^2.$$

In order to check the asymptotic properties of the estimator of  $\varphi$ , it is necessary to study to properties of the estimators of  $K$  and of  $r$ . Under regularity conditions such as the compactness of the joint distribution support and the smoothness of the density (see Darolles et al. (2002)), the estimation by boundary kernels gives the following results:

i)  $\left\| \hat{K}_n^* \hat{K}_n - K^* K \right\|^2 \sim O\left(\frac{1}{n(c_n)^p} + (c_n)^{2\rho}\right)$  where  $\rho$  is the order of the kernel and  $p$  the dimension of  $Z$ .

ii)  $\left\| \hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi \right\|^2 \sim O\left(\frac{1}{n} + (c_n)^{2\rho}\right)$

iii) A suitable choice of  $c_n$  implies

$$\sqrt{n} \left( \hat{K}_n^* \hat{r}_n - \hat{K}_n^* \hat{K}_n \varphi \right) \Longrightarrow N(0, \sigma^2 K^* K)$$

This convergence is a weak convergence in  $L_F^2(Z)$  (see Section 2.4).

Using results developed in Section 4 and in Darolles et al. (2002) it can be deduced that:

**a)** If  $\alpha_n \rightarrow 0$ ,  $\frac{c_n^{2\rho}}{\alpha_n^2} \rightarrow 0$ ,  $\frac{1}{\alpha_n^2 n c_n^\rho} \sim O(1)$  the regularized estimator  $\hat{\varphi}_n$  converge in probability to  $\varphi$  in  $L^2$  norm.

**b)** If  $\varphi \in \Phi_\beta$  ( $0 < \beta \leq 2$ ), the optimal choices of  $\alpha_n$  and  $c_n$  are:

$$\begin{aligned} \alpha_n &= k_1 n^{-\frac{1}{2\beta}} \\ c_n &= k_2 n^{-\frac{1}{2\rho}} \end{aligned}$$

and, if  $\rho$  is chosen such that  $\frac{\rho}{2} \leq \frac{\beta}{2+\beta}$ , we obtain the following bound for the rate of convergence

$$\|\hat{\varphi}_n - \varphi\| \sim O\left(n^{-\frac{\beta}{2+\beta}}\right)$$

**c)** Let us assume that  $\alpha$  is kept constant. In that case, the linear operators  $(\alpha I + K_n^* K_n)^{-1}$  and  $(\alpha I + K^* K)^{-1}$  are bounded, and using a functional version of the Slutsky theorem (see Chen and White (1992), and Section 2.4), it is immediately checked that:

$$\sqrt{n}(\hat{\varphi}_n - \varphi - b_n^\alpha) \Longrightarrow \mathcal{N}(0, \Omega), \quad (5.26)$$

where

$$b_n^\alpha = \alpha [(\alpha I + K_n^* K_n)^{-1} - (\alpha I + K^* K)^{-1}] \varphi,$$

and

$$\Omega = \sigma^2 (\alpha I + K^* K)^{-1} K^* K (\alpha I + K^* K)^{-1}.$$

Some comments may illustrate this first result:

i) The convergence obtained in (5.26) is still a functional distributional convergence in the Hilbert space  $L_F^2(Z)$ , which in particular implies the convergence of inner product  $\sqrt{n} \langle \hat{\varphi}_n - \varphi - b_n^\alpha, g \rangle$  to univariate normal distribution  $\mathcal{N}(0, \langle g, \Omega g \rangle)$ .

ii) The convergence of  $\hat{\varphi}_n$  involves two bias terms. The first bias is  $\varphi_\alpha - \varphi$ . This term is due to the regularization and does not decrease if  $\alpha$  is constant. The second one,  $\hat{\varphi}_n - \varphi_\alpha$  follows from the estimation error of  $K$ . This bias decreases to zero when  $n$  increases, but at a lower speed than  $\sqrt{n}$ .

iii) The asymptotic variance in (5.26) can be seen as the generalization of the two stage least squares asymptotic variance. An intuitive (but not correct) interpretation of this

result could be the following. If  $\alpha$  is small, the asymptotic variance is approximately  $\sigma^2(K^*K)^{-1}$ , which is the functional extension of  $\sigma^2(E(ZW')E(WW')^{-1}E(WZ'))^{-1}$ .

**d)** Let us now consider the case where  $\alpha \rightarrow 0$ . For any  $\delta \in \Phi_\beta$  ( $\beta \geq 1$ ), if  $\alpha_n$  is optimal ( $= k_1 n^{-\frac{1}{2\beta}}$ ) and if  $c_n = k_2 n^{-\left(\frac{1}{2\rho} + \varepsilon\right)}$  ( $\varepsilon > 0$ ), we have

$$\sqrt{\nu_n(\delta)} \langle \hat{\varphi}_n - \varphi, \delta \rangle - B_n \implies N(0, \sigma^2),$$

where the speed of convergence is equal to

$$\nu_n(\delta) = \frac{n}{\|K(\alpha_n I + K^*K)^{-1}\delta\|^2} \geq O\left(n^{\frac{2\beta}{2+\beta}}\right),$$

and the bias  $B_n$  is equal to  $\sqrt{\nu_n(\delta)} \langle \varphi_\alpha - \varphi, \delta \rangle$ , which in general does not vanish. If  $\delta = 1$  for example, this bias is  $O(n\alpha_n^2)$  and diverges.

The notion of  $\Phi_\beta$  permits us to rigorously define the concept of weak or strong instruments. Indeed, if  $\lambda_j$  are not zero for any  $j$ , the function  $\varphi$  is identified by Equation (5.25) and  $\hat{\varphi}_n$  is a consistent estimator. A bound for the speed of convergence of  $\hat{\varphi}_n$  is provided under the restriction that  $\varphi$  belongs to a space  $\Phi_\beta$  with  $\beta > 0$ . The condition  $\varphi \in \Phi_\beta$  means that the rate of decline of the Fourier coefficients of  $\varphi$  in the basis of  $\phi_j$  is faster than the rate of decline of the  $\lambda_j^\beta$  (which measures the dependence). In order to have asymptotic normality we need to assume that  $\beta \geq 1$ . In that case, if  $\varphi \in \Phi_\beta$ , we have asymptotic normality of inner products  $\langle \hat{\varphi}_n - \varphi, \delta \rangle$  in the vector space  $\Phi_\beta$ . Then, it is natural to say that  $W$  is a strong instrument for  $\varphi$  if  $\varphi$  is an element of a  $\Phi_\beta$  with  $\beta \geq 1$ . This may have two equivalent interpretations. Given  $Z$  and  $W$ , the set of instrumental regressions for which  $W$  is a strong instrument is  $\Phi_1$  or given  $Z$  and  $\varphi$ , any set of instruments is strong if  $\varphi$  is an element of the set  $\Phi_1$  defined using these instruments.

We may complete this short presentation with two final remarks. First, the optimal choice of  $c_n$  and  $\alpha_n$  implies that the speed of convergence and the asymptotic distribution are not affected by the fact that  $K$  is not known and is estimated. The accuracy of the estimation is governed by the estimation of the right hand side term  $K^*r$ . Secondly, the usual ‘‘curse of dimensionality’’ of nonparametric estimation appears in a complex way. The dimension of  $Z$  appears in many places but the dimension of  $W$  is less explicit. The value and the rate of decline of the  $\lambda_j$  depend on the dimension of  $W$ : Given  $Z$ , the reduction of the number of instruments implies a faster rate of decay of  $\lambda_j$  to zero and a slower rate of convergence of the estimator.

## 6. Reproducing kernel and GMM in Hilbert spaces

### 6.1. Reproducing kernel

Models based on reproducing kernels are the foundation for penalized likelihood estimation and splines (see e.g. Berlinet and Thomas-Agnan, 2004). However, it has been little used

in econometrics so far. The theory of reproducing kernels becomes very useful when the econometrician has an infinite number of moment conditions and wants to exploit all of them in an efficient way. For illustration, let  $\theta \in \mathbb{R}$  be the parameter of interest and consider an  $L \times 1$ -vector  $h$  that gives  $L$  moment conditions satisfying  $E^{\theta_0}(h(\theta)) = 0 \Leftrightarrow \theta = \theta_0$ . Let  $h_n(\theta)$  be the sample estimate of  $E^{\theta_0}(h(\theta))$ . The (optimal) generalized method of moments (GMM) estimator of  $\theta$  is the minimizer of  $h_n(\theta)' \Sigma^{-1} h_n(\theta)$  where  $\Sigma$  is the covariance matrix of  $h$ .  $h_n(\theta)' \Sigma^{-1} h_n(\theta)$  can be rewritten as  $\|\Sigma^{-1/2} h_n(\theta)\|^2$  and coincides with the norm of  $h_n(\theta)$  in a particular space called the reproducing kernel Hilbert space (RKHS). When  $h$  is finite dimensional, the computation of the GMM objective function does not raise any particular difficulty, however when  $h$  is infinite dimensional (for instance is a function) then the theory of RKHS becomes very handy. A second motivation for the introduction of the RKHS of a self-adjoint operator  $K$  is the following. Let  $T$  be such that  $K = TT^*$  then the RKHS of  $K$  corresponds to the 1-regularity space of  $T$  (denoted  $\Phi_1$  in Section 3.1).

### 6.1.1. Definitions and basic properties of RKHS

This section presents the theory of reproducing kernels, as described in Aronszajn (1950) and Parzen (1959, 1970). Let  $L_C^2(\pi) = \{\varphi : I \subset \mathbb{R}^L \rightarrow \mathbf{C} : \int_I |\varphi(s)|^2 \pi(s) ds < \infty\}$  where  $\pi$  is a pdf ( $\pi$  may have a discrete or continuous support) and denote  $\|\cdot\|$  and  $\langle \cdot, \cdot \rangle$  the norm and inner product on  $L_C^2(\pi)$ .

**Definition 6.1.** A space  $\mathcal{H}(K)$  of complex-valued functions defined on a set  $I \subset \mathbb{R}^L$  is said to be a reproducing kernel Hilbert space  $\mathcal{H}(K)$  associated with the integral operator  $K : L_C^2(\pi) \rightarrow L_C^2(\pi)$  with kernel  $k(t, s)$  if the three following conditions hold

- (i) it is a Hilbert space (with inner product denoted  $\langle \cdot, \cdot \rangle_K$ ),
  - (ii) for every  $s \in I$ ,  $k(t, s)$  as a function of  $t$  belongs to  $\mathcal{H}(K)$ ,
  - (iii) (reproducing property) for every  $s \in I$  and  $\varphi \in \mathcal{H}(K)$ ,  $\varphi(s) = \langle \varphi(\cdot), k(\cdot, s) \rangle_K$ .
- The kernel  $k$  is then called the reproducing kernel.

The following properties are listed in Aronszajn (1950):

- 1 - If the RK  $k$  exists, it is unique.
- 2 - A Hilbert space  $\mathcal{H}$  of functions defined on  $I \subset \mathbb{R}^L$  is a RKHS if and only if all functionals  $\varphi \rightarrow \varphi(s)$  for all  $\varphi \in \mathcal{H}$ ,  $s \in I$ , are bounded.
- 3 -  $K$  is a self-adjoint positive operator on  $L_C^2(\pi)$ .
- 4 - To a self-adjoint positive operator  $K$  on  $I$ , there corresponds a unique RKHS  $\mathcal{H}(K)$  of complex-valued functions.
- 5 - Every sequence of functions  $\{\varphi_n\}$  which converges weakly to  $\varphi$  in  $\mathcal{H}(K)$  (that is  $\langle \varphi_n, g \rangle_K \rightarrow \langle \varphi, g \rangle_K$  for all  $g \in \mathcal{H}(K)$ ) converges also pointwise, that is  $\lim \varphi_n(s) = \varphi(s)$ .

Note that (2) is a consequence of Riesz theorem 2.18: There exists a representer  $k$  such that for all  $\varphi \in \mathcal{H}$

$$\varphi(t) = \langle \varphi, k_t \rangle_K.$$

Let  $k_t = k(t, \cdot)$  so that  $\langle k_t, k_s \rangle_K = k(t, s)$ . (5) follows from the reproducing property. Indeed,  $\langle \varphi_n(t) - \varphi(t), k(t, s) \rangle_K = \varphi_n(s) - \varphi(s)$ .

**Example (finite dimensional case).** Let  $I = \{1, 2, \dots, L\}$ , let  $\Sigma$  be a positive definite  $L \times L$  matrix with principal element  $\sigma_{t,s}$ .  $\Sigma$  defines an inner product on  $\mathbb{R}^L$  :  $\langle \varphi, \psi \rangle_\Sigma = \varphi' \Sigma^{-1} \psi$ . Let  $(\sigma_1, \dots, \sigma_L)$  be the columns of  $\Sigma$ . Let  $\varphi = (\varphi(1), \dots, \varphi(L))'$ , then we have the reproducing property

$$\langle \varphi, \sigma_t \rangle_\Sigma = \varphi(t), \tau = 1, \dots, L$$

because  $\varphi \Sigma^{-1} \Sigma = \varphi$ . Now we diagonalize  $\Sigma$ ,  $\Sigma = PDP'$  where  $P$  is the  $m \times m$  matrix with  $(t, j)$  element  $\phi_j(t)$  ( $\phi_j$  are the orthonormal eigenvectors of  $\Sigma$ ) and  $D$  is the diagonal matrix with diagonal element  $\lambda_j$  (the eigenvalues of  $\Sigma$ ). The  $(t, s)$ th element of  $\Sigma$  can be rewritten as

$$\sigma(t, s) = \sum_{j=1}^m \lambda_j \phi_j(t) \phi_j(s).$$

We have

$$\langle \varphi, \psi \rangle_\Sigma = \varphi' \Sigma^{-1} \psi = \sum_{j=1}^m \frac{1}{\lambda_j} \langle \varphi, \phi_j \rangle \langle \psi, \phi_j \rangle$$

where  $\langle, \rangle$  is the euclidean inner product.

From this small example, we see that the norm in a RKHS can be characterized by the spectral decomposition of an operator. Let  $K$  be a positive self-adjoint compact operator with spectrum  $\{\phi_j, \lambda_j : j = 1, 2, \dots\}$ . Assume that  $\mathcal{N}(K) = 0$ . It turns out that  $\mathcal{H}(K)$  coincides with the  $1/2$ -regularization space of the operator  $K$  :

$$\mathcal{H}(K) = \left\{ \varphi : \varphi \in L^2(\pi) \text{ and } \sum_{j=1}^{\infty} \frac{|\langle \varphi, \phi_j \rangle|^2}{\lambda_j} < \infty \right\} = \Phi_{1/2}(K).$$

We can check that

(i)  $\mathcal{H}(K)$  is a Hilbert space with inner product

$$\langle \varphi, \psi \rangle_K = \sum_{j=1}^{\infty} \frac{\langle \varphi, \phi_j \rangle \overline{\langle \psi, \phi_j \rangle}}{\lambda_j}$$

and norm

$$\|\varphi\|_K^2 = \sum_{j=1}^{\infty} \frac{|\langle \varphi, \phi_j \rangle|^2}{\lambda_j}.$$

(ii)  $k(\cdot, t)$  belongs to  $\mathcal{H}(K)$

(iii)  $\langle \varphi, k(\cdot, t) \rangle_K = \varphi(t)$ .

**Proof.** (ii) follows from Mercer's formula (Theorem 2.42 (iii)) that is  $k(t, s) = \sum_{j=1}^{\infty} \lambda_j \phi_j(t) \overline{\phi_j(s)}$ . Hence  $\|k(\cdot, t)\|_K^2 = \sum_{j=1}^{\infty} |\langle \phi_j, k(\cdot, t) \rangle|^2 / \lambda_j = \sum_{j=1}^{\infty} |\lambda_j \phi_j(t)|^2 / \lambda_j = \sum_{j=1}^{\infty} \lambda_j \phi_j(t) \overline{\phi_j(t)} = k(t, t) < \infty$ . For (iii), we use again Mercer's formula.  $\langle \varphi(\cdot), k(\cdot, t) \rangle_K = \sum_{j=1}^{\infty} \langle \phi_j, k(\cdot, t) \rangle \langle \varphi, \phi_j \rangle / \lambda_j = \sum_{j=1}^{\infty} \langle \varphi, \phi_j \rangle K \phi_j(t) / \lambda_j = \sum_{j=1}^{\infty} \langle \varphi, \phi_j \rangle \phi_j(t) = \varphi(t)$ . ■

There is a link between calculating a norm in a RKHS and solving an integral equation  $K\varphi = \psi$ . We follow Nashed and Wahba (1974) to enlighten this link. We have

$$K\varphi = \sum_{j=1}^{\infty} \lambda_j \langle \varphi, \phi_j \rangle \phi_j.$$

Define  $K^{1/2}$  the square root of  $K$ :

$$K^{1/2}\varphi = \sum_{j=1}^{\infty} \sqrt{\lambda_j} \langle \varphi, \phi_j \rangle \phi_j.$$

Note that  $\mathcal{N}(K) = \mathcal{N}(K^{1/2})$ ,  $\mathcal{H}(K) = K^{1/2}(L_C^2(\pi))$ . Define  $K^{-1/2} = (K^{1/2})^\dagger$  where  $()^\dagger$  is the Moore-Penrose generalized inverse introduced in Subsection 3.1.:

$$K^\dagger\psi = \sum_{j=1}^{\infty} \frac{1}{\lambda_j} \langle \psi, \phi_j \rangle \phi_j.$$

Similarly, the generalized inverse of  $K^{1/2}$  takes the form:

$$K^{-1/2}\psi = \sum_{j=1}^{\infty} \frac{1}{\sqrt{\lambda_j}} \langle \psi, \phi_j \rangle \phi_j.$$

From Nashed and Wahba (1974), we have the relations

$$\begin{aligned} \|\varphi\|_K^2 &= \inf \{ \|p\| : p \in L_C^2(\pi) \text{ and } \varphi = K^{1/2}p \}, \\ \langle \varphi, \psi \rangle_K &= \langle K^{-1/2}\varphi, K^{-1/2}\psi \rangle, \text{ for all } \varphi, \psi \in \mathcal{H}(K). \end{aligned} \quad (6.1)$$

The following result follows from Proposition 3.6.

**Proposition 6.2.** *Let  $T : \mathcal{E} \rightarrow L_C^2(\pi)$  be an operator such that  $K = TT^*$  then*

$$\mathcal{H}(K) = \mathcal{R}(K^{1/2}) = \mathcal{R}(T^*) = \Phi_1(T).$$

Note that  $T^* : L_C^2(\pi) \rightarrow \mathcal{E}$  and  $K^{1/2} : L_C^2(\pi) \rightarrow L_C^2(\pi)$  are not equal because they take their values in different spaces.



### 6.1.2. RKHS for covariance operators of stochastic processes

In the previous section, we have seen how to characterize  $\mathcal{H}(K)$  using the spectral decomposition of  $K$ . When  $K$  is known to be the covariance kernel of a stochastic process, then  $\mathcal{H}(K)$  admits a simple representation. The main results of this section come from Parzen (1959). Consider a random element (r.e.)  $\{h(t), t \in I \subset \mathbb{R}^p\}$  defined on a probability space  $(\Omega, \mathcal{F}, P)$  and observed for all values of  $t$ . Assume  $h(t)$  is a second order random function that is  $E(|h(t)|^2) = \int_{\Omega} |h(t)|^2 dP < \infty$  for every  $t \in I$ . Let  $L_2(\Omega, \mathcal{F}, P)$  be the set of all r.v.  $U$  such that  $E|U|^2 = \int_{\Omega} |U|^2 dP < \infty$ . Define the inner product  $\langle U, V \rangle_{L_2(\Omega, \mathcal{F}, P)}$  between any two r.v.  $U$  and  $V$  of  $L_2(\Omega, \mathcal{F}, P)$  by  $\langle U, V \rangle_{L_2(\Omega, \mathcal{F}, P)} = E(U\bar{V}) = \int_{\Omega} U\bar{V} dP$ . Let  $L_2(h(t), t \in I)$  be the Hilbert space spanned by the r.e.  $\{h(t), t \in I\}$ . Define  $K$  the covariance operator with kernel  $k(t, s) = E(h(t)\bar{h}(s))$ . The following theorem implies that any symmetric nonnegative kernel can be written as a covariance kernel of a particular process.

**Theorem 6.3.**  *$K$  is a covariance operator of a r.e. if and only if  $K$  is a positive self-adjoint operator.*

The following theorem can be found in Parzen (1959) for real-valued functions. The complex case is treated in Saitoh (1997).

**Theorem 6.4.** *Let  $\{h(t), t \in I\}$  be a r.e. with mean zero and covariance kernel  $k$ . Then*

(i)  *$L_2(h(t), t \in I)$  is isometrically isomorphic or congruent to the RKHS  $\mathcal{H}(K)$ . Denote  $J : \mathcal{H}(K) \rightarrow L_2(h(t), t \in I)$  this congruence.*

(ii) *For every function  $\varphi$  in  $\mathcal{H}(K)$ ,  $J(\varphi)$  satisfies*

$$\langle J(\varphi), h(t) \rangle_{L_2(\Omega, \mathcal{F}, P)} = E\left(J(\varphi)\overline{h(t)}\right) = \langle \varphi, k(\cdot, t) \rangle_K = \varphi(t), \text{ for all } t \in I \quad (6.2)$$

where  $J(\varphi)$  is unique in  $L_2(h(t), t \in I)$  and has mean zero and variance such that

$$\|\varphi\|_K^2 = \|J(\varphi)\|_{L_2(\Omega, \mathcal{F}, P)}^2 = E(|J(\varphi)|^2).$$

Note that, by (6.2), the congruence is such that  $J(k(\cdot, t)) = h(t)$ . The r.v.  $U \in L_2(h(t), t \in I)$  corresponding to  $\varphi \in \mathcal{H}(K)$  is denoted below as  $\langle \varphi, h \rangle_K$  (or  $J(\varphi)$ ). As  $L_2(h(t), t \in I)$  and  $\mathcal{H}(K)$  are isometric, we have by Definition 2.19

$$\text{cov}[\langle \varphi, h \rangle_K, \langle \psi, h \rangle_K] = E\left[J(\varphi)\overline{J(\psi)}\right] = \langle \varphi, \psi \rangle_K$$

for every  $\varphi, \psi \in \mathcal{H}(K)$ . Note that  $\langle \varphi, h \rangle_K$  is not correct notation because  $h = \sum_j \langle h, \phi_j \rangle \phi_j$  a.s. does not belong to  $\mathcal{H}(K)$ . If it were the case, we should have  $\sum_j \langle h, \phi_j \rangle^2 / \lambda_j < \infty$  a.s.. Unfortunately  $\langle h, \phi_j \rangle$  are independent with mean 0 and variance  $\langle K\phi_j, \phi_j \rangle = \lambda_j$ . Hence,  $E\left[\sum_j \langle h, \phi_j \rangle^2 / \lambda_j\right] = \infty$  and by Kolmogorov's theorem  $\sum_j \langle h, \phi_j \rangle^2 / \lambda_j = \infty$  with

nonzero probability. The r.v.  $J(\varphi)$  itself is well-defined, only the notation  $\langle \varphi, h \rangle_K$  is not adequate; as Kailath (1971) explains, it should be regarded only as a mnemonic for finding  $J(\varphi)$  in a closed form. The rest of this section is devoted to the calculation of  $\|\varphi\|_K$ . Note that the result (6.2) is valid when  $t$  is multidimensional,  $t \in \mathbb{R}^L$ . In the next section,  $h(t)$  will be a moment function indexed by an arbitrary index parameter  $t$ .

Assume that the kernel  $k$  on  $I \times I$  can be represented as

$$k(s, t) = \int h(s, x) \overline{h(t, x)} P(dx) \quad (6.3)$$

where  $P$  is a probability measure and  $\{h(s, \cdot), s \in I\}$  is a family of functions on  $L_2(\Omega, \mathcal{F}, P)$ . By Theorem 6.4,  $\mathcal{H}(K)$  consists of functions  $\varphi$  on  $I$  of the form

$$\varphi(t) = \int \psi(x) \overline{h(t, x)} P(dx) \quad (6.4)$$

for some unique  $\psi$  in  $L_2(h(t, \cdot), t \in I)$ , the subspace of  $L_2(\Omega, \mathcal{F}, P)$  spanned by  $\{h(t, \cdot), t \in I\}$ . The RKHS norm of  $\varphi$  is given by

$$\|\varphi\|_K^2 = \|\psi\|_{L_2(\Omega, \mathcal{F}, P)}^2.$$

When calculating  $\|\varphi\|_K^2$  in practice, one looks for the solutions of (6.4). If there are several solutions, it is not always obvious to see which one is spanned by  $\{h(t, \cdot), t \in I\}$ . In this case, the right solution is the solution with minimal norm (Parzen, 1970):

$$\|\varphi\|_K^2 = \min_{\substack{\psi \text{ s.t.} \\ \varphi = \langle \psi, h \rangle_{L_2}}} \|\psi\|_{L_2(\Omega, \mathcal{F}, P)}^2.$$

Theorem 6.4 can be reinterpreted in terms of range. Let  $T$  and  $T^*$  be

$$\begin{aligned} T & : L^2(\pi) \rightarrow L_2(h(t, \cdot), t \in I) \\ \varphi & \rightarrow T\varphi(x) = \int \varphi(t) h(t, x) \pi(t) dt. \end{aligned}$$

and

$$\begin{aligned} T^* & : L_2(h(t, \cdot), t \in I) \rightarrow L^2(\pi) \\ \psi & \rightarrow T^*\psi(s) = \int \psi(x) \overline{h(s, x)} P(dx). \end{aligned}$$

To check that  $T^*$  is indeed the dual of  $T$ , it suffices to check  $\langle T\varphi, \psi \rangle_{L_2(\Omega, \mathcal{F}, P)} = \langle \varphi, T^*\psi \rangle_{L^2(\pi)}$  for  $\varphi \in L^2(\pi)$  and  $\psi(x) = h(t, x)$  as  $h(t, \cdot)$  spans  $L_2(h(t, \cdot), t \in I)$ . Using the fact that  $K = T^*T$  and Proposition 6.2, we have  $\mathcal{H}(K) = \mathcal{R}(T^*)$ , which gives Equation (6.4).

**Example.** Let  $k(t, s) = t \wedge s$ .  $k$  can be rewritten as

$$k(t, s) = \int_0^1 (t-x)_+^0 (s-x)_+^0 du$$

with

$$(s-x)_+^0 = \begin{cases} 1 & \text{if } x < s \\ 0 & \text{if } x \geq s \end{cases}.$$

It follows that  $\mathcal{H}(K)$  consists of functions  $\varphi$  of the form:

$$\begin{aligned} \varphi(t) &= \int_0^1 \psi(x) (t-x)_+^0 dx = \int_0^t \psi(x) dx, \quad 0 \leq t \leq 1 \\ \Rightarrow \psi(t) &= \varphi'(t). \end{aligned}$$

Hence, we have

$$\|\varphi\|_K^2 = \int_0^1 |\psi(x)|^2 dx = \int_0^1 |\varphi'(x)|^2 dx.$$

**Example.** Let  $k$  be defined as in (6.3) with  $h(t, x) = e^{itx}$ . Assume  $P$  admits a pdf  $f_{\theta_0}(x)$ , which is positive everywhere. Equation (6.4) is equivalent to

$$\begin{aligned} \varphi(t) &= \int \psi(x) e^{-itx} P(dx) \\ &= \int \psi(x) e^{-itx} f_{\theta_0}(x) dx. \end{aligned}$$

By the Fourier Inversion formula, we have

$$\psi(x) = \frac{1}{2\pi} \frac{1}{f_{\theta_0}(x)} \int e^{itx} \varphi(t) dt.$$

## 6.2. GMM in Hilbert spaces

First introduced by Hansen (1982), the Generalized Method of Moments (GMM) became the cornerstone of modern structural econometrics. In Hansen, the number of moment conditions is supposed to be finite. The method proposed in this section permits to deal with moment functions that take their values in finite or infinite dimensional Hilbert spaces. It was initially proposed by Carrasco and Florens (2000) and further developed in Carrasco and Florens (2001) and Carrasco, Chernov, Florens, and Ghysels (2004).

### 6.2.1. Definition and examples

Let  $\{x_i : i = 1, 2, \dots, n\}$  be an iid sample of a random vector  $X \in \mathbb{R}^p$ . The case where  $X$  is a time-series will be discussed later. The distribution of  $X$  is indexed by a parameter  $\theta \in \Theta \subset \mathbb{R}^d$ . Denote  $E^\theta$  the expectation with respect to this distribution. The unknown parameter  $\theta$  is identified from the function  $h(X; \theta)$  (called moment function) defined on  $\mathbb{R}^p \times \Theta$ , so that the following is true.

#### Identification Assumption

$$E^{\theta_0}(h(X; \theta)) = 0 \Leftrightarrow \theta = \theta_0. \quad (6.5)$$

It is assumed that  $h(X; \theta)$  takes its values in a Hilbert space  $\mathcal{H}$  with inner product  $\langle \cdot, \cdot \rangle$  and norm  $\|\cdot\|$ . When  $f = (f_1, \dots, f_L)$  and  $g = (g_1, \dots, g_L)$  are vectors of functions of  $\mathcal{H}$ , we use the convention that  $\langle f, g' \rangle$  denotes the  $L \times L$ -matrix with  $(l, m)$  element  $\langle f_l, g_m \rangle$ . Let  $B_n : \mathcal{H} \rightarrow \mathcal{H}$  be a sequence of random bounded linear operators and

$$\hat{h}_n(\theta) = \frac{1}{n} \sum_{i=1}^n h(x_i; \theta).$$

We define the GMM estimator associated with  $B_n$  as

$$\hat{\theta}_n(B_n) = \arg \min_{\theta \in \Theta} \left\| B_n \hat{h}_n(\theta) \right\|. \quad (6.6)$$

Such an estimator will in general be suboptimal; we will discuss the optimal choice of  $B_n$  later. Below, we give four examples that can be handled by the method discussed in this section. They illustrate the versatility of the method as it can deal with a finite number of moments (Example 1), a continuum (Examples 2 and 3) and a countably infinite sequence (Example 4).

**Example 1 (Traditional GMM).** Let  $h(x; \theta)$  be a vector of  $\mathbb{R}^L$ ,  $B_n$  be a  $L \times L$ -matrix and  $\|\cdot\|$  denote the Euclidean norm. The objective function to minimize is

$$\left\| B_n \hat{h}_n(\theta) \right\|^2 = \hat{h}_n(\theta)' B_n' B_n \hat{h}_n(\theta)$$

and corresponds to the usual GMM quadratic form  $\hat{h}_n(\theta)' W_n \hat{h}_n(\theta)$  with weighting matrix  $W_n = B_n' B_n$ .

**Example 2 (Continuous time process).** Suppose we observe independent replications of a continuous time process

$$X^i(t) = G(\theta, t) + u^i(t), \quad 0 \leq t \leq T, \quad i = 1, 2, \dots, n \quad (6.7)$$

where  $G$  is a known function and  $u^i = \{u^i(t) : 0 \leq t \leq T\}$  is a zero mean Gaussian process with continuous covariance function  $k(t, s) = E[u(t)u(s)]$ ,  $t, s \in [0, T]$ . Denote

$X^i = \{X^i(t) : 0 \leq t \leq T\}$ ,  $G(\theta) = \{G(\theta, t) : 0 \leq t \leq T\}$ , and  $\mathcal{H} = L^2([0, T])$ . The unknown parameter  $\theta$  is identified from the moment of the function

$$h(X^i; \theta) = X^i - G(\theta).$$

Assume  $h(X^i; \theta) \in L^2([0, T])$  with probability one. Candidates for  $B_n$  are arbitrary bounded operators on  $L^2([0, T])$  including the identity. For  $B_n f = f$ , we have

$$\|B_n \hat{h}_n(\theta)\|^2 = \int_0^T \hat{h}_n(\theta)^2 dt.$$

Estimation of Model (6.7) is discussed in Kutoyants (1984).

**Example 3 (Characteristic function).** Denote  $\psi_\theta(t) = E^\theta[e^{it'X}]$  the characteristic function of  $X$ . Inference can be based on

$$h(t, X; \theta) = e^{it'X} - \psi_\theta(t), \quad t \in \mathbb{R}^L.$$

Note that contrary to the former examples,  $h(t, X; \theta)$  is complex valued and  $|h(t, X; \theta)| \leq |e^{it'X}| + |\psi_\theta(t)| \leq 2$ . Let  $\Pi$  be a probability measure on  $\mathbb{R}^L$  and  $\mathcal{H} = L^2_C(\mathbb{R}^L, \Pi)$ . As  $h(\cdot, X; \theta)$  is bounded, it belongs to  $L^2_C(\mathbb{R}^L, \Pi)$  for any  $\Pi$ . Feuerverger and McDunnough (1981) and more recently Singleton (2001) show that an efficient estimator of  $\theta$  is obtained from  $h(\cdot, X; \theta)$  by solving an empirical counterpart of  $\int E h(t, X; \theta) \omega(t) dt = 0$  for an adequate weighting function  $\omega$ , which turns out to be a function of the pdf of  $X$ . This efficient estimator is not implementable as the pdf of  $X$  is unknown. They suggest estimating  $\theta$  by GMM using moments obtained from a discrete grid  $t = t_1, t_2, \dots, t_M$ . An alternative strategy put forward in this section is to use the full continuum of moment conditions by considering the moment function  $h$  as an element of  $\mathcal{H} = L^2_C(\mathbb{R}^L, \Pi)$ .

**Example 4 (Conditional moment restrictions).** Let  $X = (Y, Z)$ . For a known function  $\rho \in \mathbb{R}$ , we have the conditional moment restrictions

$$E^{\theta_0}[\rho(Y, Z, \theta) | Z] = 0.$$

Hence for any function  $g(Z)$ , we can construct unconditional moment restrictions

$$E^{\theta_0}[\rho(Y, Z, \theta) g(Z)] = 0.$$

Assume  $Z$  has bounded support. Chamberlain (1987) shows that the semiparametric efficiency bound can be approached by a GMM estimator based on a sequence of moment conditions using as instruments the power function of  $Z : 1, Z, Z^2, \dots, Z^L$  for a large  $L$ . Let  $\pi$  be the Poisson probability measure  $\pi(l) = e^{-1}/l!$  and  $\mathcal{H} = L^2(\mathbf{N}, \pi) = \{f : \mathbf{N} \rightarrow \mathbb{R} : \sum_{l=1}^{\infty} g(l) \pi(l) < \infty\}$ . Let

$$h(l, X; \theta) = \rho(Y, Z, \theta) Z^l, \quad l = 1, 2, \dots$$

If  $h(l, X; \theta)$  is bounded with probability one, then  $h(\cdot, X; \theta) \in L^2(\mathbf{N}, \pi)$  with probability one. Instead of using an increasing sequence of moments as suggested by Chamberlain, it is possible to handle  $h(\cdot, X; \theta)$  as a function. The efficiency of the GMM estimator based on the countably infinite number of moments  $\{h(l, X; \theta) : l \in \mathbf{N}\}$  will be discussed later.

### 6.2.2. Asymptotic properties of GMM

Let  $\mathcal{H} = L_C^2(I, \Pi) = \{f : I \rightarrow \mathbf{C} : \int_I |f(t)|^2 \Pi(dt) < \infty\}$  where  $I$  is a subset of  $\mathbb{R}^L$  for some  $L \geq 1$  and  $\Pi$  is a (possibly discrete) probability measure. This choice of  $\mathcal{H}$  is consistent with Examples 1 to 4. Under some weak assumptions,  $\sqrt{n}\hat{h}_n(\theta_0)$  converges to a Gaussian process  $\mathcal{N}(0, K)$  in  $\mathcal{H}$  where  $K$  denotes the covariance operator of  $h(X; \theta_0)$ .  $K$  is defined by

$$\begin{aligned} K & : \mathcal{H} \rightarrow \mathcal{H} \\ f & \rightarrow Kf(s) = \langle f, k(\cdot, t) \rangle = \int_I k(t, s) f(s) \Pi(ds) \end{aligned}$$

where the kernel  $k$  of  $K$  satisfies  $k(t, s) = E^{\theta_0} \left[ h(t, X; \theta_0) \overline{h(s, X; \theta_0)} \right]$  and  $k(t, s) = \overline{k(s, t)}$ . Assume moreover that  $K$  is a Hilbert Schmidt operator and hence admits a discrete spectrum. Suppose that  $B_n$  converges to a bounded linear operator  $B$  defined on  $\mathcal{H}$  and that  $\theta_0$  is the unique minimizer of  $\|BE^{\theta_0}h(X; \theta)\|$ . Then  $\hat{\theta}_n(B_n)$  is consistent and asymptotically normal. The following result is proved in Carrasco and Florens (2000).

**Proposition 6.5.** *Under Assumptions 1 to 11 of Carrasco and Florens (2000),  $\hat{\theta}_n(B_n)$  is consistent and*

$$\sqrt{n} \left( \hat{\theta}_n(B_n) - \theta_0 \right) \xrightarrow{L} \mathcal{N}(0, V)$$

with

$$\begin{aligned} V & = \langle BE^{\theta_0}(\nabla_{\theta}h), BE^{\theta_0}(\nabla_{\theta}h)' \rangle^{-1} \\ & \times \langle BE^{\theta_0}(\nabla_{\theta}h), (BKB^*)BE^{\theta_0}(\nabla_{\theta}h)' \rangle \\ & \times \langle BE^{\theta_0}(\nabla_{\theta}h), BE^{\theta_0}(\nabla_{\theta}h)' \rangle^{-1} \end{aligned}$$

where  $B^*$  is the adjoint of  $B$ .

### 6.2.3. Optimal choice of the weighting operator

Carrasco and Florens (2000) show that the asymptotic variance  $V$  given in Proposition 6.5 is minimal for  $B = K^{-1/2}$ . In that case, the asymptotic covariance becomes  $\langle K^{-1/2}E^{\theta_0}(\nabla_{\theta}h), K^{-1/2}E^{\theta_0}(\nabla_{\theta}h) \rangle^{-1}$ .

**Example 1 (continued).**  $K$  is the  $L \times L$ -covariance matrix of  $h(X; \theta)$ . Let  $K_n$  be the matrix  $\frac{1}{n} \sum_{i=1}^n h(x_i; \hat{\theta}^1) h(x_i; \hat{\theta}^1)'$  where  $\hat{\theta}^1$  is a consistent first step estimator of  $\theta$ .  $K_n$  is a consistent estimator of  $K$ . Then the objective function becomes

$$\left\langle K_n^{-1/2} \hat{h}_n(\theta), K_n^{-1/2} \hat{h}_n(\theta) \right\rangle = \hat{h}_n(\theta)' K_n^{-1} \hat{h}_n(\theta)$$

which delivers the optimal GMM estimator.

When  $\mathcal{H}$  is infinite dimensional, we have seen in Section 3.1 that the inverse of  $K$ ,  $K^{-1}$ , is not bounded. Similarly  $K^{-1/2} = (K^{1/2})^{-1}$  is not bounded on  $\mathcal{H}$  and its domain has been shown in Subsection 6.1.1 to be the subset of  $\mathcal{H}$  which coincides with the RKHS associated with  $K$  and denoted  $\mathcal{H}(K)$ .

To estimate the covariance operator  $K$ , we need a first step estimator  $\hat{\theta}^1$  that is  $\sqrt{n}$ -consistent. It may be obtained by letting  $B_n$  equal the identity in (6.6) or by using a finite number of moments. Let  $K_n$  be the operator with kernel

$$k_n(t, s) = \frac{1}{n} \sum_{i=1}^n h(t, x_i; \hat{\theta}^1) \overline{h(s, x_i; \hat{\theta}^1)}.$$

Then  $K_n$  is a consistent estimator of  $K$  and  $\|K_n - K\| = O(1/\sqrt{n})$ . As  $K^{-1}f$  is not continuous in  $f$ , we estimate  $K^{-1}$  by the Tykhonov regularized inverse of  $K_n$ :

$$(K_n^{\alpha_n})^{-1} = (\alpha_n I + K_n^2)^{-1} K_n$$

for some penalization term  $\alpha_n \geq 0$ . If  $\alpha_n > 0$ ,  $(K_n^{\alpha_n})^{-1}f$  is continuous in  $f$  but is a biased estimator of  $K^{-1}f$ . There is a trade-off between the stability of the solution and its bias. Hence, we will let  $\alpha_n$  decrease to zero at an appropriate rate. We define  $(K_n^{\alpha_n})^{-1/2} = ((K_n^{\alpha_n})^{-1})^{1/2}$ .

The optimal GMM estimator is given by

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} \left\| (K_n^{\alpha_n})^{-1/2} \hat{h}_n(\theta) \right\|.$$

Interestingly, the optimal GMM estimator minimizes the norm of  $\hat{h}_n(\theta)$  in the RKHS associated with  $K_n^{\alpha_n}$ . Under certain regularity conditions, we have

$$\left\| (K_n^{\alpha_n})^{-1/2} \hat{h}_n(\theta) \right\| \xrightarrow{P} \|E^{\theta_0}(h(\theta))\|_K.$$

A condition for applying this method is that  $E^{\theta_0}(h(\theta)) \in \mathcal{H}(K)$ . This condition can be verified using results from 6.1.1.

**Proposition 6.6.** *Under the regularity conditions of Carrasco and Florens (2000, Theorem 8),  $\hat{\theta}_n$  is consistent and*

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{L} \mathcal{N}\left(0, \langle E^{\theta_0}(\nabla_{\theta} h(\theta_0)), E^{\theta_0}(\nabla_{\theta} h(\theta_0))' \rangle_K^{-1}\right)$$

as  $n$  and  $n\alpha_n^3 \rightarrow \infty$  and  $\alpha_n \rightarrow 0$ .

The stronger condition  $n\alpha_n^3 \rightarrow \infty$  of Carrasco and Florens (2000) has been relaxed into  $n\alpha^2 \rightarrow \infty$  in Carrasco, Chernov, Florens, and Ghysels (2004). Proposition 6.6 does not indicate how to select  $\alpha_n$  in practice. A data-driven method is desirable. Carrasco and Florens (2001) propose to select the  $\alpha_n$  that minimizes the mean square error (MSE) of the GMM estimator  $\hat{\theta}_n$ . As  $\hat{\theta}_n$  is consistent for any value of  $\alpha_n$ , it is necessary to compute the higher order expansion of the MSE, which is particularly tedious. Instead of relying on an analytic expression, it may be easier to compute the MSE via bootstrap or simulations.

### 6.2.4. Implementation of GMM

There are two equivalent ways to compute the objective function

$$\left\| (K_n^{\alpha_n})^{-1/2} \hat{h}_n(\theta) \right\|^2, \quad (6.8)$$

- 1) using the spectral decomposition of  $K_n$ , or
- 2) using a simplified formula that involves only vectors and matrices.

The first method discussed in Carrasco and Florens (2000) requires calculating the eigenvalues and eigenfunctions of  $K_n$  using the method described in 2.5.3. Let  $\hat{\phi}_j$  denote the orthonormalized eigenfunctions of  $K_n$  and  $\hat{\lambda}_j$  the corresponding eigenvalues. The objective function in Equation (6.8) becomes

$$\sum_{j=1}^n \frac{\hat{\lambda}_j}{\hat{\lambda}_j^2 + \alpha_n} \left| \left\langle \hat{h}_n(\theta), \hat{\phi}_j \right\rangle \right|^2. \quad (6.9)$$

The expression (6.9) suggests a nice interpretation of the GMM estimator. Indeed, note that  $\left\langle \sqrt{n} \hat{h}_n(\theta_0), \phi_j \right\rangle$ ,  $j = 1, 2, \dots$  are asymptotically normal with mean 0 and variance  $\lambda_j$  and are independent across  $j$ . Therefore (6.9) is the regularized version of the objective function of the optimal GMM estimator based on the  $n$  moment conditions  $E[\langle h(\theta), \phi_j \rangle] = 0$ ,  $j = 1, 2, \dots, n$ .

The second method is more attractive by its simplicity. Carrasco et al. (2004) show that (6.8) can be rewritten as

$$\overline{\underline{v}(\theta)}' [\alpha_n I_n + C^2]^{-1} \underline{v}(\theta)$$

where  $C$  is a  $n \times n$ -matrix with  $(i, j)$  element  $c_{ij}$ ,  $I_n$  is the  $n \times n$  identity matrix and  $\underline{v}(\theta) = (v_1(\theta), \dots, v_n(\theta))'$  with

$$\begin{aligned} v_i(\theta) &= \int \overline{h(t, x_i; \hat{\theta}^1)}' \hat{h}_n(t; \theta) \Pi(dt) \\ c_{ij} &= \frac{1}{n} \int \overline{h(t, x_i; \hat{\theta}^1)}' h(t, x_j; \hat{\theta}^1) \Pi(dt). \end{aligned}$$

Note that the dimension of  $C$  is the same whether  $h \in \mathbb{R}$  or  $h \in \mathbb{R}^L$ .

### 6.2.5. Asymptotic Efficiency of GMM

Assume that the pdf of  $X$ ,  $f_\theta$ , is differentiable with respect to  $\theta$ . Let  $L^2(h)$  be the closure of the subspace of  $L^2(\Omega, \mathcal{F}, P)$  spanned by  $\{h(t, X_i; \theta_0) : t \in I\}$ .

**Proposition 6.7.** *Under standard regularity conditions, the GMM estimator based on  $\{h(t, x_i; \theta) : t \in I\}$  is asymptotically as efficient as the MLE if and only if*

$$\nabla_\theta \ln f_\theta(x_i; \theta_0) \in L^2(h).$$



This result is proved in Carrasco and Florens (2004) in a more general setting where  $X_i$  is Markov of order  $L$ . A similar efficiency result can be found in Hansen (1985), Tauchen (1997) and Gallant and Long (1997).

**Example 2 (continued).** Let  $K$  be the covariance operator of  $\{u(t)\}$  and  $\mathcal{H}(K)$  the RKHS associated with  $K$ . Kutoyants (1984) shows that if  $G(\theta) \in \mathcal{H}(K)$ , the likelihood ratio of the measure induced by  $X(t)$  with respect to the measure induced by  $u(t)$  equals

$$LR(\theta) = \prod_{i=1}^n \exp \left\{ \langle G(\theta), x^i \rangle_K - \frac{1}{2} \|G(\theta)\|_K^2 \right\}$$

where  $\langle G, X \rangle_K$  has been defined in Subsection 6.1.2 and denotes the element of  $L^2(X(t) : 0 \leq t \leq T)$  under the mapping  $J^{-1}$  of the function  $G(\theta)$  ( $J$  is defined in Theorem 6.4). The score function with respect to  $\theta$  is

$$\nabla_{\theta} \ln(LR(\theta)) = \left\langle \nabla_{\theta} G(\theta), \frac{1}{n} \sum_{i=1}^n (x^i - G(\theta)) \right\rangle_K.$$

For  $\theta = \theta_0$  and a single observation, the score is equal to

$$\langle \nabla_{\theta} G(\theta_0), u \rangle_K,$$

which is an element of  $L^2(u(t) : 0 \leq t \leq T) = L^2(h(X(t); \theta_0) : 0 \leq t \leq T)$ . Hence, by Proposition 6.7, the GMM estimator based on  $h(X; \theta_0)$  is asymptotically efficient. This efficiency result is corroborated by the following. The GMM objective function is

$$\|h(x; \theta)\|_K^2 = \left\langle \frac{1}{n} \sum_{i=1}^n (x^i - G(\theta)), \frac{1}{n} \sum_{i=1}^n (x^i - G(\theta)) \right\rangle_K.$$

The first order derivative equals to

$$\begin{aligned} \nabla_{\theta} \|h(x; \theta)\|_K^2 &= 2 \left\langle \nabla_{\theta} G(\theta), \frac{1}{n} \sum_{i=1}^n (x^i - G(\theta)) \right\rangle_K \\ &= 2 \nabla_{\theta} \ln(LR(\theta)). \end{aligned}$$

Therefore, the GMM estimator coincides with the MLE in this particular case as they are solutions of the same equation.

**Example 3 (continued).** Under minor conditions on the distribution of  $X_i$ , the closure of the linear span of  $\{h(t, X_i; \theta_0) : t \in \mathbb{R}^L\}$  contains all functions of  $L^2(X) = \{g : E^{\theta_0} [g(X)^2] < \infty\}$  and hence the score  $\nabla_{\theta} \ln f_{\theta}(X_i; \theta_0)$  itself. Therefore the GMM estimator is efficient. Another way to prove efficiency is to explicitly calculate the asymptotic covariance of  $\hat{\theta}_n$ . To simplify, assume that  $\theta$  is scalar. By Theorem 6.4, we have

$$\|E^{\theta_0}(\nabla_{\theta} h(\theta_0))\|_K^2 = \left\| \overline{E^{\theta_0}(\nabla_{\theta} h(\theta_0))} \right\|_K^2 = E|U|^2$$

where  $U$  satisfies

$$E^{\theta_0} \left[ \overline{U h(t; \theta_0)} \right] = \overline{E^{\theta_0} (\nabla_{\theta} h(t; \theta_0))} \text{ for all } t \in \mathbb{R}^L$$

which is equivalent to

$$E^{\theta_0} \left[ \overline{U(X)} \left( e^{it'X} - \psi_{\theta_0}(t) \right) \right] = -\nabla_{\theta} \psi_{\theta_0}(t) \text{ for all } t \in \mathbb{R}^L. \quad (6.10)$$

As  $U$  has mean zero,  $\overline{U}$  has also mean zero and we can replace (6.10) by

$$\begin{aligned} E^{\theta_0} \left[ \overline{U(X)} e^{it'X} \right] &= -\nabla_{\theta} \psi_{\theta_0}(t) \text{ for all } t \in \mathbb{R}^L \Leftrightarrow \\ \int \overline{U(x)} e^{it'x} f_{\theta_0}(x) dx &= -\nabla_{\theta} \psi_{\theta_0}(t) \text{ for all } t \in \mathbb{R}^L \Leftrightarrow \\ \overline{U(x)} f_{\theta_0}(x) &= -\frac{1}{2\pi} \int e^{-it'x} \nabla_{\theta} \psi_{\theta_0}(t) dt. \end{aligned} \quad (6.11)$$

The last equivalence follows from the Fourier inversion formula. Assuming that we can exchange the integration and derivation in the right hand side of (6.11), we obtain

$$\begin{aligned} \overline{U(x)} f_{\theta_0}(x) &= -\nabla_{\theta} f_{\theta_0}(x) \Leftrightarrow \\ U(x) &= -\nabla_{\theta} \ln f_{\theta_0}(x). \end{aligned}$$

Hence  $E^{\theta_0} |U|^2 = E^{\theta_0} [(\nabla_{\theta} \ln f_{\theta_0}(X))^2]$ . The asymptotic variance of  $\hat{\theta}_n$  coincides with the Cramer Rao efficiency bound even if, contrary to Example 3,  $\hat{\theta}_n$  differs from the MLE.

**Example 4 (continued).** As in the previous example, we intend to calculate the asymptotic covariance of  $\hat{\theta}_n$  using Theorem 6.4. We need to find  $U$  the  $p$ -vector of r.v. such that

$$\begin{aligned} E^{\theta_0} [U \rho(Y, Z; \theta_0) Z^l] &= E^{\theta_0} [\nabla_{\theta} \rho(Y, Z; \theta_0) Z^l] \text{ for all } l \in \mathbf{N}, \Leftrightarrow \\ E^{\theta_0} [E^{\theta_0} [U \rho(Y, Z; \theta_0) | Z] Z^l] &= E^{\theta_0} [E^{\theta_0} [\nabla_{\theta} \rho(Y, Z; \theta_0) | Z] Z^l] \text{ for all } l \in \mathbf{N} \end{aligned} \quad (6.12)$$

(6.12) is equivalent to

$$E^{\theta_0} [U \rho(Y, Z; \theta_0) | Z] = E^{\theta_0} [\nabla_{\theta} \rho(Y, Z; \theta_0) | Z] \quad (6.13)$$

by the completeness of polynomials (Sansone, 1959) under some mild conditions on the distribution of  $Z$ . A solution is

$$U_0 = E^{\theta_0} [\nabla_{\theta} \rho(Y, Z; \theta_0) | Z] E^{\theta_0} [\rho(Y, Z; \theta_0)^2 | Z]^{-1} \rho(Y, Z; \theta_0).$$

We have to check that this solution has minimal norm among all the solutions. Consider an arbitrary solution  $U = U_0 + U_1$ .  $U$  solution of (6.13) implies

$$E^{\theta_0} [U_1 \rho(Y, Z; \theta_0) | Z] = 0.$$

Hence  $E^{\theta_0}(UU') = E^{\theta_0}(U_0U_0') + E^{\theta_0}(U_1U_1')$  and is minimal for  $U_1 = 0$ . Then

$$\begin{aligned} & \left\| E^{\theta_0}(\nabla_{\theta} h(\theta_0)) \right\|_K^2 \\ &= E^{\theta_0}(U_0U_0') \\ &= E^{\theta_0} \left\{ E^{\theta_0}[\nabla_{\theta} \rho(Y, Z; \theta_0) | Z] E^{\theta_0}[\rho(Y, Z; \theta_0)^2 | Z]^{-1} E^{\theta_0}[\nabla_{\theta} \rho(Y, Z; \theta_0) | Z]' \right\}. \end{aligned}$$

Its inverse coincides with the semi-parametric efficiency bound derived by Chamberlain (1987).

Note that in Examples 2 and 3, the GMM estimator reaches the Cramer Rao bound asymptotically, while in Example 4 it reaches the semi-parametric efficiency bound.

### 6.2.6. Testing overidentifying restrictions

Hansen (1982) proposes a test of specification, which basically tests whether the overidentifying restrictions are close to zero. Carrasco and Florens (2000) propose the analogue to Hansen's J test in the case where there is a continuum of moment conditions. Let

$$\hat{p}_n = \sum_{j=1}^n \frac{\hat{\lambda}_j^2}{\hat{\lambda}_j^2 + \alpha_n}, \quad \hat{q}_n = 2 \sum_{j=1}^n \frac{\hat{\lambda}_j^4}{(\hat{\lambda}_j^2 + \alpha_n)^2}$$

where  $\hat{\lambda}_j$  are the eigenvalues of  $K_n$  as described earlier.

**Proposition 6.8.** *Under the assumptions of Theorem 10 of Carrasco and Florens (2000), we have*

$$\tau_n = \frac{\left\| (K_n^{\alpha_n})^{-1/2} \hat{h}_n(\hat{\theta}_n) \right\|^2 - \hat{p}_n}{\hat{q}_n} \xrightarrow{d} \mathcal{N}(0, 1)$$

as  $\alpha_n$  goes to zero and  $n\alpha_n^3$  goes to infinity.

This test can also be used for testing underidentification. Let  $\theta_0 \in \mathbb{R}$  be such that  $E[h(X, \theta_0)] = 0$ . Arellano, Hansen and Sentana (2005) show that the parameter,  $\theta_0$ , is locally unidentified if  $E[h(X, \theta)] = 0$  for all  $\theta \in \mathbb{R}$ . It results in a continuum of moment conditions indexed by  $\theta$ . Arellano et al. (2005) apply  $\tau_n$  to test for the null of underidentification.

### 6.2.7. Extension to time series

So far, the data were assumed to be iid. Now we relax this assumption. Let  $\{x_1, \dots, x_T\}$  be the observations of a time series  $\{X_t\}$  that satisfies some mixing conditions. Inference will be based on moment functions  $\{h(\tau, X_t; \theta_0)\}$  indexed by a real, possibly multidimensional index  $\tau$ .  $\{h(\tau, X_t; \theta_0)\}$  are in general autocorrelated, except in some special cases, an example of which will be discussed below.

**Example 5 (Conditional characteristic function).** Let  $Y_t$  be a (scalar) Markov process and assume that the conditional characteristic function (CF) of  $Y_{t+1}$  given  $Y_t$ ,  $\psi_\theta(\tau|Y_t) \equiv E^\theta[\exp(i\tau Y_{t+1})|Y_t]$ , is known. The following conditional moment condition holds

$$E^\theta [e^{i\tau Y_{t+1}} - \psi_\theta(\tau|Y_t) | Y_t] = 0.$$

Denote  $X_t = (Y_t, Y_{t+1})'$ . Let  $g(Y_t)$  be an instrument so that

$$h(\tau, X_t; \theta) = (e^{i\tau Y_{t+1}} - \psi_\theta(\tau|Y_t)) g(Y_t)$$

satisfies the identification condition (6.5).  $\{h(\tau, X_t; \theta)\}$  is a martingale difference sequence and is therefore uncorrelated. The use of the conditional CF is very popular in finance. Assume that  $\{Y_t, t = 1, 2, \dots, T\}$  is a discretely sampled diffusion process, then  $Y_t$  is Markov. While the conditional likelihood of  $Y_{t+1}$  given  $Y_t$  does not have a closed form expression, the conditional CF of affine diffusions is known. Hence GMM can replace MLE to estimate these models where MLE is difficult to implement. For an adequate choice of the instrument  $g(Y_t)$ , the GMM estimator is asymptotically as efficient as the MLE. The conditional CF has been recently applied to the estimation of diffusions by Singleton (2001), Chacko and Viceira (2003), and Carrasco et al. (2004). The first two papers use GMM based on a finite grid of values for  $\tau$ , whereas the last paper advocates using the full continuum of moments which permits us to achieve efficiency asymptotically.

**Example 6 (Joint characteristic function).** Assume  $Y_t$  is not Markov. In that case, the conditional CF is usually unknown. On the other hand, the joint characteristic function may be calculated explicitly (for instance when  $Y_t$  is an ARMA process with stable error, see Knight and Yu, 2002; or  $Y_t$  is the growth rate of a stochastic volatility model, see Jiang and Knight, 2002) or may be estimated via simulations (this technique is developed in Carrasco et al., 2004). Denote  $\psi_\theta(\tau) \equiv E^\theta[\exp(\tau_1 Y_t + \tau_2 Y_{t+1} + \dots + \tau_{L+1} Y_{t+L})]$  with  $\tau = (\tau_1, \dots, \tau_L)'$ , the joint CF of  $X_t \equiv (Y_t, Y_{t+1}, \dots, Y_{t+L})'$  for some integer  $L \geq 1$ . Assume that  $L$  is large enough for

$$h(\tau, X_t; \theta) = e^{i\tau' X_t} - \psi_\theta(\tau)$$

to identify the parameter  $\theta$ . Here  $\{h(\tau, X_t; \theta)\}$  are autocorrelated. Knight and Yu (2002) estimate various models by minimizing the following norm of  $h(\tau, X_t; \theta)$  :

$$\int \left( \frac{1}{T} \sum_{t=1}^T e^{i\tau' x_t} - \psi_\theta(\tau) \right)^2 e^{-\tau' \tau} d\tau.$$

This is equivalent to minimizing  $\left\| B \frac{1}{T} \sum_{t=1}^T h(\tau, X_t; \theta) \right\|^2$  with  $B = e^{-\tau' \tau / 2}$ . This choice of  $B$  is suboptimal but has the advantage of being easy to implement. The optimal weighting operator is, as before, the square root of the inverse of the covariance operator. Its estimation will be discussed shortly.

Under some mixing conditions on  $\{h(\tau, X_t; \theta_0)\}$ , the process  $\hat{h}_T(\theta_0) = \frac{1}{T} \sum_{t=1}^T h(\tau, X_t; \theta_0)$  follows a functional CLT (see Subsection 2.4.2):

$$\sqrt{T} \hat{h}_T(\theta_0) \xrightarrow{L} \mathcal{N}(0, K)$$

where the covariance operator  $K$  is an integral operator with kernel

$$k(\tau_1, \tau_2) = \sum_{j=-\infty}^{+\infty} E^{\theta_0} \left[ h(\tau_1, X_t; \theta_0) \overline{h(\tau_2, X_{t-j}; \theta_0)} \right].$$

The kernel  $k$  can be estimated using a kernel-based estimator as those described in Andrews (1991) and references therein. Let  $\omega : \mathbb{R} \rightarrow [-1, 1]$  be a kernel satisfying the conditions stated by Andrews. Let  $q$  be the largest value in  $[0, +\infty)$  for which

$$\omega_q = \lim_{u \rightarrow \infty} \frac{1 - \omega(u)}{|u|^q}$$

is finite. In the sequel, we will say that  $\omega$  is a  $q$ -kernel. Typically,  $q = 1$  for the Bartlett kernel and  $q = 2$  for Parzen, Tuckey-Hanning and quadratic spectral kernels. We define

$$\hat{k}_T(\tau_1, \tau_2) = \frac{T}{T-d} \sum_{j=-T+1}^{T-1} \omega\left(\frac{j}{S_T}\right) \hat{\Gamma}_T(j) \quad (6.14)$$

with

$$\hat{\Gamma}_T(j) = \begin{cases} \frac{1}{T} \sum_{t=j+1}^T h(\tau_1, X_t; \hat{\theta}_T^1) \overline{h(\tau_2, X_{t-j}; \hat{\theta}_T^1)}, & j \geq 0 \\ \frac{1}{T} \sum_{t=-j+1}^T h(\tau_1, X_{t+j}; \hat{\theta}_T^1) \overline{h(\tau_2, X_t; \hat{\theta}_T^1)}, & j < 0 \end{cases} \quad (6.15)$$

where  $S_T$  is some bandwidth that diverges with  $T$  and  $\hat{\theta}_T^1$  is a  $T^{1/2}$ -consistent estimator of  $\theta$ . Let  $K_T$  be the integral estimator with kernel  $\hat{k}_T$ . Under some conditions on  $\omega$  and  $\{h(\tau, X_t; \theta_0)\}$ , Carrasco et al. (2004) establish the rate of convergence of  $K_T$  to  $K$ . Assuming  $S_T^{2q+1}/T \rightarrow \gamma \in (0, +\infty)$ , we have

$$\|K_T - K\| = O_p(T^{-q/(2q+1)}).$$

The inverse of  $K$  is estimated using the regularized inverse of  $K_T$ ,  $K_T^{\alpha_T} = (K_T^2 + \alpha_T I)^{-1} K_T$  for a penalization term  $\alpha_T \geq 0$ . As before, the optimal GMM estimator is given by

$$\hat{\theta}_T = \arg \min_{\theta \in \Theta} \left\| (K_T^{\alpha_T})^{-1/2} \hat{h}_T(\theta) \right\|.$$

Carrasco et al. (2004) show the following result.

**Proposition 6.9.** Assume that  $\omega$  is a  $q$ -kernel and that  $S_T^{2q+1}/T \rightarrow \gamma \in (0, +\infty)$ . We have

$$\sqrt{T}(\hat{\theta}_T - \theta_0) \xrightarrow{L} \mathcal{N}\left(0, (\langle E^{\theta_0}(\nabla_{\theta}h), E^{\theta_0}(\nabla_{\theta}h)' \rangle_K)^{-1}\right) \quad (6.16)$$

as  $T$  and  $T^{q/(2q+1)}\alpha_T$  go to infinity and  $\alpha_T$  goes to zero.

Note that the implementation of this method requires two smoothing parameters  $\alpha_T$  and  $S_T$ . No cross-validation method for selecting these two parameters simultaneously has been derived yet. If  $\{h_t\}$  is uncorrelated, then  $K$  can be estimated using the sample average and the resulting estimator satisfies  $\|K_T - K\| = O_p(T^{-1/2})$ . When  $\{h_t\}$  are correlated, the convergence rate of  $K_T$  is slower and accordingly the rate of convergence of  $\alpha_T$  to zero is slower.

## 7. Estimating solutions of integral equations of the second kind

### 7.1. Introduction

The objective of this section is to study the properties of the solution of an integral equation of the second kind (also called Fredholm equation of the second type) defined by:

$$(I - K)\varphi = r \quad (7.1)$$

where  $\varphi$  is an element of a Hilbert space  $\mathcal{H}$ ,  $K$  is a compact operator from  $\mathcal{H}$  to  $\mathcal{H}$  and  $r$  is an element of  $\mathcal{H}$ . As in the previous sections,  $K$  and  $r$  are known functions of a data generating process characterized by its c.d.f.  $F$ , and the functional parameter of interest is the function  $\varphi$ .

In most cases,  $\mathcal{H}$  is a functional space and  $K$  is an integral operator defined by its kernel  $k$ . Equation (7.1) becomes:

$$\varphi(t) - \int k(t, s)\varphi(s)\Pi(ds) = r(t) \quad (7.2)$$

The estimated operators are often degenerate, see Subsection 2.5.1. and in that case, Equation (7.2) simplifies into:

$$\varphi(t) - \sum_{\ell=1}^L a_{\ell}(\varphi)\varepsilon_{\ell}(t) = r(t) \quad (7.3)$$

where the  $a_{\ell}(\varphi)$  are linear forms on  $\mathcal{H}$  and  $\varepsilon_{\ell}$  belongs to  $\mathcal{H}$  for any  $\ell$ .

The essential difference between equations of the first kind and of the second kind is the compactness of the operator. In (7.1),  $K$  is compact but  $I - K$  is not compact. Moreover, if  $I - K$  is one-to-one, its inverse is bounded. In that case, the inverse problem

is well-posed. Even if  $I - K$  is not one-to-one, the ill-posedness of equation (7.1) is less severe than in the first kind case because the solutions are stable in  $r$ .

In most cases,  $K$  is a self-adjoint operator (and hence  $I - K$  is also self-adjoint) but we will not restrict our presentation to this case. On the other hand, Equation (7.1) can be extended by considering an equation  $(S - K)\varphi = r$  where  $K$  is a compact operator from  $\mathcal{H}$  to  $\mathcal{E}$  (instead of  $\mathcal{H}$  to  $\mathcal{H}$ ) and  $S$  is a one-to-one bounded operator from  $\mathcal{H}$  to  $\mathcal{E}$  with a bounded inverse. Indeed,  $(S - K)\varphi = r \Leftrightarrow (I - S^{-1}K)\varphi = S^{-1}r$  where  $S^{-1}K : \mathcal{H} \rightarrow \mathcal{H}$  is compact. So that we are back to Equation (7.1), see Corollary 3.6. of Kress (1999).

This section is organized in the following way. The next paragraph recalls the main mathematical properties of the equations of the second kind. The two following paragraphs present the statistical properties of the solution in the cases of well-posed and ill-posed problems, and the last paragraph applies these results to the two examples given in Section 1.

The implementation of the estimation procedures is not discussed here because it is similar to the implementation of the estimation of a regularized equation of the first kind (see Section 3). Actually, regularizations transform first kind equations into second kind equations and the numerical methods are then formally equivalent, even though the statistical properties are fundamentally different.

## 7.2. Riesz theory and Fredholm alternative

We first briefly recall the main results about equations of the second kind as they were developed at the beginning of the 20th century by Fredholm and Riesz. The statements are given without proofs (see e.g. Kress, 1999, Chapters 3 and 4).

Let  $K$  be a compact operator from  $\mathcal{H}$  to  $\mathcal{H}$  and  $I$  be the identity on  $\mathcal{H}$  (which is compact only if  $\mathcal{H}$  is finite dimensional). Then, the operator  $I - K$  has a finite dimensional null space and its range is closed. Moreover,  $I - K$  is injective if and only if it is surjective. In that case  $I - K$  is invertible and its inverse  $(I - K)^{-1}$  is a bounded operator.

An element of the null space of  $I - K$  verifies  $K\varphi = \varphi$ , and if  $\varphi \neq 0$ , it is an eigenfunction of  $K$  associated with the eigenvalue equal to 1. Equivalently, the inverse problem (7.1) is well-posed if and only if 1 is not an eigenvalue of  $K$ . The Fredholm alternative follows from the previous results.

**Theorem 7.1 (Fredholm alternative).** *Let us consider the two equations of the second kind:*

$$(I - K)\varphi = r \tag{7.4}$$

and

$$(I - K^*)\psi = s \tag{7.5}$$

where  $K^*$  is the adjoint of  $K$ . Then:

- i) Either the two homogeneous equations  $(I - K)\varphi = 0$  and  $(I - K^*)\psi = 0$  only have the trivial solutions  $\varphi = 0$  and  $\psi = 0$ . In that case, (7.4) and (7.5) have a unique solution for any  $r$  and  $s$  in  $\mathcal{H}$
- ii) or the two homogeneous equations  $(I - K)\varphi = 0$  and  $(I - K^*)\psi = 0$  have the same finite number  $m$  of linearly independent solutions  $\varphi_j$  and  $\psi_j$  ( $j = 1, \dots, m$ ) respectively, and the solutions of (7.4) and (7.5) exist if and only if  $\langle \psi_j, r \rangle = 0$  and  $\langle \varphi_j, s \rangle = 0$  for any  $j = 1, \dots, m$ .

(ii) means that the null spaces of  $I - K$  and  $I - K^*$  are finite dimensional and have same dimensions. Moreover, the ranges of  $I - K$  and  $I - K^*$  satisfy

$$\begin{aligned}\mathcal{R}(I - K) &= \mathcal{N}(I - K^*)^\perp, \\ \mathcal{R}(I - K^*) &= \mathcal{N}(I - K)^\perp.\end{aligned}$$

### 7.3. Well-posed equations of the second kind

In this subsection, we assume that  $I - K$  is injective. In this case, the problem is well-posed and the asymptotic properties of the solution are easily deduced from the properties of the estimation of the operator  $K$  and of the right-hand side  $r$ .

The starting point of this analysis is the relation:

$$\begin{aligned}\hat{\varphi}_n - \varphi_0 &= (I - \hat{K}_n)^{-1} \hat{r}_n - (I - K)^{-1} r \\ &= (I - \hat{K}_n)^{-1} (\hat{r}_n - r) + \left[ (I - \hat{K}_n)^{-1} - (I - K)^{-1} \right] r \\ &= (I - \hat{K}_n)^{-1} \left[ \hat{r}_n - r + (\hat{K}_n - K)(I - K)^{-1} r \right] \\ &= (I - \hat{K}_n)^{-1} \left[ \hat{r}_n - r + (\hat{K}_n - K) \varphi_0 \right]\end{aligned}\tag{7.6}$$

where the third equality follows from  $A^{-1} - B^{-1} = A^{-1}(B - A)B^{-1}$ .

**Theorem 7.2.** *If*

$$i) \quad \left\| \hat{K}_n - K \right\| = o(1)$$

$$ii) \quad \left\| (\hat{r}_n + \hat{K}_n \varphi_0) - (r + K \varphi_0) \right\| = O\left(\frac{1}{a_n}\right)$$

$$\text{Then } \|\hat{\varphi}_n - \varphi_0\| = O\left(\frac{1}{a_n}\right)$$



**Proof.** As  $I - K$  is invertible and admits a continuous inverse, i) implies that  $\|(I - \hat{K}_n)^{-1}\|$  converges to  $\|(I - K)^{-1}\|$  and the result follows from (7.6). ■

In some cases  $\|r - \hat{r}_n\| = O(\frac{1}{b_n})$  and  $\|\hat{K}_n - K\| = O(\frac{1}{d_n})$ . Then  $\frac{1}{a_n} = \frac{1}{b_n} + \frac{1}{d_n}$ . In some particular examples, as will be illustrated in the last subsection, the asymptotic behavior of  $\hat{r}_n - \hat{K}_n\varphi$  is directly considered.

Asymptotic normality can be obtained from different sets of assumptions. The following theorems illustrate two kinds of asymptotic normality.

**Theorem 7.3.** *If*

$$i) \|\hat{K}_n - K\| = o(1)$$

$$ii) a_n \left( (\hat{r}_n + \hat{K}_n\varphi_0) - (r + K\varphi_0) \right) \Longrightarrow \mathcal{N}(0, \Sigma) \text{ (weak convergence in } \mathcal{H})$$

Then

$$a_n(\hat{\varphi}_n - \varphi_0) \Longrightarrow \mathcal{N}(0, (I - K)^{-1}\Sigma(I - K^*)^{-1}).$$

**Proof.** The proof follows immediately from (7.6) and Theorem 2.47. ■

**Theorem 7.4.** *We consider the case where  $\mathcal{H} = L^2(\mathbb{R}^p, \pi)$ . If*

$$i) \|\hat{K}_n - K\| = o(1)$$

$$ii) \exists a_n \text{ s.t } a_n \left[ (\hat{r}_n + \hat{K}_n\varphi_0) - (r + K\varphi_0) \right] (x) \xrightarrow{d} \mathcal{N}(0, \sigma^2(x)), \quad \forall x \in \mathbb{R}^p$$

$$iii) \exists b_n \text{ s.t } \frac{a_n}{b_n} = o(1) \text{ and}$$

$$b_n \hat{K} \left[ (\hat{r}_n + \hat{K}_n\varphi_0) - (r + K\varphi_0) \right] \Longrightarrow \mathcal{N}(0, \Omega) \text{ (weak convergence in } \mathcal{H})$$

Then

$$a_n(\hat{\varphi}_n - \varphi_0)(x) \xrightarrow{d} \mathcal{N}(0, \sigma^2(x)), \quad \forall x.$$

**Proof.** Using

$$(I - K)^{-1} = I + (I - K)^{-1}K$$

we deduce from (7.6):

$$\begin{aligned} a_n(\hat{\varphi}_n - \varphi_0)(x) &= a_n \left\{ (I - \hat{K}_n)^{-1} \left[ \hat{r}_n + \hat{K}_n\varphi_0 - r - K\varphi_0 \right] \right\} \\ &= a_n(\hat{r}_n + \hat{K}_n\varphi_0 - r - K\varphi_0)(x) \\ &+ \frac{a_n}{b_n} \left\{ b_n(I - \hat{K}_n)^{-1} \hat{K}_n(\hat{r}_n + \hat{K}_n\varphi_0 - r - K\varphi_0) \right\} (x) \end{aligned} \tag{7.7}$$

The last term in brackets converges (weakly in  $L^2$ ) to a  $\mathcal{N}(0, (I - K)^{-1}\Omega(I - K)^{-1})$  and the value of this function at any point  $x$  also converges to a normal distribution (weak convergence implies finite dimensional convergence). Then the last term in brackets is bounded and the result is verified. ■

Note that condition (iii) is satisfied as soon as premultiplying by  $K$  increases the rate of convergence of  $\hat{r}_n + \hat{K}_n\varphi$ . This is true in particular if  $K$  is an integral operator.

We illustrate these results by the following three examples. The first example is an illustrative example, while the other two are motivated by relevant econometric issues.

**Example.** Consider  $L^2(\mathbb{R}, \Pi)$  and  $(Y, Z)$  is a random element of  $\mathbb{R} \times L^2(\mathbb{R}, \Pi)$ . We study the integral equation of the second kind defined by

$$\varphi(x) + \int E^F(Z(x)Z(y))\varphi(y)\Pi(dy) = E^F(YZ(x)) \quad (7.8)$$

denoted by  $\varphi + V\varphi = r$ . Here  $K = -V$ . As the covariance operator,  $V$  is a positive operator,  $K$  is a negative operator and therefore 1 can not be an eigenvalue of  $K$ . Consequently, Equation (7.8) defines a well-posed inverse problem.

We assume that an i.i.d. sample of  $(Y, Z)$  is available and the estimated equation (7.8) defines the parameter of interest as the solution of an integral equation having the following form:

$$\varphi(x) + \frac{1}{n} \sum_{i=1}^n z_i(x) \int z_i(y)\varphi(y)\Pi(dy) = \frac{1}{n} \sum_{i=1}^n y_i z_i(x) \quad (7.9)$$

Under some standard regularity conditions, one can check that  $\|\hat{V}_n - V\| = O\left(\frac{1}{\sqrt{n}}\right)$  and that

$$\begin{aligned} & \sqrt{n} \frac{1}{n} \sum_i \left\{ z_i(\cdot) \left[ y_i - \int z_i(y)\varphi(y)\Pi(dy) \right] - E^F(YZ(\cdot)) - \int E^F(Z(\cdot)Z(y))\varphi(y)\Pi(dy) \right\} \\ & \Rightarrow \mathcal{N}(0, \Sigma) \text{ in } L^2(\mathbb{R}, \Pi). \end{aligned}$$

Suppose for instance that  $E^F(Y|Z) = \int Z(y)\varphi(y)\Pi(dy)$ . Under a homoscedasticity hypothesis, the operator  $\Sigma$  is a covariance operator with kernel  $\sigma^2 E^F(Z(x)Z(y))$  where

$$\sigma^2 = \text{Var} \left( Y - \int Z(y)\varphi(y)\Pi(dy) | Z \right).$$

Then, from Theorem 7.3,

$$\sqrt{n}(\hat{\varphi}_n - \varphi_0) \Rightarrow \mathcal{N}(0, \sigma^2(I + V)^{-1}V(I + V)^{-1}).$$

**Example. Rational expectations asset pricing models:**

Following Lucas (1978), rational expectations models characterize the pricing functional as a function  $\varphi$  of the Markov state solution of an integral equation:

$$\varphi(x) - \int a(x, y)\varphi(y) f(y|x) dy = \int a(x, y)b(y) f(y|x) dy \quad (7.10)$$

While  $f$  is the transition density of the Markov state, the function  $a$  denotes the marginal rate of substitution and  $b$  the dividend function. For the sake of expositional simplicity, we assume here that the functions  $a$  and  $b$  are both known while  $f$  is estimated nonparametrically by a kernel method. Note that if the marginal rate of substitution  $a$  involves some unknown preference parameters (subjective discount factor, risk aversion parameter), they will be estimated, for instance by GMM, with a parametric root  $n$  rate of convergence. Therefore, the nonparametric inference about  $\varphi$  (deduced from the solution of (7.10) using a kernel estimation of  $f$ ) is not contaminated by this parametric estimation; all the statistical asymptotic theory can be derived as if the preference parameters were known.

As far as kernel density estimation is concerned, it is well known that under mild conditions (see e.g. Bosq (1998)) it is possible to get the same convergence rates and the same asymptotic distribution with stationary strongly mixing stochastic processes as in the i.i.d. case.

Let us then consider a  $n$ -dimensional stationary stochastic process  $X_t$  and  $\mathcal{H}$  the space of square integrable functions of one realization of this process. In this example,  $\mathcal{H}$  is defined with respect to the true distribution. The operator  $K$  is defined by

$$K\varphi(x) = E^F(a(X_{t-1}, X_t)\varphi(X_t) | X_{t-1} = x)$$

and

$$r(x) = E^F(a(X_{t-1}, X_t)b(X_t) | X_{t-1} = x)$$

We will assume that  $K$  is compact through possibly a Hilbert-Schmidt condition (see Assumption A.1 of Section 5.5 for such a condition). A common assumption in rational expectation models is that  $K$  is a contraction mapping, due to discounting. Then, 1 is not an eigenvalue of  $K$  and (7.10) is a well-posed Fredholm integral equation.

Under these hypotheses, both numerical and statistical issues associated with the solution of (7.10) are well documented. See Rust, Traub and Wozniakowski (2002) and references therein for numerical issues. The statistical consistency of the estimator  $\hat{\varphi}_n$  obtained from the kernel estimator  $\hat{K}_n$  is deduced from Theorem 7.2 above. Assumption *i*) is satisfied because  $\hat{K}_n - K$  has the same behavior as the conditional expectation operator and

$$\begin{aligned} \hat{r}_n + \hat{K}_n\varphi - r - K\varphi &= E^{F_n}(a(X_{t-1}, X_t)(b(X_t) + \varphi(X_t)) | X_{t-1}) \\ &\quad - E^F(a(X_{t-1}, X_t)(b(X_t) + \varphi(X_t)) | X_{t-1}) \end{aligned}$$

converges at the speed  $\frac{1}{a_n} = \left(\frac{1}{nc_n^m} + c_n^4\right)^{1/2}$  if  $c_n$  is the bandwidth of the (second order) kernel estimator and  $m$  is the dimension of  $X$ .

The weak convergence follows from Theorem 7.4. Assumption ii) of Theorem 7.4 is the usual result on the normality of kernel estimation of conditional expectation. As  $K$  is an integral operator, the transformation by  $K$  increases the speed of convergence, which implies iii) of Theorem 7.4.

**Example. Partially Nonparametric forecasting model:**

This example is drawn from Linton and Mammen (2003). Nonparametric prediction of a stationary ergodic scalar random process  $X_t$  is often performed by looking for a predictor  $m(X_{t-1}, \dots, X_{t-d})$  able to minimize the mean square error of prediction:

$$E[X_t - m(X_{t-1}, \dots, X_{t-d})]^2$$

In other words, if  $m$  can be any squared integrable function, the optimal predictor is the conditional expectation

$$m_0(X_{t-1}, \dots, X_{t-d}) = E[X_t | X_{t-1}, \dots, X_{t-d}]$$

and can be estimated by kernel smoothing or any other nonparametric way of estimating a regression function. The problems with this kind of approach are twofold. First, it is often necessary to include many lagged variables and the resulting nonparametric estimation surface suffers from the well-known ‘‘curse of dimensionality’’. Second, it is hard to describe and interpret the estimated regression surface when the dimension is more than two.

A solution to deal with these problems is to think about a kind of nonparametric generalization of ARMA processes. For this purpose, let us consider semiparametric predictors of the following form

$$E[X_t | I_{t-1}] = m_\varphi(\theta, I_{t-1}) = \sum_{j=1}^{\infty} a_j(\theta) \varphi(X_{t-j}) \quad (7.11)$$

where  $\theta$  is an unknown finite dimensional vector of parameters,  $a_j(\cdot)$ ,  $j \geq 1$  are known scalar functions, and  $\varphi(\cdot)$  is the unknown functional parameter of interest. The notation

$$E[X_t | I_{t-1}] = m_\varphi(\theta, I_{t-1})$$

stresses the fact that the predictor depends on the true unknown value of the parameters  $\theta$  and  $\varphi$ , and of the information  $I_{t-1}$  available at time  $(t - 1)$  as well. This information is actually the  $\sigma$ -field generated by  $X_{t-j}$ ,  $j \geq 1$ . A typical example is

$$a_j(\theta) = \theta^{j-1} \text{ for } j \geq 1 \text{ with } 0 < \theta < 1. \quad (7.12)$$

Then the predictor defined in (7.11) is actually characterized by

$$m_\varphi(\theta, I_{t-1}) = \theta m_\varphi(\theta, I_{t-2}) + \varphi(X_{t-1}) \quad (7.13)$$

In the context of volatility modelling,  $X_t$  would denote a squared asset return over period  $[t-1, t]$  and  $m_\varphi(\theta, I_{t-1})$  the so-called squared volatility of this return as expected at the beginning of the period. Engle and Ng (1993) have studied such a partially nonparametric (PNP for short) model of volatility and called the function  $\varphi$  the “news impact function”. They proposed an estimation strategy based on piecewise linear splines. Note that the PNP model includes several popular parametric volatility models as special cases. For instance, the GARCH (1,1) model of Bollerslev (1986) corresponds to  $\varphi(x) = w + \alpha x$  while the Engle (1990) asymmetric model is obtained for  $\varphi(x) = w + \alpha(x + \delta)^2$ . More examples can be found in Linton and Mammen (2003).

The nonparametric identification and estimation of the news impact function can be derived for a given value of  $\theta$ . After that, a profile criterion can be calculated to estimate  $\theta$ . In any case, since  $\theta$  will be estimated with a parametric rate of convergence, the asymptotic distribution theory of a nonparametric estimator of  $\varphi$  is the same as if  $\theta$  were known. For the sake of notational simplicity, the dependence on unknown finite dimensional parameters  $\theta$  is no longer made explicit.

At least in the particular case (7.12)-(7.13),  $\varphi$  is easily characterized as the solution of a linear integral equation of the first kind

$$E[X_t - \theta X_{t-1} | I_{t-2}] = E[\varphi(X_{t-1}) | I_{t-2}]$$

Except for its dynamic features, this problem is completely similar to the nonparametric instrumental regression example described in Section 5.5. However, as already mentioned, problems of the second kind are often preferable since they may be well-posed. As shown by Linton and Mammen (2003) in the particular case of a PNP volatility model, it is actually possible to identify and consistently estimate the function  $\varphi$  defined as

$$\varphi = \arg \min_{\varphi} E \left[ X_t - \sum_{j=1}^{\infty} a_j \varphi(X_{t-j}) \right]^2 \quad (7.14)$$

from a well-posed linear inverse problem of the second kind. When  $\varphi$  is an element of the Hilbert space  $L^2_F(X)$ , its true unknown value is characterized by the first order conditions obtained by differentiating in the direction of any vector  $h$

$$E \left[ \left( X_t - \sum_{j=1}^{\infty} a_j \varphi(X_{t-j}) \right) \left( \sum_{l=1}^{\infty} a_l h(X_{t-l}) \right) \right] = 0$$

In other words, for any  $h$  in  $L_F^2(X)$

$$\begin{aligned}
& \sum_{j=1}^{\infty} a_j E^X [E [X_t | X_{t-j} = x] h(x)] \\
& - \sum_{j=1}^{\infty} a_j^2 E^X [\varphi(x) h(x)] \\
& - \sum_{j=1}^{\infty} \sum_{\substack{l=1 \\ l \neq j}}^{\infty} a_j a_l E^X [E [\varphi(X_{t-l}) | X_{t-j} = x] h(x)] = 0
\end{aligned} \tag{7.15}$$

where  $E^X$  denotes the expectation with respect to the stationary distribution of  $X_t$ . As the equality in (7.15) holds true for all  $h$ , it is true in particular for a complete sequence of functions of  $L_F^2(X)$ . It follows that

$$\begin{aligned}
& \sum_{j=1}^{\infty} a_j E [X_t | X_{t-j} = x] - \left( \sum_{l=1}^{\infty} a_l^2 \right) \varphi(x) \\
& - \sum_{j=1}^{\infty} \sum_{l \neq j}^{\infty} a_j a_l E [\varphi(X_{t-l}) | X_{t-j} = x] = 0
\end{aligned}$$

$P^X$  – almost surely on the values of  $x$ . Let us denote

$$r_j(X_t) = E [X_{t+j} | X_t] \quad \text{and} \quad H_k(\varphi)(X_t) = E [\varphi(X_{t+k}) | X_t].$$

Then, we have proved that the unknown function  $\varphi$  of interest must be the solution of the linear inverse problem of the second kind

$$A(\varphi, F) = (I - K) \varphi - r = 0 \tag{7.16}$$

where

$$\begin{aligned}
r &= \left( \sum_{j=1}^{\infty} a_j^2 \right)^{-1} \sum_{j=1}^{\infty} a_j r_j, \\
K &= - \left( \sum_{j=1}^{\infty} a_j^2 \right)^{-1} \sum_{j=1}^{\infty} \sum_{l \neq j}^{\infty} a_j a_l H_{j-l},
\end{aligned}$$

and, with a slight change of notation,  $F$  now characterizes the probability distribution of the stationary process  $(X_t)$ .

To study the inverse problem (7.16), it is first worth noticing that  $K$  is a self-adjoint integral operator. Indeed, while

$$K = \left( \sum_{j=1}^{\infty} a_j^2 \right)^{-1} \sum_{h=\pm 1}^{+\infty} H_h \left( \sum_{l=\max[1, 1-h]}^{+\infty} a_l a_{l+k} \right)$$

we immediately deduce from Subsection 2.5.1 that the conditional expectation operator  $H_k$  is such that

$$H_k^* = H_{-k}$$

and thus  $K = K^*$ , since

$$\sum_{l=\max[1,1-k]}^{+\infty} a_l a_{l+k} = \sum_{l=\max[1,1+k]}^{+\infty} a_l a_{l-k}$$

As noticed by Linton and Mammen (2003), this property greatly simplifies the practical implementation of the solution of the sample counterpart of equation (7.19). Even more importantly, the inverse problem (7.19) will be well-posed as soon as one maintains the following identification assumption about the news impact function  $\varphi$ .

**Assumption A.** There exists no  $\theta$  and  $\varphi \in L_F^2(X)$  with  $\varphi \neq 0$  such that  $\sum_{j=1}^{\infty} a_j(\theta) \varphi(X_{t-j}) = 0$  almost surely.

To see this, observe that Assumption A means that for any non-zero function  $\varphi$

$$0 < E \left[ \sum_{j=1}^{\infty} a_j \varphi(X_{t-j}) \right]^2,$$

that is

$$0 < \sum_{j=1}^{\infty} a_j^2 \langle \varphi, \varphi \rangle + \sum_{j=1}^{\infty} \sum_{\substack{l=1 \\ l \neq j}}^{\infty} a_l a_j \langle \varphi, H_{j-l} \varphi \rangle.$$

Therefore

$$0 < \langle \varphi, \varphi \rangle - \langle \varphi, K \varphi \rangle \tag{7.17}$$

for non zero  $\varphi$ . In other words, there is no non-zero  $\varphi$  such that

$$K \varphi = \varphi$$

and the operator  $(I - K)$  is one-to-one. Moreover, (7.17) implies that  $(I - K)$  has eigenvalues bounded from below by a positive number. Therefore, if  $K$  depends continuously on the unknown finite dimensional vector of parameters  $\theta$  and if  $\theta$  evolves in some compact set, the norm of

$(I - K)^{-1}$  will be bounded from above uniformly on  $\theta$ .

It is also worth noticing that the operator  $K$  is Hilbert-Schmidt and a fortiori compact under reasonable assumptions. As already mentioned in Subsection 2.5.1, the Hilbert-Schmidt property for the conditional expectation operator  $H_k$  is tantamount to the integrability condition

$$\int \int \left[ \frac{f_{X_t, X_{t-k}}(x, y)}{f_{X_t}(x) f_{X_t}(y)} \right]^2 f_{X_t}(x) f_{X_t}(y) dx dy < \infty$$

It amounts to saying that there is not too much dependence between  $X_t$  and  $X_{t-k}$ . This should be tightly related to the ergodicity or mixing assumptions about the stationary process  $X_t$ . Then, if all the conditional expectation operator  $H_k$ ,  $k \geq 1$  are Hilbert-Schmidt, the operator  $K$  will also be Hilbert-Schmidt insofar as

$$\sum_{j=1}^{\infty} \sum_{l \neq j} a_j^2 a_l^2 < +\infty.$$

Up to a straightforward generalization to stationary mixing processes of results only stated in the i.i.d. case, the general asymptotic theory of Theorems 7.3 and 7.4 can then be easily applied to nonparametric estimators of the news impact function  $\varphi$  based on the Fredholm equation of the second kind (7.15). An explicit formula for the asymptotic variance of  $\hat{\varphi}_n$  as well as a practical solution by implementation of matricial equations similar to those of Subsection 3.4 (without need of a Tikhonov regularization) is provided by Linton and Mammen (2003) in the particular case of volatility modelling.

However, an important difference with the i.i.d. case (see for instance assumption A.3 in Section 5.5 about instrumental variables) is that the conditional homoskedasticity assumption cannot be maintained about the conditional probability distribution of  $X_t$  given its own past. This should be particularly detrimental in the case of volatility modelling, since when  $X_t$  denotes a squared return, it will in general be even more conditionally heteroskedastic than returns themselves. Such severe conditional heteroskedasticity will likely imply a poor finite sample performance, and a large asymptotic variance of the estimator  $\hat{\varphi}_n$  defined from the inverse problem (7.15), that is from the least squares problem (7.14). Indeed,  $\hat{\varphi}_n$  is a kind of OLS estimator in infinite dimension. In order to better take into account conditional heteroskedasticity of  $X_t$  in the context of volatility modelling, Linton and Mammen (2003) propose to replace the least squares problem (7.14) by a quasi-likelihood kind of approach where the criterion to optimize is defined from the density function of a normal conditional probability distribution of returns, with variance  $m_\varphi(\theta, I_{t-1})$ . Then the difficulty is that the associated first order conditions now characterize the news impact function  $\varphi$  as solution of a nonlinear inverse problem. Linton and Mammen (2003) suggest to work with a version of this problem which is locally linearized around the previously described least squares estimator  $\hat{\varphi}_n$  (and associated consistent estimator of  $\theta$ ).

#### 7.4. Ill-posed equations of the second kind

The objective of this section is to study equations  $(I - K)\varphi = r$  where 1 is an eigenvalue of  $K$ , i.e. where  $I - K$  is not injective (or one-to-one). For simplicity, we restrict our analysis to the case where the order of multiplicity of the eigenvalue 1 is one and the operator  $K$  is self-adjoint. This implies that the dimension of the null spaces of  $I - K$  is one and using the results of Section 7.2, the space  $\mathcal{H}$  may be decomposed into

$$\mathcal{H} = \mathcal{N}(I - K) \oplus \mathcal{R}(I - K)$$



i.e.  $\mathcal{H}$  is the direct sum between the null space and the range of  $I - K$ , both closed. We denote by  $P_{\mathcal{N}}r$  the projection of  $r$  on  $\mathcal{N}(I - K)$  and by  $P_{\mathcal{R}}r$  the projection of  $r$  on the range  $\mathcal{R}(I - K)$ .

Using ii) of Theorem 7.1, a solution of  $(I - K)\varphi = r$  exists in the non injective case only if  $r$  is orthogonal to  $\mathcal{N}(I - K)$  or equivalently, if  $r$  belongs to  $\mathcal{R}(I - K)$ . In other words, a solution exists if and only if  $r = P_{\mathcal{R}}r$ . However in this case, the solution is not unique and there exists a one dimensional linear manifold of solutions. Obviously, if  $\varphi$  is a solution,  $\varphi$  plus any element of  $\mathcal{N}(I - K)$  is also a solution. This non uniqueness problem will be solved by a normalization rule which selects a unique element in the set of solutions. The normalization we adopt is

$$\langle \varphi, \phi_1 \rangle = 0 \tag{7.18}$$

where  $\phi_1$  is the eigenfunction of  $K$  corresponding to the eigenvalue equal to 1.

In most statistical applications of equations of the second kind, the  $r$  element corresponding to the true data generating process is assumed to be in the range of  $I - K$  where  $K$  is also associated with the true DGP. However, this property is no longer true if  $F$  is estimated and we need to extend the resolution of  $(I - K)\varphi = r$  to cases where  $I - K$  is not injective and  $r$  is not in the range of this operator. This extension must be done in such a way that the continuity properties of inversion are preserved.

For this purpose we consider the following generalized inverse of  $(I - K)$ . As  $K$  is a compact operator, it has a discrete spectrum  $\lambda_1 = 1, \lambda_2, \dots$  where only 0 may be an accumulation point (in particular 1 cannot be an accumulation point). The associated eigenfunctions are  $\phi_1, \phi_2, \dots$ . Then we define:

$$Lu = \sum_{j=2}^{\infty} \frac{1}{1 - \lambda_j} \langle u, \phi_j \rangle \phi_j, \quad u \in \mathcal{H} \tag{7.19}$$

Note that  $L = (I - K)^\dagger$  is the generalized inverse of  $I - K$ , introduced in Section 3. Moreover,  $L$  is continuous and therefore bounded because 1 is an isolated eigenvalue. This operator computes the unique solution of  $(I - K)\varphi = P_{\mathcal{R}}u$  satisfying the normalization rule (7.18). It can be easily verified that  $L$  satisfies:

$$\begin{aligned} LP_{\mathcal{R}} &= L = P_{\mathcal{R}}L, \\ L(I - K) &= (I - K)L = P_{\mathcal{R}}. \end{aligned} \tag{7.20}$$

We now consider estimation. For an observed sample, we obtain the estimator  $F_n$  of  $F$  (that may be built from a kernel estimator of the density) and then the estimators  $\hat{r}_n$  and  $\hat{K}_n$  of  $r$  and  $K$  respectively. Let  $\hat{\phi}_1, \hat{\phi}_2, \dots$  denote the eigenfunctions of  $\hat{K}_n$  associated with  $\hat{\lambda}_1, \hat{\lambda}_2, \dots$ . We restrict our attention to the cases where 1 is also an eigenvalue of multiplicity one of  $\hat{K}_n$  (i.e.  $\hat{\lambda}_1 = 1$ ). However,  $\hat{\phi}_1$  may be different from  $\phi_1$ .

We have to make a distinction between two cases. First, assume that the Hilbert space  $\mathcal{H}$  of reference is known and in particular the inner product is given (for example  $\mathcal{H} = L^2(\mathbb{R}^p, \Pi)$  with  $\Pi$  given). The normalization rule imposed to  $\hat{\varphi}_n$  is

$$\langle \hat{\varphi}_n, \hat{\phi}_1 \rangle = 0$$

and  $\hat{L}_n$  is the generalized inverse of  $I - \hat{K}_n$  in  $\mathcal{H}$  (which depends on the Hilbert space structure) where

$$\hat{L}_n u = \sum_{j=2}^{\infty} \frac{1}{1 - \hat{\lambda}_j} \langle u, \hat{\phi}_j \rangle \hat{\phi}_j, \quad u \in \mathcal{H}$$

Formula (7.20) applies immediately for  $F_n$ .

However, if the Hilbert space  $\mathcal{H}$  depends on  $F$  (e.g.  $\mathcal{H} = L^2(\mathbb{R}^p, F)$ ), we need to assume that  $L^2(\mathbb{R}, F_n) \subset L^2(\mathbb{R}^p, F)$ . The orthogonality condition, which defines the normalization rule (7.18) is related to  $L^2(\mathbb{R}^p, F)$  but the estimator  $\hat{\varphi}_n$  of  $\varphi$  will be normalized by

$$\langle \hat{\varphi}_n, \hat{\phi}_1 \rangle_n = 0$$

where  $\langle \cdot, \cdot \rangle_n$  denotes the inner product relative to  $F_n$ . This orthogonality is different from an orthogonality relative to  $\langle \cdot, \cdot \rangle$ . In the same way  $\hat{L}_n$  is now defined as the generalized inverse of  $I - \hat{K}_n$  with respect to the estimated Hilbert structure, i.e.

$$\hat{L}_n u = \sum_{j=2}^{\infty} \frac{1}{1 - \hat{\lambda}_j} \langle u, \hat{\phi}_j \rangle_n \hat{\phi}_j$$

and  $\hat{L}_n$  is not the generalized inverse of  $I - \hat{K}_n$  in the original space  $\mathcal{H}$ . The advantages of this definition are that  $\hat{L}_n$  may be effectively computed and satisfies the formula (7.20) where  $F_n$  replaces  $F$ . In the sequel  $P_{\mathcal{R}_n}$  denotes the projection operator on  $\mathcal{R}_n = \mathcal{R}(I - \hat{K}_n)$  for the inner product  $\langle \cdot, \cdot \rangle_n$ .

To establish consistency, we will use the following equality.

$$\begin{aligned} \hat{L}_n - L &= \hat{L}_n(\hat{K}_n - K)L \\ &+ \hat{L}_n(P_{\mathcal{R}_n} - P_{\mathcal{R}}) + (P_{\mathcal{R}_n} - P_{\mathcal{R}})L. \end{aligned} \quad (7.21)$$

It follows from (7.20) and  $\hat{L}_n - L = \hat{L}_n P_{\mathcal{R}_n} - P_{\mathcal{R}} L = \hat{L}_n(P_{\mathcal{R}_n} - P_{\mathcal{R}}) + (P_{\mathcal{R}_n} - P_{\mathcal{R}})L - P_{\mathcal{R}_n} L + \hat{L}_n P_r$  and  $\hat{L}_n(K_n - K)L = \hat{L}_n(K_n - I)L + \hat{L}_n(I - K)L = P_{\mathcal{R}_n} L + \hat{L}_n P_r$ .

The convergence property is given by the following theorem.

**Theorem 7.5.** *Let us define  $\varphi_0 = Lr$  and  $\hat{\varphi}_n = \hat{L}_n \hat{r}_n$ . If*

- i)  $\left\| \hat{K}_n - K \right\| = o(1)$
- ii)  $\|P_{\mathcal{R}_n} - P_{\mathcal{R}}\| = O\left(\frac{1}{b_n}\right)$
- iii)  $\left\| (\hat{r}_n + \hat{K}_n \varphi_0) - (r + K \varphi_0) \right\| = O\left(\frac{1}{a_n}\right)$

Then

$$\|\hat{\varphi}_n - \varphi_0\| = O\left(\frac{1}{a_n} + \frac{1}{b_n}\right).$$

**Proof.** The proof is based on:

$$\begin{aligned}\hat{\varphi}_n - \varphi_0 &= \hat{L}_n \hat{r}_n - Lr \\ &= \hat{L}_n(\hat{r}_n - r) + (\hat{L}_n - L)r \\ &= \hat{L}_n(\hat{r}_n - r) + \hat{L}_n(\hat{K}_n - K)\varphi_0 \\ &+ \hat{L}_n(P_{\mathcal{R}_n} - P_{\mathcal{R}})r + (P_{\mathcal{R}_n} - P_{\mathcal{R}})\varphi_0\end{aligned}\tag{7.22}$$

deduced from (7.21). Then

$$\begin{aligned}\|\hat{\varphi}_n - \varphi_0\| &\leq \|\hat{L}_n\| \|(\hat{r}_n + \hat{K}_n\varphi_0) - (r + K\varphi_0)\| \\ &+ (\|\hat{L}_n\| \|r\| + \|\varphi\|) \|P_{\mathcal{R}_n} - P_{\mathcal{R}}\|\end{aligned}\tag{7.23}$$

Under i) and ii)  $\|\hat{L}_n - L\| = o(1)$  from (7.21). This implies  $\|\hat{L}_n\| \rightarrow \|L\|$  and the result follows. ■

If  $\frac{a_n}{b_n} = O(1)$ , the actual speed of convergence is bounded by  $\frac{1}{a_n}$ . This will be the case in the two examples of 7.5 where  $\frac{a_n}{b_n} \rightarrow 0$ .

We consider asymptotic normality in this case. By (7.20), we have  $\hat{L}_n = P_{\mathcal{R}_n} + \hat{L}_n \hat{K}_n$ , hence:

$$\begin{aligned}\hat{\varphi}_n - \varphi_0 &= P_{\mathcal{R}_n} \left[ (\hat{r}_n + \hat{K}_n\varphi_0) - (r + K\varphi_0) \right]\end{aligned}\tag{7.24}$$

$$+ \hat{L}_n \hat{K}_n \left[ (\hat{r}_n + \hat{K}_n\varphi_0) - (r + K\varphi_0) \right]\tag{7.25}$$

$$+ \hat{L}_n(P_{\mathcal{R}_n} - P_{\mathcal{R}})r + (P_{\mathcal{R}_n} - P_{\mathcal{R}})\varphi_0\tag{7.26}$$

Let us assume that there exists a sequence  $a_n$  such that i) and ii) below are satisfied

i)  $a_n P_{\mathcal{R}_n} \left[ (\hat{r}_n + \hat{K}_n\varphi_0) - (r + K\varphi_0) \right] (x)$  has an asymptotic normal distribution,

ii)  $a_n \left[ \hat{L}_n \hat{K}_n (\hat{r}_n + \hat{K}_n\varphi_0) - r - K\varphi_0 \right] (x) \rightarrow 0$ ,  $a_n \left[ \hat{L}_n (P_{\mathcal{R}_n} - P_{\mathcal{R}}) r \right] (x) \rightarrow 0$ ,

and  $a_n [(P_{\mathcal{R}_n} - P_{\mathcal{R}})\varphi_0] (x) \rightarrow 0$  in probability.

Then the asymptotic normality of  $a_n(\hat{\varphi}_n - \varphi_0)$  is driven by the behavior of (7.24). This situation occurs in the nonparametric estimation, as illustrated in the next section.

## 7.5. Two examples: backfitting estimation in additive models and panel model

### 7.5.1. Backfitting estimation in additive models

Using the notation of Subsection 1.3.5, an additive model is defined by

$$\begin{aligned} (Y, Z, W) &\in \mathbb{R} \times \mathbb{R}^p \times \mathbb{R}^q \\ Y &= \varphi(Z) + \psi(W) + U \\ E(U|Z, W) &= 0, \end{aligned} \tag{7.27}$$

in which case (see (1.24)), the function  $\varphi$  is solution of the equation:

$$\varphi - E[E(\varphi(Z)|W)|Z] = E(Y|Z) - E[E(Y|W)|Z]$$

and  $\psi$  is the solution of an equation of the same nature obtained by a permutation of  $W$  and  $Z$ . The backfitting algorithm of Breiman and Friedman (1985), and Hastie and Tibshirani (1990) is widely used to estimate  $\varphi$  and  $\psi$  in Equation (7.27). Mammen, Linton, and Nielsen (1999) derive the asymptotic distribution of the backfitting procedure. Alternatively, Newey (1994), Tjostheim and Auestad (1994), and Linton and Nielsen (1995) propose to estimate  $\varphi$  (respectively  $\psi$ ) by integrating an estimator of  $E[Y|Z = z, W = w]$  with respect to  $w$  (respectively  $z$ ).

We focus our presentation on the estimation of  $\varphi$ . It appears as the result of a linear equation of the second kind. More precisely, we have in that case:

- $\mathcal{H}$  is the space of the square integrable functions of  $Z$  with respect to the true data generating process. This definition simplifies our presentation but an extension to different spaces is possible.
- The unknown function  $\varphi$  is an element of  $\mathcal{H}$ . Actually, asymptotic considerations will restrict the class of functions  $\varphi$  by smoothness restrictions.
- The operator  $K$  is defined by  $K\varphi = E[E(\varphi(Z)|W)|Z]$ . This operator is self adjoint and we assume its compactness. This compactness may be obtained through the Hilbert Schmidt Assumption A.1 of Subsection 5.5.
- The function  $r$  is equal to  $E(Y|Z) - E[E(Y|W)|Z]$ .

The operator  $I - K$  is not one-to-one because the constant functions belong to the null space of this operator. Indeed, the additive model (7.27) does not identify  $\varphi$  and  $\psi$ . We introduce the following assumption (see Florens, Mouchart, and Rolin (1990)), which warrants that  $\varphi$  and  $\psi$  are exactly identified up to an additive constant, or equivalently that the null space of  $I - K$  only contains the constants (meaning 1 is an eigenvalue of  $K$  of order 1).

**Identification assumption.**  $Z$  and  $W$  are measurably separated w.r.t. the distribution  $F$ , i.e. a function of  $Z$  almost surely equal to a function of  $W$  is almost surely constant.

This assumption implies that if  $\varphi_1, \varphi_2, \psi_1, \psi_2$  are such that  $E(Y|Z, W) = \varphi_1(Z) + \psi_1(W) = \varphi_2(Z) + \psi_2(W)$  then  $\varphi_1(Z) - \varphi_2(Z) = \psi_2(W) - \psi_1(W)$  which implies that  $\varphi_1 - \varphi_2$  and  $\psi_2 - \psi_1$  are a.s. constant. In terms of the null set of  $I - K$  we have:

$$\begin{aligned} K\varphi &= \varphi \\ \iff E[E(\varphi(Z)|W)|Z] &= \varphi(Z) \\ \implies E[(E[\varphi(Z)|W])^2] &= E[\varphi(Z)E(\varphi(Z)|W)] \\ &= E(\varphi^2(Z)). \end{aligned}$$

But, by Pythagore theorem

$$\begin{aligned} \varphi(Z) &= E(\varphi(Z)|W) + v \\ E(\varphi^2(Z)) &= E((E(\varphi(Z)|W))^2) + Ev^2. \end{aligned}$$

Then:

$$\begin{aligned} K\varphi &= \varphi \implies v = 0 \\ \Leftrightarrow \varphi(Z) &= E[\varphi(Z)|W]. \end{aligned}$$

Then, if  $\varphi$  is an element of the null set of  $I - K$ ,  $\varphi$  is almost surely equal to a function of  $W$  and is therefore constant.

The eigenvalues of  $K$  are real, positive and smaller than 1 except for the first one, that is  $1 = \lambda_1 > \lambda_2 > \lambda_3 > \dots$ <sup>1</sup> The eigenfunctions are such that  $\phi_1 = 1$  and the condition  $\langle \varphi, \phi_1 \rangle = 0$  means that  $\varphi$  has an expectation equal to zero. The range of  $I - K$  is the set of functions with mean equal to 0 and the projection of  $u$ ,  $P_{\mathcal{R}}u$ , equals  $u - E(u)$ .

It should be noticed that under the hypothesis of the additive model,  $r$  has zero mean and is then an element of  $\mathcal{R}(I - K)$ . Then, a unique (up to the normalization condition) solution of the structural equation  $(I - K)\varphi = r$  exists.

The estimation may be done by kernel smoothing. The joint density is estimated by

$$f_n(y, z, w) = \frac{1}{nc_n^{1+p+q}} \sum_{i=1}^n \omega\left(\frac{y - y_i}{c_n}\right) \omega\left(\frac{z - z_i}{c_n}\right) \omega\left(\frac{w - w_i}{c_n}\right) \quad (7.28)$$

and  $F_n$  is the c.d.f. associated to  $f_n$ . The estimated  $\hat{K}_n$  operator verifies:

$$(\hat{K}_n\varphi)(z) = \int \varphi(u) \hat{a}_n(u, z) du \quad (7.29)$$

where

$$\hat{a}_n(u, z) = \int \frac{\hat{f}_n(\cdot, u, w) \hat{f}_n(\cdot, z, w)}{\hat{f}_n(\cdot, \cdot, w) \hat{f}_n(\cdot, z, \cdot)} dw.$$

---

<sup>1</sup>Actually  $K = T^*T$  when  $T\varphi = E(\varphi|W)$  and  $T^*\psi = E(\psi|Z)$  when  $\psi$  is a function of  $W$ . The eigenvalues of  $K$  correspond to the squared singular values of the  $T$  and  $T^*$  defined in Section 2.

The operator  $\hat{K}_n$  must be an operator from  $\mathcal{H}$  to  $\mathcal{H}$  (it is by construction an operator from  $L_Z^2(F_n)$  into  $L_Z^2(F_n)$ ). Therefore,  $\frac{\omega\left(\frac{z-z_\ell}{c_n}\right)}{\sum_\ell \omega\left(\frac{z-z_\ell}{c_n}\right)}$  must be square integrable w.r.t.  $F$ .

The estimation of  $r$  by  $\hat{r}_n$  verifies

$$\hat{r}_n(z) = \frac{1}{\sum_{\ell=1}^n \omega\left(\frac{z-z_\ell}{c_n}\right)} \sum_{\ell=1}^n \left( y_\ell - \sum_{i=1}^n y_i \omega_{\ell i} \right) \omega\left(\frac{z-z_\ell}{c_n}\right)$$

where  $\omega_{\ell i} = \frac{\omega\left(\frac{w_\ell - w_i}{c_n}\right)}{\sum_{j=1}^n \omega\left(\frac{w_\ell - w_j}{c_n}\right)}$ .

The operator  $\hat{K}_n$  also has 1 as the greatest eigenvalue corresponding to an eigenfunction equal to 1. Since  $F_n$  is a mixture of probabilities for which  $Z$  and  $W$  are independent, the measurable separability between  $Z$  and  $W$  is fulfilled. Then, the null set of  $I - \hat{K}_n$  reduces a.s. (w.r.t.  $F_n$ ) to constant functions. The generalized inverse of an operator depends on the inner product of the Hilbert space because it is defined as the function  $\varphi$  of minimal norm which minimizes the norm of  $\hat{K}_n \varphi - \hat{r}_n$ . The generalized inverse in the space  $L_Z^2(F)$  cannot be used for the estimation because it depends on the actual unknown  $F$ . Then we construct  $\hat{L}_n$  as the generalized inverse in  $L_Z^2(F_n)$  of  $I - \hat{K}_n$ . The practical computation of  $\hat{L}_n$  can be done by computing the  $n$  eigenvalues of  $\hat{K}_n$ ,  $\hat{\lambda}_1 = 1, \dots, \hat{\lambda}_n$  and the  $n$  eigenfunctions  $\hat{\phi}_1 = 1, \hat{\phi}_2, \dots, \hat{\phi}_n$ . Then

$$\hat{L}_n u = \sum_{j=2}^n \frac{1}{1 - \hat{\lambda}_j} \left\{ \int u(z) \hat{\phi}_j(z) \hat{f}_n(z) dz \right\} \hat{\phi}_j$$

It can be easily checked that property (7.20) is verified where  $P_{\mathcal{R}_n}$  is the projection (w.r.t.  $F_n$ ) on the orthogonal of the constant function. This operator subtracts from any function its empirical mean, which is computed through the smoothed density:

$$P_{\mathcal{R}_n} u = u - \frac{1}{nc_n^p} \sum_i \int u(z) \omega\left(\frac{z-z_i}{c_n}\right) dz$$

The right hand side of the equation  $(I - \hat{K}_n)\varphi = \hat{r}_n$  has a mean equal to 0 (w.r.t.  $F_n$ ). Hence, this equation has a unique solution  $\hat{\varphi}_n = \hat{L}_n \varphi_0$  which satisfies the normalization condition  $\frac{1}{nc_n^p} \sum_i \int \hat{\varphi}_n(z) \omega\left(\frac{z-z_i}{c_n}\right) dz = 0$ .

The general results of Section 7.4 apply. First, we check that the conditions i) to iii) of Theorem 7.5 are fulfilled.

- i) Under very general assumptions,  $\|\hat{K}_n - K\| \rightarrow 0$  in probability.

ii) We have to check the properties of  $P_{\mathcal{R}_n} - P_{\mathcal{R}}$

$$(P_{\mathcal{R}_n} - P_{\mathcal{R}})\varphi = \frac{1}{nc_n^p} \sum_i \int \varphi(z)\omega\left(\frac{z - z_i}{c_n}\right) dz - \int \varphi(z)f(z)dz.$$

The asymptotic behavior of  $\|(P_{\mathcal{R}_n} - P_{\mathcal{R}})\varphi\|^2 = \left| \frac{1}{nc_n^p} \sum_{i=1}^n \int \varphi(z)\omega\left(\frac{z - z_i}{c_n}\right) dz - E(\varphi) \right|^2$  is the same as the asymptotic behavior of the expectation of this positive random variable:

$$E \left( \frac{1}{nc_n^p} \sum_{i=1}^n \int \varphi(z)\omega\left(\frac{z - z_i}{c_n}\right) dz - E(\varphi) \right)^2.$$

Standard computation on this expression shows that this mean square error is  $O\left(\frac{1}{n} + c_n^{2\min(d,d')}\right) \|\varphi\|^2$ , where  $d$  is the smoothness degree of  $\varphi$  and  $d'$  the order of the kernel.

iii) The last term we have to consider is actually not computable but its asymptotic behavior is easily characterized. We simplify the notation by denoting  $E^{F_n}(\cdot|\cdot)$  the estimation of a conditional expectation. The term we have to consider is

$$\begin{aligned} (\hat{r}_n + \hat{K}_n\varphi) - (r + K\varphi) &= E^{F_n}(Y|Z) - E^{F_n}(E^{F_n}(Y|W)|Z) + E^{F_n}(E^{F_n}(\varphi(Z)|W)|Z) \\ &\quad - E^F(Y|Z) + E^F(E^F(Y|W)|Z) - E^F(E^F(\varphi(Z)|W)|Z) \\ &= E^{F_n}(Y - E^F(Y|W) + E^F(\varphi(Z)|W)|Z) \\ &\quad - E^F(Y - E^F(Y|W) + E^F(\varphi(Z)|W)|Z) \\ &\quad - R \end{aligned}$$

where  $R = E^F \{ E^{F_n}(Y - \varphi(Z)|W) - E^F(Y - \varphi(Z)|W) \}$ . Moreover

$$E^F(Y|W) = E^F(\varphi(Z)|W) + \psi(W).$$

Then

$$\begin{aligned} (\hat{r}_n + \hat{K}_n\varphi) - (r + K\varphi) &= E^{F_n}(Y - \psi(W)|Z) - E^F(Y - \psi(W)|Z) \\ &\quad - R. \end{aligned}$$

The  $R$  term converges to zero at a faster rate than the first part of the r.h.s. of this equation and can be neglected. We have seen in the other parts of this chapter that

$$\|E^{F_n}(Y - \psi(W)|Z) - E^F(Y - \psi(W)|Z)\|^2 = O\left(\frac{1}{nc_n^p} + c_n^{2\rho}\right)$$

where  $\rho$  depends on the regularity assumptions. Therefore, Condition iii) of Theorem 7.5 is fulfilled.

From Theorem 7.5, it follows that  $\|\hat{\varphi}_n - \varphi_0\| \rightarrow 0$  in probability and that  $\|\hat{\varphi}_n - \varphi_0\| = 0 \left( \frac{1}{\sqrt{nc_n^p}} + c_n^\rho \right)$ .

The pointwise asymptotic normality of  $\sqrt{nc_n^p}(\hat{\varphi}_n(z) - \varphi_0(z))$  can now be established. We apply the formulas (7.24) to (7.26) and Theorem 7.4.

- 1) First, consider (7.26). Under a suitable condition on  $c_n$  (typically  $nc_n^{\rho+2\min(d,r)} \rightarrow 0$ ), we have:

$$\sqrt{nc_n^p} \left\{ \hat{L}_n(P_{\mathcal{R}_n} - P_{\mathcal{R}})r + (P_{\mathcal{R}_n} - P_{\mathcal{R}})\varphi_0 \right\} \rightarrow 0 \text{ in probability.}$$

- 2) Second, consider (7.25). Using the same argument as in Theorem 7.4, a suitable choice of  $c_n$  implies that

$$\sqrt{nc_n^p} \hat{L}_n \hat{K}_n \left[ (\hat{r}_n + \hat{K}_n \varphi_0) - (r + K \varphi_0) \right] \rightarrow 0. \quad (7.30)$$

Actually, while  $E^{F_n}(Y - \psi(W)|Z) - E^F(Y - \psi(W)|Z)$  only converges pointwise at a nonparametric speed, the transformation by the operator  $\hat{K}_n$  converts this convergence into a functional convergence at a parametric speed. Then

$$\sqrt{nc_n^p} \left\| \hat{K}_n \left[ E^{F_n}(Y - \psi(W)|Z) - E^F(Y - \psi(W)|Z) \right] \right\| \rightarrow 0.$$

Moreover,  $\hat{L}_n$  converges in norm to  $L$ , which is a bounded operator. Hence, the result of (7.30) follows.

- 3) The term (7.24) remains. The convergence of  $\sqrt{nc_n^p}(\varphi_{F_n}(z) - \varphi_F(z))$  is then identical to the convergence of

$$\begin{aligned} & \sqrt{nc_n^p} P_{\mathcal{R}_n} \left[ E^{F_n}(Y - \psi(W)|Z = z) - E^F(Y - \psi(W)|Z = z) \right] \\ &= \sqrt{nc_n^p} \left[ E^{F_n}(Y - \psi(W)|Z = z) - E^F(Y - \psi(W)|Z = z) \right] \\ & \quad - \frac{1}{n} \sum_i (y_i - \psi(w_i)) - \frac{1}{nc_n^p} \sum_i \int \int (y - \psi(w)) f(y, w|Z = z) \omega \left( \frac{z - z_i}{c_n} \right) dz dw \Big]. \end{aligned}$$

It can easily be checked that the difference between the two sample means converge to zero at a higher speed than  $\sqrt{nc_n^p}$  and these two last terms can be neglected. Then using standard results on nonparametric estimation, we obtain:

$$\sqrt{nc_n^p}(\varphi_{F_n}(z) - \varphi_F(z)) \xrightarrow{d} \mathcal{N} \left( 0, \text{Var}(Y - \psi(W)|Z = z) \frac{\int \omega(u)^2 du}{f_Z(z)} \right)$$

where the 0 mean of the asymptotic distribution is obtained thanks to a suitable choice of the bandwidth, which needs to converge to 0 faster than the optimal speed.



Note that the estimator of  $\varphi$  has the same properties as the oracle estimator based on the knowledge of  $\psi$ . This attractive feature was proved by Mammen, Linton, and Nielsen (1999) using different tools.

### 7.5.2. Estimation of the bias function in a measurement error equation

We have introduced in Example 1.3.6, Section 1, the measurement error model:

$$\begin{cases} Y_1 = \eta + \varphi(Z_1) + U_1 & Y_1, Y_2 \in \mathbb{R} \\ Y_2 = \eta + \varphi(Z_2) + U_2 & Z_1, Z_2 \in \mathbb{R}^p \end{cases}$$

where  $\eta, U_i$  are random unknown elements and  $Y_1$  and  $Y_2$  are two measurements of  $\eta$  contaminated by a bias term depending on observable elements  $Z_1$  and  $Z_2$ . The unobservable component  $\eta$  is eliminated by differencing and we get the model under consideration :

$$Y = \varphi(Z_2) - \varphi(Z_1) + U \tag{7.31}$$

when  $Y = Y_2 - Y_1$  and  $E(Y|Z_1, Z_2) = \varphi(Z_2) - \varphi(Z_1)$ . We assume that i.i.d. observations of  $(Y, Z_1, Z_2)$  are available. Moreover, the order of measurements is arbitrary or equivalently  $(Y_1, Y_2, Z_1, Z_2)$  is distributed identically to  $(Y_2, Y_1, Z_2, Z_1)$ . This implies that  $(Y, Z_1, Z_2)$  and  $(-Y, Z_2, Z_1)$  have the same distribution. In particular,  $Z_1$  and  $Z_2$  are identically distributed.

- The reference space  $\mathcal{H}$  is the space of random variables defined on  $\mathbb{R}^p$  that are square integrable with respect to the true marginal distribution on  $Z_1$  (or  $Z_2$ ). We are in a case where the Hilbert space structure depends on the unknown distribution.
- The function  $\varphi$  is an element of  $\mathcal{H}$  but this set has to be reduced by the smoothness condition in order to obtain asymptotic properties of the estimation procedure.
- The operator  $K$  is the conditional expectation operator

$$\begin{aligned} (K\varphi)(z) &= E^F(\varphi(Z_2)|Z_1 = z) \\ &= E^F(\varphi(Z_1)|Z_2 = z) \end{aligned}$$

from  $\mathcal{H}$  to  $\mathcal{H}$ . The two conditional expectations are equal because  $(Z_1, Z_2)$  and  $(Z_2, Z_1)$  are identically distributed (by the exchangeability property). This operator is self-adjoint and we suppose that  $K$  is compact. This property may be deduced as in previous cases from a Hilbert Schmidt argument.

Equation (7.31) introduces an overidentification property because it constrains the conditional expectation of  $Y$  given  $Z_1$  and  $Z_2$ . In order to define  $\varphi$  for any  $F$  (and in particular for the estimated one), the parameter  $\varphi$  is now defined as the solution of the minimization problem:

$$\varphi = \arg \min_{\varphi} E(Y - \varphi(Z_2) + \varphi(Z_1))^2$$

or, equivalently as the solution of the first-order conditions:

$$E^F [\varphi(Z_2) | Z_1 = z] - \varphi(z) = E(Y | Z_1 = z)$$

because  $(Y, Z_1, Z_2)$  and  $(-Y, Z_2, Z_1)$  are identically distributed.

The integral equation, which defines the function of interest,  $\varphi$ , may be denoted by

$$(I - K) \varphi = r$$

where  $r = E(Y | Z_2 = z) = -E(Y | Z_1 = z)$ . As in the additive models, this inverse problem is ill-posed because  $I - K$  is not one-to-one. Indeed, 1 is the greatest eigenvalue of  $K$  and the eigenfunctions associated with 1 are the constant functions. We need an extra assumption to warrant that the order of multiplicity is one, or in more statistical terms, that  $\varphi$  is identified up to a constant. This property is obtained if  $Z_1$  and  $Z_2$  are measurably separated, i.e. if the functions of  $Z_1$  almost surely equal to some functions of  $Z_2$ , are almost surely constant.

Then the normalization rule is

$$\langle \varphi, \phi_1 \rangle = 0$$

where  $\phi_1$  is constant. This normalization is then equivalent to

$$E^F(\varphi) = 0.$$

If  $F$  is estimated using a standard kernel procedure, the estimated  $F_n$  does not satisfy in general, the exchangeability assumption ( $(Y, Z_1, Z_2)$  and  $(-Y, Z_2, Z_1)$  are identically distributed). A simple way to incorporate this constraint is to estimate  $F$  using a sample of size  $2n$  by adding to the original sample  $(y_i, z_{1i}, z_{2i})_{i=1, \dots, n}$  a new sample  $(-y_i, z_{2i}, z_{1i})_{i=1, \dots, n}$ . For simplicity, we do not follow this method here and consider an estimation of  $F$ , which does not verify the exchangeability. In that case,  $\hat{r}_n$  is not in general an element of  $\mathcal{R}(I - \hat{K}_n)$ , and the estimator  $\hat{\varphi}_n$  is defined as the unique solution of

$$(I - \hat{K}_n) \varphi = P_{\mathcal{R}_n} \hat{r}_n,$$

which satisfies the normalization rule

$$E^{F_n}(\varphi) = 0.$$

Equivalently, we have seen that the functional equation  $(I - \hat{K}_n) \varphi = \hat{r}_n$  reduces to a  $n$  dimensional linear system, which is solved by a generalized inversion. The asymptotic properties of this procedure follow immediately from the theorems of Section 7.4 and are obtained identically to the case of additive models.

## References

- [1] Ai, C. and X. Chen (2003) “Efficient Estimation of Models with Conditional Moment Restrictions Containing Unknown Functions”, *Econometrica*, 71, 1795-1843.
- [2] Ait-Sahalia, Y., L.P. Hansen, and J.A. Scheinkman (2004) “Operator Methods for Continuous-Time Markov Processes”, forthcoming in the *Handbook of Financial Econometrics*, edited by L.P. Hansen and Y. Ait-Sahalia, North Holland.
- [3] Arellano, M., L. Hansen, and E. Sentana (2005) “Underidentification?”, mimeo, CEMFI.
- [4] Aronszajn, N. (1950) “Theory of Reproducing Kernels”, *Transactions of the American Mathematical Society*, Vol. 68, 3, 337-404.
- [5] Bai, J. and S. Ng (2002) “Determining the Number of Factors in Approximate Factor Models”, *Econometrica*, 70, 191-221.
- [6] Basman, R.L. (1957), “A Generalized Classical Method of Linear Estimation of Coefficients in a Structural Equations”, *Econometrica*, 25, 77-83.
- [7] Berlinet, A. and C. Thomas-Agnan (2004) *Reproducing Kernel Hilbert Spaces in Probability and Statistics*, Kluwer Academic Publishers, Boston.
- [8] Blundell, R., X. Chen, and D. Kristensen (2003) “Semi-Nonparametric IV Estimation of Shape-Invariant Engel Curves”, Cemmap working paper CWP 15/03, University College London.
- [9] Blundell, R. and J. Powell (2003) “Endogeneity in Nonparametric and Semiparametric Regression Models”, in *Advances in Economics and Econometrics*, Vol. 2, eds by M. Dewatripont, L.P. Hansen and S.J. Turnovsky, Cambridge University Press, 312-357.
- [10] Bollerslev, T. (1986), “Generalized Autoregressive Conditional Heteroskedasticity”, *Journal of Econometrics* 31, 307-327.
- [11] Bosq, D. (1998) *Nonparametric Statistics for Stochastic Processes. Estimation and Prediction*, Lecture Notes in Statistics, 110. Springer-Verlag, NewYork.
- [12] Bosq, D. (2000) *Linear processes in function spaces. Theory and applications*, Lecture Notes in Statistics, 149. Springer-Verlag, NewYork.
- [13] Breiman, L. and J.H. Friedman (1985) “Estimating Optimal Transformations for Multiple Regression and Correlation”, *Journal of American Statistical Association*, 80, 580-619.

- [14] Carrasco, M., M. Chernov, J.-P. Florens, and E. Ghysels (2004) “Efficient estimation of jump diffusions and general dynamic models with a continuum of moment conditions”, mimeo, University of Rochester.
- [15] Carrasco, M. and J.-P. Florens (2000) “Generalization of GMM to a continuum of moment conditions”, *Econometric Theory*, 16, 797-834.
- [16] Carrasco, M. and J.-P. Florens (2001) “Efficient GMM Estimation Using the Empirical Characteristic Function”, mimeo, University of Rochester.
- [17] Carrasco, M. and J.-P. Florens (2002) “Spectral method for deconvolving a density”, mimeo, University of Rochester.
- [18] Carrasco, M. and J.-P. Florens (2004) “On the Asymptotic Efficiency of GMM”, mimeo, University of Rochester.
- [19] Carroll, R. and P. Hall (1988) “Optimal Rates of Convergence for Deconvolving a Density”, *Journal of American Statistical Association*, 83, No.404, 1184-1186.
- [20] Carroll, R., A. Van Rooij, and F. Ruymgaart (1991) “Theoretical Aspects of Ill-posed Problems in Statistics”, *Acta Applicandae Mathematicae*, 24, 113-140.
- [21] Chacko, G. and L. Viceira (2003) “Spectral GMM estimation of continuous-time processes”, *Journal of Econometrics*, 116, 259-292.
- [22] Chen, X., L.P. Hansen and J. Scheinkman (1998) “Shape-preserving Estimation of Diffusions”, mimeo, University of Chicago.
- [23] Chen, X., and H. White (1992), “Central Limit and Functional Central Limit Theorems for Hilbert Space-Valued Dependent Processes”, Working Paper, University of San Diego.
- [24] Chen, X. and H. White (1996) “Law of Large Numbers for Hilbert Space-Valued mixingales with Applications”, *Econometric Theory*, 12, 284-304.
- [25] Chen, X. and H. White (1998) “Central Limit and Functional Central Limit Theorems for Hilbert Space-Valued Dependent Processes”, *Econometric Theory*, 14, 260-284.
- [26] Darolles, S., J.-P. Florens, and C. Gourieroux (2004) “Kernel Based Nonlinear Canonical Analysis and Time Reversibility”, *Journal of Econometrics*, 119, 323-353.
- [27] Darolles, S., J.-P. Florens, and E. Renault (1998), “Nonlinear Principal Components and Inference on a Conditional Expectation Operator with Applications to Markov Processes”, presented in Paris-Berlin conference 1998, Garchy, France.

- [28] Darolles, S., J.-P. Florens, and E. Renault (2002), “Nonparametric Instrumental Regression”, Working paper 05-2002, CRDE.
- [29] Das, M. (2005) “Instrumental variables estimators of nonparametric models with discrete endogenous regressors”, *Journal of Econometrics*, 124, 335-361.
- [30] Dautray, R. and J.-L. Lions (1988) *Analyse mathématique et calcul numérique pour les sciences et les techniques*. Vol. 5. Spectre des opérateurs, Masson, Paris.
- [31] Davidson, J. (1994) *Stochastic Limit Theory*, Oxford University Press, Oxford.
- [32] Debnath, L. and P. Mikusinski (1999) *Introduction to Hilbert Spaces with Applications*, Academic Press. San Diego.
- [33] Dunford, N. and J. Schwartz (1988) *Linear Operators, Part II: Spectral Theory*, Wiley, New York.
- [34] Engl. H. W., M. Hanke, and A. Neubauer (1996) *Regularization of Inverse Problems*, Kluwer Academic Publishers.
- [35] Engle R.F., (1990), “Discussion: Stock Market Volatility and the Crash of '87”, *Review of Financial Studies* 3, 103-106.
- [36] Engle, R.F., D.F. Hendry and J.F. Richard, (1983), “Exogeneity”, *Econometrica*, 51 (2) 277-304.
- [37] Engle, R.F., and V.K. Ng (1993), “Measuring and Testing the Impact of News on Volatility”, *The Journal of Finance* XLVIII, 1749-1778.
- [38] Fan, J. (1993) “Adaptively local one-dimensional subproblems with application to a deconvolution problem”, *The Annals of Statistics*, 21, 600-610.
- [39] Feuerverger, A. and P. McDunnough (1981), “On the Efficiency of Empirical Characteristic Function Procedures”, *Journal of the Royal Statistical Society, Series B*, 43, 20-27.
- [40] Florens, J.-P. (2003) “Inverse Problems in Structural Econometrics: The Example of Instrumental Variables”, in *Advances in Economics and Econometrics*, Vol. 2, eds by M. Dewatripont, L.P. Hansen and S.J. Turnovsky, Cambridge University Press, 284-311.
- [41] Florens, J.-P. (2005) “Engogeneity in nonseparable models. Application to treatment models where the outcomes are durations”, mimeo, University of Toulouse.
- [42] Florens, J.-P., J. Heckman, C. Meghir and E. Vytlacil (2002), “Instrumental Variables, Local Instrumental Variables and Control Functions”, Manuscript, University of Toulouse.

- [43] Florens, J.-P. and Malavolti (2002) “Instrumental Regression with Discrete Variables”, mimeo University of Toulouse, presented at ESEM 2002, Venice.
- [44] Florens, J.-P. and M. Mouchart (1985), “Conditioning in Dynamic Models”, *Journal of Time Series Analysis*, 53 (1), 15-35.
- [45] Florens, J.-P., M. Mouchart, and J.F. Richard (1974), “Bayesian Inference in Error-in-variables Models”, *Journal of Multivariate Analysis*, 4, 419-432.
- [46] Florens, J.-P., M. Mouchart, and J.F. Richard (1987), “Dynamic Error-in-variables Models and Limited Information Analysis”, *Annales d’Economie et Statistiques*, 6/7, 289-310.
- [47] Florens, J.-P., M. Mouchart, and J.-M. Rolin (1990) *Elements of Bayesian Statistics*, Dekker, New York.
- [48] Florens, J.-P., C. Protopopescu, and J.F. Richard, (1997), “Identification and Estimation of a Class of Game Theoretic Models”, GREMAQ, University of Toulouse.
- [49] Forni, M., M. Hallin, M. Lippi, and L. Reichlin (2000) “The generalized dynamic factor model: identification and estimation”, *Review of Economic and Statistics*, 82, 4, 540-552.
- [50] Forni, M. and L. Reichlin (1998) “Let’s Get Real: A Factor Analytical Approach to Disaggregated Business Cycle Dynamics”, *Review of Economic Studies*, 65, 453-473.
- [51] Gallant, A. R. and J. R. Long (1997) “Estimating Stochastic Differential Equations Efficiently by Minimum Chi-squared”, *Biometrika*, 84, 125-141.
- [52] Gaspar, P. and J.-P. Florens, (1998), “Estimation of the Sea State Bias in Radar Altimeter Measurements of Sea Level: Results from a Nonparametric Method”, *Journal of Geophysical Research*, 103 (15), 803-814.
- [53] Guerre, E., I. Perrigne, and Q. Vuong, (2000), “Optimal Nonparametric Estimation of First-Price Auctions”, *Econometrica*, 68 (3), 525-574.
- [54] Groetsch, C. (1993) *Inverse Problems in Mathematical Sciences*, Vieweg Mathematics for Scientists and Engineers, Wiesbaden.
- [55] Hall, P. and J. Horowitz (2004) “Nonparametric Methods for Inference in the Presence of Instrumental Variables”, mimeo, Northwestern University.
- [56] Hansen, L.P., (1982), “Large Sample Properties of Generalized Method of Moments Estimators”, *Econometrica*, 50, 1029-1054.
- [57] Hansen, L.P. (1985) “A method for calculating bounds on the asymptotic covariance matrices of generalized method of moments estimators”, *Journal of Econometrics*, 30, 203-238.

- [58] Hardle, W. and O. Linton (1994) “Applied Nonparametric Methods”, *Handbook of Econometrics*, Vol. IV, edited by R.F. Engle and D.L. McFadden, North Holland, Amsterdam.
- [59] Hastie, T.J. and R.J. Tibshirani (1990), *Generalized Additive Models*, Chapman and Hall, London.
- [60] Hausman, J., (1981), “Exact Consumer’s Surplus and Deadweight Loss”, *American Economic Review*, 71, 662-676.
- [61] Hausman, J. (1985), “The Econometrics of Nonlinear Budget sets”, *Econometrica*, 53, 1255-1282.
- [62] Hausman, J. and W.K. Newey, (1995) “Nonparametric Estimation of Exact Consumers Surplus and Deadweight Loss”, *Econometrica*, 63, 1445-1476.
- [63] Heckman, J., H. Ichimura, J. Smith, and P. Todd (1998), “Characterizing Selection Bias Using Experimental Data”, *Econometrica*, 66, 1017-1098.
- [64] Heckman, J., and E. Vytlacil (2000), “Local Instrumental Variables”, in *Nonlinear Statistical Modeling: Proceedings of the Thirteenth International Symposium in Economic Theory and Econometrics: Essays in Honor of Takeshi Amemiya*, ed. by C. Hsiao, K. Morimune, and J. Powells. Cambridge: Cambridge University Press, 1-46.
- [65] Hoerl, A. E. and R. W. Kennard (1970) “Ridge Regression: Biased Estimation of Nonorthogonal Problems”, *Technometrics*, 12, 55-67.
- [66] Horowitz, J. (1999) “Semiparametric estimation of a proportional hazard model with unobserved heterogeneity”, *Econometrica*, 67, 1001-1028.
- [67] Imbens, G., and J. Angrist (1994), “Identification and Estimation of Local Average Treatment Effects”, *Econometrica*, 62, 467-476.
- [68] Jiang, G. and J. Knight (2002) “Estimation of Continuous Time Processes Via the Empirical Characteristic Function”, *Journal of Business & Economic Statistics*, 20, 198-212.
- [69] Judge, G., W. Griffiths, R. C. Hill, H. Lutkepohl, and T-C. Lee (1980) *The Theory and Practice of Econometrics*, John Wiley and Sons, New York.
- [70] Kargin, V. and A. Onatski (2004) “Dynamics of Interest Rate Curve by Functional Auto-Regression”, mimeo Columbia University, presented at the CIRANO and CIREQ Conference on Operator Methods (Montreal, November 2004).
- [71] Kitamura, Y. and M. Stutzer (1997), “An Information Theoretic Alternative to Generalized Method of Moments Estimation”, *Econometrica*, 65, 4, 861-874.

- [72] Knight, J. L. and J. Yu (2002) “Empirical Characteristic Function in Time Series Estimation”, *Econometric Theory*, 18, 691-721.
- [73] Kress, R. (1999), *Linear Integral Equations*, Springer.
- [74] Kutoyants, Yu. (1984), *Parameter estimation for stochastic processes*, Heldermann Verlag, Berlin.
- [75] Lancaster, H. (1968), “The Structure of Bivariate Distributions”, *Annals of Mathematical Statistics*, 29, 719-736.
- [76] Linton, O. and E. Mammen (2003), “Estimating Semiparametric ARCH( $\infty$ ) models by kernel smoothing methods”, forthcoming in *Econometrica*.
- [77] Linton, O. and J.P. Nielsen (1995) “A Kernel Method of Estimating Structured Nonparametric regression Based on Marginal Integration”, *Biometrika*, 82, 93-100.
- [78] Loubes, J.M. and A. Vanhems (2001), “Differential Equation and Endogeneity”, Discussion Paper, GREMAQ, University of Toulouse, presented at ESEM 2002, Venice.
- [79] Loubes, J.M. and A. Vanhems (2003), “Saturation Spaces for Regularization Methods in Inverse Problems”, Discussion Paper, GREMAQ, University of Toulouse, presented at ESEM 2003, Stockholm.
- [80] Lucas, R. (1978) “Asset Prices in an Exchange Economy”, *Econometrica*, 46, 1429-1446.
- [81] Luenberger, D. G. (1969) *Optimization by Vector Space Methods*, Wiley, New York.
- [82] Malinvaud, E. (1970), *Methodes Statistiques de l'Econometrie*, Dunod, Paris.
- [83] Mammen, E., O. Linton, and J. Nielsen (1999) “The existence and asymptotic properties of a backfitting projection algorithm under weak conditions”, *The Annals of Statistics*, 27, 1443-1490.
- [84] Nashed, N. Z. and G. Wahba (1974) “Generalized inverses in reproducing kernel spaces: An approach to regularization of linear operator equations”, *SIAM Journal of Mathematical Analysis*, 5, 974-987.
- [85] Natterer (1984) “Error bounds for Tikhonov regularization in Hilbert scales”, *Applicable Analysis*, 18, 29-37.
- [86] Newey, W. (1994) “Kernel Estimation of Partial Means”, *Econometric Theory*, 10, 233-253.
- [87] Newey, W., and J. Powell (2003), “Instrumental Variables for Nonparametric Models”, *Econometrica*, 71, 1565-1578.



- [88] Newey, W., Powell, J., and F. Vella (1999), “Nonparametric Estimation of Triangular Simultaneous Equations Models”, *Econometrica*, 67, 565-604.
- [89] Owen, A. (2001) *Empirical likelihood*, Monographs on Statistics and Applied Probability, vol. 92. Chapman and Hall, London.
- [90] Pagan, A. and A. Ullah (1999), *Nonparametric Econometrics*, Cambridge University Press.
- [91] Parzen, E. (1959) “Statistical Inference on time series by Hilbert Space Methods,I.”, Technical Report No.23, Applied Mathematics and Statistics Laboratory, Stanford. Reprinted in (1967) *Time series analysis papers*, Holden-Day, San Francisco.
- [92] Parzen, E. (1970) “Statistical Inference on time series by RKHS methods”, Proc. 12th Biennial Canadian Mathematical Seminar, R. Pyke, ed. American Mathematical Society, Providence.
- [93] Politis, D. and J. Romano (1994) “Limit theorems for weakly dependent Hilbert space valued random variables with application to the stationary bootstrap”, *Statistica Sinica*, 4, 451-476.
- [94] Polyanin, A. and A. Manzhirov (1998) *Handbook of Integral Equations*, CRC Press, Boca Raton, Florida.
- [95] Qin, J. and J. Lawless, (1994), “Empirical Likelihood and General Estimating Equations”, *The Annals of Statistics*, 22, 1, 300-325.
- [96] Reiersol, O. (1941), “Confluence Analysis of Lag Moments and other Methods of Confluence Analysis”, *Econometrica*, 9, 1-24.
- [97] Reiersol, O. (1945), “Confluence Analysis by Means of Instrumental Sets of Variables”, *Arkiv for Matematik, Astronomie och Fysik*, 32A, 119.
- [98] Ross, S. (1976) “The Arbitrage Theory of Capital Asset Pricing”, *Journal of Finance*, 13, 341-360.
- [99] Rust, J., J. F. Traub, and H. Wozniakowski (2002) “Is There a Curse of Dimensionality for Contraction Fixed Points in the Worst Case?”, *Econometrica*, 70, 285-330.
- [100] Ruymgaart, F. (2001) “A short introduction to inverse statistical inference”, lecture given at the conference “L’Odyssée de la Statistique”, Institut Henri Poincaré, Paris.
- [101] Saitoh, S. (1997) *Integral transforms, reproducing kernels and their applications*, Longman.
- [102] Sansone, G. (1959) *Orthogonal Functions*, Dover Publications, New York.

- [103] Sargan, J.D. (1958), “The Estimation of Economic Relationship using Instrumental Variables”, *Econometrica*, 26, 393-415.
- [104] Schaumburg, E. (2004) “Estimation of Markov Processes of Levy Type Generators”, mimeo, Kellogg School of Management.
- [105] Singleton, K. (2001) “Estimation of Affine Pricing Models Using the Empirical Characteristic Function”, *Journal of Econometrics*, 102, 111-141.
- [106] Stefanski, L. and R. Carroll (1990) “Deconvoluting Kernel Density Estimators”, *Statistics*, 2, 169-184.
- [107] Stock, J. and M. Watson (1998) “Diffusion Indexes”, NBER working paper 6702.
- [108] Stock, J. and M. Watson (2002) “Macroeconomic Forecasting Using Diffusion Indexes”, *Journal of Business and Economic Statistics*, 20, 147-162.
- [109] Tauchen, G. (1997) “New Minimum Chi-Square Methods in Empirical Finance”, in *Advances in Econometrics, Seventh World Congress*, eds. D. Kreps and K. Wallis, Cambridge University Press, Cambridge.
- [110] Tauchen, G. and R. Hussey (1991) “Quadrature-Based Methods for Obtaining Approximate Solutions to Nonlinear Asset Pricing Models”, *Econometrica*, 59, 371-396.
- [111] Tautenhahn, U. (1996) “Error estimates for regularization methods in Hilbert scales”, *SIAM Journal of Numerical Analysis*, 33, 2120-2130.
- [112] Theil, H.(1953), “Repeated Least-Squares Applied to Complete Equations System”, The Hague: Central Planning Bureau (mimeo).
- [113] Tjostheim, D. and B. Auestad (1994) “Nonparametric Identification of Nonlinear Time Series Projections”, *Journal of American Statistical Association*, 89, 1398-1409.
- [114] Vanhems, A. (2000), “Nonparametric Solutions to Random Ordinary Differential Equations of First Orders”, GREMAQ-University of Toulouse.
- [115] Van Rooij and F. Ruymgaart (1991) “Regularized Deconvolution on the Circle and the Sphere”, in *Nonparametric Functional Estimation and Related Topics*, edited by G. Roussas, 679-690, Kluwer Academic Publishers, the Netherlands.
- [116] Van Rooij, A., F. Ruymgaart (1999) “On Inverse Estimation”, in *Asymptotics, Nonparametrics, and Time Series*, 579-613, Dekker, NY.
- [117] Van Rooij, A., F. Ruymgaart, and W. Van Zwet (2000) “Asymptotic Efficiency of Inverse Estimators”, *Theory of Probability and its Applications*, 44, 4, 722-738.
- [118] Vapnik A.C.M. (1998), *Statistical Learning Theory*, Wiley, New York.

- [119] Wahba, G. (1973) “Convergence Rates of Certain Approximate Solutions to Fredholm Integral Equations of the First Kind”, *Journal of Approximation Theory*, 7, 167-185.