

IDENTIFYING SOURCES OF INEFFICIENCY IN HEALTHCARE*

Amitabh Chandra
Harvard and the NBER

Douglas O. Staiger
Dartmouth and the NBER

Version: December 18th 2016

Abstract

In medicine, variations in treatment rates across providers serving similar patient populations are usually interpreted as some hospitals delivering too much or too little care. An economic interpretation is that variation in treatment rates reflects comparative advantage. One interpretation pushes towards standardization while the other towards learning the sources of comparative advantage. We build a simple economic model which provides an empirical framework to separate these explanations. It generates the key insight is that allocative inefficiency—doing too much or too little-- affects treatment effects conditional on the propensity to receive reperfusion, even when comparative-advantage is present. Estimating this model with data on reperfusion treatments for heart attack patients, we find that (1) consistent with the medical literature, a substantial proportion of hospitals are overusing treatment; (2) consistent with the economic literatures, variation in comparative-advantage, with the variation being the same order of magnitude as the average treatment effect; (3) that one mechanism for overuse is that smaller hospitals have imprecise information about their comparative-effectiveness and misperceive their ability to deliver treatment. Consequently, they overuse treatment relative to their actual expertise. A stylized welfare-calculation suggests that eliminating allocative inefficiency would reduce mortality by about 5 percent, which is more than double the effectiveness of treatment.

* Contact: Amitabh.Chandra@Harvard.EDU and Douglas.Staiger@Dartmouth.EDU. This research was funded by the National Institute of Aging (NIA) P01 AG19783-02. We thank Josh Angrist, Janet Currie, Joe Doyle, Amy Finkelstein, Jonathan Skinner, Heidi Williams for comments that have greatly improved our paper. We obtained access to the proprietary data used in this paper through a data use agreement between the Centers for Medicare and Medicaid Services (CMS) and Dartmouth Medical School. Readers wishing to use these data must obtain them from CMS. Programs and output are available from the authors.

A large and influential literature in economics and medicine has measured the association between treatment rates and risk-adjusted outcomes [Skinner (2013); Institute of Medicine (2013)]. But variation in treatment rates could arise from two different mechanisms. The conventional interpretation in the medical literature is that there is a correct amount of use, so that variation across providers in risk-adjusted treatment rates is evidence of allocative inefficiency: some providers are using too much care and others are using too little. This interpretation of variation leads to an emphasis on guidelines and developing and disseminating information on cost-effectiveness of care. An alternative interpretation argues that the ability to deliver treatment varies across providers, so that hospitals who can obtain higher benefits from intensive treatments deliver more treatment because of comparative advantage. This interpretation leads to an emphasis on understanding the sources of variation in hospital-specific skill, and efforts to improve quality, instead of trying to standardize care.

We develop a simple economic framework that can distinguish between these explanations and shed light on the mechanisms behind them. In our model, hospitals choose intensive treatment or non-intensive treatment for each patient, based on patient characteristics, the hospitals comparative-advantage at performing this treatment, and the minimum benefit threshold for treatment. The minimum threshold should be zero if hospitals only maximize health, but may be negative if hospitals are maximizing something other than health or if are overconfident about their ability to deliver treatment. The model generates two results: that risk-adjusted treatment rates capture comparative advantage in treatment and allocative efficiency, and are therefore not sufficient to separate the two. Second, the model demonstrates that it is possible to identify allocative inefficiency in the presence of comparative-advantage by using an outcomes test where we compare outcomes for patients receiving treatment with the same propensity to receive treatment. The propensity to receive treatment captures patient characteristics and hospital-effects that capture comparative advantage and allocative inefficiency. If one hospital overuses treatment relative to another, the benefit from treatment, conditional on this propensity, will be lower at hospitals that overuse treatment because of working into less appropriate patients.

We begin by setting the clinical stage for the theory and empirical work, and focus on heart-attacks and their treatment—reperfusion therapy. While our model can be applied to any healthcare setting, heart-attack treatments have several features that make them particularly suited for the analysis: outcomes are easily measured and agreed upon, and questions about overuse, underuse and comparative-advantage are central to treatment decisions. Estimating our model with data for elderly patients following a heart-attack, we find strong evidence of allocative inefficiency, with a substantial proportion of hospitals overusing reperfusion therapy among patients who are harmed by the treatment. This is consistent with the medical literatures interpretation of variations. However, we also we find substantial variation in hospitals ability to perform treatment, with the variation across patients and hospitals in the

survival benefit from reperfusion being the same order of magnitude as the average treatment effect of reperfusion. This is evidence of economic explanations of variations in treatment-rates.

Our framework allows us to identify mechanisms that could lead to the allocative inefficiency that we observe in the data. One possibility, motivated by Currie and MacLeod (2017) is that allocative inefficiency would arise if hospital's had imperfect information and misperceived their ability to deliver treatment. In this mechanism, allocative inefficiency arises because hospitals base treatment decisions on their incorrect perception of the return to treatment in their patients, rather than on the true return to treatment. Given the general lack of systematic performance feedback and small samples of their own treated patients to observe, it is quite plausible that hospitals and physicians will have inaccurate beliefs about their own treatment effectiveness. We find evidence in favor of this mechanism, with smaller hospitals having particularly imprecise information about their own treatment effectiveness. Another explanation is that hospitals are optimizing something other than the survival of a given patient, e.g. over-treating for financial gain (particularly in for-profit hospitals) or because of benefits to future patients through learning-by-doing (particularly in teaching hospitals). This type of mechanism would suggest that allocative inefficiency would be related to hospital characteristics such as ownership, teaching status, etc. We find little *prima facie* evidence for this hypothesis—overuse is not correlated with a hospital's for-profit status or other characteristics such as being a teaching-hospital.

The paper proceeds as follows: Section II develops the theoretical model underlying our analysis, and links it to our estimation strategy, paying particular attention to how the theoretical model will be evaluated using the CCP data. Section III presents results that rely on transparent sources of identification such as graphs. At hospitals with high risk-adjusted reperfusion rates, we are able to identify substantial subsets of treated patients who are harmed by overuse of this treatment, suggesting that allocative inefficiency in the form of overuse is a real concern. These results motivate a more parametric approach in Section IV that allows for a richer set of parameters to be estimated: we are able to estimate not only hospital-level differences in thresholds, but also estimate their comparative-advantage parameters, and study their covariance. We find a large role for comparative advantage in treatment, with the variation in comparative-advantage being quantitatively as large as the direct effect of reperfusion. We also find that hospitals that have a high-comparative advantage at delivering reperfusion therapy are bad at producing survival for patients that don't receive reperfusion therapy (i.e, comparative-advantage for reperfusion and the ability to perform medical management are negatively correlated) —indeed, that may be one reason why it makes sense for them to treat more patients with reperfusion therapy. This is a different view of hospital performance than thinking that hospitals with a higher benefit from reperfusion are also better at medical management. We use our framework to explore mechanisms and find that smaller hospitals *incorrectly* believe that they receive higher benefit from performing the treatment, resulting in

overuse among those hospitals that did not have comparative-advantage in performing the treatment. This is a different than hospitals maximizing financial objectives or being risk-averse. Indeed, we find evidence that higher rates of allocative inefficiency aren't correlated with non-profit status and other hospital characteristics. We conclude by performing a highly stylized calculation of the welfare-loss from variation in treatment rates.

I. Heart-Attacks: Biology, Treatments, and Data

A. Heart-Attack Biology and Treatments

Heart attacks (more precisely, acute myocardial infarction (AMI)) occur when the heart-muscle (the myocardium) does not receive sufficient oxygen, because of a blockage in one of the coronary arteries which supply blood to the heart. The blockage is typically caused by a blood clot that occurs because of coagulation induced by the rupture of atherosclerotic plaque inside the coronary arteries, and must be reperfused rapidly. There are two ways to give patients reperfusion: first, through thrombolytics, which are also known as fibrinolytics, are administered intravenously and break down blood clots by pharmacological means (these drugs include tissue plasminogen activators, streptokinase and urokinase). Reperfusion can be performed through angioplasty (where a balloon on a catheter is inflated inside the blocked coronary artery to restore blood flow). Following the clinical literature, we define a patient to have received reperfusion if any of these therapies was provided within 12 hours of the heart attack. In our data from the mid-1990s, over 90 percent of patients receiving reperfusion received thrombolytics.

We focus our empirical work on the treatment of AMI for a number of reasons. First, cardiovascular disease, of which heart attacks are the primary manifestation, is the leading cause of death in the US. A perusal of the leading medical journals would indicate that heart attack treatments are constantly being refined, and a large body of trial evidence points to significant therapeutic gains from many of these treatments. In this context, variation in treatments across hospitals may directly translate into lost lives, and there is a rich tradition of studying variation across hospitals in treatments and outcomes after heart attacks.

Second, because of what is known about heart attack treatments from randomized controlled trials, the benefits from reperfusion which is intensive therapy for heart-attacks, we are able to assess whether our regression estimates of the benefits from reperfusion are comparable to those found in the medical literature, or whether they are confounded by selection-bias. We focus on reperfusion, where our use of chart data allows us to replicate the RCT evidence that is summarized by the Fibrinolytic Therapy Trialists' Collaborative Group (1994). Chart data provides comprehensive documentation on the patient's condition at the time that the treatment decision is made, and therefore minimizes the possibility that unobserved clinical factors related to a patient's survival are correlated with treatment.

Third, because mortality post-AMI is high (mortality rates at 30 days are nearly 20 percent), a well-defined endpoint is available to test the efficacy of heart attack treatments. Moreover during the acute phase of the heart attack the therapeutic emphasis is on maximizing survival, which is achieved by timely reperfusion, and hospital staff (not patients and their families) make treatment decisions. This would not be true if we focused on treatment variation for more chronic conditions such as diabetes, chronic obstructive pulmonary disease, or arthritis where because of the importance of quality-of-life there would be considerable disagreement on how to measure productivity.

Fourth, heart attacks are an acute condition for which virtually all patients are hospitalized at a nearby hospital and receive some medical care. This may not be true of more chronic conditions such as diabetes or heart-failure where many patients aren't diagnosed and selection into the sample confounds the analysis.

B. Data

Because acute myocardial infarction is both common and serious, it has been the topic of intense scientific and clinical interest. One effort to incorporate evidence-based practice guidelines into the care of heart attack patients, begun in 1992, is the Health Care Financing Administration's Health Care Quality Improvement Initiative Cooperative Cardiovascular Project (CCP). Information about patients admitted to hospitals for treatment of heart attacks in 1994/1995 was obtained from clinical records. The CCP is considerably superior to administrative/claims data of the type used by McClellan et al. (1994) as it collects chart data on the patients—detailed information is provided on laboratory tests, enzyme levels, the location of the myocardial infarction, and the condition of the patient at the time of admission. Detailed clinical data were abstracted from each patient's chart using a standard protocol. Further details about the CCP data are available in Marciniak et al. (1998), O'Connor et al. (1999), and in the appendix to this paper. The choice of sample and variables is identical to what we used and described in Barnato et al. (2005) and Chandra and Staiger (2007, 2010).

II. Theory and Estimation

A Roy model of patient treatment choice guides our empirical work. We assume that a hospital must choose between two treatment options for every patient: whether to offer reperfusion or not. Treatment is provided to each patient whenever a patient's expected benefit from the treatment exceeds a minimal threshold. In our framework, there are two ways in which a patient's hospital could affect treatment. First, because of comparative-advantage, the benefit of treatment for a given patient may vary across hospitals, reflecting each hospital's expertise in providing the treatment. Second, because of allocative efficiency, the minimum threshold for receiving care may vary across hospitals. From the

patient's point of view, treatment should be provided whenever the expected benefit from treatment exceeds zero. Therefore, there is underuse of the treatment in hospitals that set a minimum benefit threshold above zero, and overuse in hospitals that set a minimum threshold below zero.

To formalize this, let Y_{ih}^1 represent the survival for patient i at hospital h if the patient receives the treatment (reperfusion) and let Y_{ih}^0 represent the survival if the patient does not receive the treatment, but otherwise receives usual medical care. We focus on the health benefits of the treatment, which in our setting is survival, but in other settings would include any reduction in mortality or morbidity that was expected from the treatment, e.g. the impact of the treatment on Quality Adjusted Live Years (QALYs).¹ Treatment decisions are based on expected survival given the information available to the provider at the time of treatment (within 12 hours of admission). If treated medically, a patient's expected survival $E(Y_{ih}^0)$ depends on the hospital's general level of expertise α_h^0 , observable patient characteristics X_{ih} such as age, medical history and lab results, and other unmeasured factors affecting baseline mortality v_{ih}^0 that are observed by the healthcare provider but not by the econometrician. If treated intensively, a patient's expected survival $E(Y_{ih}^1)$ depends on a similar set of factors representing the hospital's expertise at providing reperfusion α_h^1 , patient characteristics (which may have a different relationship to survival when patients receive the treatment), and other unmeasured factors v_{ih}^1 that affect the expected benefits of reperfusion. The presence of two productivity parameters α_h^1 and α_h^0 , allows us to model hospital specific benefits at both forms of medicine—and medical (non-intensive) and reperfusion (intensive).

Actual (realized) survival if treated medically or with reperfusion is equal to expected survival plus a random error term $(\varepsilon_{ih}^0, \varepsilon_{ih}^1)$, which yields survival equations of the following form:

$$(1a) \quad Y_{ih}^0 = E(Y_{ih}^0) + \varepsilon_{ih}^0 = \alpha_h^0 + X_i\beta_h^0 + v_{ih}^0 + \varepsilon_{ih}^0$$

$$(1b) \quad Y_{ih}^1 = E(Y_{ih}^1) + \varepsilon_{ih}^1 = \alpha_h^1 + X_i\beta_h^1 + v_{ih}^1 + \varepsilon_{ih}^1$$

The benefit, or gain, or return, from reperfusion treatment for patient i in hospital h , Y_{ih}^Δ given by:

$$(1c) \quad Y_{ih}^\Delta = \alpha_h^\Delta + X_i\beta_h^\Delta + v_{ih}^\Delta + \varepsilon_{ih}^\Delta,$$

$$\text{where } \alpha_h^\Delta = \alpha_h^1 - \alpha_h^0, \quad \beta_h^\Delta = \beta_h^1 - \beta_h^0, \quad v_{ih}^\Delta = v_{ih}^1 - v_{ih}^0 \text{ and } \varepsilon_{ih}^\Delta = \varepsilon_{ih}^1 - \varepsilon_{ih}^0$$

And similarly the expected benefit from reperfusion at the time of choosing treatment is given by:

$$(1d) \quad E(Y_{ih}^\Delta) = \alpha_h^\Delta + X_i\beta_h^\Delta + v_{ih}^\Delta$$

¹ For now, we abstract from the problem that hospitals should stop treatment prior to achieving zero marginal benefits—that is, that providers should maximize benefits net of costs. In practice, the cost of treating heart-attacks is small relative to the survival benefit but we will return to this point later in the paper. In other settings, such as oncology, this may not be true.

In Equation (1d), α_h^Δ represents the hospital-specific benefit in providing reperfusion. One could think of α_h^0 as representing a hospital's Total Factor Productivity (TFP)— because increases in it reflect improvements that are unrelated to intensive treatments such as reperfusion or surgery [Garber and Skinner (2012) and Syverson (2001)]. Efforts to increase α_h^0 are efforts to increase productive efficiency— increasing the fraction receiving beta-blockers or improving patient safety are examples. The higher the α_h^0 , the lower the benefit from reperfusion, for a fixed level of α_h^1 . Because α_h^Δ represents the difference between the ability to perform two treatments, we call it comparative-advantage at reperfusion.. Hospitals may have comparative advantage in providing reperfusion because of either being particularly good at reperfusion treatment or being particularly bad at caring for patients medically. In the above equations, we have also allowed to hospital-level variation in how patient characteristics affect outcomes through the β_h^Δ term.

Each patient receives treatment if the expected benefit from treatment exceeds a minimal threshold τ_h , where the threshold varies across hospitals due to incentives or information. Since $E(Y_{ih}^\Delta)$ captures the total expected benefit to the patient of providing treatment, then the optimal decision from the patient's perspective would let $\tau_h=0$ and provide treatment whenever the expected benefits to the patient exceed zero. There is underuse if $\tau_h > 0$, since patients with positive benefits are under the threshold and do not receive treatment. There is overuse if $\tau_h < 0$, since patients with negative benefits (who would do better without treatment) are above the threshold and receive treatment.

Figure 1A illustrates the intuition behind a Roy model of treatment at the hospital level. The two lines denote patient survival if a hospital treats a given patient medically (intercept is α_h^0) or using reperfusion (intercept is α_h^1) as a function of patient characteristics (i.e. patient X's) on the x-axis. To simplify exposition, we have suppressed the distribution of unobservables that come out of the plane. In reality, as well as in our model and empirical work, providers observe these unmeasured characteristics and use them to determine treatment. Expertise at medical and intensive care is captured by the intercepts α_h^0 and α_h^1 respectively, with comparative advantage being the difference between them. Allocative efficiency means that reperfusion should be performed to the point that the marginal patient receiving it receives zero benefit, that $\tau_h = 0$, everyone to the right of the point of intersection should be treated intensively and to the left of it medically.

First, consider the role of comparative-advantage in explaining treatment rates: *ceteris paribus*, a hospital that is better at reperfusion would have a higher intercept for reperfusion α_h^1 , which would increase the fraction of patients receiving reperfusion at that hospital. A hospital may also have a relative advantage at reperfusion because it is worse at medical management. Either would increase $\alpha_h^\Delta = \alpha_h^1 - \alpha_h^0$ and also increase the fraction of all patients being reperfused. Next, consider allocative inefficiency by

a hospital that overtreats patients with reperfusion therapy, and stops to the left of point of intersection. This harms patients and lowers the average benefit from reperfusion amongst all patients receiving reperfusion. Underuse of reperfusion happens when patients who are appropriate for reperfusion don't receive it—a possibility that increases the benefits of reperfusion amongst patients receiving it.

This figure provides four pieces of intuition. First, knowledge of comparative-advantage doesn't tell us where it originates from—it could arise from low α_h^0 , a high α_h^1 , or both. Second, allocative inefficiency may arise from overuse $\tau_h < 0$ (a willingness to perform reperfusion even if the marginal benefit is negative) or underuse $\tau_h > 0$ (leaving benefits on the table because of stopping reperfusion too soon). Third, how a patient is treated depends on patient characteristics, the hospital's comparative advantage at delivering reperfusion (α_h^Δ), and the level of allocative efficiency at the hospital (τ_h): all determine the propensity to be reperfused for all patients at a hospital. This brings us to the fourth insight: that the presence of comparative-advantage doesn't imply anything for the presence of allocative efficiency or vice versa. Another way to see this is to note that risk-adjusted hospital treatment rates capture both mechanisms—high risk-adjusted rates may arise because of high levels of hospital-specific benefits at performing the procedure or a very low threshold for performing the treatment-- and do not, by themselves, isolate the source of variation even with perfect risk-adjustment.

Figure 1B illustrates how a hospital that misperceives its comparative-advantage from reperfusion and believes it to be higher than it is, through overconfidence or imperfect knowledge about its comparative advantage, would overuse reperfusion. It could also be that hospital overuses reperfusion because it is risk-averse, or because it is maximizing something other than health. These are alternative mechanisms that we explore in Section IV (we find evidence for the misperception mechanism). Regardless of the mechanism for allocative inefficiency, they cause a welfare loss whose magnitude is illustrated by area of the triangle in the figure. The height of the triangle is the threshold, and the base is the threshold multiplied by how much the threshold increases the probability of receiving reperfusion. At the end of the paper, we aggregate the area of these triangles to estimate the welfare loss from allocative inefficiency.

We now specify our model more completely, paying particular attention to how one can identify the different sources of inefficiency. The probability of receiving treatment is the probability that expected benefits exceed the minimum threshold:

$$(2) \quad \Pr(\text{Treatment}_{ih} = 1) = \Pr(E(Y_{ih}^\Delta) > \tau_h) = \Pr(\alpha_h^\Delta + X_i\beta^\Delta + v_{ih}^\Delta > \tau_h) = \Pr(-v_{ih}^\Delta < I_{ih}),$$

where $I_{ih} = X_{ih}\beta^\Delta + (\alpha_h^\Delta - \tau_h)$

In the terminology of Heckman, Urzua and Vytlacil (2006), our model allows for *essential* heterogeneity where the decision to provide treatment to each patient is made with knowledge of their

idiosyncratic response to treatment (v_{ih}^Δ). Equation (2) can be estimated with a logit or OLS regression of treatment on patient characteristics and hospital effects. The hospital effect $\theta_h = \alpha_h^\Delta - \tau_h$, which means that a hospital may be more likely to provide treatment because of greater comparative advantage at delivering treatment, $\alpha_h^\Delta > 0$, or using a lower benefit threshold for providing care ($\tau_h < 0$), which reflects overuse.² Another way to say this is that hospital-specific abilities to perform treatment and potential overuse/underuse affect a patient's propensity to receive care in the same way that patient characteristics affect a patient's propensity for care. Even if treatment rates were the same across hospitals, there could still be overuse or underuse if, say, hospitals with greater comparative advantage set a correspondingly higher threshold for providing care. Thus, because variation in treatment rates across hospitals confounds variation in hospital specific benefits and hospital treatment thresholds, such variation cannot by itself say anything about overuse or underuse.

A. Identifying Allocative Inefficiency

We now demonstrate that allocative efficiency can be identified separately from comparative-advantage if we can estimate the treatment effect for those patients receiving treatment. The treatment-on-the-treated parameter is the average gain from treatment amongst those who were given treatment, and can be obtained by conditioning the expression for Y_{ih}^Δ (equation 1c) on the condition for receiving treatment (equation 2):

$$(3) \quad E(Y_{ih}^\Delta | \text{Treatment}_{ih} = 1) = E(Y_{ih}^\Delta | -v_{ih}^\Delta < I_{ih}) = X_i \beta^\Delta + \alpha_h^\Delta + E(v_{ih}^\Delta | -v_{ih}^\Delta < I_{ih})$$

Noting that $X_i \beta^\Delta + \alpha_h^\Delta = I_{ih} + \tau_h$, we can rewrite Equation (3) as:

$$(4) \quad E(Y_{ih}^\Delta | \text{Treatment}_{ih} = 1) = \tau_h + g(I_{ih})$$

where $g(I_{ih}) = I_{ih} + E(v_{ih}^\Delta | -v_{ih}^\Delta < I_{ih})$

Here, $g(I)$ is an unknown function of the propensity to receive treatment and the conditional expectation $E(v_{ih}^\Delta | -v_{ih}^\Delta < I_{ih})$. Under a 'single-index' assumption this conditional expectation is only a function of the index, which means that $g(I)$ is only a function of I . In the empirical work we will show evidence that hospitals do not differ in the weights that they place on the extremely rich set of observable characteristics suggesting that this is not a first-order channel in our setting. The assumption would be violated if the distribution of patient level idiosyncratic gains had a hospital level component, such as the variance of v_{ih}^Δ varying across hospitals; hospitals with larger variance in unobservables will have larger treatment on

² Technically, logit models estimate $(X\beta + \theta)/\sigma_\varepsilon$, where σ_ε is the SD of the unobservables. For now, we assume that $\sigma_\varepsilon=1$ (a standard assumption) but we will return to this point in the results section when we try to recover estimates of τ_h .

the treated effects (which will look like having a higher treatment threshold and underuse) because of the higher values of conditional error term. This is a key assumption in our work, and we will discuss supporting evidence for it in the empirical work.

Equation (4) states that if hospitals differ in their minimum threshold to deliver care (τ_h), differences in benefit for two patients with the same propensity to receive treatments, identify these differences. Note that the propensity to receive treatment depends on the hospital treatment rate and that includes both the presence of productive and allocative efficiency. By conditioning on this propensity and examining differences in benefit across hospitals, we can isolate differences in hospital thresholds.

We note that our model does not, by itself, uncover mechanisms for overuse or underuse—we will investigate these later. It is possible that overuse occurs because providers are worried about malpractice, because they're maximizing something other than health, because they incorrectly believe that they're better at offering intensive treatment, or because they inaccurately assess patients are more appropriate for treatment than they actually are (perceiving a rightward shift in the distribution of patients X_s).

B. Graphical Intuition

The graphical intuition for our model can be seen in Figure 2. The expected benefit from treatment, $E(Y_{ih}^\Delta | Treatment_{ih} = 1) = E(Y_{ih}^\Delta | E(Y_{ih}^\Delta) > \tau_h)$, is given on the vertical axis, while the propensity of being treated (which depends on I) is given on the horizontal axis; as in the theory section, the threshold is set at zero meaning that providers should treat patients until the point of zero-marginal benefit. The top curve in Figure 2 represents the treatment-on-the-treated effect for a patient with a given propensity that is treated in a hospital with a high minimum threshold for treatment, *i.e.* it represents $E(Y_{ih}^\Delta | E(Y_{ih}^\Delta) > \tau_h) = \tau_h + g(I_{ih})$. The lower curve represents the same thing for a hospital with a low minimum threshold. Treatment-on-the-treated approaches the minimum threshold (τ_{high} or τ_{low}) for a patient with a low propensity of being treated (small value of I), since no patient is ever treated with a benefit below this threshold. For a patient with a high propensity of being treated (large value of I), truncation becomes irrelevant and the treatment-on-the-treated effect asymptotes to the unconditional benefit of treatment. However, conditional on a patient's propensity, the treatment effect is always higher in the hospital with the higher threshold.

The graphs allow us to make two observations that reinforce the theoretical model: First, we can identify overuse and underuse by focusing on patients with the lowest probability of receiving treatment. In these patients, there is overuse when the treatment effect for the lowest propensity patients is negative, and underuse when the treatment effect for the lowest propensity patients remains positive.

Second, differences in comparative-advantage at performing reperfusion would show up as a movement along the curves – higher comparative advantage at reperfusion α_h^Δ increases the propensity of patients to be treated, and therefore the treatment effect, but does not affect treatment effects conditional on propensity. As an example, suppose that there are two hospitals A and B, where hospital A has $\alpha_h^\Delta = 0$ and $\tau_h = 0$, and hospital B has $\alpha_h^\Delta = 10$ and $\tau_h = 0$. Comparing patients with the same propensity to receive reperfusion at hospitals A and B would generate the same treatment effects even though younger patients (or more generally, clinically less appropriate patients) in hospital A are being compared to older patients in hospital B. In other words, being treated at a hospital with higher α_h^Δ is just like having more $X_i\beta^\Delta$. This thinking may cause one to think that we should be matching comparing patients with identical X's across hospitals (as opposed to similar propensities to be treated) and concluding that there is overuse if hospitals with higher treatment rates have lower returns and underuse if they have lower returns. This intuition is tempting but incorrect. It certainly works if all the variation across hospitals in risk-adjusted treatment rates reflects only comparative-advantage or only thresholds, but breaks down if both are present at the same hospital. To see this, consider hospital C, which has $\alpha_h^\Delta = 5$ and $\tau_h = -5$. Hospitals B and C both have $\theta_h = 10$. A patient with the same X's would do unambiguously better at hospital B than C (because hospital B has higher α_h^Δ and lower τ_h) and so one would incorrectly conclude that variation was driven by higher comparative at hospital B, when it is also driven by a lower threshold at hospital C. The only way to isolate the difference is thresholds is to compare patients with the same propensity, $X_i\beta^\Delta + \theta_h$. This is because increasing both the mean of the distribution and the truncation point by tau also increases the conditional benefit of reperfusion by tau. In summary, the key difference between identifying comparative advantage from overuse/underuse is that differences in hospital comparative advantage have an impact on treatment effects by shifting the propensity to be treated, while differences in the minimum threshold have an impact on treatment effects conditional on the propensity to receive reperfusion.

C. Identification

In a potential outcomes framework, the equation relating the level of survival to treatment is:

$$(5) \quad Y_{ih} = Y_{ih}^0 + Y_{ih}^\Delta \text{Treatment}_{ih} \\ = \alpha_h^0 + X_{ih}\beta^0 + Y_{ih}^\Delta \text{Treatment}_{ih} + (v_{ih}^0 + \varepsilon_{ih}^0)$$

Here, survival for the ih patient at hospital h depends on a hospital effect that captures medical expertise α_h^0 , patient risk adjusters X_{ih} and a patient-specific treatment effect Y_{ih}^Δ . Regression estimates of this equation identify the treatment-on-the-treated effect: $E(Y_{ih}^\Delta | \text{Treatment} = 1)$, if the receipt of

treatment is uncorrelated with the unobservable characteristics of patients who were not reperfused (v_{ih}^0 and ε_{ih}^0). This treatment-on-the-treated effect is the same as Equation 4—and identifies overuse (if negative) and underuse (if positive).

It is important to see that we are not assuming that receipt of treatment is uncorrelated with the gain from treatment, which is the conventional assumption required to estimate average treatment effects. Indeed, we explicitly allow for ‘selection on gains’ where providers use information on the gain, $v_{ih}^1 - v_{ih}^0$, to determine treatment [Wooldridge (2002, p.606)]. This is likely in our setting: whether a patient gets treated medically or surgically depends on factors such as the experience of the doctor and team, as well as the capacity of the hospital at the moment of the patients arrival. Providers observe these factors, which are idiosyncratic to every patient situation, and act on them. This is a weaker set of identifying assumptions than the conventional ‘selection on observables’ model where one assumes that conditional on rich observable patient characteristics, patients receive treatments randomly. In terms of the analogy to random assignment, we are assuming that conditional on X’s, patients are randomly assigned to hospitals, but within hospitals doctors triage them according to the benefit that they would receive from treatment. In other words, the random assignment that we require is that treatment is random with respect to v_{ih}^0 (baseline unobservables), which is plausible given the rich covariates that we have. We are *not* assuming that patients are randomly allocated to hospitals, and within hospitals that they are randomly assigned to treatment (conditional on X’s).

To support this assumption we estimate a simple logit model for the impact of reperfusion on 30-day survival, controlling for the patient risk-adjusters using the CCP data. We compared the estimates from this model to those obtained from clinical trials to evaluate the plausibility this assumption. A summary of nine trials was published in the journal *Lancet* by the Fibrinolytic Therapy Trialists' Collaborative Group (FTTCG, 1994). This was the same time-period as the CCP data and each trial evaluated fibrinolytic therapy in heart-attack patients. Across these nine trials, reperfusion within 12 hours reduced 35-day mortality from 11.5% to 9.6%, which implies that the treatment on the treated effect of reperfusion on the log-odds of mortality is 0.20. In our CCP data, a logit model controlling for the CCP risk-adjusters estimates an identical effect, with reperfusion increasing the log-odds of survival (equivalently reducing the logodds of mortality) by 0.206 (S.E. = 0.023). We take this evidence as supporting the case that we can estimate unbiased estimates of the treatment on treated effect.

Our test for allocative efficiency requires comparing the treatment on the treated parameter across hospitals, while holding the propensity to receive treatment constant. The propensity to receive treatment, I , is obtained from a logit model of treatment receipt at the individual level, and we can control for it using indicator variables for the 100 percentiles of I , which allows $g(I)$ to have any shape.

$$(6) \quad Y_{ih} = \alpha_h^0 + X_{ih}\beta^0 + (\tau_h + g(I_{ih}))Treatment_{ih} + v_{ih}^0 + \varepsilon_{ih}^0$$

$$Y_{ih} = \alpha_h^0 + X_{ih}\beta^0 + \tau_h Treatment_{ih} + \sum_{p=1}^{100} 1(g_{p-1} < I_{ih} < g_p) * Treatment_{ih} + v_{ih}^0 + \varepsilon_{ih}^0$$

Here, the hospital-specific treatment effects identify τ_h (the coefficients will be negative at hospitals that overuse and positive at hospitals that underuse). The hospital fixed-effects in this regression identify hospital TFP (α_h^0). The coefficients on the indicator variables for the percentiles of $g(I)$ provide a test for whether the benefit of reperfusion therapy is increasing the propensity to receive such treatment—as would be the case if a Roy Model of treatment allocation was at work, as opposed to model where providers select patients randomly or without regard to benefits.

Noting that the risk-adjusted hospital reperfusion rate is $\theta_h = \alpha_h^\Delta - \tau_h$, we can manipulate Equation 6 to produce another test for allocative inefficiency:

$$(7a) \quad Y_{ih} = \alpha_h^0 + X_{ih}\beta^0 + \theta_h Treatment_{ih} + \sum_{p=1}^{100} 1(g_{p-1} < I_{ih} < g_p) * Treatment_{ih} + v_{ih}^0 + \varepsilon_{ih}^0$$

If hospital-level reperfusion rates were entirely driven by comparative advantage in treatment, we would get a coefficient on 0 on hospital reperfusion rates, and a simple t-test on the interaction of the hospital treatment effect with receipt of treatment allows us to test this null.³ Equation 7a will therefore be the principal estimating equation in our analysis.

One further simplification that we experiment with, and find evidence for, is that the benefit from treatment is linear in I (as in a first-order Taylor expansion: $g(I) = \lambda_0 + \lambda_1 I$). There was no reason to expect this, for the theory only predicts monotonic relationship, but we do find that estimates from a simple linear specification are very similar to those that allow $g(I)$ to have a completely flexible form. With this simplification, the test for allocative efficiency becomes:

$$(7b) \quad Y_{ih} = \alpha_h^0 + X_i\beta^0 + \theta_h Treatment_{ih} + \lambda_0 Treatment_{ih} + \lambda_1 Treatment_{ih} * I_{ih} + v_{ih}^0 + \varepsilon_{ih}^0$$

We will exploit this specification later in the paper, where we impose additional parametric structure to recover hospital measures of α_h^Δ and α_h^0 . In this specification, if we demean I to have a value of 0, the coefficient λ_0 captures the average effect of reperfusion.

III. Results

In Table 1 we report some basic characteristics of our sample overall, and by whether the patient received reperfusion within 12 hours of admission to the hospital. In our sample, 19% of patients received reperfusion within 12 hours of admission for a heart attack. Overall, 81% of patients were still alive 30

³ The approach in Equation 7a does not require one to include Hospital*Treatment indicator variables, but provides a simpler, but powerful, test for understanding the importance of allocative inefficiency in medical care. The simplicity halves the number of covariates that need to be estimated: there are 4690 hospital in our sample, and Equation 6 would require including $4690 \times 2 = 9380$ fixed-effects (one for each α_h^0 and one for each hospital specific treatment term).

days after admission, but survival was higher for patients receiving reperfusion (86%) than for patients who did not receive reperfusion (80%). However, much of the difference in survival between these two groups was due to differences in underlying health and pre-existing conditions, rather than the result of reperfusion. Patients receiving reperfusion were younger, and much less likely to have pre-existing conditions such as congestive heart failure, hypertension, diabetes, and dementia.

A. Graphical Results

We first present some simple graphical results. In Figure 3 we plot the survival benefit from reperfusion against a patient's treatment propensity index. The graphs report the treatment benefit at each point in the distribution of the propensity to receive care using our preferred specification with 100 indicator variables for the percentiles of the propensity index (Equation 6). We separate the panels for hospital's in the lowest tercile and highest tercile of the estimated hospital effect ($\hat{\theta}_h$). The estimation of hospital-effects is described in Appendix-B. As noted in this Appendix, we are using hospital random-effects which allow for empirical bayes-shrinkage because their correlation with hospital fixed-effects is over 0.999. Both the hospital effect, which is simply the risk-adjusted reperfusion rate, and the patient's treatment propensity are obtained from logit estimation of equation (2). We report the treatment effect using a local-regression, using an epinochov kernel (triangular, or linearly declining weights) that included 30% of the sample on either side.

Both plots show a strong upward slope, with higher benefit from treatment for patients with a higher propensity to receive reperfusion—and exactly mirrors the theoretical illustration in Figure 2. But at every propensity, the benefits of reperfusion are lower in the top-tercile hospitals. At the lowest propensity levels, the survival benefits from reperfusion are significantly negative for the top-tercile hospitals, suggesting that there is overuse among these hospitals. In the bottom-tercile hospitals, the estimate survival benefits from reperfusion for the lowest propensity patients are less negative and not significantly different from zero, which is consistent with appropriate use of reperfusion in these hospitals. Finally, we note that plots are also linear in log-odds despite the non-parametric nature of the estimation—later in the paper, this will allow us to use logit models and control for the propensity linearly.

In the left-hand panel of Figure 3 we plot the estimated survival benefit from reperfusion and 95% CI against the hospital effect from the propensity equation, once again, controlling non-parametrically for the propensity index. The right-hand panel is the analogous plot but estimated only for low-propensity patients whose propensity index implied that they had below a 20% probability of receiving reperfusion. Both plots show a clear downward slope, with lower benefit from treatment for patients treated by hospitals with higher random effects in the propensity equation. Among all patients

(the left-hand plot), the estimated survival benefit from reperfusion is positive for all hospitals, although it is small and not significant in hospitals with the highest treatment rates (those 2 standard deviations above average, with $\hat{\theta}_h=0.6$). In contrast, among the lowest propensity patients (the right-hand plot), only hospital's with the lowest treatment rates are estimated to have survival benefits from reperfusion that are near to zero. The estimated survival benefit from reperfusion is negative and significant in hospitals with the highest treatment rates, suggesting that there is overuse in these hospitals and, as a result, we were able to identify substantial subsets of patients who were harmed by reperfusion treatment.

B. Regression Results

Table 2 reports the estimation of the key estimating equation as described in Equation 8 and 9. The regressions include hospital fixed-effects, as the theory tells us to condition on them to control for hospital TFP (α_h^0). To help with interpretation, we have normed the propensity-index so that a value of 0 refers to the average patient receiving reperfusion. Thus, the coefficient on reperfusion is an estimate of the effect of reperfusion on an average patient receiving reperfusion. The first three columns report OLS estimates and the last two, logit estimates where the coefficients are odds ratios. The second and third columns of Table 2 report estimation of Equations 7 using two different approaches to control for the propensity index. In the first (which corresponds to Equation 7b), we restrict the index to be a simple linear function—an approach that was justified in the graphical analysis in Figure 3, in the second we control for it with 100 indicator variables (this corresponds to Equation 7a).

Column (1) is not a test of our theory, but it is included to demonstrate that the benefit of reperfusion is increasing in the propensity to receive reperfusion, and consequently, that a Roy-model of triage describes provider decision making. The coefficient on the interaction of reperfusion with the propensity index is positive and highly significant, implying that the treatment effect of reperfusion on survival is increasing in the patient's propensity index, $\lambda_1 > 0$, as predicted by our model. The coefficient on this interaction implies that an increase in the propensity index of one (about one standard deviation of the propensity index in the treated population) is associated with roughly a doubling of the treatment effect. Thus, it appears that hospitals are choosing patients for treatment based on the benefit of the treatment, and the heterogeneity in the treatment effect is large relative to the average treatment effect.

In column 2, we estimate equation 7b and include hospital fixed-effects, a linear control for the index, and an interaction of the propensity with the hospital risk-adjusted treatment rate. The propensity and hospital effects are demeaned, so the regression identifies the effect of reperfusion for the typical patient and at the typical hospital in our sample. As noted earlier, if hospital-level reperfusion rates were entirely driven by comparative advantage in treatment, we would get an OLS coefficient on 0 on this

variable and a coefficient of -1 if the variation was driven entirely by differences in thresholds. We can safely reject both numbers, and conclude that both factors are at play. The coefficient is similar in column 3, where we non-parametrically control for the interaction of reperfusion with a set of 100 dummies for each propensity percentile. This means that the non-parametric control for $g(I)$ produces the same estimates as a simple linear control. Columns 2 and 3 imply that conditional on a patient's propensity, the treatment effect is smaller in hospitals with a high propensity to treat meaning that more aggressive hospitals with lower minimum thresholds for treatment treat more patients and have lower benefits to treatment. The estimated coefficients suggest that a one standard deviation increase in the hospital effect from the propensity equation (about 0.3) lowers the return to reperfusion by about .06-.09. The last two columns of Table 2 are logit analogs to the earlier OLS regressions.

A different way to approach the problem is to estimate Equation 6 and directly obtain estimates of hospital specific thresholds. Table 4 reports results from using a mixed-effects logit and including the propensity (I) as a linear control. The mixed-effects logit allows us to include random effects for the hospital medical management terms (which identify α_h^0), and allow each hospital to have its own random-coefficient for the returns to treatment. Including hospital-random effects is much easier than hospital fixed effects and side-steps the challenges of fixed-effects estimation for small hospitals. This structure may appear to be restrictive relative to the simpler logit models in Table 2, but the restrictions do not change the estimated effect of reperfusion or reperfusion*propensity relative to the simple models. The mixed-logit allows us to recover the hospital-level thresholds, hospital-level medical management (TFP in our parlance), and correlate them. We find that the SD of hospital thresholds is large—and of the same magnitude as the direct effect of reperfusion. There is also considerable variation in hospitals ability to deliver medical-management as seen by the SD of α_h^0 . The two quantities are negatively correlated with the quality of medical management, meaning that hospitals with higher thresholds (conservative hospitals that do less) are worse at medical management. Later in the paper we find that this stems from hospitals that are worse at medical management not being able to tell, especially if they're small, that the benefits from doing more reperfusion are actually high for them.

In Table 4 we investigate the sensitivity of our results to alternative survival windows. The purpose of using 7-day survival was to examine whether the patterns noted above are evident soon after admission, and reflect decisions about how the heart-attack was initially treated. If they do not appear 7-days, the concern would be that we're picking up the effect of later treatments—for example, the quality of post-discharge care. At 7-days relative to 30 days, we expect the effect of the treatment to be even more tightly linked to a patient's propensity to receive it and that it is exactly what we find in Panel A. This relationship is half as strong for 1-year survival relative to 7 day survival (Panel B), and represents the importance of post-discharge factors in affecting 1-year survival. In both panels, the benefits of

reperfusion fall in hospitals that do more of it which is consistent with more intensive hospital working into less appropriate patients.

Our analysis relies on three key assumptions: (1) that hospitals triage patients according to a Roy-model; (2) the ‘single-index’ assumption, that the distribution of unobservables does not have a hospital specific component; (3) that we are able to estimate a ‘treatment on the treated’ parameter. We presented evidence on the first and last assumptions. The third would can be understood by noting that logit models estimate $\Pr(Treatment) = (X\beta + \theta)/\sigma_v$, where σ_v is the SD of the unobservables. If hospitals vary in σ_v , then estimates of β from more intensive or less intensive hospitals will be different. But this is testable—and we estimated separate propensity equations (as in equation 2) by intensity of hospital and found that their predictions are correlated 0.9987, suggesting that differences in variances are a first-order source of bias. We acknowledge that we’re not able to reject a world in which the variances vary across hospitals and where the coefficients scale up to perfectly offset the increased variance, but we’re unaware of an economic story for why this should happen.

To summarize the evidence so far, we have shown that (i) patients with higher appropriateness receive higher benefits from treatment (ii) that this relationship is approximately linear—which is why simpler linear-controls for the propensity to receive care do as well as non-parametric controls for the propensity to receive care (iii) that less appropriate patients are harmed in high-intensity hospitals—which is consistent with overuse (iv) that higher intensity hospitals have lower average treatment benefits which is consistent with working into less appropriate patients. Because of the simultaneous present of variation in thresholds and variation in treatment rates, we now turn to a more parametric framework to estimate both quantities.

C. Identifying Hospital Productivity Parameters

It is possible to directly estimate the minimum treatment threshold τ_h using estimates of the joint distribution of the hospital effects in the treatment equation and survival equation. This approach requires us to use the linear-approximation for $g(I) = \lambda_0 + \lambda_1 I$ instead of the non-parametric control. The linear assumption was justified by the figures, and in Table 2 where we showed very similar results from this restriction to a fully non-parametric approach. To pursue this approach, we note that since $I_{ih} = X_i\beta + \theta_h$, equation 7b can we rewritten as:

$$(7c) Y_{ih} = \alpha_h^0 + X_{ih}\beta^0 + \lambda_0 Treat_{ih} + (\lambda_1 \theta_h + \tau_h) Treat_{ih} + \lambda_1 Treat_{ih} * X_i\beta + v_{ih}^0 + \varepsilon_{ih}^0$$

Equation 7c is a hierarchical logistic model with patients nested within hospitals. It has a number of assumptions: (i) random coefficients at the hospital level that identify $\lambda_1 \theta_h + \tau_h$; (ii) We estimate the treatment propensity equation and the survival equation jointly, treating the hospital-effect in the

propensity equation (which identifies $\hat{\theta}_h$) and the hospital-specific effect of reperfusion (which identifies α_h^0) in the survival equation as jointly normal random effects. The remaining parameters determining the effect of reperfusion (λ_0, λ_1) and the variance and covariance of the hospital-level random coefficients were estimated by maximum likelihood using xtmelogit in Stata. Knowledge of $(\lambda_1 \theta_h + \tau_h)$, λ_1 and θ_h allows to back out a transformed estimate of τ_h . All the reported standard errors are conditional on the first-stage estimates $X_{ih}\hat{\beta}, X_{ih}\hat{\beta}^0$, but any adjustment for using these generated regressors is likely to be second-order because of the large samples used to estimate the patient-level coefficients.

While these may seem like a number of assumptions we are able to verify that this model is able to replicate the results and magnitudes from simpler models. For example, the coefficient on the benefit for reperfusion from the above table is 0.27 (Table 5) compared to 0.31 in the simpler model using a simpler logit model in Table 3 (both coefficients are in log-odds). The benefit of reperfusion increases with the index with similar magnitudes in both models—0.292 in Table 3 vs. 0.276 from the hierarchical logit in Table 5. The threshold and quality of medical management are correlated -0.331 in the simpler model and are correlated -0.321 in more complex model. This reassures us that the estimates are not entirely a consequence of the structure that we have imposed.

The estimates in Table 5 suggest that there is considerable variation across hospitals in τ_h , the minimum threshold for treatment (Std. Dev. = 0.327). Consistent with the evidence presented in Table 2, there is a negative correlation (-0.341) between τ_h and the reperfusion intercept θ_h , suggesting that some of the variation in treatment rates across hospitals is associated with variation in the treatment threshold (mostly overuse, as suggested by figures 2-3). The reperfusion intercept is barely correlated with hospital comparative advantage, suggesting that risk-adjusted treatment rates aren't picking up expertise as much as variation in the treatment thresholds. Thus, most of the variation across hospitals in the observed treatment effect (with Std. Dev. = 0.307) is the result of variation in the treatment threshold rather than comparative advantage. The negative correlation (-0.321) between τ_h and the survival intercept α_h^0 which represents hospital general-productivity suggests that hospitals with better survival rates when not using the treatment tend to set too low a treatment threshold and overuse the treatment. This is consistent with a model in which hospitals that are highly skilled at caring for patients without reperfusion overestimate the benefits of treatment and overuse reperfusion (or that patients who come to these hospitals demand intensive treatment because of similarly erroneous judgement).

One criticism of our approach is that distribution of α_h^Δ depends on σ – the scale parameter in the treatment propensity logit, which we have so far ignored by assuming that it is one (as in all our earlier tables). So it is important to know the scale factor in order to know specific magnitudes as opposed

to making directional statements about correlations. Under normality, the scale parameter would be equal to the standard deviation of the unobservable factors determining benefit from treatment (and would have the notation (σ_v)). While we cannot estimate it directly, we used a range of values for σ to recalculate the standard deviation of α_h^Δ and its correlation with τ_h . These are presented in Figure 5. The left hand panel plots estimates of the standard deviation of comparative advantage α_h^Δ for values of σ from 0.01 to 3, while the right hand panel plots estimates of the correlation of α_h^Δ with τ_h for the same range of σ . Interestingly, the estimates from Table 5 bound the standard deviation of α_h^Δ to be above 0.3. Thus, our estimates imply that the variation across hospitals in comparative advantage is at least as large as the variation across hospitals in the treatment threshold (SD=.327) and possibly much larger. For $\sigma < 1$, corresponding to relatively less variation in unobservable differences across patients in the benefits from treatment, our estimates imply similar amounts of variation in τ_h and α_h^Δ , and that the two are strongly positively correlated (between 0.4 and 1). In this case, a hospital's minimum treatment threshold τ_h is positively correlated with a hospital's comparative advantage α_h^Δ , meaning that hospitals with low comparative advantage tend to have low thresholds and overuse reperfusion. Interestingly, such a correlation would arise if all hospital's *incorrectly* believed that they had high comparative advantage in performing the treatment, resulting in overuse among those hospitals that actually did not have a high comparative advantage in performing the treatment. This is the mechanism that we examine in the next section.

IV. Mechanisms

As noted earlier, there are two broad mechanisms that could lead to allocative inefficiency. First, hospitals may be over-treating for financial gain (particularly in for-profit hospitals) or because of benefits to future patients through learning-by-doing (particularly in teaching hospitals). This type of mechanism would suggest that allocative inefficiency (τ) would be related to hospital characteristics such as ownership, teaching status, etc. To investigate this hypothesis, we interacted treatment with a number of hospital-characteristics such as ownership, teaching status and size (Table 6). Overuse at for-profit hospitals or teaching hospitals, or at hospitals with characteristics that are included in the table would mean that the return to treatment would be lower at such facilities. In Column 1 we establish that these factors do predict variation in the use of reperfusion, with for-profit and high-volume hospitals doing more reperfusion, and teaching hospitals and hospitals with high DSH (Disproportionate Share, a proxy for serving a poverty population) doing less reperfusion. In Column 2, however, we show that there is no evidence that these characteristics are associated with the return to treatment, conditional on the patient's

propensity to receive treatment. A joint-test on all the *Treatment*Hospital Characteristics* interactions can't reject the null-hypothesis that these variables are jointly zero (chi-squared statistic=2.96; p-value=0.96). Yet, since overtreatment is clearly evidence from the earlier exhibits, we need another mechanism for why it happens.

A second mechanism for allocative inefficiency is that the hospital had imperfect information and misperceived their comparative advantage.⁴ Given the general lack of systematic performance feedback and small samples of their own treated patients to observe, it is quite plausible that hospitals and physicians will have inaccurate beliefs about their own comparative advantage. Put differently, there is no reason to think that physicians or hospitals know their α_h^Δ perfectly—it's the difference of two parameters (α_h^0 and α_h^1) and both are probably measured with error. In this mechanism, θ represents a hospital's belief about their comparative advantage and τ represents a hospital's misperception (or prediction error) of their own comparative advantage.

More formally, we can reinterpret our empirical model in the following way. Suppose that a hospital does not know its comparative advantage, but instead has a belief about its comparative advantage which is given by θ . Based on this belief, they treat patients if the expected benefit of treatment is positive. Thus, patients are treated based on beliefs (if $\theta_h + X_i\beta^\Delta + v_{ih}^\Delta > 0$) rather than based on actual comparative advantage (if $\alpha_h^\Delta + X_i\beta^\Delta + v_{ih}^\Delta > 0$). Let τ_h represent the difference between a hospital's actual comparative advantage and their beliefs about it, so that $\tau_h = \alpha_h^\Delta - \theta_h$ is the hospital's prediction error (and therefore $\theta_h = \alpha_h^\Delta - \tau_h$, as in our empirical model). Thus, this framework interprets τ_h as arising from an inaccurate belief about α , rather than assuming that hospitals know α and consciously set $\tau \neq 0$ to achieve other objectives. A negative τ_h implies that the hospital over-estimated their comparative advantage and, as a result, treated some patients who were in fact harmed by the treatment. In this reframing, the key question is how hospitals form their beliefs.

Suppose that each hospital receives a noisy signal of their comparative advantage (S), where $S = \alpha + \omega$ and the noise (ω) is independent of α with variance σ_ω^2 (we have suppressed the subscripts & superscripts to simplify notation). Based on this signal, the hospital forms a prediction of its comparative advantage (θ). If the hospital knew the reliability of the signal ($r = \sigma_\alpha^2 / (\sigma_\alpha^2 + \sigma_\omega^2)$, where σ_α^2 is the variance of α_h^Δ across hospitals), then the optimal prediction of α given S is the posterior mean, given by $E(\alpha_h^\Delta | S) = r * S$. More generally, we assume that hospitals may not know the reliability of the signal, and form their prediction using $\theta = w * S$, where $w \neq r$. Incorrectly weighting the signal generates additional variation in the prediction error (τ) which leads to greater allocative inefficiency. Even if hospital beliefs are optimal

⁴ We are grateful to Janet Currie for suggesting this interpretation and alerting us to related work in Currie and MacLeod (2017).

given S (i.e., $w=r$), there will be allocative inefficiency ($\tau \neq 0$) because hospitals have imperfect information, and this information only predicts a fraction (r) of the true variation in comparative advantage.

This simple framework delivers a number of strong empirical implications. First, because the error in the signal is assumed to be independent of the hospital's actual productivity and comparative advantage, this framework constrains the number of parameters in the structural model to six, which allows us to identify the scale parameter and the variance in comparative advantage. More specifically, our empirical model from Table 5 estimated 6 reduced-form moments (the variances and covariances of θ , τ , and α_0) which are a function of the 6 unknown structural parameters in this framework (the variance and covariance of α and α_0 , the reliability of the signal r , the weight placed on the signal w , and the scale parameter from the logit σ).⁵ Therefore, we can derive estimates of the unknown structural parameters for this model using minimum chi-squared estimation (Wooldridge 2010, p.442-446). Minimum chi-squared estimation chooses the structural parameters that provide a best fit of the reduced-form estimates (in a weighted least squares sense, using the standard errors & covariance of the reduced form estimates to form weights). In the just-identified case the resulting structural parameters fit the reduced form moments exactly because they are just a function of the reduced form parameters. Restrictions on the structural parameters can be tested based on how they affect the structural models ability to fit the reduced form estimates (through a chi-squared goodness of fit statistic).

Just-identified estimates of the structural parameters for this model are provided in the first column of Table 7. There is substantial variation in comparative advantage (standard deviation of $\alpha = 0.317$), with the variation across hospitals being as large as the average treatment effect. The signal that hospital's receive about their comparative advantage is estimated to have very low reliability ($r=.065$), but hospital's place more weight on the signal than is optimal, with $w=0.154$.

If $w=r$ in this framework then beliefs are optimal. When we impose $w=r$ (column 2 of Table 7), we are over-identified (estimating 5 parameters from 6 moments) and can use the chi-squared goodness of fit statistic to test the restriction (Wooldridge, 2010, pp. 444-445). This statistic rejects the hypothesis that $w=r$ (chi-squared with 1 df = 10.4, $p=.001$). In other words, the constrained model with $w=r$ implies reduced-form variances and correlations of θ , τ , and α_0 that are significantly different from those estimated in Table 5. More specifically, if $w=r$ then the hospital's prediction (θ) is optimal and, therefore, should be uncorrelated with the prediction error (τ). The fact that we estimated a significant negative correlation of -0.34 between θ and τ in table 5 implies that hospitals' predictions are not optimal and they are overweighting the noisy signal ($w>r$), i.e. they over-react to the signal. One might not expect hospitals

⁵ See Appendix-C for derivation of the equations stating the reduced-form moments in terms of the structural parameters.

to have the information necessary to form optimal weights – in particular, they may not know how much true variation in comparative advantage there is across hospitals, and are acting as if they are using an over-diffuse prior (placing too much weight on their own signal, and not shrinking enough to a prior mean).

Finally, if hospitals learn based on their experience with patients, then one would expect that low-volume hospitals would have less reliable signals of their comparative advantage than high-volume hospitals. In the third column of table 7 we fit our model to reduced-form moments estimated separately for low, medium and high-volume hospitals (6 moments for each group, for a total of 18 moments), allowing the reliability parameter to vary across the 3 groups but otherwise constraining the remainder of the model parameters to be equal across the 3 groups (8 parameters total). As expected, the reliability of the signal is estimated to be highest for the high-volume hospitals and lowest for the low-volume hospitals. Moreover, the goodness of fit statistic cannot reject our model (chi-squared with 10 df = 12.7, p=.24) suggesting that this simple model provides an adequate fit of the data. In other words, the model estimated in column 3 of Table 7 implies reduced-form variances and correlations of θ , τ , and α_0 that are not significantly different from the unconstrained reduced-form estimates for low, medium, and high-volume hospitals. Assuming that reliability of the signal is the same for high, medium and low-volume hospitals (final column of table 7) is strongly rejected (chi-squared with 2 df = 49.7, p<.0001) and such a model is strongly rejected by the goodness-of-fit test (chi-squared with 12 df = 62.4, p<.0001).

V. Conclusion

Using a Roy model of treatment to motivate our empirical framework, we find significant evidence of allocative inefficiency across hospitals ($\tau \neq 0$). We can use our results to construct a stylized estimate of the welfare loss generated by this allocative inefficiency, along the lines suggested by Phelps (2000).

Returning to the intuition from Figure 1b, the effect of a non-zero τ generates a standard welfare loss

triangle, where $Welfare\ Loss_h = \frac{1}{2} \cdot \overbrace{(\tau_h)}^{height} \overbrace{(\tau_h)}^{base} \frac{dPr(Treatment)}{d\tau}$. This welfare loss is the average reduction in (logodds) survival across all patients at the hospital due to allocative inefficiency.⁶

The expected welfare loss across all hospitals is therefore given by $\frac{1}{2} \cdot dPr(Treatment)/d\tau_h E(\tau_h^2)$

To get an estimate of $dPr(Treatment)/d\tau$, which is a change in the propensity to receive treatment for a small increase in τ , we took a tiny change of 0.01 in τ_h , divided it by our estimate of the scale factor (σ_v) of 0.44 to turn it into how much change that would create in the hospitals risk-adjusted treatment rate θ_h . This yielded an estimate of $dPr(Treatment)/d\tau = -0.26$. To estimate $E(\tau_h^2)$ note that $E(\tau_h^2) =$

⁶ The welfare loss is measured in the same units as τ (logodds of survival in our estimates) and is the welfare loss per patient because we use the probability of treatment rather than total number treated.

$Var(\tau_h) + [E(\tau_h)]^2$. From the hierarchical-logit model in Table 5, we estimated $SD(\tau_h)=0.33$.⁷ If we assume that there is no allocative inefficiency on average ($E(\tau_h) = 0$), then the welfare loss from variation in allocative inefficiency as $(1/2)*(0.33^2)*-0.26=-0.014$, i.e. the allocative inefficiency across hospitals results in an average reduction in the logodds of survival per patient of .014. The overall benefit from treatment is the benefit among the treated (0.20 in log-odds), who comprise 20 percent of the patient population for a total benefit of 0.04. This means that we could increase the effectiveness of treatment by about a third if we removed the allocative inefficiency across hospitals. There is additional welfare loss if the mean of τ is not equal to zero, e.g. if there is systematic overuse across all hospitals ($E(\tau_h) < 0$). This part of the welfare calculation is far more speculative, but a good guess about systematic overuse across all hospitals comes from the average treatment effect among low propensity patients, which is about -.1 from Figure 4. Thus, the additional welfare loss from systematic overuse would be $(1/2)*(-0.1^2)*-0.26=-0.0013$. This calculation suggests that systematic overuse is a minor concern relative to the welfare loss from the overall variation.

In addition to the welfare loss from allocative inefficiency, we also found evidence of substantial variation in comparative advantage across hospitals, with the benefits from treatment being much higher in some hospitals than others. This variation in the benefits from treatment implies that “one size fits all” policies such as strict treatment guidelines are incorrect, since hospitals with greater comparative advantage at a treatment should use it more among their patients. Moreover, our evidence suggests that much of the allocative inefficiency that we observe is due to hospitals having imperfect information and misperceiving their comparative advantage. This is a different mechanism than explaining variations by appealing to medical malpractice or financial entrepreneurship by providers (Gawande, 2009). Thus, rather than reducing treatment variation across hospitals, better information about treatment effect heterogeneity across hospitals is key to improving patient welfare. We don’t know if these findings and conclusions generalize to settings beyond the treatment of heart attack patients, but our framework is general and can be applied to a variety of settings.

Finally, our work suggests new directions for research on productivity in healthcare. A large literature in economics and medicine has studied variation in how patients are treated, and much of the emphasis has been on the importance of risk-adjustment, and whether higher risk-adjusted treatment rates generate superior outcomes [Fisher et al (2003a, 2003b), Doyle (2001), Skinner (2011), Doyle (2011), Finkelstein, Gentzkow, Williams (2016), Doyle, Graves, Gruber, Kleiner (2015)]. This literature interprets variation in risk-adjusted treatment rates as isolating allocative inefficiency. Another research cluster, influenced by the productivity literature (Syverson, 2011), has emphasized productivity differences in healthcare [Chandra and Staiger (2007), Chandra et al (2016), Skinner and Staiger (2015), Currie and MacLeod

⁷ The mixed-logit model had yielded a very similar estimate of 0.31(see Table 3)

(2017)], but ignored the possibility that of variation in treatment thresholds across hospitals. Our work unites these literatures and provides a way to disentangle allocative inefficiency and productive inefficiency when they're both present. Moreover, our results highlight the key role that variation in productive efficiency (comparative advantage) plays in generating variation, and how lack of information about this variation generates welfare loss. Thus, future work should explore the sources of variation in productive efficiency and how patients and providers learn about and respond to variation in productive efficiency (Chandra et al., 2016).

References

- Baicker Katherine, Chandra Amitabh, Skinner Jonathan. 2012. Saving Money or Just Saving Lives? Improving the Productivity of the U.S. Health Care Spending. *Annual Review of Economics* 2012;4: 33- 56.
- Chandra, Amitabh, Amy Finkelstein, Adam Sacarny and Chad Syverson. 2016. "Health Care Exceptionalism? Performance and Allocation in the US Health Care Sector." *American Economic Review*, 106(8): 2110-44.
- Chandra, Amitabh and Douglas O. Staiger. 2007. "Productivity Spillovers in Healthcare: Evidence from the Treatment of Heart Attacks." *Journal of Political Economy*.
- Chandra, Amitabh and Jonathan Skinner. 2012. "Technology Growth and Expenditure Growth in Health Care." *Journal of Economic Literature*, 50(3): 645-80.
- Currie, Janet and MacLeod, W. Bentley. 2017. Diagnosing Expertise: Human Capital, Decision Making and Performance Among Physicians. *Journal of Labor Economics*, forthcoming.
- Doyle Jr, Joseph J. 2001. "Returns to local-area healthcare spending: evidence from health shocks to patients far from home." *American economic journal. Applied economics* 3.3: 221.
- Doyle Jr, Joseph J., John A. Graves, Jonathan Gruber, and Samuel A. Kleiner. 2015. "Measuring Returns to Hospital Care: Evidence from Ambulance Referral Patterns," *Journal of Political Economy* 123, no. 1 (February): 170-214.
- Finkelstein, Amy and Gentzkow, Matthew and Williams, Heidi L. 2016. Sources of Geographic Variation in Health Care: Evidence from Patient Migration. *Quarterly Journal of Economics*:
- Fibrinolytic Therapy Trialists' Collaborative Group. 1994. "Indications for Fibrinolytic Therapy in Suspected Acute Myocardial Infarction: Collaborative Overview of Early Mortality and Major Morbidity Results from all Randomized Trials of More Than 1000 patients." *Lancet* 343(8893): 311-22.
- Fisher Elliott, et al. The Implications Of Regional Variations In Medicare Spending. Part 1: The Content, Quality, And Accessibility Of Care. *Annals of Internal Medicine*. 2003a February 18;138(4):273–87.
- Fisher Elliott, et al. The implications of regional variations in Medicare spending. Part 2: Health outcomes and satisfaction with care. *Annals of Internal Medicine*. 2003b February 18;138(4):288–98.
- Garber, Alan M. and Jonathan Skinner. 2008. "Is American Health Care Uniquely Inefficient?" *Journal of Economic Perspectives*, 22(4): 27-50.
- Gawande, Atul. 2009. "The Cost Conundrum," *The New Yorker*, June 1. Available at <http://www.newyorker.com/magazine/2009/06/01/the-cost-conundrum>
- Heckman JH, S Urzua, E Vytlacil. 2006. "Understanding Instrumental Variables in Models with Essential Heterogeneity." *The Review of Economics and Statistics* 88(3):389-432.
- Marciniak, TA, EF Ellerbeck, MJ Radford, et al. 1998. "Improving the Quality of Care for Medicare Patients with Acute Myocardial Infarction: Results from the Cooperative Cardiovascular Project." *Journal of the American Medical Association* 279:1351-7.
- McClellan, Mark, Barbara J. McNeil, Joseph P. Newhouse. 1994. "Does More Intensive Treatment of Acute Myocardial Infarction in the Elderly Reduce Mortality? Analysis Using Instrumental Variables." *Journal of the American Medical Association* 272(11): 859-866.

- O'Connor, GT, HB Quinton, ND Traven, et al. 1999. "Geographic Variation in the Treatment of Acute Myocardial Infarction: The Cooperative Cardiovascular Project." *Journal of the American Medical Association* 281: 627-33.
- Phelps, Charles E. 2000. Information Diffusion and Best Practice Adoption. *Handbook of Health Economics*, Volume 1, Elsevier: 223-264.
- Skinner, Jonathan. 2011. "Causes and Consequences of Regional Variations in Health Care," in *Handbook of Health Economics*, Volume 2, Elsevier: 45-94.
- Syverson, Chad. 2011. "What Determines Productivity?" *Journal of Economic Literature*, 49(2): 326-65.
- Wooldridge, J. M. 2010. *Econometric Analysis of Cross Section and Panel Data*. 2nd ed. Cambridge, Massachusetts: MIT Press.

Appendix-A

Construction of CCP Estimation Sample

The CCP used bills submitted by acute care hospitals (UB-92 claims form data) and contained in the Medicare National Claims History File to identify all Medicare discharges with an International Classification of Diseases, Ninth Revision, Clinical Modification (ICD-9-CM) principal diagnosis of 410 (myocardial infarction), excluding those with a fifth digit of 2, which designates a subsequent episode of care. The study randomly sampled all Medicare beneficiaries with acute myocardial infarction in 50 states between February 1994 and July 1995, and in the remaining 5 states between August and November, 1995 (Alabama, Connecticut, Iowa, and Wisconsin) or April and November 1995 (Minnesota); for details see O'Connor et al. (1999). Among patients with multiple myocardial infarction (MIs) during the study period, only the first AMI was examined. The Claims History File does not reliably include bills for all of the approximately 12% of Medicare beneficiaries insured through managed care risk contracts, but the sample was representative of the Medicare fee-for-service (FFS) patient population in the United States in the mid-1990s. After sampling, the CCP collected hospital charts for each patient and sent these to a study center where trained chart abstracters abstracted clinical data. Abstracted information included elements of the medical history, physical examination, and data from laboratory and diagnostic testing, in addition to documentation of administered treatments. The CCP monitored the reliability of the data by monthly random reabstractions. Details of data collection and quality control have been reported previously in Marciniak et al. (1998). For our analyses, we delete patients who were transferred from another hospital, nursing home or emergency room since these patients may already have received care that would be unmeasured in the CCP. We transformed continuous physiologic variables into categorical variables (e.g., systolic BP < 100 mm Hg or \geq 100 mm Hg, creatinine <1.5, 1.5-2.0 or >2.0 mg/dL) and included dummy variables for missing data. Our choice of variables was based on those selected by Fisher et al. (2003a,b) and Barnato et al. (2005). With the exception of two variables that are both measured by blood-tests, albumin and bilirubin (where the rates of missing data were 24 percent), we do not have a lot of missing data (rates were less than 3 percent). Included in our model are the following risk-adjusters:

Age, Race, Sex (full interactions)
previous revascularization (1=y)
hx old mi (1=y)
hx chf (1=y)
history of dementia
hx diabetes (1=y)
hx hypertension (1=y)
hx leukemia (1=y)
hx ef <= 40 (1=y)
hx metastatic ca (1=y)
hx non-metastatic ca (1=y)
hx pvd (1=y)
hx copd (1=y)
hx angina (ref=no)

hx angina missing (ref=no)
hx terminal illness (1=y)
current smoker
atrial fibrillation on admission
cpr on presentation
indicator mi = anterior
indicator mi = inferior
indicator mi = other
heart block on admission
chf on presentation
hypotensive on admission
hypotensive missing
shock on presentation
peak ck missing
peak ck gt 1000

no-ambulatory (ref=independent)
ambulatory with assistance
ambulatory status missing
albumin low(ref>=3.0)
albumin missing(ref>=3.0)
bilirubin high(ref<1.2)
bilirubin missing(ref<1.2)
creat 1.5-<2.0(ref=<1.5)
creat >=2.0(ref=<1.5)
creat missing(ref=<1.5)
hematocrit low(ref=>30)
hematocrit missing(ref=>30)
ideal for CATH (ACC/AHA criteria)

Appendix-B

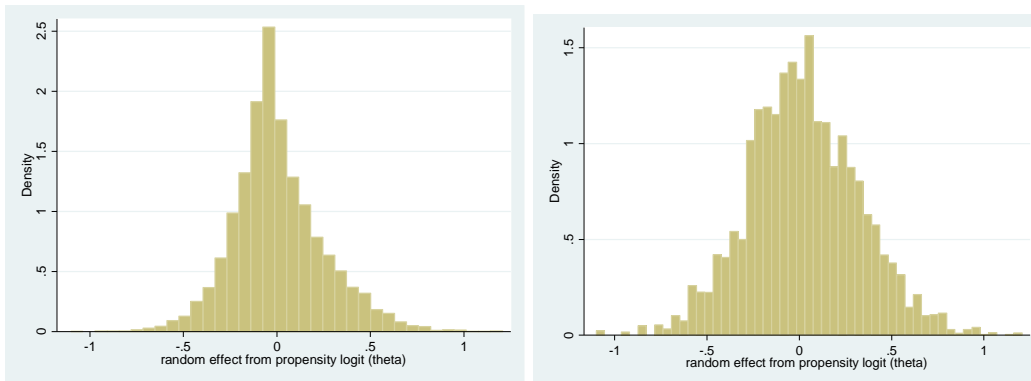
Estimation of Propensity to Receive Reperfusion

We compared results from fixed-effects and random-effects logits predicting reperfusion as a function of the full list of patient risk-adjusters and a random hospital-level intercept (Equation 5), and examined the sensitivity of different approaches to estimating the slope parameters for this equation, ranging from OLS with hospital fixed-effects, conditional logit with hospital fixed effects, logit with random effects, and mixed-logit. If equation (2) is estimated using xtmelogit in Stata, then we obtain Bayesian posterior estimates of hospital random-effects, commonly referred to in the literature as shrinkage estimates or, in linear models, best linear unbiased predictions. As the table below notes, the choice between OLS and logit and between random effects and fixed effects is not significant for the estimation of hospital effects, but the use of shrinkage is because of the substantial number of small hospitals in our sample.

Model: $\Pr(\widehat{Treatment}_{ih}) = F(\widehat{I}_{ih}) = F(X_i\hat{\beta} + \hat{\theta}_h)$	Correlation of Patient Characteristics ($X_i\hat{\beta}$)	Correlation of Hospital Effects ($\hat{\theta}_h$)
Fixed Effects OLS and Fixed-Effects Logit	0.9745	0.9997
Fixed Effects OLS and Random-Effects Logit (Unshrunk)	0.9997	0.9997
Fixed Effects OLS and Random-Effects (Shrunk)	0.9732	0.9998

The results from the shrunken random-effects were used to form posterior estimates of the hospital random effects $\hat{\theta}_h$ and an estimate of the propensity index for each patient $\hat{I}_{ih} = X_{ih}\hat{\beta} + \hat{\theta}_h$. The coefficients on the patient-level variables are consistent with the medical literature, with reperfusion being less likely among patients with pre-existing conditions and who are older, and also depending on the location and severity of the heart attack. The estimated standard deviation of the hospital effect is 0.44 (Std. Err. = 0.01), which implies that a one standard deviation in the hospital effect increases the logodds of receiving reperfusion by 0.44, which would increase an average patients probability of receiving reperfusion from 19% to 26%. Thus, there is sizable variation across hospitals in the rate at which they provide reperfusion to observationally similar patients. The model is able to predict much of the hospital-level variation, with the posterior prediction of each hospital's effect on reperfusion having a standard deviation of 0.30 in our data.

Histogram of hospital-risk adjusted treatment rates are below (Panel A is hospital level and Panel B is patient weighted)



Appendix-C

Equations used in minimum distance estimation

This appendix describes the equations used in minimum distance estimation that state the reduced-form estimates in terms of the parameters of the structural model. The mixed-logit model from Table 5 estimates 6 reduced-form moments (the variances and covariances of θ , τ , and α_0) and their associated variance covariance matrix. Call this 1x6 vector of reduced form parameters $\beta = (\sigma_\theta^2, \sigma_\tau^2, \sigma_{\alpha_0}^2, \sigma_{\theta\tau}, \sigma_{\theta\alpha_0}, \sigma_{\tau\alpha_0})$ and let $\hat{\beta}$ be the vector of estimates of these parameters and $V = Var(\hat{\beta})$ be the associated 6x6 variance matrix of these estimates. Our structural model has 6 unknown structural parameters: the variance and covariance of α and α_0 , the variance of the noise in the signal ω , the weight placed on the signal w , and the scale parameter from the logit σ_v . Call this 1x6 vector $\delta = (\sigma_{\alpha^\Delta}^2, \sigma_{\alpha_0}^2, \sigma_{\alpha^\Delta\alpha_0}, \sigma_\omega^2, w, \sigma_v)$. We can state the reduced form parameters as a function of the structural parameters (as shown below) so that $\beta = f(\delta)$. Then minimum distance estimates of δ minimize the objective function $(\hat{\beta} - f(\delta))V^{-1}(\hat{\beta} - f(\delta))'$. In the just-identified case the resulting structural parameters fit the reduced form moments exactly ($\hat{\beta} = f(\hat{\delta})$). Restrictions on the structural parameters can be tested based on how they affect the structural models ability to fit the reduced form estimates, using the fact that the objective function has a chi-squared distribution with degrees of freedom equal to the degree of over-identification (the difference between the dimension of $\hat{\beta}$ and the dimension of δ). Fitting the model to reduced form estimates from low, medium and high-volume hospitals is done similarly, where $\hat{\beta}$ stacks the estimates from the three samples into a 1x18 vector and V is 18x18 with the variance covariance matrix for estimates from each of the samples along the diagonal and zeros everywhere else.

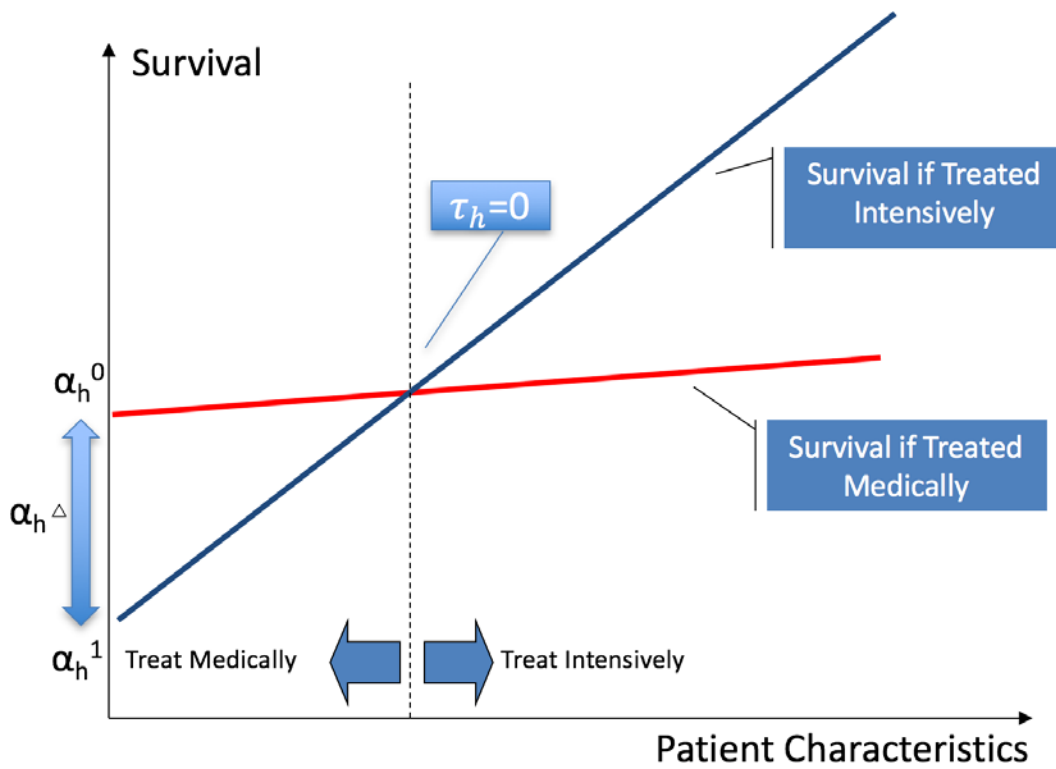
The structural model interprets our original model parameters (θ , τ , and α_0) as follows:

1. α_0 is unchanged
2. $\theta^{**} = wS = w(\alpha^\Delta + \omega)$ which is the hospital's prediction of its comparative advantage given its signal. In the treatment propensity logit we estimate $\theta = \frac{\theta^*}{\sigma_v} = \frac{1}{\sigma_v}w(\alpha^\Delta + \omega)$
3. $\tau = \alpha^\Delta - \theta^* = (1 - w)\alpha^\Delta - w\omega$

Using these definitions, it is straightforward to derive the following relationships between the reduced form estimates and the parameters of the structural model:

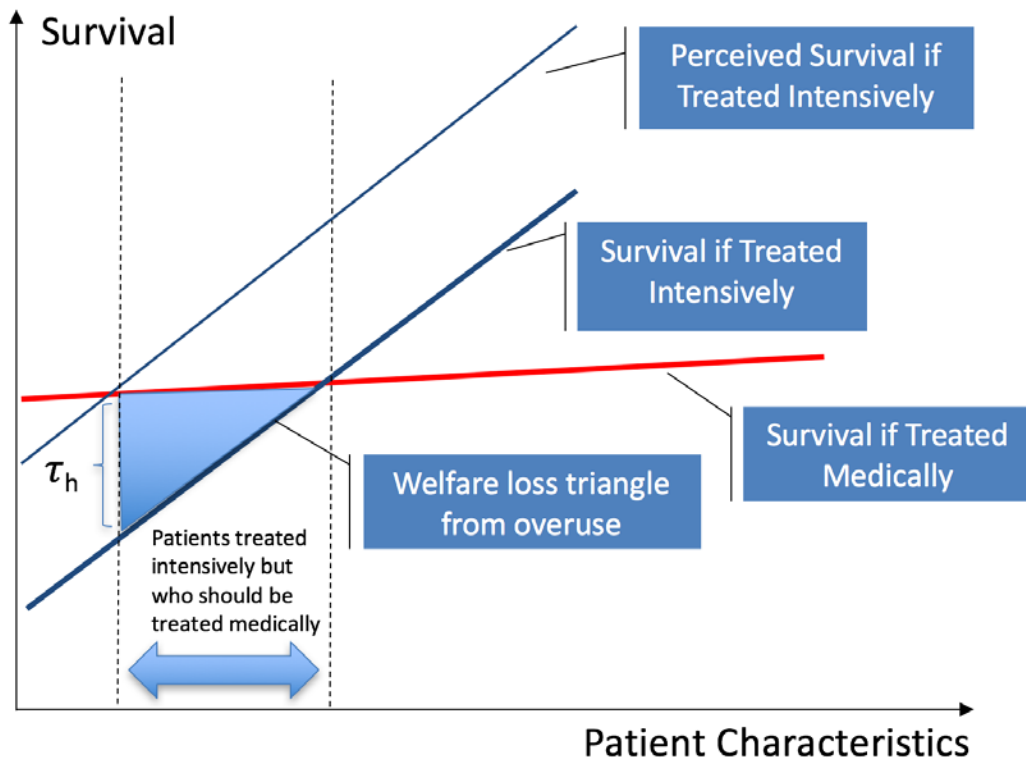
- i. $\sigma_\theta^2 = Var\left[\frac{1}{\sigma_v}w(\alpha^\Delta + \omega)\right] = \frac{1}{\sigma_v^2}w^2(\sigma_{\alpha^\Delta}^2 + \sigma_\omega^2)$
- ii. $\sigma_\tau^2 = Var[(1 - w)\alpha^\Delta - w\omega] = (1 - w)^2\sigma_{\alpha^\Delta}^2 + w^2\sigma_\omega^2$
- iii. $\sigma_{\alpha_0}^2 = Var[\alpha_0] = \sigma_{\alpha_0}^2$
- iv. $\sigma_{\theta\tau} = Cov\left[\frac{1}{\sigma_v}w(\alpha^\Delta + \omega), (1 - w)\alpha^\Delta - w\omega\right] = \frac{1}{\sigma_v}w(1 - w)\sigma_{\alpha^\Delta}^2 - \frac{1}{\sigma_v}w^2\sigma_\omega^2$
- v. $\sigma_{\theta\alpha_0} = Cov\left[\frac{1}{\sigma_v}w(\alpha^\Delta + \omega), \alpha_0\right] = \frac{1}{\sigma_v}w\sigma_{\alpha^\Delta\alpha_0}$
- vi. $\sigma_{\tau\alpha_0} = Cov[(1 - w)\alpha^\Delta - w\omega, \alpha_0] = (1 - w)\sigma_{\alpha^\Delta\alpha_0}$

Figure 1A: A Roy model of Treatment at the Hospital level



The two lines denote patient survival if a hospital treats a given patient medically (intercept is α_h^0) or using reperfusion (intercept is α_h^1) as a function of patient characteristics (i.e. patient X's) on the x-axis. We have suppressed the distribution of unobservables that come out of the plane. Expertise at medical and intensive care is captured by the intercepts α_h^0 and α_h^1 respectively, with comparative advantage being the difference between them. Allocative efficiency means that reperfusion should be performed to the point that the marginal patient receiving it receives zero benefit, that $\tau_h = 0$.

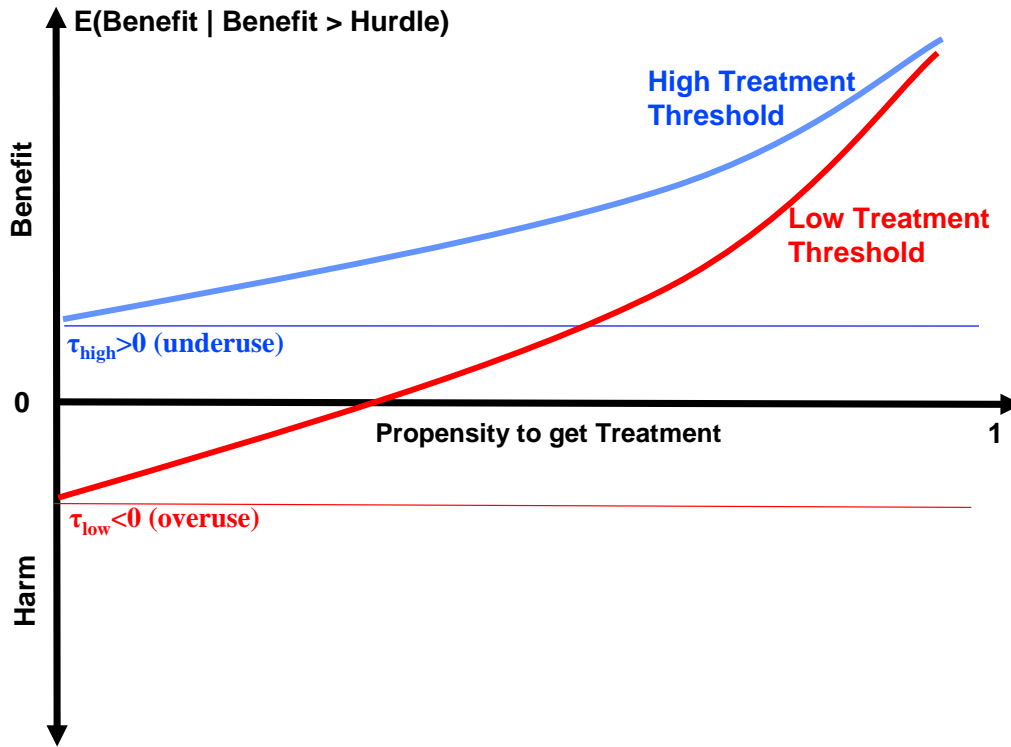
Figure 1B: A Roy model of Treatment at the Hospital level with Allocative Inefficiency



The figure illustrates the presence of allocative inefficiency. Here, perceptions about comparative advantage at delivering intensive treatment mean that more patients are treated with intensively than is optimal. As drawn, the hospital overuses intensive treatment and uses a negative threshold (τ_h). It is also possible that some hospital overuse intensive treatment because of maximizing something other than survival. A welfare loss triangle emerges.

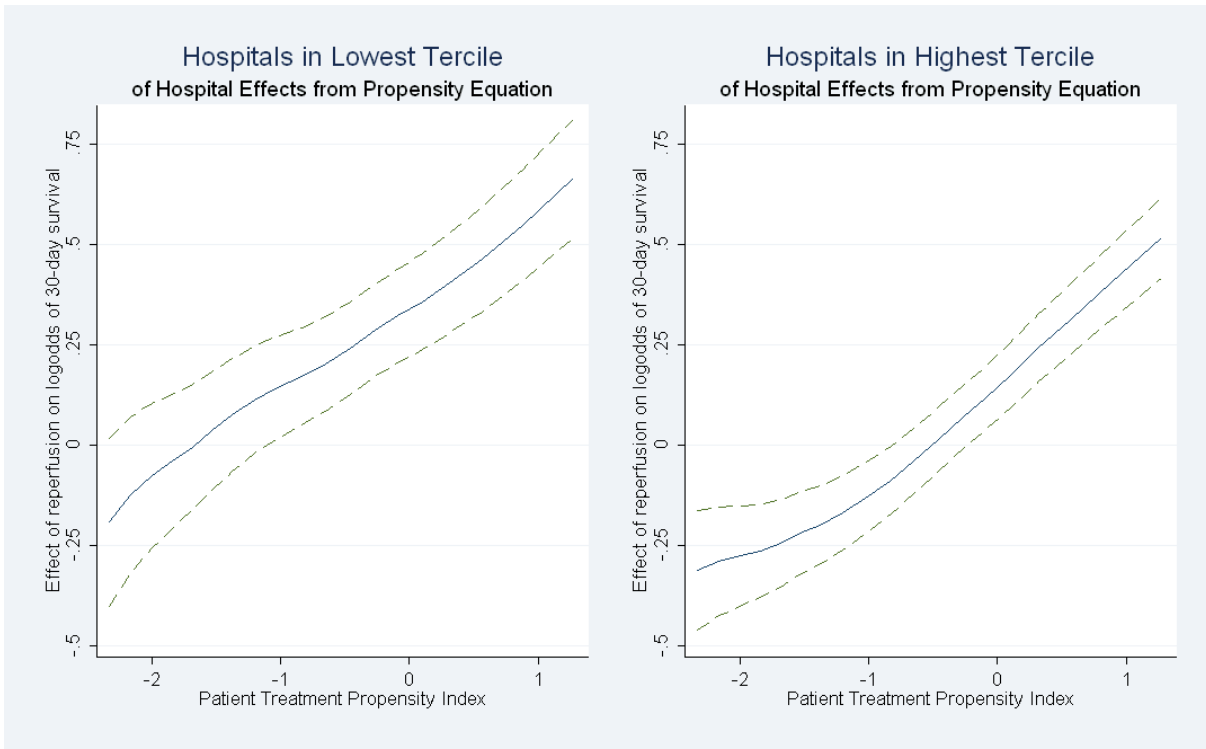
Figure 2: Distinguishing underuse and overuse using the propensity to receive treatment.

Underuse and Overuse



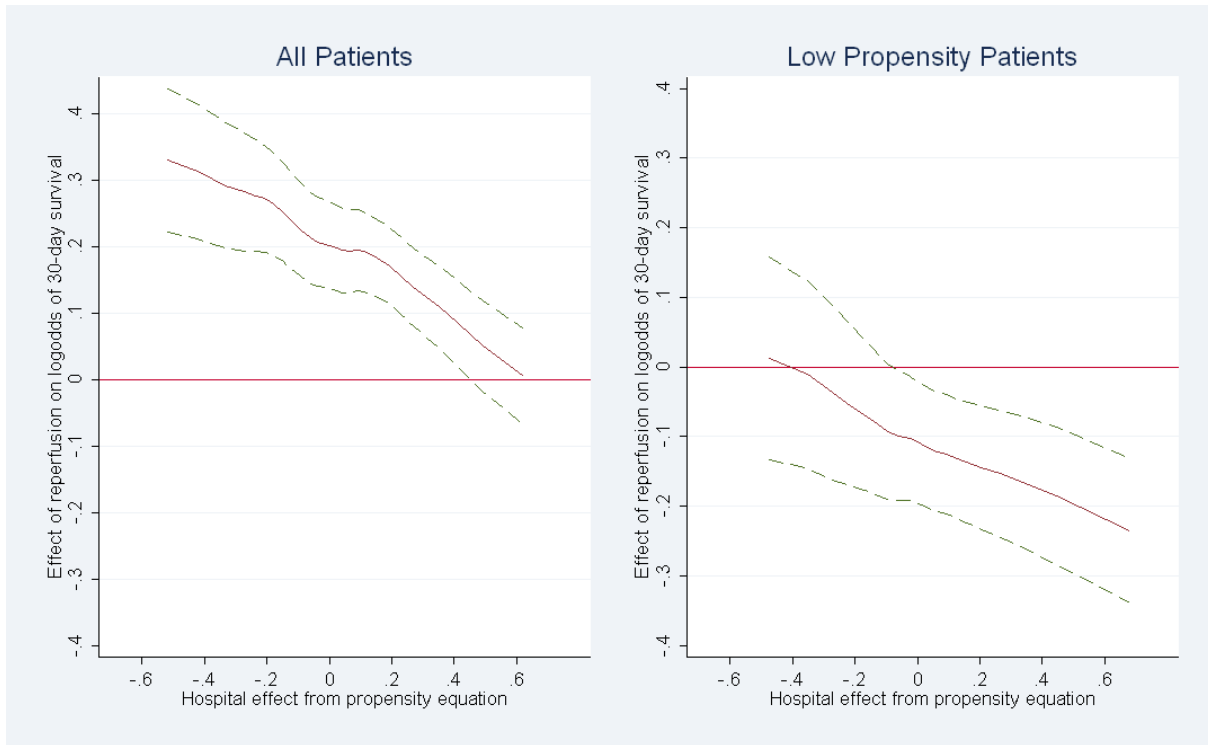
The figure illustrates the relationship between the expected benefit from treatment, $E(Y_{ih}^{\Delta} \mid Y_{ih}^{\Delta} > \tau_h)$, on the vertical axis, and the propensity index I on the horizontal axis. The propensity to receive treatment depends on patient characteristics and a hospital's assessment of its hospital-specific benefit from treatment. The curves represent the treatment-on-the-treated effect for a patient with index I , and approach the minimum threshold (τ) for a patient with a low propensity of being treated. The top curve represents a hospital with a high treatment threshold (underuse) and the bottom curve represents a hospital with a low treatment threshold (overuse).

Figure 3: Survival Benefit from Reperfusion According to Patient’s Treatment Propensity, Low-Treatment-Rate (Left) and High-Treatment-Rate (Right) Hospitals.



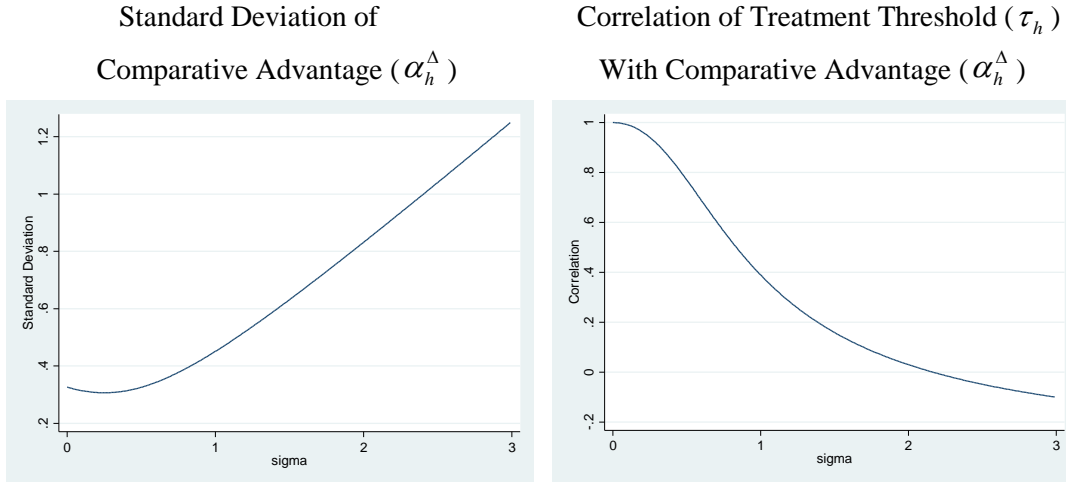
The figures plot the estimated survival benefit (and 95% confidence intervals) from reperfusion against a patient’s treatment propensity index for hospital’s in the lowest (left-hand side) and highest (right-hand side) terciles of the estimated hospital effect from the propensity equation. Propensity Equation is $\Pr(\text{Reperfusion}) = F(Xb + \text{Hospital Effect})$ and is estimated using a logit model; see Appendix-B. Propensity index refers to the logit index $(XB + \text{Hospital Effect})$. It is demeaned to the average value of patients receiving reperfusion. All models include all CCP risk-adjusters.

Figure 4: Survival Benefit from Reperfusion by to Hospital Effect from Treatment Propensity, All patients (Left) and Low-propensity patients (Right).



The left-hand panel plots the estimated survival benefit from reperfusion (and 95% confidence interval) against the hospital effect from the propensity equation using a locally-weighted logit model to estimate the reperfusion effect (controlling non-parametrically for the propensity index as was done in column 3 of Table 2). The right-hand panel is the analogous plot estimated only for low-propensity patients whose propensity index implied that they had below a 20% probability of receiving reperfusion. Propensity Equation is $\Pr(\text{Reperfusion})=F(Xb+\text{Hospital Effect})$ and is estimated using a logit model; see Appendix-B. Propensity index refers to the logit index $(XB+\text{Hospital Effect})$. It is demeaned to the average value of patients receiving reperfusion. All models include all CCP risk-adjusters.

Figure 5: Estimates of the Standard Deviation of Comparative Advantage (α_h^Δ) and its Correlation with a Hospital's Treatment Threshold (τ_h) for a Range of Values of the Scale Parameter (σ).



We used the estimates from Table 3 to calculate estimates of the standard deviation of α_h^Δ and its correlation with τ_h for a range of values for σ . The left hand panel plots estimates of the standard deviation of comparative advantage α_h^Δ for values of σ from 0.01 to 3, while the right hand panel plots estimates of the correlation of α_h^Δ with τ_h for the same range of σ .

Table 1: Patient Characteristics, Full Sample and by Reperfusion

Variable	Full Sample	Received Reperfusion w/in 12 hours	No Reperfusion w/in 12 hours
Survival 30 days post-AMI	81%	86%	80%
Reperfusion within 12 hours	19%	100%	0%
Age	77	73	77
Previous diagnoses:			
Congestive Heart Failure	22%	7%	25%
Hypertension	62%	56%	63%
Diabetes	30%	23%	32%
Dementia	6%	2%	7%
Number of observations	138,957	25,876	113,081

Note: Full-list of variables is in Appendix-A.

Table 2: Effect of Reperfusion on 30-day Survival, OLS and Logit Estimates

	OLS (1)	OLS (2)	OLS (3)	Logit (4)	Logit (5)
Reperfusion	0.039 (0.003)	0.043 (0.003)	non-parametric	0.328 (0.027)	non-parametric
Reperfusion*Propensity index	0.040 (0.002)	0.042 (0.002)	non-parametric	0.291 (0.018)	non-parametric
Reperfusion* Hospital Treatment Rate (θ)		-0.031 (0.009)	-0.037 (0.009)	-0.211 (0.076)	-0.254 (0.077)
Hospital Fixed-Effects	Yes	Yes	Yes	Yes	Yes
Control for Propensity Index	None	Linear	Non-Parametric	Linear	Non-Parametric

Note: Dependent variable is the whether patient survived to 30 days. Reperfusion measures receipt of reperfusion therapy within 12 hours of admission. OLS coefficients are percentage-point changes in survival and logit coefficients are odds ratios. Propensity Equation is $\Pr(\text{Reperfusion})=F(\text{XB}+\text{Hospital Effect})$ and is estimated using a logit model; see Appendix-B. Propensity index refers to the logit index ($\text{XB}+\text{Hospital Effect}$). It is demeaned to the average value of patients receiving reperfusion. All models include all CCP risk-adjusters. Columns 2 and 4 include linear controls for propensity-index. Columns 3 and 5 includes 100 percentiles of propensity-index interacted with the receipt of Reperfusion.

Table 3: Effect of Reperfusion on 30-day Survival, Mixed-Logit Estimates

	(1)	(2)
Reperfusion	0.297 (0.0218)	0.314 (0.0243)
Reperfusion*Propensity index	0.289 (0.0169)	0.292 (0.0171)
Std dev of hospital intercept (α^0)	0.188 (0.0151)	0.198 (0.0168)
Hospital Level Random-Intercept (α^0)	Yes	Yes
Hospital Level Random Coefficient on Reperfusion (τ)	No	Yes
Std dev of hospital coefficient on reperfusion (identifies τ ; hospital level thresholds)		0.313 (0.0557)
corr(hospital level intercept, coefficient on reperfusion) (identifies corr (α^0 , τ))		-0.331 (0.154)
Number of Hospitals	4,690	4,690
Control for Propensity Index	Linear	Linear

Note: Coefficients are log-odds. Propensity Equation is $\Pr(\text{Reperfusion})=F(Xb+\text{Hospital Effect})$ and is estimated using a logit model; see Appendix. Propensity index refers to the logit index $(XB+\text{Hospital Effect})$. It is demeaned to the average value of patients receiving reperfusion. All models include all CCP risk-adjusters. Columns 2 and 4 include linear controls for propensity-index as in equation 7a. Sample-size in every regression is 138,957.

Table 4: Effect of Reperfusion on 7-day and 365-day Survival, Logit Estimates

			Conditional on Propensity	Conditional on Propensity
	(1)	(2)	(3)	(4)
Panel A: 7 Day Survival				
Reperfusion	0.172 (0.029)	0.084 (0.031)	0.233 (0.031)	non-parametric
Reperfusion * propensity index	0.348 (0.020)		0.362 (0.020)	non-parametric
Reperfusion * Hospital Treatment Rate (θ)		-0.166 (0.081)	-0.368 (0.084)	-0.511 (0.087)
Control for Propensity Index	None	None	Linear	Non-Parametric
Panel B: 365 Day Survival				
Reperfusion	0.374 (0.021)	0.330 (0.023)	0.403 (0.024)	non-parametric
Reperfusion * propensity index	0.173 (0.016)		0.181 (0.016)	non-parametric
Reperfusion * Hospital Treatment Rate (θ)		-0.068 (0.064)	-0.177 (0.066)	-0.351 (0.068)
Control for Propensity Index	None	None	Linear	Non-Parametric

Note: Table is analogous to Table 2. Coefficients are log-odds. Propensity Equation is $\Pr(\text{Reperfusion})=F(\text{XB}+\text{Hospital Effect})$ and is estimated using a logit model; see Appendix. Propensity index refers to the logit index $(\text{XB}+\text{Hospital Effect})$. It is demeaned to the average value of patients receiving reperfusion. Column 3 reports equation 7c and Column 4 reports equation 7b. All models include all CCP risk-adjusters. Sample-size in every regression is 138,957.

Table 5: Effect of Reperfusion on 30 day Survival, Heirarchical-Logit Estimates

Reperfusion Equation:

Std. Dev. Of Hospital Reperfusion Rate (θ)	0.442 (0.013)
---	------------------

30-day Survival Equation:

Reperfusion	0.265 (0.026)
-------------	------------------

Reperfusion * Propensity Index	0.276 (0.018)
--------------------------------	------------------

<u>Hospital-level intercept (α_0; general productivity)</u> Standard Deviation	0.199 (0.017)
---	------------------

Correlation with Hospital Reperfusion Rate (θ)	-0.100 (0.073)
---	-------------------

<u>Hospital-level Return to reperfusion ($\lambda_1\theta + \tau$)</u> Standard deviation	0.307 (0.056)
---	------------------

Correlation with Hospital Reperfusion Rate (θ)	0.035 (0.112)
---	------------------

Correlation with General Productivity (α^0)	-0.381 (0.151)
--	-------------------

Transformed Estimates:

<u>Hospital minimum treatment threshold (τ)</u> Standard deviation	0.327 (0.055)
---	------------------

Correlation with Hospital Reperfusion Rate (θ)	-0.341 -(0.106)
---	--------------------

Correlation with General Productivity (α^0)	-0.321 (0.150)
--	-------------------

Correlation with Hospital-level Return to reperfusion ($\lambda_1\theta + \tau$)	0.928 (0.026)
--	------------------

Table 5: Effect of Reperfusion on 30--day Survival, Hierarchical-Logit Estimates

Reperfusion Equation:	
Std. Dev. Of Hospital Reperfusion Rate (θ)	0.442 (0.013)
30-day Survival Equation:	
Reperfusion	0.265 (0.026)
Reperfusion * Propensity Index	0.276 (0.018)
<u>Hospital-level intercept (α_0; general productivity)</u>	
Standard Deviation	0.199 (0.017)
Correlation with Hospital Reperfusion Rate (θ)	-0.100 (0.073)
<u>Hospital-level Return to reperfusion ($\lambda_1\theta + \tau$)</u>	
Standard deviation	0.307 (0.056)
Correlation with Hospital Reperfusion Rate (θ)	0.035 (0.112)
Correlation with General Productivity (α^0)	-0.381 (0.151)
Transformed Estimates:	
<u>Hospital minimum treatment threshold (τ)</u>	
Standard deviation	0.327 (0.055)
Correlation with Hospital Reperfusion Rate (θ)	-0.341 (-0.106)
Correlation with General Productivity (α^0)	-0.321 (0.150)
Correlation with Hospital-level Return to reperfusion ($\lambda_1\theta + \tau$)	0.928 (0.026)

Note: Table reports estimates from hierarchical logit, where the propensity to receive treatment is estimated simultaneously with the survival equation. See text and discussion of equation 7c for details.

**Table 6: Variation in use of Reperfusion and Return to Reperfusion,
By type of Hospital**

	(1) Reperfusion		(2) Survival
Reperfusion		Reperfusion	0.205 (0.195)
Reperfusion*Propensity		Reperfusion*Propensity	0.289 (0.0171)
Church Operated Hospital	0.0593 (0.0320)	Reperfusion*Church Operated Hospital	-0.0119 (0.0596)
For-Profit Hospital	0.0957 (0.0384)	Reperfusion*For-Profit Hospital	0.0341 (0.0750)
Government (Federal and Non-Federal)	0.0117 (0.0357)	Reperfusion*Government Hospital	0.0230 (0.0719)
ln (Discharge Volume)	0.132 (0.0208)	Reperfusion*ln (Discharge Volume)	0.0214 (0.0429)
Major Teaching Hospital	-0.115 (0.0447)	Reperfusion*Major Teaching Hospital	0.0352 (0.0816)
Minor Teaching Hospital	-0.0458 (0.0360)	Reperfusion*Minor Teaching Hospital	-0.00622 (0.0675)
Percent of DSH Patients	-0.306 (0.0987)	Reperfusion*Percent of DSH Patients	-0.156 (0.200)
ln (Beds)	-0.0178 (0.0253)	Reperfusion*ln (Beds)	0.00949 (0.0516)
Resident to Bed Ratio	0.0156 (0.153)	Reperfusion*Resident to Bed Ratio	-0.230 (0.296)
Constant	-0.360 (0.0948)	Constant	0.0984 (0.0757)
Propensity	Linear		Linear
Hospital Random Effects	Yes		Yes
Observations	138,957		138,957
Number of Hospitals	4,690		4,690

Both columns report coefficients from mixed-logits that allow for random coefficients. Coefficients are log-odds. In Column (1) Reperfusion is regressed on hospital characteristics. Omitted characteristics is a non-profit hospital. In Column (2) 30-day survival is regressed on hospital-characteristics and hospital characteristics interacted with treatment. Column two reports the interaction effects. A test of joint-significance on these interactions yielded a chi-square statistic of 2.96, p=.097.

Table 7: Minimum Chi-Squared Estimates of Structural Parameters

	Reduced Form Estimates Pooling All Hospitals		Reduced Form Estimates Separately by Hospital Volume	
	Just- identified	Constrain w=r	Different Reliability by Hospital Volume	Same Reliability by Hospital Volume
Std. Dev (α_0)	0.198 (0.0167)	0.204 (0.017)	0.200 (0.016)	0.200 (0.016)
Std. Dev (α)	0.317 (0.058)	0.407 (0.059)	0.337 (0.057)	0.336 (0.057)
Corr(α, α_0)	-0.390 (0.145)	-0.438 (0.130)	-0.457 (0.148)	-0.438 (0.154)
sigma	0.435 (0.152)	0.367 (0.218)	0.431 (0.117)	0.441 (0.135)
w (weight)	0.154 (0.169)	constrained	0.119 (0.107)	0.114 (0.127)
r (reliability)	0.065 (0.106)	0.155 (0.162)		0.040 (0.069)
r (big Hospitals)			0.069 (0.093)	
r(Medium Hospitals)			0.047 (0.063)	
r(Small Hospitals)			0.019 (0.026)	
# moments being fit	6	6	18	18
Degrees of freedom	0	1	10	12
Chi-Squared statistic (p-value)	NA	10.4 (p=.001)	12.7 (p=.24)	62.4 (p<.001)

The first two columns fit 6 reduced-form moments estimated from our empirical model (the variances and covariances of θ , τ , and α_0) as a function of the unknown structural parameters in this framework (the variance and covariance of α and α_0 , the reliability of the signal r , the weight placed on the signal w , and the scale parameter from the logit σ). The reduced-form moments were estimated pooling all hospitals. The unknown structural parameters were estimated using minimum chi-squared methods. The last two columns fit our model to reduced-form moments estimated separately for low (20 or fewer patients), medium (21-80 patients) and high (81 or more patients) volume hospitals - 6 moments for each group, for a total of 18 moments.