

Cream skimming: Theory and evidence from hospital transfers and capacity utilization

by

Ou Yang* Marc K. Chan* Terence C. Cheng** Jongsay Yong*

*University of Melbourne

**University of Adelaide

Revised version: 11 March 2020

Acknowledgement

This research was funded by an NHMRC Partnership Grant (Grant ID: 567217) and in collaboration with the Victorian Department of Health and Human Services (formerly the Victorian Department of Health). The views expressed herein are those of the authors and do not necessarily reflect the views of the Victorian Department of Health.

Abstract

The paper examines cream skimming behaviour by studying hospital transfers in a mixed public-private hospital system. A key innovation is the use of capacity utilization to identify cream skimming. We develop a dynamic model with uncertain patient arrivals and hospital capacity constraints to clarify the conditions under which a profit maximizing hospital will engage in patient selection by transferring ‘hard’ patients—those with severe/complex conditions—to free up capacity to accommodate ‘easy’ patients with few severe/complex conditions. Given finite capacity, public hospitals are strictly less likely to transfer patients than profit-motivated private hospitals at the same level of capacity. We test implications of the model using hospital administrative data from Victoria, Australia, and find empirical support for the cream skimming predictions of the model.

Keywords: Cream skimming; Hospital transfers; Capacity utilization; Mixed private-public hospital system.

JEL Classifications: D25; I11; L33

1 Introduction

The proliferation of prospective payment arrangements, and more recently bundled payment models, has highlighted issues surrounding cream skimming or cherry picking, whereby providers maximize financial gains by selecting patients with the best risks. Incentives for this behavior arise when providers are reimbursed a fixed and predetermined amount for care services through diagnosis-related groups (DRGs) payments. Providers stand to reap financial gains by selectively treating patients with lower than average cost within a DRG (Ellis, 1998). A number of risk adjustment mechanisms have been developed to reduce the incentive to cream skim (Glazer and McGuire, 2000; van de Ven and Ellis, 2000). More recently, under bundled payments, risk adjustment models are proposed to differentiate patients and procedures according to their expected outcomes and financial risks (McLawhorn et al., 2018). However, no risk adjustment model is perfect. Accounting for cream skimming is therefore an important consideration in the design of payment models.

This paper examines cream skimming by hospitals in the context of the mixed public-private hospital system of Australia. We consider two types of hospitals: a private hospital that maximize profits and a public (government) hospital that maximize output or volume.¹ We examine the decision of hospitals in transferring patients and propose a novel approach to identifying cream skimming by linking the likelihood of transfers with a hospital’s capacity utilization. We develop a dynamic model with uncertain patient arrivals and hospital capacity constraints to clarify the conditions under which a profit maximizing hospital will engage in patient selection by transferring ‘hard’ patients—those with more severe or complex conditions—to free up capacity to accommodate ‘easy’ patients who have less severe or complex conditions. Given finite capacity, the hospital is more likely to transfer hard patients as capacity constraints bind.

The model produces a new result: although public hospitals do not engage in selecting

¹We will show later that the model can easily accommodate hospitals with mixed objectives that maximize a combination of output and profits.

patients for financial gains, a forward-looking welfare or output maximizing public hospital will still transfer hard patients when its capacity utilization is high. This result, driven by a dynamic “turnover” incentive by public hospitals to maximize output because of capacity constraints, is new to the literature. However, the model also shows that, given the same level of capacity utilization, the likelihood of public hospitals transferring hard patients is strictly dominated by the likelihood of profit-motivated private hospitals transferring hard patients. That is, public hospitals are strictly less likely to transfer hard patients than profit-motivated private hospitals. This result is key in identifying cream skimming in our empirical work.

Our theoretical model provides a structural interpretation which helps to inform the identification strategy for the empirical analysis. The model presents several empirically testable implications which we examine using hospital administrative data from the state of Victoria, Australia. First, complex patients are more likely to be transferred than non-complex patients as capacity utilization rises, and this is true for patients in both public and private hospitals. Second, private hospitals are more likely to transfer complex patients than public hospitals. We implement the empirical analysis at two levels of aggregation: a hospital level analysis on the number and proportion of transfers of high complexity patients (relative to non-complex patients), and a patient level examination on the probability of transfers versus discharge of high complexity patients relative to non-complex patients. The latter allows for patient characteristics to be taken into account. The empirical findings lend support to the theoretical predictions of cream skimming, and the evidence from the hospital level analysis is especially strong.

This study contributes to the theoretical and empirical literature on cream skimming in health care in several ways. Theoretical studies on payment systems and incentives for cream skimming have its origins on the seminal article by Ellis (1998) which, together with much of subsequent follow-up work (e.g. Barros, 2003), involve the use of a static framework. We formulate a richer and dynamic model enabling us to identify cream skimming behaviour with an important aspect of provider heterogeneity that has

attracted little attention, namely providers' capacity.

Our theoretical work resonates with a recent study on the long-term care (LTC) market in the US by Hackmann and Pohl (2018). Their work, developed independently from ours, adopts a different dynamic modelling approach by examining endogenous discharges of a relatively homogeneous population of nursing home care residents, whereas we examine endogenous transfers of heterogeneous patients of different case complexity. A distinguishing feature of our model is that providers endogenize the effect of current transfers on future capacity utilization. As a result, we show that an output-maximizing public hospital would also transfer complex patients to optimize the use of limited capacity. Two recent studies (Einav et al., 2018; Eliason et al., 2018) also examine provider incentives in the US LTC market using dynamic structural models, although cream skimming is not a feature in these studies, and identification in both studies rest on a payment schedule discontinuity imposed from an exogenous policy design. We add to this literature by examining the provider's incentive to cream skim in an inpatient care context, where no exogenous policy events can be exploited. Our use of capacity utilization to identify cream skimming behaviour adds a novel feature to the literature.

Our study also differs from the preceding three studies on the US LTC market in that they focus on hospital discharges (to home or to another hospital), whereas our focus is on hospital transfers, because cream skimming is more likely to occur in transfers between hospitals than discharges to home. Although hospitals (and doctors) may discharge patients to home earlier than clinically optimal, the extent they could do so is bounded, legally by the prospect of lawsuits and morally by their duties of care. Given Australia's mixed public-private hospital system and the prevalence of dual-practice doctors, private hospitals have often been accused of cream skimming by transferring very sick or complex patients to public hospitals, or turning them away (at the point of admission) while enticing low complexity and profitable patients via transfers from public to private hospitals (O'Loughlin, 2002; Cheng et al., 2015).

Due to patient heterogeneity, empirically identifying cream skimming has been a chal-

lenge. Providers engaging in cream skimming often claim that they act on the best interest of patients. A provider who refuses to admit a critically ill patient, for example, may claim that the patient would be better served elsewhere at another, better equipped or better staffed, facility. Disentangling this from cream skimming motive is often difficult. We propose an approach that links cream skimming to capacity utilization, thereby providing a feasible approach for identifying cream skimming in empirical settings absent exogenous policy events. To our knowledge, this is the first study using capacity utilization and the case complexity of patients to overcome the observational equivalence between cream skimming and non-cream skimming behavior.

This study also contributes to the empirical literature on cream skimming by health care providers, for which the body of evidence is relatively scant. Existing studies show that changes in remuneration systems can create incentives for cream skimming, and that incentives vary with ownership types depending on the extent hospitals can appropriate profit. For example, Duggan (2000) examines a policy change that provided financial incentives for hospitals to treat low income patients in California and finds that both private-for-profit and not-for-profit hospitals cream-skimmed profitable patients, leaving unprofitable ones to public hospitals. Brown et al. (2014) find that while the introduction of risk-adjusted payments led to Medicare Advantage (MA) plans to enroll individuals with high risk scores, these plans selected individuals who have low costs conditional on their risk scores.

Outside of the US, the introduction of DRG-based activity-based funding in Norway, where reimbursement within a DRG is fixed regardless of whether patients receive same-day or overnight treatment, led to hospitals prioritising patients with the shortest length of stay and lowest costs (Martinussen and Hagen, 2009). In Italy, where hospitals are reimbursed through a prospective payment system, Berta et al. (2010) finds that private hospitals cream skim at a much higher intensity compared with public or not-for-profit hospitals. In an empirical study of cream skimming behavior of hospitals in Australia, Cheng et al. (2015) find that patients with severe conditions are more likely to be

transferred from private to public hospitals and have higher length of stay and overall costs. However, as remarked earlier, empirically identifying cream skimming is difficult since it is observationally equivalent to providers acting, or claiming to act, on the best interest of patients.

Our findings have implications for policies encouraging greater private sector participation in health care markets and related market-driven reforms such as those implemented in the UK (Barros and Siciliani, 2011; Cooper et al., 2018). Our results also highlight the importance of accounting for cream skimming in designing payment models, e.g., the use of financial incentives and risk sharing in value-based payment schemes (Eijkenaar, 2013; Scott et al., 2018), and on attempts to increase market transparency by publicizing hospital performance via public performance reporting mechanisms (Dranove et al., 2003; Dranove and Jin, 2010; Chen and Meinecke, 2012; Mak, 2017).

2 The Model

We set up a simple model to illustrate how a hospital can optimize the use of its limited capacity through transferring patients. We consider two types of hospitals: a private hospital maximizing an expected stream of profits, and a public (government) hospital that maximizes societal welfare by maximizing an expected stream of output (or volume). We show later in Section 2.3 that the model can be easily extended to accommodate hospitals with mixed objectives maximizing a combination of output and profits.

The hospital objective functions we assumed are consistent with the conventional view. For example, Iversen (1997) considers a public hospital as one owned by the local government who allocate a fixed budget to the hospital to maximize output, as measured by the number of patients treated. Later studies by Ma (2004) and Grassi and Ma (2012) consider a public provider that maximizes social welfare by providing the volume of services up to a capacity limit determined by the available budget. Private hospitals are traditionally modelled as profit maximizers (Sloan, 2000). Models of mixed

objectives have also been used to characterize the behaviour of private not-for-profit hospitals with a broader set of societal or altruistic objectives as a point of differentiation from for-profit hospitals focusing solely on revenue- or profit-maximization. In considering not-for-profit firms, Lakdawalla and Philipson (2006), for example, specify a model where the utility of a firm's owner depends on the firm's output, input, and the firm owner's own good consumption. The firm owner may choose to tradeoff profits for higher or better quality output. In Gaynor and Vogt (2003), not-for-profit hospitals are assumed to maximize both quantity of output and profits.

While we follow the conventional view in our choice of hospital objective functions, we note the lack of agreement in the literature on what public and private hospitals maximize. Some studies assume public hospitals maximize profits, much like private hospitals do. Others, in contrast, assume that all health care providers, both public and private, are motivated by both profits and altruism, and maximize a weighted sum of patient benefits, output and profits; for a review, see Barros and Siciliani (2011). Our results remain qualitatively unchanged as long as private hospitals are assumed to place greater weights on profit than output in comparison to public hospitals.

2.1 Private Hospitals

We assume that private hospitals maximize future streams of profits. In each period, patients arrive at the hospital and are admitted for treatment. The hospital discharges patients at a rate of θ each period. When a patient is discharged, the hospital receives a profit π . We assume that there are two types of patients, 'easy' patients (denoted by subscript e) and 'hard' patients (denoted by subscript h). We assume that $\theta^e > \theta^h$ and $\pi^e \geq \pi^h$, that is, easy patients are subject to a higher discharge rate each period and entail a (weakly) higher profit than hard patients. While it is more reasonable to assume $\pi^e > \pi^h$, we will discuss the special case of $\pi^e = \pi^h$ later. The assumption on discharge rate implies that easy patients have lower average length of stay than hard patients. Thus the exogeneity in discharge rates reflects the technological differences in

treatment between easy and hard patients, under the premise that the implicit cost of manipulating discharges is prohibitive relative to transfers.²

Suppose the hospital faces an infinite horizon. Let z_t^e and z_t^h denote the number of respectively easy and hard patients arriving at the hospital in period t . Both z_t^e and z_t^h are random and unknown prior to period t . Let x_t^e and x_t^h denote the *stock* of respectively easy and hard patients in the hospital at the beginning of period t . For simplicity, we abstract from capacity and transfer considerations for now and will introduce them later. Numerous further extensions are possible, and we consider several in Section 2.3, including hospitals with mixed profit-output objectives, strategic interactions between hospitals, destinations of transferred patients and the market equilibrium, new patients reacting to transfer rates, heterogeneous discharge rates, and accommodating resource-related incentives.

Without capacity constraints, all patients arriving at the hospital are admitted and treated. The stock of patients evolves according to the following laws of motion:

$$\begin{aligned}x_{t+1}^e &= (1 - \theta^e)x_t^e + z_{t+1}^e, \\x_{t+1}^h &= (1 - \theta^h)x_t^h + z_{t+1}^h.\end{aligned}$$

The stock of patients next period is the sum of undischarged patients this period and the new patients who arrive next period. The profit function of the hospital in period t is:

$$\Pi(x_t^e, x_t^h) = \pi^e \theta^e x_t^e + \pi^h \theta^h x_t^h. \quad (1)$$

Denote the mean number of arrivals of easy and hard patients by $E(z_t^e) = \bar{z}^e$ and $E(z_t^h) = \bar{z}^h$. Using the laws of motion, it is straightforward to show that the steady-state stock of easy and hard patients are respectively $\bar{x}^e = \frac{\bar{z}^e}{\theta^e}$ and $\bar{x}^h = \frac{\bar{z}^h}{\theta^h}$. Therefore, the steady-state profit is $\pi^e \theta^e \bar{x}^e + \pi^h \theta^h \bar{x}^h = \pi^e \bar{z}^e + \pi^h \bar{z}^h$.

²For example, the legal consequences may be disastrous if the hospital discharges a patient prematurely and results in death.

We now introduce capacity and endogenous transfers in the hospital's dynamic optimization problem. Suppose the hospital has a finite capacity limit K per period in serving the total stock of patients, x_t :³

$$x_t = x_t^e + x_t^h \leq K.$$

In each period t , the hospital decides the optimal number of hard patients to be transferred, n_t^h , given the existing stock of easy and hard patients, x_t^e and x_t^h . Transferring patients entails an immediate cost $C(n_t^h)$, a twice-differentiable, increasing and convex function in n_t^h . These assumptions are made partly for analytical tractability. However, increasing and convex costs may reflect the fact that it can become disproportionately more difficult and costly to transfer patients due to legal, reputation or regulation costs, or to successfully transfer a large number of patients possibly due to implicit search costs (we return to this issue later). Transferring hard patients will free up the capacity and reduce the stock of hard patients next period; in certain cases, it will increase the stock of easy patients next period (discussed below). The hospital chooses n_t^h to maximize the sum of the current net profit and the expected discounted stream of net profits in the future.

The laws of motion are more complicated in the presence of the capacity constraint. Let $\tilde{x}_t = (1 - \theta^e)x_t^e + (1 - \theta^h)x_t^h$ denote the total stock of undischarged patients at the beginning of period t prior to the transfer of hard patients in period t and the arrival of new patients in period $t + 1$. The hospital cannot admit all new patients when the total number of new patient arrivals, $z_{t+1} := z_{t+1}^e + z_{t+1}^h$, exceed the remaining capacity, that is, $z_{t+1} > K - (\tilde{x}_t - n_t^h)$. In this case, we assume that the hospital rations the admissions via randomization within the pool of all new arrivals (e.g., patient type is not revealed

³A more general form of the constraint is $\frac{x_t^e}{\rho} + x_t^h \leq K$, where $\rho \geq 1$, e.g., a hard patient may use more resources than an easy patient. When $\rho > 1$, our theoretical results still holds because the contrast between easy and hard patients becomes even starker.

prior to admission) to fill all the remaining capacity.⁴ The admission rule is:

$$\tilde{z}_{t+1}^e = \begin{cases} \frac{z_{t+1}^e}{z_t} (K - \tilde{x}_t + n_t^h) & \text{if } z_{t+1} > K - \tilde{x}_t + n_t^h, \\ z_{t+1}^e & \text{otherwise.} \end{cases} \quad (2)$$

$$\tilde{z}_{t+1}^h = \begin{cases} \frac{z_{t+1}^h}{z_t} (K - \tilde{x}_t + n_t^h) & \text{if } z_{t+1} > K - \tilde{x}_t + n_t^h, \\ z_{t+1}^h & \text{otherwise.} \end{cases} \quad (3)$$

The laws of motion are:

$$x_{t+1}^e = (1 - \theta^e)x_t^e + \tilde{z}_{t+1}^e, \quad (4)$$

$$x_{t+1}^h = (1 - \theta^h)x_t^h + \tilde{z}_{t+1}^h - n_t^h, \quad (5)$$

The hospital's (net) profit function is:

$$\Pi(x_t^e, x_t^h) = \pi^e \theta^e x_t^e + \pi^h \theta^h x_t^h - C(n_t^h). \quad (6)$$

The Bellman equation for the profit maximization problem is:

$$V(x_t^e, x_t^h) = \max_{x_t^h \geq n_t^h \geq 0} \left(\pi^e \theta^e x_t^e + \pi^h \theta^h x_t^h - C(n_t^h) + \beta \int \int_{z_{t+1}^e, z_{t+1}^h} V(x_{t+1}^e, x_{t+1}^h) dF(z_{t+1}^e, z_{t+1}^h) \right), \quad (7)$$

subject to the laws of motion (4)–(5). The value function $V(\cdot)$ depends on state variables known in period t (x_t^e and x_t^h). Given the current stocks of patients, the hospital chooses the optimal number of transfer of hard patients, n_t^{h*} , such that it maximizes the sum of current net profit and the discounted expected value in the future, with discount factor β . The solution is obtained by the following first order condition:

$$-\frac{dC}{dn_t^h} + \beta \frac{dEV}{dn_t^h} = 0, \quad (8)$$

where $EV := \int \int_{z_{t+1}^e, z_{t+1}^h} V(x_{t+1}^e, x_{t+1}^h) dF(z_{t+1}^e, z_{t+1}^h)$ is the expected value function and strict equality holds when $x_t^h > n_t^{h*} > 0$. The marginal expected future value can be rewritten as:

$$\frac{dEV}{dn_t^h} = \underbrace{E\left(\frac{\partial V(x_{t+1}^e, x_{t+1}^h)}{\partial x_{t+1}^h} \frac{\partial x_{t+1}^h}{\partial n_t^h}\right)}_{\leq 0} + \underbrace{E\left(\frac{V(x_{t+1}^e, x_{t+1}^h)}{\partial x_{t+1}^e} \frac{\partial x_{t+1}^e}{\partial n_t^h}\right)}_{\geq 0}. \quad (9)$$

⁴Under the law of large numbers, $\frac{z_t^e}{z_t}$ of the admitted patients will belong to the easy type and $\frac{z_t^h}{z_t}$ will belong to the hard type. We assume that this holds true for any z_t .

To simplify the discussion, suppose $\frac{z_t^e}{z_t} = p_e$ and $\frac{z_t^h}{z_t} = p_h$, where $0 < p_e, p_h < 1$ are constants such that $p_e + p_h = 1$. This simplifies the computation of EV as the only shock that needs to be integrated out is z . The model can be generalized to incorporate variations in the risk mix, i.e., proportion of hard versus easy patient arrivals. In this case there will be two shocks to integrate out in EV , i.e., z and the risk mix. The result will be qualitatively similar.

Assumptions. For clarity, we list the assumptions as follows:

- (A1) $\theta^e > \theta^h$: easy patients are subject to a higher discharge rate than hard patients.
- (A2) $\pi^e \geq \pi^h$: easy patients entail a higher profit than hard patients.
- (A3) The cost function $C(n_t^h)$ is a twice-differentiable, increasing and convex function in n_t^h . In addition, $C(0) = 0$.
- (A4) The total number of patient arrivals, z , follows a twice-differentiable cumulative distribution function $F(z)$, which has continuous support over $\mathbb{R}_{\geq 0}$.
- (A5) The number of easy and hard patient arrivals are constant proportions of z : $\frac{z_t^e}{z_t} = p_e$ and $\frac{z_t^h}{z_t} = p_h$, where $0 < p_e, p_h < 1$ are constants such that $p_e + p_h = 1$.

We have the following convenient result (All proofs are collected in Appendix A):

Lemma 1 *Suppose assumption A5 holds. If $z_{t+1} > K - \tilde{x}_t + n_t^h$, then $\frac{\partial x_{t+1}^h}{\partial n_t^h} = -p_e$ and $\frac{\partial x_{t+1}^e}{\partial n_t^h} = p_e$; otherwise, $\frac{\partial x_{t+1}^h}{\partial n_t^h} = -1$ and $\frac{\partial x_{t+1}^e}{\partial n_t^h} = 0$.*

Therefore, transferring hard patients today increases the stock of easy patients next period only when too many patients arrive relative to the remaining capacity (i.e., the hospital is at full capacity next period). Note that, due to the admission rule, the result holds for all n_t^h , not just at the optimal level.

Define $A(n_t^h; x_t^e, x_t^h) := \{z_{t+1} | z_{t+1} > K - \tilde{x}_t + n_t^h\}$, that is, $A(n_t^h; x_t^e, x_t^h) = (K - \tilde{x}_t + n_t^h, \infty)$, an open interval, represents the case in which the hospital is at full capacity in period

$t + 1$. By Lemma 1, we can rearrange the right hand side of (9) as follows:

$$\begin{aligned} \frac{dEV}{dn_t^h} &= \int_{z_{t+1} \in A(n_t^h)} \left(\frac{\partial V}{\partial x_{t+1}^h}(-p_e) + \frac{\partial V}{\partial x_{t+1}^e} p_e \right) dF(z_{t+1}) + \\ &\quad \int_{z_{t+1} \notin A(n_t^h)} \left(\frac{\partial V}{\partial x_{t+1}^h}(-1) + \frac{\partial V}{\partial x_{t+1}^e}(0) \right) dF(z_{t+1}) \\ &= p^e \int_{z_{t+1} \in A(n_t^h)} \left(\frac{\partial V}{\partial x_{t+1}^e} - \frac{\partial V}{\partial x_{t+1}^h} \right) dF(z_{t+1}) - \int_{z_{t+1} \notin A(n_t^h)} \frac{\partial V}{\partial x_{t+1}^h} dF(z_{t+1}), \quad (10) \end{aligned}$$

Putting changes in the limit of integration aside, the magnitude of $\frac{dEV}{dn_t^h}$ depends on the expected gain when the hospital is at full capacity next period (the first term on r.h.s.) versus the expected loss otherwise (the second term on r.h.s.).⁵ However, importantly, we have $\frac{\partial \inf A}{\partial n_t^h} \leq 0$ by the definition of $A(n_t^h; x_t^e, x_t^h)$, therefore transferring more patients today will reduce the probability that the hospital is at full capacity next period.⁶ Therefore, $\frac{dEV}{dn_t^h}$ decreases as n_t^h increases, or $\frac{d^2 EV}{d(n_t^h)^2} < 0$.

Proposition 1 *Suppose Assumptions A1 to A5 hold. Then there exists a unique solution n_t^{h*} to the profit maximization problem in (7).*

Remark 1 *Although the model does not feature transfers of easy patients, the hospital will never endogenously transfer any easy patients even if this feature exists. According to (10), the sole benefit of transferring hard patients now is that it may free up future capacity which increases the proportion of easy patients. This benefit has to be large enough to override the forgone profit that could accrue from the transferred hard patients. Using a similar logic, transferring easy patients now may free up future capacity which increases the proportion of hard patients; this has no benefit at all because hard patients entail lower profits than easy patients, all else being equal. Thus the hospital will never make endogenous transfers of easy patients.*

⁵Due to the linear profit function in (6), it follows that $V(\cdot)$ is a linear function of x^e and x^h , and $\frac{\partial V}{\partial x^e} > \frac{\partial V}{\partial x^h}$.

⁶Note again that according to the admission rule, the decision to transfer is made before arrivals are known, so that whether the capacity constraint is binding or not is stochastic.

Remark 2 *The model can be easily extended to incorporate exogenous transfers of easy (and hard) patients, which reflect the technology of the hospital. To illustrate, let $s^e \geq 0$ and $s^h \geq 0$ be the exogenous transfer rates of easy and hard patients, respectively. Then, the stock of patients will evolve according to $x_{t+1}^e = (1 - \theta^e - s^e)x_t^e + z_{t+1}^e$ and $x_{t+1}^h = (1 - \theta^h - s^h)x_t^h + z_{t+1}^h$, while the profit function will remain unchanged. The number of transfers of easy patients is $s^e x_t^e$. The number of transfers of hard patients is $s^h x_t^h$ plus an endogenously determined amount given by the solution of the reformulated dynamic optimization problem.*

We next derive the comparative statics. Given a fixed capacity limit K , as the stock of patients (x_t) increases, available capacity becomes smaller, i.e., capacity utilization becomes higher. We will use the term capacity utilization and the stock of patients synonymously. The following result shows that, in the interior solution, the optimal transfer of hard patients, n_t^{h*} , is strictly increasing as capacity utilization rises. Note that the optimal transfer function may be concave or convex in capacity utilization. If the cost function $C(n_t^h)$ is highly convex in n_t^h , it is reasonable to expect this function to be concave.

Proposition 2 *Suppose Assumptions A1 to A5 hold. Then $\frac{dn_t^{h*}}{dx_t}, \frac{dn_t^{h*}}{dx_t^h}, \frac{dn_t^{h*}}{dx_t^e} \geq 0$ and, in addition, $\frac{dn_t^{h*}}{dx_t}, \frac{dn_t^{h*}}{dx_t^h}, \frac{dn_t^{h*}}{dx_t^e} > 0$ where $\frac{dn_t^{h*}}{dx_t^h} > \frac{dn_t^{h*}}{dx_t^e}$ when n_t^{h*} is an interior solution, i.e., $x_t^h > n_t^{h*} > 0$.*

2.2 Public Hospitals

Following Iversen (1997), we assume that public hospitals care about output (or volume), as measured by the number of discharged patients, i.e., $\theta^e x_t^e + \theta^h x_t^h$. A forward-looking public hospital maximizes an expected stream of output. Note that we have implicitly assumed that the output from discharged hard and easy patients are comparable, e.g., we consider hard and easy patients within the same disease category,⁷ rather than hard

⁷In the empirical analysis we control for disease category by using dummies of diagnostic categories.

patients from, say, cancer treatment versus easy patients from flu shots.⁸ The output function in period t is:

$$\Pi^b(x_t^e, x_t^h) = \pi^b(\theta^e x_t^e + \theta^h x_t^h), \quad (11)$$

where $\pi^b = 1$ without loss of generality. Note that this is a special case of the profit function in equation (1); π^b can be interpreted as the amount of social welfare generated by each discharged patient. The following corollary summarizes the behavior of private and public hospitals:

Corollary 1 *Suppose Assumptions A1 to A5 hold. Consider a private hospital (v) with profit function $\Pi^v(x_t^e, x_t^h) = \pi^e \theta^e x_t^e + \pi^h \theta^h x_t^h$ where $\pi^e > \pi^h$, and a public hospital (b) with output function $\Pi^b(x_t^e, x_t^h) = \theta^e x_t^e + \theta^h x_t^h$. Suppose both hospitals are otherwise identical. Given the current stocks of patients (x_t), let $n_t^{(v)h*}$ and $n_t^{(b)h*}$ denote the optimal transfer of hard patients by the private hospital and public hospital, respectively. We have the following results:*

C1.1 The private hospital transfers more hard patients than the public hospital, all else being equal: $n_t^{(v)h} \geq n_t^{(b)h*}$ and, in particular, $n_t^{(v)h*} > n_t^{(b)h*}$ when $n_t^{(b)h*}$ is an interior solution, i.e., $x_t^h > n_t^{(b)h*} > 0$.*

C1.2 Both the private and public hospitals increase transfers of hard patients when capacity utilization increases: Proposition 2 holds for both hospitals.

Remark 3 *The forward-looking public hospital anticipates that keeping a hard patient will crowd out more easy patients next period because hard patients stay longer than easy patients. To maximize the expected stream of output, it will transfer more hard patients as the capacity utilization increases (the “turnover” incentive). However, all else being equal, the public hospital will transfer fewer hard patients than the private hospital due to the absence of a profit motive. In this sense, we can distinguish cream skimming from*

⁸It is also likely that cancer patients have a higher value of treatment that is offset by higher resource use; for details, see Section 2.3 on the extension of the model to consider resource usage.

non-cream skimming behavior. Recall that cream skimming refers to hospitals selecting patients and the selection is driven by profit motives. Even though our results suggest that both the public and private hospitals select patients by transferring hard patients, the public hospital's behavior is non-cream skimming due to the absence of profit motives.

2.3 Possible Extensions

We discuss below some possible extensions to our model.

Hospitals with Mixed Objectives

Our model can readily accommodate hospitals that maximize a combination of output and profits. It is not unreasonable to suggest that most private hospitals can be so characterized, with some placing more weights on profits than others. The weighted output-profit function in period t is:

$$\Pi^m(x_t^e, x_t^h) = w(\pi^e \theta^e x_t^e + \pi^h \theta^h x_t^h) + (1 - w)(\theta^e x_t^e + \theta^h x_t^h) \quad (12)$$

$$= \underbrace{[w\pi^e + (1 - w)]}_{\tilde{\pi}^e(w, \pi^e)} \theta^e x^e + \underbrace{[w\pi^h + (1 - w)]}_{\tilde{\pi}^h(w, \pi^h)} \theta^h x^h, \quad (13)$$

where $0 < w < 1$ is the weight assigned to the profit motive and $1 - w$ is the weight assigned to the output motive. This is analogous to the profit function in (1), where the relative profits π^e and π^h are replaced by $\tilde{\pi}^e$ and $\tilde{\pi}^h$, respectively.

We now consider optimal transfers by this hospital. Given $\pi^e > \pi^h \geq 1$ (upon normalization), by simple arithmetic, we have $\tilde{\pi}^e > \tilde{\pi}^h \geq 1$ and $\frac{\pi^e}{\pi^h} > \frac{\tilde{\pi}^e}{\tilde{\pi}^h} > 1$ for any $0 < w < 1$. Hence we can view a hospital with mixed objectives as a private hospital with “muted” relative profits between easy and hard patients. In addition, it is straightforward to show that $\frac{\tilde{\pi}^e}{\tilde{\pi}^h}$ is an increasing function of w . Therefore, using a similar proof to Corollary 1, we establish the following result:

Corollary 2 *All else being equal, a mixed-objective hospital transfers fewer (more) hard patients than a private (public) hospital. Moreover, it increases transfers of hard patients*

when (i) capacity utilization increases, or (ii) the weight assigned to the profit motive increases.

In the empirical analysis, anonymized hospital identities in the administrative data prevent us from distinguishing between for-profit from not-for-profit private hospitals. Hence we examine these two types of private hospitals as a single group. Our estimate of cream-skimming behaviour will be a weighted average of the true effects of for-profit, and not-for-profit, private hospitals.

Strategic interactions

The model can be extended to consider intra- and inter-hospital strategic interactions. However, numerous possibilities exist and a detailed treatment is beyond the scope of this paper. As an illustration, we consider intra-hospital interactions in the context of our empirical application. A distinguishing feature of our empirical application is public hospitals in Australia admit both public and private patients (see Section 3 below). We consider this possibility by letting there be two managers or decision makers, one is in charge of private patients and the other is in charge of public patients. The former is profit maximizing while the latter welfare maximizing. They share a common capacity K and as such, entails spillover effects. Strategic interactions exist between the two managers. Relative to the case of complete internalization, we expect that each manager will have a smaller incentive to transfer his/her own hard patients, because the future ‘cost’ of reaching full capacity is partially borne by the other manager. However, we still expect that each manager will transfer more hard patients as capacity utilization rises. A sketch of the model and its Markov perfect equilibrium solution can be found in Appendix S ⁹.

Using a similar technique, we can also consider inter-hospital strategic interactions. We argue that the main results from the baseline model still hold. The intuition is as follows. Consider two hospitals, one private and the other public. Transfers from one hospital goes to the other hospital. If both hospitals treat patient arrival as exogenous

⁹Appendix S is included as supplementary material.

(e.g., transfers constitute a tiny proportion of arrivals), then the solution will be the same as in the baseline model. If both hospitals treat patient arrival as endogenous instead, then the best response transfer of each hospital will be an increasing function of the other hospital's transfers; given a higher transfer from the other hospital, the hospital closer to the capacity limit will transfer more patients as the best response. In addition, the private hospital's optimal transfers will be more responsive to the public hospital's transfers than vice versa. In equilibrium, both hospitals will transfer more patients relative to the exogenous transfer case. However, transfers will still increase when capacity utilization increases, and the private hospital will still transfer more patients than the public hospital (given the same capacity utilization).

Hospitals accepting transferred patients

The baseline model is one of partial equilibrium analysis in that the arrival of easy and hard patients is taken as exogenous by the hospital. This assumption is reasonable if, as in our empirical application, transferred patients constitute a small proportion of patient arrivals, so that transfers from other hospitals have minimal impact on the hospital's optimization problem. Our model also assumes that whenever a hospital decides to transfer patients, they will be admitted elsewhere. At the market level, this implicitly assumes that the system is never saturated, i.e., at any given point in time vacancies exist in some hospitals. For hospitals experiencing low capacity utilization, it can be optimal to accept transferred patients. However, we do posit that the cost of transferring patients is increasing and convex in n_t^h , suggesting that the search and placement costs will be rising disproportionately as more patients are transferred. A market-wide implication is that some patients may be transferred multiple times, but will eventually be admitted into a hospital.

New patients' reactions to transfer rates

The baseline model assumes that the flow of new patients is exogenous. However, if new patients have knowledge about the transfer rate and/or capacity utilization of a hospital and they dislike the risk of being transferred, they may rationally adjust their

behaviour by avoiding that hospital. Theoretically, from the hospital’s point of view, this adjustment may create two additional effects when more patients are transferred now. First, the transfer rate is lifted and new patients will avoid the hospital, which relieves the capacity pressure in the future. Second, the hospital’s profit/output is hurt in the future due to reduced inflows. Both will reduce the hospital’s incentive to transfer, but we still expect transfers to increase when capacity utilization increases. Formally, patients’ behaviour are endogenous to the hospital’s action too, which calls for a model with hospital(s) and patients as players engaging in strategic interactions. Solving for an equilibrium will be complicated. In practice, there are institutional and informational constraints, e.g., patients are unlikely to have information about transfer rates or capacity utilization of hospitals.

Heterogeneous discharge rates

The base model assumes a uniform discharge rate for each patient type, regardless of the patient’s length of stay since admission. While this assumption avoids the curse of dimensionality in the dynamic programming model, it is unlikely to be correct in practice, e.g., the hospital may limit the maximum lengths of stay at $\max s^e$ and $\max s^h$ for easy and hard patients. The distributions of hard and easy patients’ length of stay will have to be taken into account in solving the dynamic programming program, which is analytically intractable. Nevertheless, if we assume no “crossover” in maximum stay, e.g., $\max s^e \leq \max s^h$, then our base model is a reasonable approximation to this extended model and the results will be qualitatively similar.

Resource-related incentives

The condition of hard patients may be more likely to worsen during hospitalization, taking up more hospital resources along the way. The hospital may transfer these patients because it may not have sufficient resources to accommodate their needs at full capacity. The base model partially reflects this differential resource use, as hard patients have lower discharge rates (and hence longer stays), implying that they take up more hospital resources in the steady state. Indeed, this generates the dynamic “turnover” incentive

to transfer hard patients in an output-maximizing public hospital. To extend the model, the technology can be modified to reflect deterioration in the quality of care: (1) the discharge rate of hard patients (θ^h) is a decreasing function of capacity utilization, (2) hard patients die at a rate of $\theta^{h,d}$, which is increasing in capacity utilization, and (3) there is an extra cost function $C^d(\theta^{h,d}x_t^h)$, increasing and convex in the number of hard patients who die. We expect the results to remain qualitative similar, as they magnify the “turnover” incentive (e.g., patients who die cannot be part of the output).

2.4 Empirical Implications

The following are the key empirical implications of the model. First, a forward-looking hospital, whether it be motivated by profits or social welfare, will tend to transfer patients with high complexity (i.e., hard patients) as capacity utilization increases. Thus, the tendency to transfer complex patients as capacity utilization rises is necessary but insufficient evidence for establishing cream skinning. Second, all else being equal, public hospitals have weaker incentives to transfer complex patients than private hospitals, since for public hospitals the turnover incentive still exists but the profit incentive is absent. Thus, private hospitals engaging in cream skinning will tend to transfer more complex patients than public hospitals. We translate these implications into two testable hypotheses:

H1 Compared to non-complex patients, patients with higher complexity are more likely to get transferred as capacity utilization rises.

H2 Private hospitals are more likely to transfer high complexity patients (relative to non-complex patients) than public hospitals (see C1.1).

In the empirical analysis below we will test hypothesis *H1* at four discrete levels of capacity utilization. We say that hypothesis *H1* is rejected by the data if high complexity patients are no more likely to get transferred than non-complex patients at successively higher levels of capacity utilization. For hypothesis *H2*, it is rejected if the rates of

transfers of high complexity patients (relative to non-complex patients) from private hospitals are no higher than those from public hospitals. We will formulate these tests in precise terms after we describe the empirical model below.

3 Empirical Analysis

The Australian hospital system is a mixed system involving the coexistence of public and private hospitals. Public hospitals are wholly government-owned, managed, and funded. Private hospitals comprise both for-profit and not-for-profit institutions, and are remunerated under prospectively determined per-episode payments similar to the DRG payment system used in public hospitals. The payments depend on patient casemix complexity and are negotiated directly between providers and insurers. The data do not allow us to distinguish not-for-profit from for-profit private hospitals, since hospital identities are anonymized and only public and private hospital types can be identified. In Australia, not-for-profit private hospitals are financed in a similar fashion as for-profit private hospitals, with both relying on funding sources from private insurance funds, patient co-payments and government subsidies in the form of Medicare rebates. As mentioned in Section 2.3, the cream-skimming behaviour we identify will be a combination of the true effects of for-profit and not-for-profit private hospitals.

For public hospitals, a large proportion (more than 80%) of their revenue is generated through the DRG payment system. However, unlike private hospitals where financial viability is paramount, public hospitals face a soft-budget constraint and the government may for political reasons intervene to provide additional funding if public hospitals exhaust their budgets.

The empirical analysis is conducted at two levels of aggregation: hospital level and individual patient level. For the hospital level analysis, admissions are aggregated up to hospital level and on a weekly basis. We examine the aggregate transfer of complex patients relative to non-complex patients each week and relate transfers to capacity

utilization during that week and hospital type (public or private hospitals). For the disaggregate patient level analysis, we examine the probability of transfers of individual patients given their characteristics, including their case complexity and patient type (public or private patients), and the characteristics of hospitals, including their capacity utilization levels. The patient level analysis allows for individual characteristics such as gender and age to be included as controls. A further advantage is the patient level analysis can account for the presence of private patients in public hospitals.

A distinguishing feature of the Australian health care system that has a bearing on our empirical analysis is the presence of private patients in public hospitals. Patients with private hospital insurance can choose to be admitted into public hospitals as private patients. Doing so reduces the waiting time for surgery, and also allows patients to choose their doctors and enjoy better amenities (e.g., a private room). Public hospitals do admit private patients and they derive higher revenue from private patients. In recent years private patients account for about 15–20% of all admissions in public hospitals.

For the hospital level analysis, given that hospital type is the focus, we drop all private patients in public hospitals when aggregating the admission data of public hospitals. However, private patients were accounted for in measuring capacity utilization. For the patient level analysis, the focus is on patient type, hence private patients in public hospitals are included in the analysis. At the patient level, we also carry out a subgroup analysis focusing on public hospitals that admit both public and private patients (see Section 3.4 below).

For the empirical analysis, we derive a working definition of patient complexity using the Charlson comorbidity score (Charlson et al., 1994; Sundararajan et al., 2004). The Charlson score is a good predictor of mortality and correlates well with resource use. It is used as a summary measure of patient health conditions that are readily observable to hospitals. The Charlson score is based on a scoring system that ranges from zero to a theoretical maximum of 30. The typical distribution is heavily right skewed with few patients scored above 6. We define three categories of patients based on the Charlson

comorbidity index: c_0 , c_1 and c_2 , corresponding to patients with no, low and high complexity. The majority of patients have a Charlson score of zero, i.e., no comorbidity, and are grouped under c_0 . Patients with a Charlson score of 1 or 2 are low complexity and classified as c_1 . The remaining are high complexity patients with a Charlson score of 3 or higher, and are classified as c_2 . The focus of the analysis will be on patients with high complexity, c_2 , *vis-a-vis* those with no complexity, c_0 .

Another key variable for the empirical analysis is capacity utilization. The theory model predicts different transfer patterns as capacity utilization rises, but is silent on the exact level of utilization at which changes are expected to occur. We construct a capacity utilization measure and examine its empirical distribution; for tractability we define four discrete categories of utilization levels which we describe below.

3.1 Data and descriptive statistics

The data are extracted from the Victorian Admitted Episodes Dataset (VAED), an administrative database containing all hospital admission episodes, public and private, in the state of Victoria, Australia. The database contains admission episode-level data which include information on patient demographics such as age and gender, clinical details such as relevant diagnoses and comorbidity, administrative details such as the date of admission and discharge, and whether the patient was admitted as a private or public patient. The data used for this study cover a twelve-year period from 2000/01 to 2011/12 and contain all admission episodes occurring in all hospitals in Victoria during the period. In total there were 24.9 million admission episodes recorded in some 300 hospitals. Private patients account for about 40% of the admission episodes.

For the purpose of our analysis, the following admission episodes are excluded: 892,477 (3.59%) episodes of which the patient type is neither public nor private (e.g., Department of Veterans' Affairs health card holders and Transport Accident Commission claimants), 210,652 (0.85%) episodes during which the patient died (hence with no possibility of

transfers), 1,932,561 (7.77%) episodes of which the patients are younger than 18 years or older than 100 years, 3,701,771 (14.88%) episodes involving chemotherapy, dialysis, or radiotherapy¹⁰, and 4,523,335 (18.19%) episodes with missing fields required for the analysis. Lastly, for the convergence of the logistic estimation, we also exclude hospitals with fewer than 10 transfers during the entire 12-year period, and this results in the exclusion of 1,461,923 (5.88%) episodes.

After removing all excluded episodes, we end up with a sample of 12.1 million admission episodes. The admission episodes were classified into 425 three-digit DRGs (out of a total of 452 possible three-digit DRGs), and occurring in 180 hospitals, of which 64 are private hospitals.

Hospital transfers are identified by the separation and admission dates in the data. A patient is said to be transferred if she is separated from a hospital and admitted into another hospital on the same day. Hospital capacity is not observed in our data. In previous studies capacity is often measured using physical infrastructure such as floor space and number of beds. However, these measures are not meaningful for our purpose, since we intend to measure capacity use at the time a patient is transferred. The measure must capture any short-term, e.g., weekly, fluctuations in the hospital's ability to deliver services. In reality these short-term fluctuations will be closely related to factors such as work shift arrangements and the number of doctors, nurses and allied health workers available at the time. Since we have no information on weekly fluctuations in workforce, we approximate hospital capacity using the maximum number of patients that a hospital can accommodate in a given week. Note that our capacity measure is computed using all data before applying any sample exclusion criteria.

In measuring a hospital's capacity, we think of the hospital as consisting of multiple specialty departments, which we approximate using major diagnostic categories (MDC).¹¹

¹⁰Chemotherapy, dialysis and radiotherapy patients are excluded primarily because of inconsistent coding across hospitals—in some cases these were coded as in-patient whereas in other cases as out-patient episodes. As a result, the occurrence of transfers cannot be reliably identified. In any case, most of these are same-day procedure and transfers rarely, if ever, occur.

¹¹A MDC corresponds to a single body system or aetiology, broadly reflecting the specialty providing

However, we posit that decisions about transfers are made at the hospital level, not departmental level.¹² Specifically, we assume that hospital management to some extent can reallocate resources across departments when some department is facing capacity pressure; e.g., nursing staff can be reassigned, and empty ward space in adjacent departments can be redeployed. Therefore it is plausible that hospitals make transfer decisions after making an overall assessment of available capacity across all departments.

We construct an overall capacity utilization measure built from utilization rates of MDCs in a hospital. Let K_{jmt}^* denote the capacity (unobserved) of hospital j for MDC m during time t . We assume that K_{jmt}^* is constant during some time period t_0 to t_1 , e.g., within a quarter. For brevity we omit the hospital subscript j below.

Let M be the set of all MDCs. For a given week t and MDC $m \in M$, we compute the weekly utilization, denoted by U_{mt} , which is the stock of patients in MDC m staying in the hospital during week t . For each MDC m , we calculate U_{mt} as the sum of the number of patients admitted under the MDC prior to the week and remain in the hospital (i.e., not transferred, discharged or died) at the end of the week, and patients admitted under the MDC during the week whether separated during the week or not. We compute the maximum weekly utilization for MDC m during the period as $U_{mt}^{\max} = \max\{U_{m,t_0}, \dots, U_{m,t_1}\}$. We can think of this quantity as an approximate capacity for the department specializing in MDC m .

For each MDC m , the observed capacity is assumed to be some fraction of the true capacity during the period, i.e., $K_{mt} = \delta K_{mt}^*$, where $\delta \in (0, 1]$. For convenience, we express capacity utilization as a rate: $V_{mt} = U_{mt}/K_{mt}^* = \delta U_{mt}/K_{mt}$. An overall hospital level capacity utilization measure can be defined as a weighted average of capacity utilization rates across all MDCs: $\bar{V}_t = \sum_{m \in M} \omega_m V_{mt} = \delta \sum_{m \in M} \omega_m U_{mt}/K_{mt}$, where the weight, ω_m , represents the relative size/importance of a MDC in the hospital, is defined

care. The Australian refined diagnosis related groups (AR-DRG) classify all diagnoses into one of 23 MDCs, which we regard as an approximation of specialty departments in hospitals.

¹²Our treatment of departmental capacity is similar to Chan and Zeng (2018), although there it is unclear whether optimizing fee decisions were made at the departmental or hospital level.

as the share of a MDC to the hospital's total number of separations during the period (quarterly in our case), i.e., $\omega_m = \sum_{t=t_0}^{t_1} S_{mt} / \sum_{m \in M} \sum_{t=t_0}^{t_1} S_{mt}$, where S_{mt} denotes the total number of separations in MDC m during week t .¹³

In a static framework, U_{mt} reflects the quantity of services that equalize demand and supply, as a result \bar{V}_t is subject to both demand and supply shocks. Our dynamic model provides a structural interpretation which helps to inform the identification strategy for the empirical analysis. We maintain the key (yet plausible) assumption that these shocks do not directly affect the differential likelihood of transfers of complex versus non-complex patients except through capacity utilization; specifically, these shocks affect the evolution of the state variable (capacity utilization), which in turn affects the differential likelihood of future transfers. One possible complication is variations in new patient arrivals (i.e., demand shocks) may, for whatever reason, be correlated with variations in the arrivals of complex versus non-complex patients. We guard against this by controlling for the composition of patient stock in our empirical model. In addition, due to the dynamic structure of our model, the current state (capacity utilization) is determined by the hospital's past states and decisions, yet the current demand shock is unanticipated by the hospital in prior periods.¹⁴

Negative supply shocks, such as a nursing staff strike or sudden resignation of key staff, will negatively affect the ability of the hospital to provide services and increase capacity utilization. This will, from our theoretical results, increase the transfer of high complexity patients relative to non-complex patients. However, the observed capacity in our empirical estimation does not change since by construction, U_{mt}^{\max} remains constant during the given period. This implies that, with such shocks, our observed measure of capacity utilization will be lower than its true value. Because the transfer of high complexity patients becomes more likely, the estimated effect of capacity utilization will

¹³For the empirical implementation below, δ is treated as a constant to be absorbed by the regression coefficients.

¹⁴Our setting is similar to Chan (2018), who discusses an empirical application in a different context but derived from a similar dynamic structural model in which current state variables depend on previous state variables and shocks.

be conservative as it is biased toward zero. Thus our empirical analysis will tend to underestimate the effect of cream skimming.

Table 1 shows some descriptive statistics of selected variables at both patient and hospital levels of aggregation. Out of the total 12.1 million episodes, 386,048 episodes (3.2%) involve transfers. The average age of patients in the sample was 54 years, and about 43% were male patients. Private patients accounted for nearly 42% of the sample. In terms of case complexity, most episodes (more than 81%) were non-complex, 13% had low complexity, and only 5% were with high complexity. Rural patients account for close to 5% of the total episodes, while same-day patients account for nearly half of all episodes. Patients in rare DRGs, defined as DRGs with fewer than 5,000 episodes during the 12-year sample period, account for 1.4% of total episodes.

At the hospital level, the number of transfers of high-complexity patients, net of non-complex patients, is on average -1.9 per week (i.e, non-complex patients outnumber high complexity patients), as a percentage to total weekly separations, it is -1.1%, while the average weekly capacity utilization rate is 56%. Private hospitals account for about 34% of all hospitals, and on a per hospital basis, the percentage to total weekly separations of rare DRG episodes is 1.3%, rural patient is 15% and same-day episodes is 39%. The weekly number of patients are 274 on average, and the number of non-, low- and high-complexity patients per week average respectively 190, 57, and 27 patients.

Figure 1 shows the distribution of the capacity utilization measure, which is somewhat skewed. The median capacity utilization is at 60% while the 25th and 75th percentiles are respectively 39% and 75%. To allow for potential nonlinear effects and nonlinear interactions with patient complexity and hospital (or patient) type, we divide V_j into discrete categories so that nonlinear effects can be reflected in the coefficients of the dummy variables and the interaction terms. For the estimation we use a capacity utilization variable with four discrete levels, denoted $w_{jk} \in \{w_{j1}, w_{j2}, w_{j3}, w_{j4}\}$, corresponding to $V_j \in (0, 0.70]$, $V_j \in (0.70, 0.80]$, $V_j \in (0.80, 0.85]$, $V_j \in (0.85, 1.00]$. The ‘cut-points’ correspond roughly to the 70th, 90th and 95th percentiles of the distribution of V_j . As

Table 1: Descriptive statistics by level of aggregation

	Variable	Mean	Std dev	Min	Max
Patient level	Transferred	0.032	0.175	0	1
	Age	54.3	19.7	18	99
	Male	0.426	0.494	0	1
	Married	0.627	0.484	0	1
	Private patient	0.418	0.493	0	1
	Admitted via ED	0.309	0.462	0	1
	Complexity category				
	Non complex (c_0)	0.816	0.387	0	1
	Low complexity (c_1)	0.134	0.341	0	1
	High complexity (c_2)	0.050	0.218	0	1
	Rural patient	0.049	0.217	0	1
	Rare DRG patient	0.014	0.116	0	1
	Same-day patient	0.486	0.500	0	1
	Hospital level	Diff. in number of transfers of high-complexity v non-complex patients (D_j)	-1.903	4.992	-57
Diff. in proportion of transfers of high-complexity v non-complex patients (d_j)		-0.011	0.047	-1	1
Capacity utilization (V)		0.560	0.230	0	1
Private hospital		0.344	0.475	0	1
Proportion rare DRGs*		0.013	0.034	0	1
Proportion rural patients		0.154	0.334	0	1
Proportion same-day episodes		0.391	0.252	0	1
Weekly patient stock, all patients		274.3	408.6	0	2,797
Weekly non-complex patient stock		190.5	264.3	0	1,903
Weekly low complexity patient stock		56.7	124.0	0	1,471
Weekly high complexity patient stock		27.1	51.6	0	381
No. admission episodes					12,149,040
No. hospitals					180
No. 3-digit DRGs					425
No. MDCs				23	

*Rare DRGs are defined as DRGs with fewer than 5,000 episodes during the entire 12-year sample period.

Additional variables used in the hospital-level estimation include financial year dummies and month dummies and in the patient-level estimation include MDC dummies, SEIFA quintiles, financial year, month and day of separation dummies.

a robustness check, we repeat the analysis using eight discrete categories; results are largely unchanged and are reported in Appendix S.

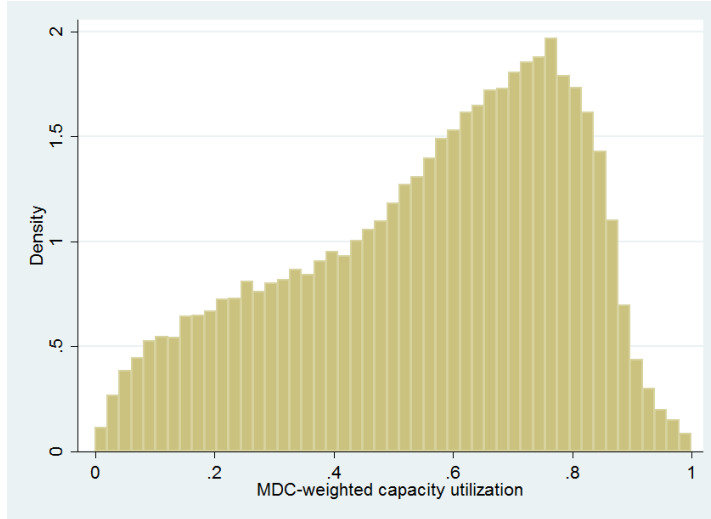


Figure 1: MDC-weighted weekly capacity utilization

3.2 Hospital level analysis

We begin by specifying the number of high complexity patients net of non-complex patients transferred from hospital j during week t , $D_{jt} = n_{jt}^h - n_{jt}^e$, as a function of the stock of high complexity patients, non-complex patients, and capacity utilization. In the theoretical discussion hospitals do not transfer easy patients, although we remarked that our results could be extended to allow for an exogenous transfer rate of easy (and hard) patients (see Remark 2 above). In practice, the transfer of easy patients is not zero—patients are transferred for many reasons, not all of which are initiated by hospitals. Since more than 80% of the patients are non-complex, it is therefore necessary to consider the magnitude of transfers of high complexity patients relative to that of non-complex patients.

As an alternative to the number of patients transferred, we also examine the proportion of patients transferred. Let Q_{jt} be the total number of separations (i.e., discharges and transfers) from hospital j during week t . The proportion of transfers is calculated as $d_{jt} = D_{jt}/Q_{jt}$.

Assuming a linear form, we specify:

$$E[y_{jt} | X_{jt}] = \alpha + \sum_k \beta_{1k} w_{jtk} + \beta_2 v_j + \sum_k \beta_{3k} v_j \times w_{jtk} + \kappa_{jt} \gamma + \mu_j + \nu_t, \quad (14)$$

where $y_j = D_j$ or d_j , $X_{jt} = [w_{jt}, v_j, \kappa_{jt}]$ is the vector of key explanatory variables, with $w_{jt} \equiv [w_{jt1}, \dots, w_{jt4}]$ denoting the vector of capacity utilization dummy variables, $v_j = 1$ if hospital j is private, otherwise $v_j = 0$ for public hospital, κ_{jt} a vector of other time-varying or time-invariant covariates, and μ_j and ν_t the hospital effects and time effects respectively.¹⁵

We say that hypothesis $H1$ is rejected by the data if: (i) one or more marginal effects of capacity utilization are negative, i.e., $\frac{\Delta y_j}{\Delta w_{jk}} \leq 0$ for some $k = 2, 3, 4$ with $\frac{\Delta y_j}{\Delta w_{jk}} < 0$ for at least one k ; or $\frac{\Delta y_j}{\Delta w_{jk}} = 0$ for all k . These conditions imply $H1$ is not rejected if none of the marginal effects is negative and one or more is strictly positive, i.e., $\frac{\Delta y_j}{\Delta w_{jk}} \geq 0 \forall k$ with $\frac{\Delta y_j}{\Delta w_{jk}} > 0$ for at least one k . Note that the marginal effects are evaluated with reference to the lower adjacent capacity utilization level, i.e., the marginal effect of w_{jk} is with reference to $w_{j,k-1}$ for $k = 2, 3, 4$. With respect to hypothesis $H2$, the rejection criterion is simply to reject $H2$ if $\frac{\Delta y_j}{\Delta v_j} \leq 0$.

We use the correlated random effects (CRE) to estimate equation (14), which is a viable alternative to fixed effects estimation in capturing individual unit heterogeneity (Wooldridge, 2019). Conventional approaches for estimating fixed effects models (e.g. the ‘within estimator’) are not feasible because the key variable, hospital type (v_j), is time invariant. Under the CRE, individual heterogeneity term μ_j is expressed as a function of time averages of time-varying covariates, i.e. $\mu_j = \eta_0 + \bar{X}_j \eta_1 + \varepsilon_j$. Here, $\bar{X}_j = T^{-1} \sum_{t=1}^T X_{jt}$. It is usually assumed that $E[\varepsilon_j | X_j] = 0$.

¹⁵An alternative approach is to aggregate the data by hospital-week-complexity and estimate the proportion (or number) of transfers on a complexity dummy interacted with capacity and hospital type as two separate covariates on the right-hand side. This approach embodies a triple-differencing specification. It can be shown mathematically that by taking the difference between the transfers of high complexity and non-complex patients, the specification will reduce to (14), and both specifications are based on the same set of identifying conditions. Proofs are available upon request. We are grateful to an anonymous reviewer for suggesting this alternative approach.

In addition to the CRE model, we further show that our results are robust with respect to unobservable hospital effects using two restricted forms of fixed effects estimation. First, under the hospital level analysis, we omit all time-invariant variables and estimate the fixed effects models (without interacting hospital type with capacity utilization). Second, at the patient level we consider a subsample of public and private patients admitted to public hospitals only. Here we imagine that each public hospital has two divisions—public and private division—run by managers with different objective functions. The results, discussed in Section 3.4, are broadly consistent with our main results obtained from the CRE models.

For the CRE models, the unit of observations is hospital-week and in total there are 95,131 observations. The explanatory variables, in addition to hospital type and capacity utilization variables (and their interactions), include the weekly stock of all patients, weekly stocks of no- and high-complexity patients¹⁶, proportion of admissions in rare DRGs, proportion of rural patients and proportion of same-day episodes. Also included are group means of time-varying variables as Mundlak adjustment terms, month-of-the-year and year dummies.

The marginal effect estimates from the CRE estimation are reported in Table 2; coefficient estimates are reported in Appendix B. The estimates of the marginal effects of the coefficients β_{1k} are positive for all k , and indicate that the number (D) and proportion (d) of net transfers are increasing with higher levels of capacity utilization. Specifically, as capacity utilization is varied from w_1 to w_2 , w_2 to w_3 , and w_3 to w_4 , the number of transfers of high complexity patients net of non-complex patients, increase by respectively, 0.457, 0.248, and 0.144 per week. The analysis of the net transfer proportions show the same result: the proportions of transfers of high complexity patients (net of non-complex patients) increase by 0.17, 0.09 and 0.04 percentage points for higher levels of capacity utilization w_2 to w_4 .

¹⁶Weekly stock proportions of no- and high-complexity patients (in total patients) are used in regression of d_j .

The positive estimate of the marginal effect on β_2 indicates that private hospitals are more likely to transfer high complexity patients than public hospitals. Compared to public hospitals, private hospitals transfer 1.09 or 0.64 percentage points more high complexity patients net of non-complex patients. While the marginal effects may appear small in magnitude, the effects are in fact large relative to the means of the net number and proportion of transfers, which are -1.9 and -1.1% respectively.

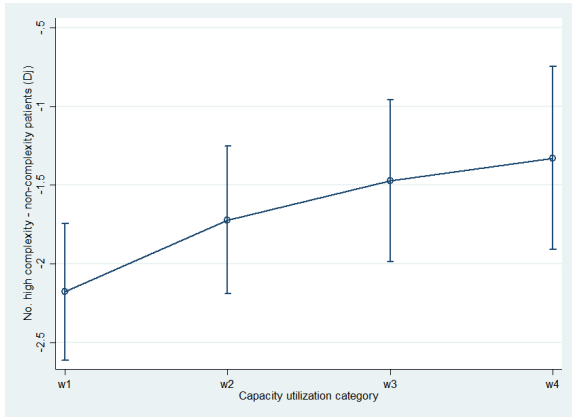
Table 2: Marginal effects, hospital-level CRE estimation

	D_j		d_j	
	MEff.	<i>p</i> -val	MEff.	<i>p</i> -val
Capacity utilization level				
w_2 (v. w_1)	0.4570 (0.0861)	$p < 0.001$	0.0017 (0.0008)	0.029
w_3 (v. w_2)	0.2483 (0.0967)	0.010	0.0009 (0.0004)	0.023
w_4 (v. w_3)	0.1442 (0.1452)	0.321	0.0004 (0.0008)	0.613
Private (v. public) hospital	1.0915 (0.4812)	0.023	0.0064 (0.0024)	0.007
No. of observations	95,131		95,131	
No. of hospitals	180		180	

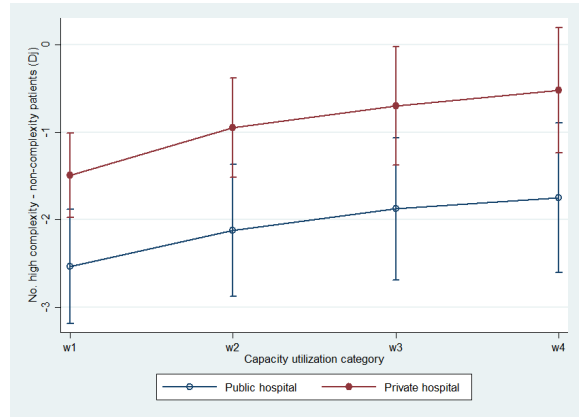
Notes: (i) D_j denotes number of high-complexity (net of non-complex) patients transferred. (ii) d_j denotes proportion of high-complexity (net of no-complexity) patients transferred. (iii) Robust standard errors are obtained by clustering (by hospital). (iv) Figures in parentheses are standard errors obtained by delta method.

The estimates reported in Table 2 suggest that both hypotheses $H1$ and $H2$ are supported by the data. All marginal effects of capacity utilization are positive with at least one statistically significant at 5% or higher, while the marginal effects of private hospital are positive and statistically significant at 5% or higher.

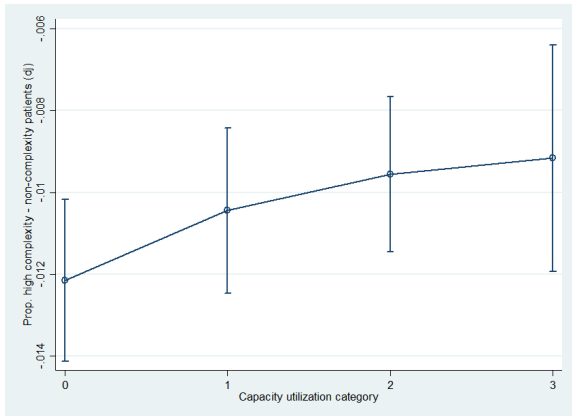
To visualize the results, we compute average predicted values of the dependent variables by varying the capacity utilization levels and separately for public and private hospitals, while holding all other explanatory variables at their sample values. The results are shown in Figure 2. Panels (a) and (b) show the difference in the number of transfers between high complexity and non-complex patients, while panels (c) and (d) show the proportionate difference in transfers.



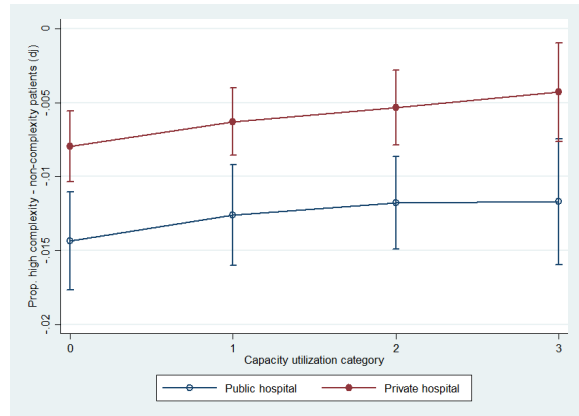
(a) D_j by capacity utilization



(b) D_j by cap. uti. & pub./priv. hospital



(c) d_j by capacity utilization



(d) d_j by cap. uti. & pub./priv. hospital

Note: D_j denotes number of high-complexity (net of non-complex) patients transferred; d_j denotes proportion of high complexity (net of non-complex) patients transferred.

Figure 2: Average predicted number and proportion of transfers

The figures reveal distinct patterns in the transfers of high complexity relative to non-complex patients. First, the left-hand panels show the number and proportion of transfers of high complexity patients (relative to non-complex patients) increase monotonically as capacity utilization rises. Second, there is a persistent gap in transfers, in number and proportion, across all capacity utilization levels between private and public hospitals, as shown by the right-hand panels. Note that in Panel (b), although the confidence intervals for public and private hospitals appear to overlap at every capacity utilization level, the difference on average across all utilization levels is statistically significant with a p -value of 0.023, as reported in Table 2.

Hospital level fixed effects estimation

To show that our results are robust to unobservable hospital-specific heterogeneity, we estimate a variant of (14) where all time-invariant covariates are omitted. Hospital fixed effects are accounted for using a set of hospital dummies. To test hypothesis $H1$, we compute the marginal effects of capacity utilization as before. However, testing hypothesis $H2$ requires the marginal effect of the private hospital dummy which we omit as this variable is time invariant. To obtain the difference in net transfers by hospital type, we first use the estimated coefficients to calculate the predicted outcomes for each public and private hospital in our sample. We then use the difference in the average predicted values for public and private hospitals as the estimate of the marginal effect of hospital type. The results, reported in Table 3, are in line with our earlier findings that both hypotheses $H1$ and $H2$ are supported by the data.

3.3 Patient level analysis

A disadvantage of the hospital level analysis is that individual patient characteristics cannot be accounted for. To do so, we conduct a patient level analysis by examining the probability of a patient getting transferred conditional on his or her characteristics and

Table 3: Marginal effects, hospital-level FE models (no interactions)

	D_j		d_j	
	MEff.	p -val	MEff.	p -val
Capacity utilization level				
w_2 (v. w_1)	0.4864 (0.0798)	$p < 0.001$	0.0017 (0.0006)	0.007
w_3 (v. w_2)	0.2526 (0.0957)	0.009	0.0009 (0.0004)	0.030
w_4 (v. w_3)	0.1491 (0.1421)	0.295	0.0003 (0.0007)	0.669
Private (v. public) Hospital	1.1719 (0.0097)	$p < 0.001$	0.0055 (0.0001)	$p < 0.001$
No. of observations	95,131		95,131	
No. of hospitals	180		180	

Notes: (i) D_j denotes number of high-complexity (net of no-complexity) patients transferred. (ii) d_j denotes proportion of high-complexity (net of no-complexity) patients transferred. (iii) Robust standard errors are obtained by clustering (by hospital). (iv) Figures in parentheses are standard errors obtained by delta method.

the characteristics of the transferring hospital.

We consider two possible outcomes for a patient i admitted to hospital j : the patient is either discharged or transferred.¹⁷ We postulate that the probability of transfer of patient i from hospital j depends on the characteristics of the patient and the transferring hospital. The key specification is on the interaction between patient complexity, patient type and levels of capacity utilization. Given that public hospitals in Australia also admit private patients, it is natural for the patient level analysis to focus on patient type rather than hospital type. We also perform fixed effects estimation using the subset of public hospitals that also admit private patients.

Let T_{ij} be a binary variable denoting transfer if $T_{ij} = 1$, otherwise $T_{ij} = 0$ indicates that the patient is discharged but not transferred. As before, $c_{ij} \in \{c_0, c_1, c_2\}$ denote the case complexity of patient i which is constructed based on the Charlson comorbidity score, and $w_j \in \{w_1, w_2, w_3, w_4\}$ is the capacity utilization measure for hospital j faced by patient i in the week of her separation and is given in four discrete categories as before.

¹⁷The possibility that the patient died in hospital is not considered since the sample excludes patients who died.

Omitting the time subscript, we postulate that

$$\Pr(T_{ij} = 1 \mid X_{ij}) = g(v_{ij}, c_{ij}, w_j, Z_{ij}), \quad (15)$$

where v_{ij} is the patient type (public or private) of patient i and Z_{ij} is other relevant characteristics of patient i such as gender and age. We focus on the interaction between c_{ij} , v_{ij} and w_j .

Note that v_{ij} now denotes patient type, not hospital type. As mentioned earlier, patients with private health insurance can elect to be admitted into public hospitals as private patients. For patients in private hospitals, patient type and hospital type are synonymous since private hospitals only admit private patients.¹⁸ This consideration rules out estimating hospital fixed effects, since within-hospital variation in patient types are found only in public hospitals which admit both private and public patients. As a supplement to the main analysis, we perform a subgroup analysis of public hospitals that admit both public and private patients using hospital fixed effects and discuss the results in Section 3.4.

For the estimation of (15) using the full sample, we specify a logistic model with triple interactions of patient type, complexity, and capacity utilization. Cluster-robust standard errors are obtained by clustering at the hospital-MDC unit level. Marginal effects of these three variables are reported in Table 4. For the purpose of testing hypotheses $H1$ and $H2$, we compute estimates of the change in marginal effects as patient complexity is varied from no complexity ($c_{ij} = 0$) to high complexity ($c_{ij} = 2$), that is, we compute the quantities Λ_{w_j} and Λ_v , defined as

$$\Lambda_{w_j} \equiv \frac{\Delta \Pr(T_{ij} = 1)}{\Delta w_j} \Big|_{c_{ij}=2} - \frac{\Delta \Pr(T_{ij} = 1)}{\Delta w_j} \Big|_{c_{ij}=0}, \text{ for } j = 2, 3, 4.$$

and

$$\Lambda_v \equiv \frac{\Delta \Pr(T_{ij} = 1)}{\Delta v_{ij}} \Big|_{c_{ij}=2} - \frac{\Delta \Pr(T_{ij} = 1)}{\Delta v_{ij}} \Big|_{c_{ij}=0}.$$

Note that in computing Λ_{w_j} , w_{j-1} is set as the reference category.

¹⁸Public patients can be admitted to private hospitals under some special arrangements but these cases are rare and are omitted from the sample.

Similar to the case of hospital level analysis, we say that hypothesis $H1$ is rejected if one or more Λ_{w_j} is significantly negative, or none is significantly positive, and hypothesis $H2$ is rejected if Λ_v is insignificant from zero or significantly negative. The estimates of Λ_{w_j} and Λ_v are reported in Table 4. We list all coefficient estimates and standard errors in Appendix B.

Table 4: Estimates of marginal effects and changes in marginal effects, logistic estimation

Dep. var.: Transfer (T_{ij})	Marg. eff.	std. err.	p -val
Marginal effect estimates			
Patient complexity (Ref: Non-complex (c_0))			
Low complexity (c_1)	0.0064	(0.0008)	$p < 0.001$
High complexity (c_2)	0.0104	(0.0014)	$p < 0.001$
Capacity utilization level (Ref: w_1)			
w_2	0.0057	(0.0014)	$p < 0.001$
w_3	0.0047	(0.0019)	0.0117
w_4	0.0017	(0.0020)	0.3955
Private patient (Ref: public patient)	-0.0097	(0.0016)	$p < 0.001$
Estimated changes in marginal effects			
Λ_{w_2}	-0.0009	(0.0015)	0.533
Λ_{w_3}	0.0050	(0.0016)	0.002
Λ_{w_4}	0.0009	(0.0019)	0.649
Λ_v	0.0110	(0.0023)	$p < 0.001$
No. observations	12,072,395		

Notes: (i) Robust standard errors are obtained by clustering at hospital-MDC unit level. (ii) Standard errors are obtained by delta method.

We make several observations about Table 4. First, complex patients are more likely to be transferred than non-complex patients, and among complex patients, patients with high complexity are almost twice as likely to be transferred compared to patients with low complexity. Second, transfers become more likely as capacity utilization rises, although the increase appears to plateau at higher levels of capacity utilization. Third, conditional on patient complexity, private patients are less likely to be transferred than public patients.

The estimated changes in marginal effects, Λ_{w_j} and Λ_v , suggest that hypotheses $H1$ and $H2$ are supported by the data. We find that high complexity patients are more likely to be transferred at higher levels of capacity utilization than non-complex patients. This is

most evident for capacity level w_3 where the estimated effect is positive and statistically significant. The estimated effects for capacity levels w_2 and w_4 , in contrast, are not statistically significant. We further find that high complexity private patients are more likely to be transferred than public patients of the same complexity level. This is shown in the estimate of Λ_v which is positive and statistically significant. Overall, the estimates we find are economically large given that the unconditional mean transfer rate is 3.2%.

To visualize the change in probabilities as capacity utilization varies, we show in Figure 3 the difference in predicted probabilities of transfers between high complexity and non-complex patients. Figure 3(a) shows the difference in transfer probabilities is generally rising at higher levels of capacity utilization. Figure 3(b) shows that, although high complexity patients of both types are more likely to be transferred than non-complex patients, the differences are bigger for private than public patients.

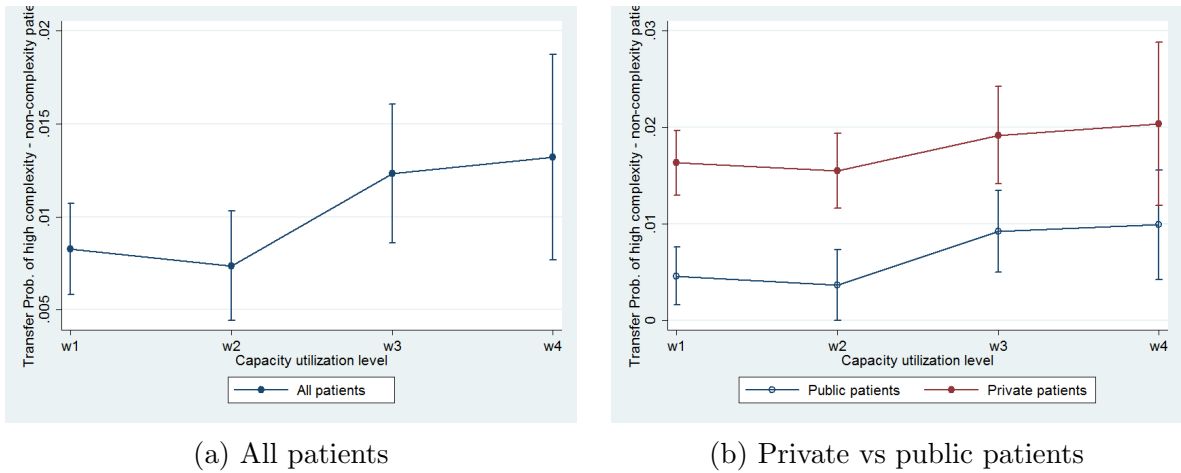


Figure 3: Difference in predicted transfer probability between high complexity and non-complex patients

3.4 Public and private patients in public hospitals

We carry out a subgroup analysis focusing on patients admitted to public hospitals that admit both public and private patients. For this subgroup, we can employ hospital fixed effects in the logistic estimation while keeping patient type as an explanatory variable.

Hospitals in this group are assumed to wear two hats—profit-maximizing when admitting private patients and output maximizing when admitting public patients. However, there could be complex interactions between the two motives. As discussed in Section 2.3, a possible extension of the theory is to assume different managers in charge of private and public patients in the same hospital.

Table 5 presents the estimated marginal effects. The key results are substantively similar to those from the full sample, with a few notable exceptions. As before, complex patients are more likely to be transferred. However, unlike the previous results, the probabilities of transfers do not vary with capacity utilization, and the marginal effect of private patients, which was negative under the full sample, turns out to be positive and significant. This suggests that private patients in public hospitals are more likely to get transferred than public patients. Despite these differences, the estimated changes in marginal effects, Λ_{w_j} and Λ_v , remain broadly similar. Although some estimates change signs, these are not statistically significant from zero. Testing of hypotheses *H1* and *H2* produces the same conclusion as before.

Table 5: Estimates of marginal effects and changes in marginal effects, public and private patients in public hospitals, logistic estimation

Dep. var.: Transfer (T_{ij})	Marg. eff.	std. err.	<i>p</i> -val
Marginal effect estimates			
Patient complexity (Ref: Non-complex (c_0))			
Low complexity (c_1)	0.0074	(0.0012)	$p < 0.001$
High complexity (c_2)	0.0115	(0.0020)	$p < 0.001$
Capacity utilization level (Ref: w_1)			
w_2	-0.0006	(0.0005)	0.214
w_3	-0.0001	(0.0007)	0.872
w_4	-0.0002	(0.0010)	0.844
Private patient (Ref: public patient)	0.0047	(0.0023)	0.037
Estimated changes in marginal effects			
Λ_{w_2}	0.0007	(0.0019)	0.730
Λ_{w_3}	0.0049	(0.0020)	0.012
Λ_{w_4}	-0.0005	(0.0025)	0.830
Λ_v	0.0100	(0.0035)	0.004
No. observations	7,730,057		

Note: (i) Robust standard errors are obtained by clustering at hospital-MDC unit level. (ii) Standard errors are obtained by delta method.

4 Concluding Remarks

This study examines cream skimming in a model of hospital transfers, in which the desire to optimize the use of limited capacity drives cream skimming. The model predicts that private hospitals have a stronger incentive than public hospitals to transfer complex patients to free up limited capacity. Although public hospitals have the same tendency to do so, their incentive to maximize social welfare is strictly weaker than the profit motive of private hospitals. We show that cream skimming leads to the following testable predictions. First, complex patients are more likely to be transferred relative to non-complex patients as capacity utilization rises. Second, private hospitals are strictly more likely to transfer complex patients (relative to non-complex patients) than public hospitals. The empirical analysis tests the theoretical predictions at two levels of aggregation using hospital administrative data from Australia. The hospital level analysis examines the weekly transfers from public and private hospitals of complex relative to non-complex patients. In addition, transfers at the patient level is examined by relating the probability of transfers (*vis-a-vis* discharges) to patient complexity, capacity utilization and private patient status, while allowing for additional patient and hospital characteristics. The empirical results lend support to the theoretical predictions of cream skimming in Australia's hospital system.

This paper contributes to the theoretical and empirical literature on cream skimming, and to the broader literature on how financial incentives affect the behavior of health care providers. Our results have implications for countries aiming to increase private sector participation in health care markets. Constraints in public sector budgets coupled with rapidly rising health care costs have prompted many countries to look to the private sector. In the UK, market-based reforms were introduced in the mid 2000s to improve competition by allowing patients to choose their hospital care providers, and by providing information on the quality (e.g., risk-adjusted mortality rates, waiting times, etc.) of providers (Gaynor et al., 2013, 2016). Competition for patients has resulted in new entries of for-profit surgical centres. Although efficiency may improve from greater

private participation (e.g. Siciliani et al., 2013), the potential of cream skimming has also increased, which could erode the gain from private sector involvement by leaving public hospitals with increasingly sicker and costlier mix of patients (Cooper et al., 2018).

For Australia's mixed public-private hospital system, cream skimming has long been a policy concern (O'Loughlin, 2002). The Australian government has for many years been subsidizing private hospital insurance with a purported aim of encouraging private hospital admissions to alleviate the pressure on public hospitals. However, the likely occurrence of cream skimming means that private hospitals get to enjoy higher profits while public hospitals are likely to be saddled with complex and difficult to treat patients, worsening the financial situation of the public sector.

References

- Barros, P.P, 2003. Cream-skimming, incentives for efficiency and payment system. *Journal of Health Economics* 22, 419–443.
- Barros, P.P, Siciliani, L, 2011. Public and Private Sector Interface. In: Pauly, M.V, McGuire, T.G, Barros, P.P (Eds.), *Handbook of Health Economics*, vol. 2. Elsevier B.V., 927–1001.
- Berta, P, Callea, G, Martini, G, Vittadini, G, 2010. The effects of upcoding, cream skimming and readmissions on the Italian hospitals efficiency: A population-based investigation. *Economic Modelling* 27, 812–821.
- Brown, J, Duggan, M, Kuziemko, I, Woolston, W, 2014. How Does Risk Selection Respond to Risk Adjustment? New Evidence from the Medicare Advantage Program. *American Economic Review* 104, 3335–3364.
- Chan, M.K, 2018. Measuring the effects of welfare time limits. *Journal of Human Resources* 53, 232–271.
- Chan, M.K, Zeng, G, 2018. Unintended consequences of supply-side cost control? Evidence from China’s new cooperative medical scheme. *Journal of Health Economics* 61, 27–46.
- Charlson, M, Szatrowski, T.P, Peterson, J, Gold, J, 1994. Validation of a combined comorbidity index. *Journal of Clinical Epidemiology* 47, 1245–1251.
- Chen, Y, Meinecke, J, 2012. Do healthcare report cards cause providers to select patients and raise quality of care? *Health Economics* 21, 33–55.
- Cheng, T.C, Haisken-DeNew, J.P, Yong, J, 2015. Cream skimming and hospital transfers in a mixed public-private system. *Social Science and Medicine* 132, 156–164.
- Cooper, Z, Gibbons, S, Skellern, M, 2018. Does competition from private surgical centres improve public hospitals’ performance? Evidence from the English National Health Service. *Journal of Public Economics* 166, 63–80.
- Dranove, D, Jin, G.Z, 2010. Quality Disclosure and Certification: Theory and Practice. *Journal of Economic Literature* 48, 935–963.
- Dranove, D, Kessler, D, McClellan, M, Satterthwaite, M, 2003. Is more information better? The effects of ”Report Cards” on health care providers. *Journal of Political Economy* 111, 555–588.
- Duggan, M.G, 2000. Hospital Ownership and Public Medical Spending. *The Quarterly Journal of Economics* 115, 1343–1373.

- Eijkenaar, F, 2013. Key issues in the design of pay for performance programs. *European Journal of Health Economics* 14, 117–131.
- Einav, L, Finkelstein, A, Mahoney, N, 2018. Provider Incentives and Healthcare Costs: Evidence From Long-Term Care Hospitals. *Econometrica* 86, 2161–2219.
- Eliason, P.J, Grieco, P.L, McDevitt, R.C, Roberts, J.W, 2018. Strategic patient discharge: The case of long-term care hospitals. *American Economic Review* 108, 3232–3265.
- Ellis, R.P, 1998. Creaming, skimping and dumping: Provider competition on the intensive and extensive margins. *Journal of Health Economics* 17, 537–555.
- Gaynor, M, Moreno-Serra, R, Propper, C, 2013. Death by Market Power: Reform, Competition, and Patient Outcomes in the National Health Service. *American Economic Journal: Economic Policy* 5, 134–166.
- Gaynor, M, Propper, C, Seiler, S, 2016. Free to choose? Reform, choice, and consideration sets in the english national health service. *American Economic Review* 106, 3521–3557.
- Gaynor, M, Vogt, W.B, 2003. Competition among Hospitals. *RAND Journal of Economics* 34, 764–785.
- Glazer, J, McGuire, T.G, 2000. Optimal risk adjustment in markets with adverse selection: An application to managed care. *American Economic Review* 90, 1055–1071.
- Grassi, S, Ma, C.T.A, 2012. Public Sector Rationing and Private Sector Selection. *Journal of Public Economic Theory* 14, 1–34.
- Hackmann, M.B, Pohl, R.V, 2018. Patient vs. Provider Incentives in Long Term Care. NBER Working Paper No. w25178. Available at SSRN: <https://ssrn.com/abstract=3270758>.
- Iversen, T, 1997. The effect of a private sector on the waiting time in a national health service. *Journal of Health Economics* 16, 381–396.
- Lakdawalla, D, Philipson, T, 2006. The nonprofit sector and industry performance. *Journal of Public Economics* 90, 1681–1698.
- Ma, C.T.A, 2004. Public rationing and private cost incentives. *Journal of Public Economics* 88, 333–352.
- Mak, H.Y, 2017. Provider performance reports and consumer welfare. *RAND Journal of Economics* 48, 250–280.

- Martinussen, P.E, Hagen, T.P, 2009. Reimbursement systems, organisational forms and patient selection: Evidence from day surgery in Norway. *Health Economics, Policy and Law* 4, 139–158.
- McLawhorn, A.S, Schairer, W.W, Schwarzkopf, R, Halsey, D.A, Iorio, R, Padgett, D.E, 2018. Alternative Payment Models Should Risk-Adjust for Conversion Total Hip Arthroplasty: A Propensity Score-Matched Study. *Journal of Arthroplasty* 33, 2025–2030.
- O’Loughlin, M.A, 2002. Conflicting interests in private hospital care. *Australian Health Review* 25, 106–117.
- Scott, A, Liu, M, Yong, J, 2018. Financial Incentives to Encourage Value-Based Health Care. *Medical Care Research and Review* 75, 3–32.
- Siciliani, L, Sivey, P, Street, A, 2013. Differences in length of stay for hip replacement between public hospitals, specialised treatment centres and private providers: Selection or efficiency? *Health Economics* 22, 234–242.
- Sloan, F.A, 2000. Not-for-profit ownership and hospital behavior. In: Culyer, A.J, Newhouse, J.P (Eds.), *Handbook of Health Economics*, vol. 1. Elsevier, 1141–1174.
- Sundararajan, V, Henderson, T, Perry, C, Muggivan, A, Quan, H, Ghali, W.A, 2004. New ICD-10 version of the Charlson comorbidity index predicted in-hospital mortality. *Journal of Clinical Epidemiology* 57, 1288–1294.
- van de Ven, P.W, Ellis, R.P, 2000. Risk adjustment in competitive health plan markets. In: Culyer, A.J, Newhouse, J.P (Eds.), *Handbook of Health Economics*, vol. 1. Elsevier, 755–845.
- Wooldridge, J.M, 2019. Correlated random effects models with unbalanced panels. *Journal of Econometrics* 211, 137–150.

Appendix A

Proofs

Proof of Lemma 1

Proof. If $z_{t+1} > K - \tilde{x}_t + n_t^h$, then by the laws of motion we have

$$x_{t+1}^h = (1 - \theta^h)x_t^h - n_t^h + p_h(K - \tilde{x}_t + n_t^h) = (1 - \theta^h)x_t^h + p_h(K - \tilde{x}_t) - p_e n_t^h \text{ and}$$

$$x_{t+1}^e = (1 - \theta^e)x_t^e + p_e(K - \tilde{x}_t + n_t^h) = (1 - \theta^e)x_t^e + p_e(K - \tilde{x}_t) + p_e n_t^h. \text{ Thus}$$

$$\frac{\partial x_{t+1}^h}{\partial n_t^h} = -p_e \text{ and } \frac{\partial x_{t+1}^e}{\partial n_t^h} = p_e. \text{ If } z_{t+1} \leq K - \tilde{x}_t + n_t^h, \text{ then } x_{t+1}^h = (1 - \theta^h)x_t^h - n_t^h + z_{t+1}^h$$

$$\text{and } x_{t+1}^e = (1 - \theta^e)x_t^e + z_{t+1}^e. \text{ Thus } \frac{\partial x_{t+1}^h}{\partial n_t^h} = -1 \text{ and } \frac{\partial x_{t+1}^e}{\partial n_t^h} = 0. \blacksquare$$

Proof of Proposition 1

Proof. Note that both $\frac{dEV}{dn_t^h}$ and $\frac{dC}{dn_t^h}$ are continuous in n_t^h , and $\frac{d^2EV}{d(n_t^h)^2} < 0$ and

$$\frac{d^2C}{d(n_t^h)^2} \geq 0. \text{ If } \frac{dEV}{dn_t^h} \Big|_{n_t^h=0} < \frac{dC}{dn_t^h} \Big|_{n_t^h=0}, \text{ then } n_t^{h*} = 0. \text{ If } \frac{dEV}{dn_t^h} \Big|_{n_t^h=x_t^h} > \frac{dC}{dn_t^h} \Big|_{n_t^h=x_t^h}, \text{ then}$$

$$n_t^{h*} = x_t^h. \text{ Otherwise, there exists a unique } x_t^h > n_t^{h*} > 0 \text{ such that } \frac{dEV}{dn_t^h} \Big|_{n_t^{h*}} = \frac{dC}{dn_t^h} \Big|_{n_t^{h*}}. \blacksquare$$

Proof of Proposition 2

Proof. We will show that the marginal future value, $\frac{dEV}{dn_t^h}$, increases as x_t^e or x_t^h

increases. Once this is established, it follows that n_t^{h*} is an increasing function of x_t^e and x_t^h because the marginal cost does not depend on x_t^e and x_t^h . From (9) and (10) and Lemma 1, the only channel by which x_t^e and x_t^h can affect $\frac{dEV}{dn_t^h}$ is via affecting the

integration space $A(n_t^h; x_t^e, x_t^h)$. The integrands do not depend on x_t^e and x_t^h because

the derivatives $\frac{\partial V}{\partial x_{t+1}^e}, \frac{\partial V}{\partial x_{t+1}^h}, \frac{\partial x_{t+1}^e}{\partial n_t^h}, \frac{\partial x_{t+1}^h}{\partial n_t^h}$ are not functions of x_t^e and x_t^h . By the definition

of $A(\cdot)$ and assumption A1, $\frac{\partial A}{\partial x_t^h} > \frac{\partial A}{\partial x_t^e} > 0$. Thus, $\frac{\partial}{\partial x_t^h}(\frac{dEV}{dn_t^h}) > \frac{\partial}{\partial x_t^e}(\frac{dEV}{dn_t^h}) > 0$. It follows

that $\frac{dn_t^{h*}}{dx_t^h}, \frac{dn_t^{h*}}{dx_t^e}, \frac{dn_t^{h*}}{dx_t^e} \geq 0$. When n_t^{h*} is an interior solution, the first order condition in

$$(8) \text{ implies that } \frac{dn_t^{h*}}{dx_t^h}, \frac{dn_t^{h*}}{dx_t^e}, \frac{dn_t^{h*}}{dx_t^e} > 0 \text{ and } \frac{dn_t^{h*}}{dx_t^h} > \frac{dn_t^{h*}}{dx_t^e}. \blacksquare$$

Proof of Corollary 1

Proof. The profit and social welfare functions can be rescaled as $\frac{\pi^e \theta^e}{\pi^h \theta^h} x_t^e + x_t^h$ and $\frac{\theta^e}{\theta^h} x_t^e + x_t^h$, respectively. The rescaling does not affect the optimal patient transfer decisions, as is evident from the value function in (7) and the marginal future value in

(10). Because $\frac{\pi^e \theta^e}{\pi^h \theta^h} > \frac{\theta^e}{\theta^h}$, the marginal future value of transferring hard patients is larger in the private hospital than the public hospital, all else being equal (the

integrand $\frac{\partial V}{\partial x_{t+1}^e} - \frac{\partial V}{\partial x_{t+1}^h}$ is larger in the private hospital). Therefore, C1.1 follows in a similar argument to that in Proposition 2. Next, C1.2 follows directly from

Proposition 2. \blacksquare

Appendix B Estimation Results

This Appendix presents a listing of coefficient estimates and standard errors of our base models discussed in the text. Table B1 shows the estimated coefficients and standard errors of the correlated random effects models estimated at the hospital-week level. For the patient level analysis, Table B2 shows the coefficient estimates and standard errors of the logistic regression using the full sample.

Table B1: Correlated random effects models, hospital level estimation

Dep. var:	D_j			d_j		
	Coeff.	Std. err.	p -val	Coeff.	Std. err.	p -val
Capacity utilization (Ref: w_1)						
w_2	0.4123	0.1171	$p < 0.001$	0.0017	0.0011	0.1291
w_3	0.6597	0.1906	$p < 0.001$	0.0026	0.0011	0.0173
w_4	0.7853	0.2355	$p < 0.001$	0.0026	0.0017	0.1114
Private Hospital (Ref: Public)	1.0435	0.4429	0.0185	0.0064	0.0025	0.0092
Private \times Capacity utilization (Ref: Public $\times w_1$)						
Private $\times w_2$	0.1299	0.1557	0.4042	-0.0001	0.0013	0.9633
Private $\times w_3$	0.1327	0.2420	0.5835	0.0000	0.0014	0.9829
Private $\times w_4$	0.1867	0.2788	0.5030	0.0010	0.0023	0.6504
Weekly patient stock, all patients	-3.5872	0.5549	$p < 0.001$	-0.0023	0.0012	0.0544
Weekly non-complex patient stock*	-0.4648	0.2747	0.0907	-0.0211	0.0027	$p < 0.001$
Weekly high complexity patient stock*	0.5087	0.1745	0.0035	0.0408	0.0067	$p < 0.001$
Proportion rare DRGs	-0.2598	0.1325	0.0499	-0.0259	0.0103	0.0121
Proportion rural patients	0.1442	0.0950	0.1292	0.0007	0.0064	0.9095
Proportion same-day episodes	0.9302	0.1789	$p < 0.001$	0.0207	0.0031	$p < 0.001$
Group mean of Capacity utilization	1.5089	1.7160	0.3792	-0.0234	0.0114	0.0398
Average weekly patient stock, all patients	6.7592	2.7346	0.0134	0.0007	0.0015	0.6275
Average weekly non-complex patient stock*	-5.2747	2.4250	0.0296	-0.0029	0.0238	0.9016
Average weekly high complexity patient stock*	-0.9674	0.6955	0.1642	-0.0332	0.0342	0.3311
Average proportion rare DRGs	17.7318	18.2134	0.3303	0.2872	0.1219	0.0185
Average proportion rural patients	0.4440	0.2512	0.0771	0.0043	0.0066	0.5123
Average proportion same-day episodes	0.3988	0.8874	0.6531	0.0107	0.0135	0.4276
Constant	-1.0343	0.4362	0.0177	0.0020	0.0170	0.9048
No. of observations		95,131			95,131	
No. of hospitals		180			180	

*Weekly proportions of no- and high-complexity patients (to total stock) are used in the regression of d_j , while the stocks of no- and high-complexity patients are used in the regression of D_j .

Notes: (i) D_j denotes number of high-complexity (net of non-complex) patients transferred. (ii) d_j denotes proportion of high-complexity (net of non-complex) patients transferred. (iii) Robust standard errors are obtained by clustering (by hospital). (iv) Included in the estimation but not shown are year dummies and month dummies.

Table B2: Logistic model, patient level estimation

Dep. var.: Transfer ($T_{ij} = 0$ or 1)	Coeff.	Std. err.	<i>p</i> -val
Age Group (Ref: 45 or below)			
Age 46 to 55	0.3475	0.0377	$p < 0.001$
Age 56 to 65	0.6366	0.0481	$p < 0.001$
Age 66 to 75	0.9742	0.0638	$p < 0.001$
Age 76 to 85	1.2381	0.0752	$p < 0.001$
Age 86 or above	1.2905	0.0816	$p < 0.001$
Male (Ref: Female)	-0.0677	0.0197	$p < 0.001$
Marital Status (Ref: Never married)			
Married or de facto	-0.0758	0.0172	$p < 0.001$
Divorced or separated	-0.0298	0.0140	0.0327
SEIFA index decile (Ref: 1st decile)			
2nd decile	-0.2444	0.1184	0.0389
3rd decile	-0.0233	0.1187	0.8447
4th decile	-0.2206	0.1072	0.0397
5th decile	-0.3428	0.0804	$p < 0.001$
6th decile	-0.4934	0.0892	$p < 0.001$
7th decile	0.0613	0.0852	0.4720
8th decile	0.2191	0.0852	0.0101
9th decile	-0.0302	0.0696	0.6645
10th decile	0.1413	0.0794	0.0752
Capacity utilization (Ref: w_1)			
w_2	0.2430	0.0692	$p < 0.001$
w_3	0.1705	0.0885	0.0539
w_4	0.0490	0.1021	0.6312
Patient complexity (Ref: Non-complex, c_0)			
Low complexity (c_1)	0.1402	0.0591	0.0177
High complexity (c_2)	0.1655	0.0557	0.0030
Capacity utilization \times Complexity (Ref: $w_1 \times c_0$)			
$w_2 \times c_1$	0.0480	0.0591	0.4165
$w_2 \times c_2$	-0.0540	0.0610	0.3763
$w_3 \times c_1$	0.0364	0.0695	0.6003
$w_3 \times c_2$	0.1140	0.0783	0.1455
$w_4 \times c_1$	-0.0226	0.0827	0.7843
$w_4 \times c_2$	0.1581	0.1018	0.1203
Private patient (Ref: Public patient)	-0.4518	0.0802	$p < 0.001$
Capacity utilization \times Private patient (Ref: $w_1 \times$ Public patient)			
$w_2 \times$ Private patient	-0.1355	0.0852	0.1117
$w_3 \times$ Private patient	-0.0374	0.1135	0.7417
$w_4 \times$ Private patient	0.0322	0.1611	0.8417
Patient complexity \times Private patient (Ref: $c_0 \times$ Public patient)			
$c_1 \times$ Private patient	0.3011	0.0669	$p < 0.001$
$c_2 \times$ Private patient	0.5132	0.0801	$p < 0.001$
Capacity utilization \times Complexity \times Private patient (Ref: $w_1 \times c_0 \times$ Public)			
$w_2 \times c_1 \times$ Private patient	-0.0829	0.0778	0.2864
$w_2 \times c_2 \times$ Private patient	-0.0144	0.0874	0.8689
$w_3 \times c_1 \times$ Private patient	-0.0498	0.0970	0.6076
$w_3 \times c_2 \times$ Private patient	-0.0825	0.1131	0.4658
$w_4 \times c_1 \times$ Private patient	-0.0054	0.1440	0.9703

... continued on next page

Table B2 . . . continued from previous page

Dep. var.: Transfer ($T_{ij} = 0$ or 1)	Coeff.	Std. err.	<i>p</i> -val
$w_4 \times c_2 \times$ Private patient	-0.0704	0.1727	0.6833
Weekly patient stock, all patient	-0.0633	0.1080	0.5579
Weekly no complexity patient stock	0.3076	0.0885	$p < 0.001$
Weekly high complexity patient stock	-0.1104	0.0665	0.0967
Admitted via ED (Ref: Not admitted via ED)	1.0646	0.0550	$p < 0.001$
No. obs		12,072,395	

Note: (i) Robust standard errors are obtained by clustering (by hospital \times MDC). (ii) Included in the estimation but not shown are year dummies, month dummies, day dummies, hospital dummies and MDC dummies.

Appendix S (Supplementary Material)

S1 A two-agent model of a public hospital with mixed public and private patients

Suppose the hospital has two managers, v and b . Manager v is in charge of private patients, and maximizes profit by deciding how many hard private patients to transfer ($n_t^{h(v)}$) in each period. Manager b is in charge of public patients; she maximizes social welfare by deciding how many hard public patients to transfer ($n_t^{h(b)}$) in each period. The hospital has an overall capacity K . There are T periods.

Denote the stock of easy private (public) and hard private (public) patients by $x_t^{e(v)}$ ($x_t^{e(b)}$) and $x_t^{h(v)}$ ($x_t^{h(b)}$), respectively. The hospital's capacity constraint and the laws of motion are:

$$x_t^{e(v)} + x_t^{e(b)} + x_t^{h(v)} + x_t^{h(b)} \leq K \quad (\text{S16})$$

$$x_{t+1}^{e(v)} = (1 - \theta^e)x_t^{e(v)} + \tilde{z}_{t+1}^{e(v)} \quad (\text{S17})$$

$$x_{t+1}^{e(b)} = (1 - \theta^e)x_t^{e(b)} + \tilde{z}_{t+1}^{e(b)} \quad (\text{S18})$$

$$x_{t+1}^{h(v)} = (1 - \theta^h)x_t^{h(v)} + \tilde{z}_{t+1}^{h(v)} - n_t^{h(v)} \quad (\text{S19})$$

$$x_{t+1}^{h(b)} = (1 - \theta^h)x_t^{h(b)} + \tilde{z}_{t+1}^{h(b)} - n_t^{h(b)}, \quad (\text{S20})$$

where $\theta^e > \theta^h$ (see Assumption A1). Patient arrivals, which are potentially rationed, are denoted by variables $\tilde{z}_{t+1}^{(\cdot)}$ (discussed below).¹⁹ Manager v 's objective is to maximize profits. Her payoff in period t is:

$$\pi^e \theta^e x_t^{e(v)} + \pi^h \theta^h x_t^{h(v)} - C(n_t^{h(v)}), \quad (\text{S21})$$

where $\pi^e > \pi^h$ (see Assumption A2). Manager b 's objective is to maximize welfare. Her payoff in period t is:

$$\pi^b (\theta^e x_t^{e(b)} + \theta^h x_t^{h(b)}) - C(n_t^{h(b)}), \quad (\text{S22})$$

where $\pi^b > 0$. The cost function $C(\cdot)$ satisfies Assumption A3.

We consider the Markov strategies of each manager, $\kappa_t^{(v)}, \kappa_t^{(b)} : \chi \rightarrow D$, which are mappings from the state space χ to the choice set D . For each manager, her continuation strategy profile in period t is defined as her set of strategies from periods t to T . Their profiles are denoted by $\kappa_t^{(v)+} := \{\kappa_a^{(v)}\}_{a=t}^T$ and $\kappa_t^{(b)+} := \{\kappa_a^{(b)}\}_{a=t}^T$, respectively. Denote the state variables by $\mathbf{x}_t := (x_t^{e(v)}, x_t^{e(b)}, x_t^{h(v)}, x_t^{h(b)})$ and the raw patient arrivals by $\mathbf{z}_t := (z_t^{e(v)}, z_t^{e(b)}, z_t^{h(v)}, z_t^{h(b)})$.

¹⁹As in Section 2, we assume proportional rationing when the capacity limit is reached.

Given manager b 's continuation strategy profile, manager v 's value function is:

$$V_t^{(v)}(\mathbf{x}_t; \kappa_t^{(b)+}) = \max_{n_t^{h(v)}} \left(\pi^e \theta^e x_t^{e(v)} + \pi^h \theta^h x_t^{h(v)} - C(n_t^{h(v)}) + \beta \iiint_{\mathbf{z}_{t+1}} V_{t+1}^{(v)}(\mathbf{x}_{t+1}; \kappa_{t+1}^{(b)+}) dF(\mathbf{z}_{t+1}) \right). \quad (\text{S23})$$

Given manager v 's continuation strategy profile, manager b 's value function is:

$$V_t^{(b)}(\mathbf{x}_t; \kappa_t^{(v)+}) = \max_{n_t^{h(b)}} \left(\pi^b (\theta^e x_t^{e(b)} + \theta^h x_t^{h(b)}) - C(n_t^{h(b)}) + \beta \iiint_{\mathbf{z}_{t+1}} V_{t+1}^{(b)}(\mathbf{x}_{t+1}; \kappa_{t+1}^{(v)+}) dF(\mathbf{z}_{t+1}) \right). \quad (\text{S24})$$

We assume that the total number of patient arrivals $z_t := z_t^{e(v)} + z_t^{e(b)} + z_t^{h(v)} + z_t^{h(b)}$ satisfies Assumption A4, the number of easy and hard patient arrivals satisfy Assumption A5, and half of the easy and hard patient arrivals are private patients, i.e., $\frac{z_t^{e(v)}}{z_t^{e(v)} + z_t^{e(b)}} = \frac{z_t^{h(v)}}{z_t^{h(v)} + z_t^{h(b)}} = \frac{1}{2}$. Then, the only shock that needs to be integrated out is z .

The solution(s) of the model $(\kappa_1^{*(v)+}, \kappa_1^{*(b)+})$ represents a Markov perfect equilibrium, which is generally difficult to compute in a similar class of models. However, the following key features of the model greatly simplify the solution: (1) the discharge rate (θ) between private and public patients is the same given patient complexity, and (2) transfers made by either party are “perfectly substitutable” in the sense that they have the same effect on the aggregate stocks of hard and easy patients next period. Because the capacity limit K is common, the expected future benefit of transferring patients now depends only on these aggregate stocks next period and how close they are to K . As a result, keeping the other party's strategy profile fixed, a manager's dynamic optimization problem is no different than that given in Section 2. This allows us to use backward recursion to compute the unique Markov perfect equilibrium of the model.²⁰

We now discuss the intuition of the solution, which follows from mathematical induction. Consider the problem in period $t < T$, given the optimal continuation strategy profiles in $t + 1$ $(\kappa_{t+1}^{*(v)+}, \kappa_{t+1}^{*(b)+})$. Denote the aggregate stock of hard patients by $x_t^h := x_t^{h(v)} + x_t^{h(b)}$. Denote the best response functions of manager v and b by $n_t^{*h(v)}(n_t^{h(b)}; x_t^h, \cdot)$ and $n_t^{*h(b)}(n_t^{h(v)}; x_t^h, \cdot)$, respectively. To derive these functions, we proceed with the following arguments: (#1) The optimal transfer by manager v , conditional on *no* transfer by manager b , is given by $n_t^{*h(v)}(0; x_t^h, \cdot)$. As in Propositions 1 and 2, it is strictly increasing in x_t^h beyond a threshold $\underline{x}_t^{h(v)}$: $n_t^{*h(v)}(0; x_t^h, \cdot) = 0$ for $x_t^h < \underline{x}_t^{h(v)}$ and $1 > \frac{dn_t^{*h(v)}(0; x_t^h, \cdot)}{dx_t^h} > 0$ for $x_t^h \geq \underline{x}_t^{h(v)}$. (#2) Due to the property of perfect substitutability discussed above, we have $n_t^{*h(v)}(q; x_t^h, \cdot) = n_t^{*h(v)}(0; x_t^h - q, \cdot)$ for $q \geq 0$. (#3) Similarly, for manager b , $n_t^{*h(b)}(0; x_t^h, \cdot) = 0$ for $x_t^h < \underline{x}_t^{h(b)}$ and $1 > \frac{dn_t^{*h(b)}(0; x_t^h, \cdot)}{dx_t^h} > 0$

²⁰Note that in period T , no one will transfer patients under any state. The value functions are computed accordingly and carried over to the decision problems in $T - 1$, and so on.

for $x_t^h \geq \underline{x}_t^{h(b)}$. We have $n_t^{*h(b)}(q; x_t^h, \cdot) = n_t^{*h(b)}(0; x_t^h - q, \cdot)$ for $q \geq 0$. (#4) As in Corollary 1, the thresholds and the slopes differ: $\underline{x}_t^{h(v)} < \underline{x}_t^{h(b)}$ and $\frac{dn_t^{*h(v)}(0; x_t^h, \cdot)}{dx_t^h} \geq \frac{dn_t^{*h(b)}(0; x_t^h, \cdot)}{dx_t^h}$. Thus, the best response functions differ between both managers.

Figure S1 illustrates the solution. At stock level x_t^h , line AB'' is the best response function of manager v , while line $A''B$ is the best response function of manager b . Point C represents the optimal action pair $(n_t^{*h(b)}, n_t^{*h(v)})$ at stock level x_t^h . When the stock is increased to $x_t^h + 1$, the best response functions are shifted outward to line $A'B'''$ and line $A'''B'$, respectively. The new solution is given by point C' .

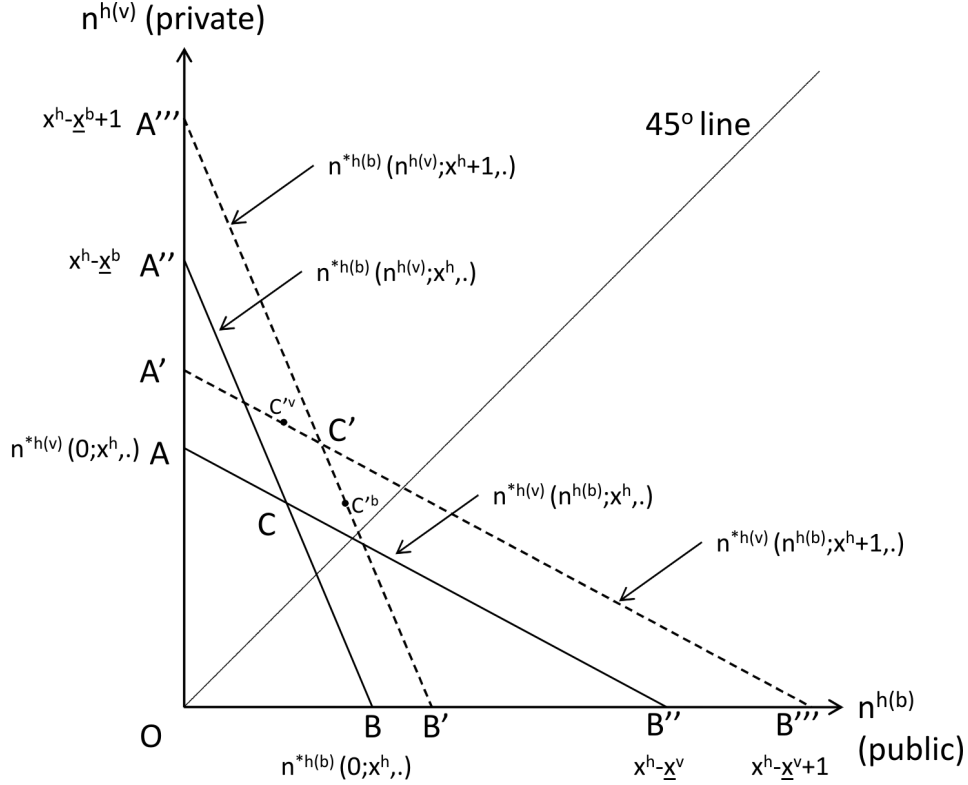


Figure S1: Illustration of Markov perfect solution

The solution C satisfies the following properties:

1. $n_t^{*h(v)} > n_t^{*h(b)}$: Manager v transfers more hard patients than manager b . This result is analogous to C1.1 in Corollary 1, which is based on the single-agent model.
2. When the stock x_t^h increases, both managers increase their transfers of hard patients. This result is analogous to C1.2 in Corollary 1.
3. When the stock x_t^h increases, due to strategic interactions, each manager has an incentive to free-ride on the increase in transfers by the other manager. Manager

v would have increased transfers by CC'^v if she took manager b 's decision as fixed. Manager b would have increased transfers by CC'^b if she took manager v 's decision as fixed.

Proof of results. Assume that the best response functions are in linear form – AB'' : $n^{h(v)} = a_1 - b_1 n^{h(b)}$; $A''B$: $n^{h(v)} = a_2 - b_2 n^{h(b)}$, where $a_2 > a_1$. By arguments #1 and #3 above, $0 < b_1, \frac{1}{b_2} < 1$ (i.e., $b_2 > 1 > b_1 > 0$). By argument #4, $\frac{1}{b_2} \leq b_1$ (i.e., $b_1 b_2 \geq 1$).

Proof of property 1: if $OB = OA$ and $OB'' = OA''$, then by symmetry, point C will lie on the 45 degree line. However, by argument #4, $OB < OA$ because $n^{*h(b)}(0; x^h, \cdot) < n_t^{*h(v)}(0; x^h, \cdot)$ and $OA'' < OB''$ because $x^h - \underline{x}^{h(b)} < x^h - \underline{x}^{h(v)}$.

Therefore, the intersection C must lie on the left of the 45 degree line.

Proof of property 2: By simple arithmetic, point C gives $n^{*h(b)} = \frac{a_2 - a_1}{b_2 - b_1}$ and $n^{*h(v)} = \frac{a_1 b_2 - a_2 b_1}{b_2 - b_1}$. When the stock increases from x_t^h to $x_t^h + 1$, the best response functions become $A'B'''$: $n^{h(v)} = a_1 + b_1 - b_1 n^{h(b)}$; $A'''B'$: $n^{h(v)} = a_2 + 1 - b_2 n^{h(b)}$ by arguments #2 and #3. Point C' gives $n^{**h(b)} = \frac{a_2 - a_1}{b_2 - b_1} + \frac{1 - b_1}{b_2 - b_1} = n^{*h(b)} + \frac{1 - b_1}{b_2 - b_1} > n^{*h(b)}$. In addition, $n^{**h(v)} = a_2 + 1 - b_2 n^{**h(b)} = n^{*h(v)} + 1 - \frac{1 - b_1}{b_2 - b_1} = n^{*h(v)} + \frac{b_2 - 1}{b_2 - b_1} > n^{*h(v)}$.

Proof of property 3: this is evident from the positions of C , C' , C'^v and C'^b .

S2 Estimation results

This section contain additional estimation results of models discussed in the text. Table S1 shows the coefficients of hospital-level fixed effects estimation in which all terms involving time invariant covariates are omitted, including any such interaction terms. Table S2 presents the coefficients of patient level logistic estimation where the sample is restricted to patients in public hospitals that admit both public and private patients. Hospital fixed effects are accounted for using hospital dummies.

Table S1: Fixed effects models, hospital level estimation (no interactions)

Dep. var:	D_j			d_j		
	Coeff.	Std. err.	p -val	Coeff.	Std. err.	p -val
Capacity utilization (Ref: w_1)						
w_2	0.4864	0.0798	$p < 0.001$	0.0017	0.0006	0.0070
w_3	0.7389	0.1521	$p < 0.001$	0.0026	0.0008	0.0011
w_4	0.8881	0.2054	$p < 0.001$	0.0030	0.0013	0.0221
Weekly patient stock, all patients	-3.5905	0.5556	$p < 0.001$	-0.0023	0.0012	0.0584
Weekly non-complex patient stock*	-0.4638	0.2747	0.0932	-0.0211	0.0028	$p < 0.001$
Weekly high complexity patient stock*	0.5059	0.1741	0.0041	0.0408	0.0068	$p < 0.001$
Proportion rare DRGs	-0.2507	0.1358	0.0666	-0.0259	0.0104	0.0132
Proportion rural patients	0.1437	0.0949	0.1316	0.0007	0.0064	0.9086
Proportion same-day episodes	0.9295	0.1789	$p < 0.001$	0.0207	0.0031	$p < 0.001$
Constant	0.5611	0.1052	$p < 0.001$	0.0023	0.0012	0.9048
No. of observations		95,131			95,131	
No. of hospitals		180			180	

*Weekly proportions of no- and high-complexity patients (to total stock) are used in the regression of d_j , while the stocks of no- and high-complexity patients are used in the regression of D_j .

Notes: (i) D_j denotes number of high-complexity (net of non-complex) patients transferred. (ii) d_j denotes proportion of high-complexity (net of non-complex) patients transferred. (iii) Robust standard errors are obtained by clustering (by hospital). (iv) Included in the estimation but not shown are year dummies and month dummies.

Table S2: Logistic model, patient level sub-sample estimation

Dep. var.: Transfer ($T_{ij} = 0$ or 1)	Coeff.	Std. err.	p -val
Age Group (Ref: 45 or below)			
Age 46 to 55	0.3602	0.0364	$p < 0.001$
Age 56 to 65	0.6257	0.0467	$p < 0.001$
Age 66 to 75	0.8891	0.0634	$p < 0.001$
Age 76 to 85	1.1283	0.0791	$p < 0.001$
Age 86 or above	1.2245	0.0884	$p < 0.001$
Male (Ref: Female)	-0.0210	0.0223	0.3466
Marital Status (Ref: Never married)			
Married or de facto	-0.0279	0.0154	0.0692
Divorced or separated	-0.0291	0.0151	0.0538
SEIFA index decile (Ref: 1st decile)			
2nd decile	0.1273	0.1338	0.3413
3rd decile	0.2044	0.0549	$p < 0.001$
4th decile	-0.1659	0.0850	0.0510

... continued on next page

Table S2 ... continued from previous page

Dep. var.: Transfer ($T_{ij} = 0$ or 1)	Coeff.	Std. err.	<i>p</i> -val
5th decile	-0.0772	0.0539	0.1521
6th decile	0.0192	0.0509	0.7053
7th decile	0.0366	0.0473	0.4389
8th decile	0.1959	0.0530	$p < 0.001$
9th decile	0.0838	0.0406	0.0392
10th decile	0.2772	0.0522	$p < 0.001$
Capacity utilization (Ref: w_1)			
w_2	-0.0527	0.0230	0.0218
w_3	-0.0372	0.0318	0.2418
w_4	-0.0091	0.0401	0.8203
Patient complexity (Ref: Non-complex, c_0)			
Low complexity (c_1)	0.2145	0.0555	$p < 0.001$
High complexity (c_2)	0.1967	0.0475	$p < 0.001$
Capacity utilization \times Complexity (Ref: $w_1 \times c_0$)			
$w_2 \times c_1$	0.0592	0.0519	0.2539
$w_2 \times c_2$	0.0422	0.0503	0.4014
$w_3 \times c_1$	0.0090	0.0642	0.8887
$w_3 \times c_2$	0.1523	0.0707	0.0313
$w_4 \times c_1$	-0.1002	0.0730	0.1700
$w_4 \times c_2$	0.1334	0.0969	0.1685
Private patient (Ref: Public patient)			
-0.0732	0.0606	0.2272	
Capacity utilization \times Private patient (Ref: $w_1 \times$ Public patient)			
$w_2 \times$ Private patient	0.2474	0.1449	0.0877
$w_3 \times$ Private patient	0.2031	0.0576	$p < 0.001$
$w_4 \times$ Private patient	0.1760	0.1021	0.0847
Patient complexity \times Private patient (Ref: $c_0 \times$ Public patient)			
$c_1 \times$ Private patient	0.0869	0.0717	0.2256
$c_2 \times$ Private patient	0.3470	0.0778	$p < 0.001$
Capacity utilization \times Complexity \times Private patient (Ref: $w_1 \times c_0 \times$ Public)			
$w_2 \times c_1 \times$ Private patient	-0.1818	0.1593	0.2538
$w_2 \times c_2 \times$ Private patient	-0.2334	0.1644	0.1557
$w_3 \times c_1 \times$ Private patient	-0.1030	0.0778	0.1856
$w_3 \times c_2 \times$ Private patient	-0.1449	0.0872	0.0967
$w_4 \times c_1 \times$ Private patient	-0.0259	0.0951	0.7856
$w_4 \times c_2 \times$ Private patient	-0.1362	0.1163	0.2415
Weekly patient stock, all patients	0.1540	0.0933	0.0989
Weekly no complexity patient stock	-0.1362	0.0335	$p < 0.001$
Weekly high complexity patient stock	0.0871	0.0351	0.0132
Admitted via ED (Ref: Not admitted via ED)	1.1352	0.0686	$p < 0.001$
No. obs		7,730,057	

Note: (i) Robust standard errors are obtained by clustering (by hospital \times MDC). (ii) Included in the estimation but not shown are year dummies, month dummies, day dummies, hospital dummies and MDC dummies.

S3 Hospital level analysis with eight capacity utilization categories

In estimating the base models, we define a discrete capacity utilization variable with four categories of utilization levels. Here we divide capacity utilization into eight discrete categories to verify the results are not dependent on the specific way of discretizing capacity utilization in our data.

Define w_{j1}, \dots, w_{j8} as binary variables, where

$$\begin{aligned}w_{j1} &= 1 \text{ if } V_j \leq 0.60, \\w_{j2} &= 1 \text{ if } V_j \in (0.60, 0.70], \\w_{j3} &= 1 \text{ if } V_j \in (0.70, 0.75], \\w_{j4} &= 1 \text{ if } V_j \in (0.75, 0.80], \\w_{j5} &= 1 \text{ if } V_j \in (0.80, 0.825], \\w_{j6} &= 1 \text{ if } V_j \in (0.825, 0.850], \\w_{j7} &= 1 \text{ if } V_j \in (0.850, 0.875], \\w_{j8} &= 1 \text{ if } V_j > 0.875;\end{aligned}$$

otherwise $w_{jk} = 0$, $k = 1, \dots, 8$. We repeat both hospital and patient level analyses using the capacity utilization variable with eight discrete levels.

Table S3 shows the coefficient estimates and robust standard errors of the estimation where the dependent variables are the difference in the number and proportion of high complexity and non-complexity patients. The marginal effect estimates are reported in Table S4. The results closely mimic those in the main body of the paper, thus our findings are not dependent on specific values of the capacity utilization rates chosen.

Table S3: Hospital level CRE analysis, eight capacity utilization categories

Dep. var:	D_j			d_j		
	Coeff.	Std. err.	p -val	Coeff.	Std. err.	p -val
Capacity utilization (Ref: w_1)						
w_2	0.2733	0.0728	$p < 0.001$	0.0022	0.0011	0.0419
w_3	0.5871	0.1223	$p < 0.001$	0.0031	0.0017	0.0741
w_4	0.6285	0.2096	0.0027	0.0033	0.0018	0.0618
w_5	0.9162	0.2470	$p < 0.001$	0.0042	0.0017	0.0134
w_6	0.8480	0.2404	$p < 0.001$	0.0042	0.0017	0.0136
w_7	1.0148	0.2804	$p < 0.001$	0.0048	0.0020	0.0129
w_8	1.0478	0.2787	$p < 0.001$	0.0038	0.0024	0.1191
Private Hospital (Ref: Public Hospital)	0.9587	0.4165	0.0214	0.0065	0.0026	0.0128
Private \times Capacity utilization (Ref: Public $\times w_1$)						
Private $\times w_2$	0.1777	0.1065	0.0953	-0.0004	0.0012	0.7060
Private $\times w_3$	0.1520	0.1641	0.3543	-0.0006	0.0019	0.7432
Private $\times w_4$	0.2975	0.2793	0.2867	-0.0002	0.0020	0.9279
Private $\times w_5$	0.1389	0.3197	0.6640	-0.0006	0.0020	0.7656
Private $\times w_6$	0.2894	0.3090	0.3491	-0.0002	0.0020	0.9238
Private $\times w_7$	0.0761	0.3350	0.8202	-0.0011	0.0022	0.5998
Private $\times w_8$	0.3152	0.3387	0.3521	0.0017	0.0032	0.5955
Weekly patient stock, all patients	-3.6702	0.5503	$p < 0.001$	-0.0028	0.0012	0.0241
Weekly no complexity patient stock*	-0.4677	0.2663	0.0790	-0.0211	0.0027	$p < 0.001$
Weekly high complexity patient stock*	0.5079	0.1757	0.0038	0.0409	0.0067	$p < 0.001$
Proportion rare DRGs	-0.2362	0.1303	0.0698	-0.0258	0.0103	0.0124
Proportion rural patients	0.1426	0.0925	0.1233	0.0007	0.0064	0.9130
Proportion same-day episodes	0.8770	0.1771	$p < 0.001$	0.0204	0.0031	$p < 0.001$
Group mean of Capacity utilization	1.0246	1.7322	0.5542	-0.0264	0.0118	0.0247
Average weekly patient stock, all patients	6.8589	2.7216	0.0117	0.0011	0.0015	0.4461
Average weekly no complexity patient stock*	-5.2857	2.4190	0.0289	-0.0025	0.0238	0.9179
Average weekly high complexity patient stock*	-0.9703	0.6958	0.1631	-0.0331	0.0342	0.3336
Average proportion rare DRGs	17.7217	18.1136	0.3279	0.2876	0.1224	0.0188
Average proportion rural patients	0.4366	0.2461	0.0760	0.0043	0.0065	0.5087
Average proportion same-day episodes	0.5129	0.8837	0.5616	0.0112	0.0135	0.4085
Constant	-0.8858	0.4359	0.0422	0.0025	0.0170	0.8809
No. of observations		95,131			95,131	
No. of hospitals		180			180	

*Weekly proportions of no- and high-complexity patients (to total stock) are used in the regression of d_j , while the stocks of no- and high-complexity patients are used in the regression of D_j .

Notes: (i) D_j denotes number of high-complexity (net of non-complex) patients transferred. (ii) d_j denotes proportion of high-complexity (net of non-complex) patients transferred. (iii) Robust standard errors are obtained by clustering (by hospital). (iv) Included in the estimation but not shown are year dummies and month dummies.

Table S4: Marginal effects, Hospital level CRE analysis, eight capacity

	D_j		d_j	
	MEff.	p -val	MEff.	p -val
Capacity utilization level				
w_2 (v. w_1)	0.3343 (0.0559)	$p < 0.001$	0.0020 (0.0007)	0.006
w_3 (v. w_2)	0.3050 (0.0571)	$p < 0.001$	0.0009 (0.0006)	0.144
w_4 (v. w_3)	0.0913 (0.0846)	0.280	0.0004 (0.0003)	0.242
w_5 (v. w_4)	0.2333 (0.0899)	0.010	0.0008 (0.0004)	0.070
w_6 (v. w_5)	-0.0165 (0.0932)	0.859	0.0002 (0.0004)	0.664
w_7 (v. w_6)	0.0935 (0.1176)	0.427	0.0003 (0.0005)	0.569
w_8 (v. w_7)	0.1151 (0.1351)	0.394	-0.0001 (0.0010)	0.902
Private (v. public) Hospital	1.0622 (0.4773)	0.026	0.0064 (0.0024)	0.008
No. of observations	95,131		95,131	
No. of hospitals	180		180	

Notes: (i) D_j denotes number of high-complexity (net of no-complexity) patients transferred. (ii) d_j denotes proportion of high-complexity (net of no-complexity) patients transferred. (iii) Robust standard errors are obtained by clustering (by hospital). (iv) Figures in parentheses are standard errors obtained by delta method.