

NBER WORKING PAPER SERIES

FIRST DO NO HARM? DOCTOR DECISION MAKING AND PATIENT OUTCOMES

Janet Currie
W. Bentley MacLeod
Kate Musen

Working Paper 32788
<http://www.nber.org/papers/w32788>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
August 2024

We would like to thank Jonathan Gruber, Amanda Kowalski, David Romer, and 4 anonymous referees for helpful comments. Kate Musen gratefully acknowledges support from the National Science Foundation (Grant Number DGE2036197). The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2024 by Janet Currie, W. Bentley MacLeod, and Kate Musen. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

First Do No Harm? Doctor Decision Making and Patient Outcomes
Janet Currie, W. Bentley MacLeod, and Kate Musen
NBER Working Paper No. 32788
August 2024
JEL No. I11,I12

ABSTRACT

Doctors facing similar patients often make different treatment choices. These decisions can have important effects on patient health and health care spending. This paper seeks to organize the recent economics literature on physician decision making using a simple model that incorporates doctor diagnostic and procedural skills, differences in beliefs and patient populations, and incentives. Economic considerations that affect the quality of decision making include training, experience, peer effects, financial incentives and time constraints. We also consider interventions aimed at improving decision making including provision of informational, heuristics and guidelines, and the use of technologies including electronic medical records and algorithmic decision tools. Our review suggests that we have learned a great deal about specific factors that influence doctor decision making but that our knowledge of how to apply that knowledge to improve health care is still quite limited.

Janet Currie
Department of Economics
Center for Health and Wellbeing
185A Julis Romo Rabinowitz Building
Princeton University
Princeton, NJ 08544
and NBER
jcurrie@princeton.edu

Kate Musen
Department of Economics
Columbia University
420 W 118th St.
New York, NY 10027
khm2128@columbia.edu

W. Bentley MacLeod
Department of Economics
Princeton University
214 Robertson Hall, 20 Prospect St.
Princeton, NJ 08544-1013
and NBER
wbmacleod@wbmacleod.net

First Do No Harm? Doctor Decision Making and Patient Outcomes*

Janet Currie[†], W. Bentley MacLeod[‡] and Kate Musen[§]

July 27, 2024

Abstract

Doctors facing similar patients often make different treatment choices. These decisions can have important effects on patient health and health care spending. This paper seeks to organize the recent economics literature on physician decision making using a simple model that incorporates doctor diagnostic and procedural skills, differences in beliefs and patient populations, and incentives. Economic considerations that affect the quality of decision making include training, experience, peer effects, financial incentives and time constraints. We also consider interventions aimed at improving decision making including provision of informational, heuristics and guidelines, and the use of technologies including electronic medical records and algorithmic decision tools. Our review suggests that we have learned a great deal about specific factors that influence doctor decision making but that our knowledge of how to apply that knowledge to improve health care is still quite limited.

1 Introduction

Doctors facing similar patients often make different treatment choices, and these can have large consequences for patient outcomes and health care spending. A rapidly growing literature focuses on understanding the sources of this variation. We are all health care consumers who want our doctors to give us good advice and make sound choices, so the question of what drives doctor decision making is of intrinsic interest. At a more macro level, health care accounts for almost 20% of U.S. GDP and many observers feel that much of that spending is misdirected, wasted, or even harmful (Chandra and Skinner (2012), Cutler (2014)). Badinski, Finkelstein, Gentzkow, and Hull (2023) use data from Medicare patients and physicians who move to show that roughly a third of regional differences in healthcare utilization in Americans over 65 is explained by differences in the average physician treatment intensity. A third reason to study the doctor decision making is that doctors share many features with other experts such as lawyers, top managers, or even professors, so some insights from the literature on doctor decision making to understanding other types of experts.

This paper seeks to organize the recent literature (since 2010) on physician decision making by looking at it through the lens of a model that has several key elements. First, doctors care about patients, but they

*We would like to thank Jonathan Gruber, Amanda Kowalski, David Romer, and 4 anonymous referees for helpful comments.

[†]Princeton University and NBER

[‡]Princeton University, Columbia University and NBER

[§]Columbia University

are influenced their beliefs about appropriate care, time constraints, and profit motives, all of which can vary across doctors. Hence doctors are imperfect agents from the point of view of patients since they care about other considerations in addition to patient utility. Second, doctors' skill levels vary. We distinguish between skill involved in deciding what to do (diagnosis), and procedural skill, defined as skilled execution of a given decision. Third, patients care about medical outcomes, as well as other factors including quality of life and out-of-pocket costs. Both doctors and patients may have strong beliefs about treatments (e.g. doctors may have been trained to think that a procedure is necessary, and patients may believe, for example, that vaccines are harmful). All of these factors mean that patients with identical conditions can end up being treated differently.

Table 1 describes a number of studies demonstrating that physicians often treat similar patients so differently that they can be said to have distinct "practice styles" (see Table 1). For example, Berndt, Gibbons, Kolotilin, and Taub (2015) study concentration in the way that doctors prescribe anti-psychotics and show that typically, two thirds of a doctor's prescriptions are for the same drug, and that crucially, doctors have different favorite drugs. Cutler, Skinner, Stern, and Wennberg (2019) use Medicare claims data to identify "cowboys" who recommend aggressive treatments that go beyond clinical guidelines and "comforters" who recommend palliative care for severely ill patients. Focusing on heart attack (i.e. acute myocardial infarction or AMI) patients, they find that a one standard deviation increase in the share of doctors who are cowboys leads to a 13% increase in annual spending, whereas a one standard deviation increase in the share of comforters leads to a small decrease in annual spending. Notably, neither share is associated with changes in survival probabilities. Fadlon and van Parys (2020) look at patients who switch providers after their primary care physician retired or moved away. They find that changing to a provider who spends more on primary care increases spending on primary care, which they interpret as evidence of distinct practice styles. Although their design looks at exogenous separations between patients and providers, patients choose their new provider, which may lead to sorting of patients with different preferences. Ahammer and Schober (2020) show similar results in the Austrian context. Marquardt (2021) examines variation in diagnoses of ADHD and finds that a one standard deviation increase in physician "intensity" (measured as the intercept in a doctor-specific regression) increases the probability a patient is diagnosed by 22.45 percent.

The model outlined in the next section builds on work in four of the papers shown in Table 1, Abaluck et al. (2016), Currie, MacLeod, and Van Parys (2016) Currie and MacLeod (2017), and Chan et al. (2022) to provide a framework to think about alternative reasons for the observed variation in physician decision making, and about interventions that have been suggested in an effort to improve outcomes. The literature on health disparities, discussed in section 3, shows that an individual physician may vary treatment based upon characteristics of the patient that are unrelated to their health status, illustrating the role that idiosyncratic physician preferences play in treatment decisions. Economic considerations that affect the quality of decision making include financial incentives, experience, training, peer effects, and time constraints, are discussed in section 4. Another branch of the literature asks whether decision making can be improved through informational interventions, guidelines, or the use of technology including algorithmic decision tools. These studies are discussed in section 5 of the paper.

Understandably, most of the studies we review focus on the role of a single explanatory factor, although this often requires strong assumptions about the other factors. Our first objective is to make these assumptions more explicit. Second, we try to connect aspects of the decision process that are typically studied in

isolation, such as the relationship between doctor skill and thresholds for choosing aggressive procedures. Third, we offer an empirical assessment of what we have learned to date about doctor decision making, and suggestions for further research.

2 A simple model of physician behavior and patient outcomes

This section sketches a simple model of physician decision making. A more formal model and proofs are relegated to the Appendix. Consider a physician’s choice of two treatments, a non-intensive treatment (NI) and an intensive treatment (I). For example, Chandra and Staiger (2007) consider heart patients where the choice is cardiac catheterization (the intensive procedure) vs. medical (i.e. drug) management. Currie and MacLeod (2017) study childbirth, where a vaginal delivery is the noninvasive procedure and a C-section is the invasive procedure. In Abaluck et al. (2016), the “invasive” (or at least more expensive) procedure is to test a patient for a pulmonary embolism and the alternative is not to test.

Consider patient $i \in \mathcal{N}$ who seeks treatment from doctor $j \in J$. The doctor chooses between a non-intensive or a more intensive treatment, denoted by $t \in \{NI, I\}$. Assume that there is a best medical choice for the patient given by their *unobserved* state $\alpha_i \in \{0, 1\}$. If $\alpha_i = 0$, then the non-intensive treatment is preferred, while $\alpha_i = 1$ implies that the intensive treatment is more appropriate. Let the fraction of patients in \mathcal{N} for whom $\alpha_i = 0$ be given by $p_0 \in (0, 1)$, while a fraction $p_1 = 1 - p_0$ are in state $\alpha_i = 1$. Doctor j cannot perfectly observe the patient’s state, rather the doctor observes a noisy signal:

$$T_{ij} = \alpha_i + \epsilon/\gamma_j, \tag{1}$$

where $\epsilon \sim N(0, 1)$ and γ_j is the diagnostic skill of the doctor. An increase in diagnostic skill implies a more precise assessment of a person’s state. Although diagnostic skill is often ignored by economists, the National Academy of Sciences (Balogh, Miller, and Ball (2015)) notes that diagnostic errors are frequent, affecting 5 percent of American outpatients annually, contributing to between 6 and 17 percent of hospital adverse events, and ultimately leading to 10 percent of patient deaths. Diagnostic errors are also a leading cause of successful medical malpractice cases.

The signal T_{ij} is increasing in α_i so it follows that the optimal diagnostic rule for the treatment $t_{ij} \in \{NI, I\}$ takes the form:

$$t_{ij} = \begin{cases} I, & T_{ij} \geq \tau_j, \\ NI, & T_{ij} < \tau_j, \end{cases}$$

where τ_j is the doctor’s *decision threshold* for deciding when to implement the intensive treatment. Increasing the threshold reduces the probability that the intensive treatment is chosen. The determination of the optimal threshold is discussed in the next subsection.

The quality of diagnosis is measured by the likelihood that a patient is assigned to the correct medical treatment. In this framework there are two measures of performance that correspond to whether patients correctly or incorrectly receive the intensive treatment. The first is the probability that a patient of type $\alpha = 1$ receives the appropriate treatment. The second measure is the probability that a patient type $\alpha = 0$ receives the inappropriate intensive treatment. Since there is uncertainty in the doctor’s mind regarding the

true state, increasing the probability of the type 1 patients getting the intensive treatment will mechanically have the negative consequence of increasing the probability that patients of type 0 get the inappropriate intensive treatment.

This trade off is illustrated in Figure 1 showing a plot of the probability of appropriate versus inappropriate intensive treatment for different levels of diagnostic skill γ_j . As γ_j increases, the frontier moves up and left. The top left corner represents perfect diagnosis - the patient receives the intensive treatment if and only if they are of type $\alpha = 1$. Conversely, as γ_j approaches zero, the frontier approaches the dashed 45 degree line. The decision threshold τ_j defines a point on the diagnostic frontier. As τ_j increases, the doctor has a higher threshold for performing the intensive procedure so the probability of intensive treatment falls. This decline is indicated by a move to the left along a given frontier.¹

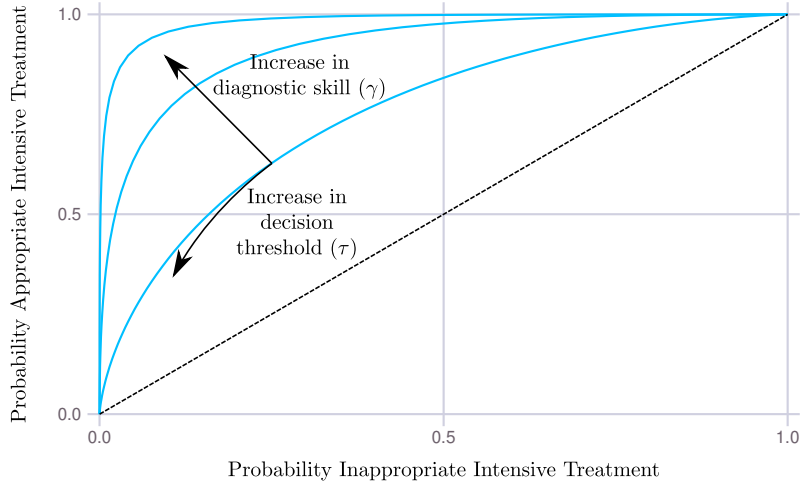


Figure 1: Effect of Diagnostic Performance

In order to determine a doctor's optimal decision threshold, we need to consider the doctor's utility:

$$U_{\alpha t j} = u_{\alpha t j} + \lambda_{t j}, \quad (2)$$

where $u_{\alpha t j}$ is the expected medical benefit to a patient of type $\alpha \in \{0, 1\}$ getting treatment $t \in \{NI, N\}$ from doctor j . The outcome $u_{\alpha t j}$ can differ by doctor, depending on the doctor's *procedural skill*. Additional factors that affect treatment, such as doctor payments for administering the treatment, are captured by $\lambda_{t j}$. Suppose for the moment $\lambda_{t j} = 0$ so that the doctor only cares about the medical benefit to the patient.

If the patient is type $\alpha_i = 0$, then non-intensive treatment is preferred ($u_{0NIj} > u_{0Ij}$), while for type $\alpha_i = 1$ intensive treatment is preferred ($u_{1Ij} > u_{1NIj}$). Let $\Delta_{1Ij} = \{u_{1Ij} - u_{1NIj}\} > 0$ and $\Delta_{0NIj} = \{u_{0NIj} - u_{0Ij}\} > 0$ be the increases in utility for patients getting the appropriate treatment. The doctor's *ex ante* beliefs regarding the appropriate treatment for a patient in this pool of potential patients is given by:

$$p_{1j} = \Pr[\alpha = 1].$$

¹This curve is taken from the machine learning literature where it is called the *receiver-operator* curve, or ROC curve (see Fawcett (2006)). This terminology comes from the use of this curve during World War II to describe the performance of radar systems.

The belief that the probability that $\alpha_i = 0$ is $p_{0j} = 1 - p_{1j}$.

Given this set up, the optimal threshold τ_{ij}^* is:

$$\tau_{ij}^* = \frac{1}{2} + b_{ij}^*/\gamma_j^2, \quad (3)$$

where $b_{ij}^* \equiv (\ln(\Delta_{0NIj}/\Delta_{1Ij}) + \ln(p_{0j}/p_{1j}))$ is the *optimal threshold shifter*.²

Equation (3) shows that the optimal decision threshold depends on diagnostic skill, γ_j , the relative effectiveness of non-intensive and intensive treatments for the two types of patients, $\Delta_{0NIj}/\Delta_{1Ij}$, and the doctor's beliefs about the relative proportion of patient types, p_{0j}/p_{1j} , in the population. For example, a doctor who believes that most patients need non-intensive treatment will adopt a higher decision threshold for use of the intensive treatment than a doctor who believes the reverse. If the relative benefit from intensive treatment is higher, doctors will adopt a *lower* decision threshold resulting in more use of the intensive procedure. Greater diagnostic skill makes the doctor's beliefs about the distribution of patient types in the population and the expected relative benefits of the procedures less important, because in the limit a doctor with perfect diagnostic skill would choose the best procedure for the patient. In effect, as diagnostic skill falls, physicians choose the treatment that they believe is *ex-ante* optimal for most patients, and more patients receive the same treatment.³

These results are illustrated in Figure (2). It shows outcomes for two doctor types with different practice styles:

- A cautious doctor (C), or “comforter” in the Cutler et al. (2019) terminology, is one who is more likely to give a non-intensive treatment. In this case the shift parameter is $b_{iC} = \log\left(\frac{\Delta_{0NIC}}{\Delta_{1IC}} \times \frac{p_{0C}}{p_{1C}}\right) > 0$, and the optimal decision threshold is at the point where the slope, $\frac{\Delta_{0NIC}}{\Delta_{1IC}} \times \frac{p_{0C}}{p_{1C}} > 1$, is tangent to the diagnostic frontier. The points τ_{CH}^* , τ_{CM}^* and τ_{CL}^* , correspond to cautious doctors with high, medium and low diagnostic skills respectively.
- An aggressive doctor (A), or “cowboy” in the Cutler et al. (2019) terminology, is one who is more likely to do the intensive treatment. In this case the shift parameter is $b_{iA} = \log\left(\frac{\Delta_{0NIA}}{\Delta_{1IA}} \times \frac{p_{0A}}{p_{1A}}\right) < 0$, and the optimal decision threshold is at the point where the slope, $\frac{\Delta_{0NIC}}{\Delta_{1IC}} \times \frac{p_{0C}}{p_{1C}} < 1$, is tangent to the diagnostic frontier. The points τ_{AH}^* , τ_{AM}^* and τ_{AL}^* correspond to doctors with high, medium and low diagnostic skill respectively.

The figure shows that even if doctors base their decisions on what is best for the patient, it is still the case that *ex ante* beliefs about the probability that the non-intensive treatment is appropriate (p_{0j}/p_{1j}) affect their choices.

At this point, we can re-introduce the pecuniary concerns that are captured with the λ_{tj} term in the doctor's utility function. Let $\delta_j = \lambda_{Ij} - \lambda_{NIj}$ denote the doctor's pecuniary preference for intensive treatment. If $\delta_j \in (\Delta_{0NIj}, -\Delta_{1Ij})$ then doctor j chooses a decision threshold to satisfy:

$$\tau_{ij}^0 = \frac{1}{2} + b_{ij}^0/\gamma_j^2, \quad (4)$$

²See Propositions 1 and 2 in the Appendix.

³See Proposition 3 in the Appendix.

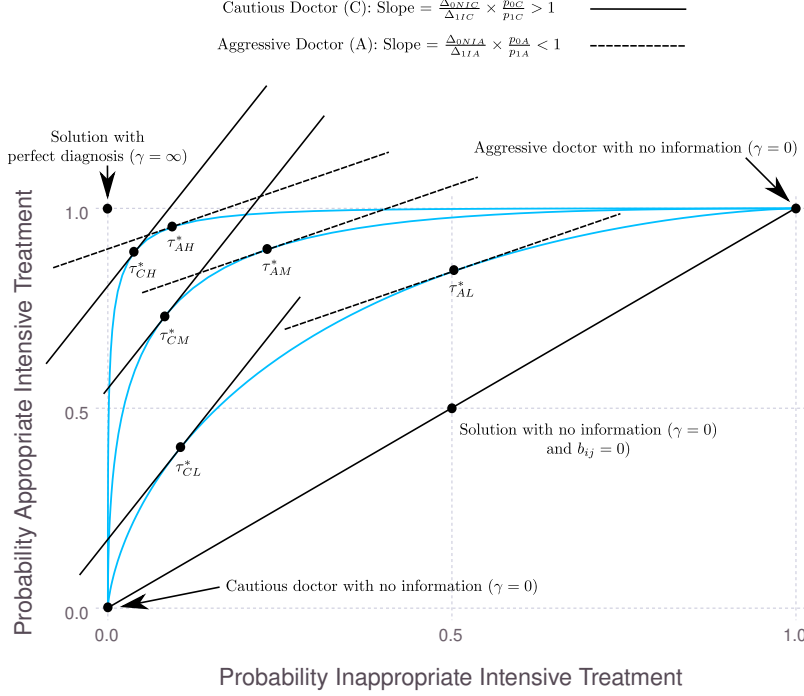


Figure 2: Optimal Diagnostic Rule

where $b_{ij}^0 \equiv (\ln(\Delta_{0NIj} - \delta_j) - \ln(\Delta_{1Ij} + \delta_j) + \ln(p_{0j}/p_{1j}))$ is now the *optimal threshold shifter*. If $\delta_j > \Delta_{0NIj}$ then the doctor always chooses the intensive procedure, while if $\delta_j < \Delta_{1Ij}$, the doctor always chooses the non-intensive procedure. The proof is in the Appendix.

2.1 Identifying Doctor Diagnostic Thresholds, Diagnostic Skill, and Procedural Skill From Data

In this section we consider conditions under which the econometrician can separately identify diagnostic skill, decision thresholds, and procedural skill. Data often includes information about the primary treatment choice ($t_{ij} \in \{NI, I\}$), some measures of patient outcomes following treatment, and some information about patient type from medical records. Let \vec{x}_i be the observable patient characteristics. We illustrate three approaches to solving the identification problem used respectively in Abaluck et al. (2016), Chan et al. (2022), and Currie and MacLeod (2017). Abaluck et al. (2016) focus on the doctor's diagnostic rule, τ , while both Currie and MacLeod (2017) and Chan et al. (2022) seek to measure both the diagnostic rule, τ , and diagnostic skill, γ . In addition, Currie and MacLeod (2017) seeks to measure the doctor's procedural skill.

Abaluck et al. (2016) study doctors treating patients who may have a life-threatening pulmonary embolism (PE). A near definitive diagnosis can be made with a computerized tomography (CT) scan, but doctors may be ordering too many CT scans since CT scans are expensive and expose patients to potentially harmful

radiation.⁴ The doctor forms an estimate of the person’s likelihood of having a PE:

$$\begin{aligned} p_{1i} &= \Pr [\alpha_i = 1 | \vec{x}_i, j], \\ &= \Pr [\vec{x}_i \beta + d_j + \eta_{ij} > 0], \end{aligned}$$

where \vec{x}_i is a vector of observed patient characteristics, d_j is a doctor fixed effect, and η_{ij} reflects unobserved characteristics of the patient. Their data come from Medicare claims, Medicare being the public health insurance program that covers most U.S. elderly. Abaluck et al. (2016) ask whether the doctor’s decision rule varies after controlling for the characteristics of the population they are treating. The statistic used to allocate patients to a CT scan is:

$$\begin{aligned} T_{ij} &= p_{1i} - \tau_j, \\ &= \vec{x}_i \beta + d_j + \eta_{ij} - \tau_j. \end{aligned}$$

Hence, a patient is tested if and only if the expected probability of having PE, p_{1i} , is greater than τ_j . Since the statistic is positively correlated with patient outcomes, self-selection implies:

$$E \{ \alpha_i | T_{ij} > 0 \} > E \{ \alpha_i \} > E \{ \alpha_i | T_{ij} < 0 \},$$

where $E \{ \alpha_i | T_{ij} > 0 \}$ is the probability of a PE in the population of tested individuals. Since some doctors may be treating sicker patient populations than others, the positive test rate, $E \{ \alpha_i | T_{ij} > 0 \}$, can vary even if all the doctors have the same decision threshold.

Abaluck et al. (2016) provide a clever solution to this selection problem. Essentially they rank all the patients in terms of their appropriateness for the procedure using $\vec{x}_i \beta$. Then they infer that doctors whose least appropriate patient (the marginal patient) is sicker, must have higher testing thresholds. A nice thing about their setting is that since a PE is a serious medical condition, people with a missed diagnosis will likely return to the hospital. Hence, they can infer the patient’s true type after observing the treatment decision. With this information they estimate a model of the probability that a patient has a PE given a vector of observable patient characteristics. The β ’s from this regression can be thought of as the true weights as long as there are not too many important omitted variables that the doctor can see but the econometrician cannot see. They compare these β ’s to those obtained from a model in which the dependent variable is whether a test was ordered. A comparison of the two sets of coefficients implies that, on average, doctors are using the wrong weights when deciding whether to order a test. However, this section of the paper assumes that all doctors use the same weights, i.e. that they all have similar diagnostic skills where those skills are approximated by the weights. Hence, by construction, variation in doctor behavior in their model comes only from differences in their thresholds and patient pools.

Chan, Gentzkow, and Yu (2022) explicitly consider the effect of diagnostic skill in a group of radiologists who must decide, on the basis of a chest x-ray, whether a patient has pneumonia. Their goal is to identify both the decision threshold, τ_j , (they call this doctor preferences) and diagnostic skill, γ_j . A valuable feature of their setting is that cases are approximately randomly assigned to radiologists as they arrive at the hospital,

⁴The authors note that the downstream cancer risk from radiation exposure is less of a concern in the elderly population they study.

so on average they all see similar patient pools. Building on work that uses the random assignment of judges to defendants for the determination of bail (Kling (2006), Arnold, Dobbie, and Hull (2022)),⁵ they propose a procedure for estimating PI_j , the probability that a patient correctly receives the intensive treatment, and PNI_j , the probability that a patient incorrectly receives the intensive treatment, for each doctor. They observe that from these measures one can identify both the decision threshold and diagnostic skill. This result follows from the definition of these quantities:

$$\begin{aligned}
 PI(\tau_j, \gamma_j) &\equiv \Pr[T_{ij} \geq \tau_j | \alpha_i = 1], \\
 &= \Pr[1 + \epsilon/\gamma_j \geq \tau_j], \\
 &= F(\gamma_j(1 - \tau_j)),
 \end{aligned} \tag{5}$$

where $F(\cdot)$ is the Normal cumulative probability distribution, and

$$\begin{aligned}
 PNI(\tau_j, \gamma_j) &\equiv \Pr[T_{ij} \geq \tau_j | \alpha = 0] \\
 &= \Pr[\epsilon/\gamma_j \geq \tau_j] \\
 &= F(-\gamma_j\tau_j).
 \end{aligned} \tag{6}$$

Hence, given $PI_j \in (0, 1)$, $PNI_j \in (0, 1)$ and $PI > PNI$ there is a unique solution for $\tau_j \in (-\infty, \infty)$ and $\gamma_j > 0$ solving (5-6). See the Appendix for details.

Like Figures (1-2), Figure (3), taken from Chan, Gentzkow, and Yu (2022), illustrates the relationship between appropriate and inappropriate testing. Each point corresponds to the true positive and the false positive rate of a radiologist. If doctors only varied in terms of their decision thresholds then all the points would lie on the same curve. Similarly, if all the doctors differed only in terms of diagnostic skill, then the points would follow a line such as that connecting the points τ_{AH}^* , τ_{AM}^* and τ_{AL}^* in Figure 2. Instead, these data suggest a great deal of variation in diagnostic skill as well as some variation in thresholds.

In addition to the random assignment of cases to doctors, another valuable feature of Chan et al. (2022)'s setting is that in the case of a radiologist interpreting an x-ray image all variation in outcomes is due only to diagnostic skill. However, in many other medical situations such as surgery, there is a meaningful distinction between deciding when an intensive procedure is appropriate, and actually performing the intensive procedure. Currie and MacLeod (2017) discuss doctor thresholds for intensive procedures, diagnostic skill, and procedural skill in the context of child birth. The doctor's decision is between a vaginal delivery (the non-intensive treatment), and cesarean section (CS: the intensive treatment). The doctor deciding on the CS will normally also perform it. Procedural skill will be reflected in the relative returns from treatment, $\Delta_{0NIj}/\Delta_{1Ij}$. Doctors who are better at performing vaginal deliveries will have a higher Δ_{0NI} , while better surgeons have a higher Δ_{1Ij} .

As in Abaluck et al. (2016) and Chan et al. (2022), one can use the vector of observed patient characteristics, \vec{x}_i to estimate the patient's appropriateness for the intensive procedure and treat this estimated propensity as an index of appropriateness, i.e. the medical benefit of the procedure. Call this index $\rho(\vec{x}_i)$.

⁵See Rambachan (2024) for a recent extension of these identification results.

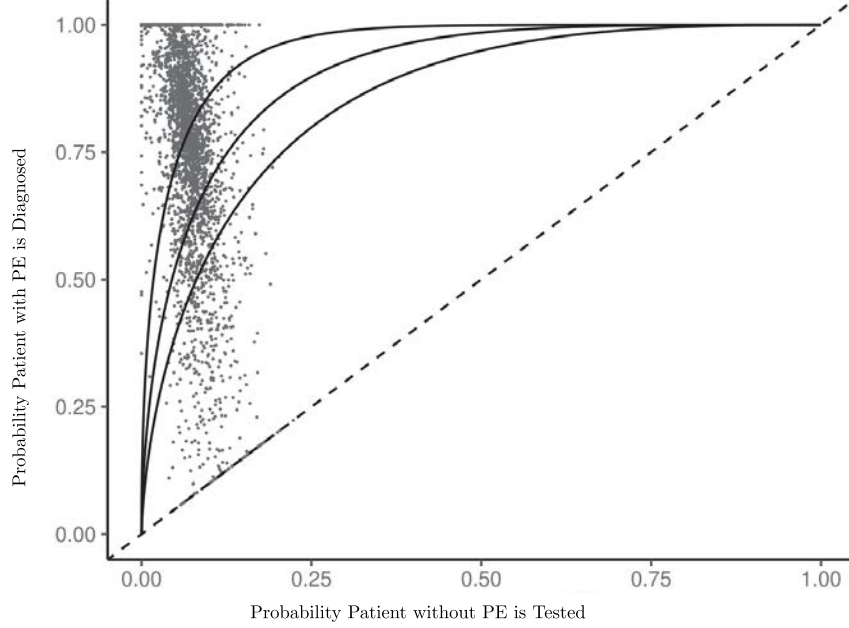


Figure 3: Distribution of Decision Thresholds and Diagnostic Skill for Radiologists (Modified version of figure V of Chan et al. (2022))

Note: Each point represents one radiologist.

It is assumed that there is a monotonic ordering of patients from the least appropriate for the intensive procedure to the most appropriate. Assume that $p_{1i} = \rho(\vec{x}_i)$, and that the doctor gets the signal given by (1) which is used to update these prior beliefs. Empirically, the ranking of patients is very stable across several different ways of estimating $\rho(\vec{x}_i)$ and patients who do not get a C-section have much lower estimated propensities on average than those who do.

Currie and MacLeod (2017) show that the doctor’s estimated probability of performing an intensive procedure is:

$$\Pr [t_{ij} = I | j, \vec{x}_i] = (PI_j - PNI_j) \rho(\vec{x}_i) + PNI_j. \quad (7)$$

The slope term, $\theta_j = (PI_j - PNI_j)$ is a doctor-specific measure that increases with doctor diagnostic skill:

$$\frac{d\theta_j}{d\gamma_j} > 0,$$

where $\Pr [t_i = I] = \rho(\beta\vec{x}_i)$ is the estimated probability of intensive treatment. Hence, doctors who have better diagnostic skills are more responsive to the measure of patient appropriateness for the procedure, $\rho(\vec{x}_i)$. See Proposition 4 in the Appendix.

Accordingly, Currie and MacLeod (2017) estimate a doctor-specific model of procedure choice as a function of $\rho(\vec{x}_i)$ and focus on each doctor’s estimated slope as a measure of diagnostic skill. Intuitively, a doctor with lower diagnostic skill has a noisier signal of the patient’s condition, and hence is less sensitive to the appropriateness measure. A doctor with poor diagnostic skill will be less likely to correctly match the pro-

cedure to the patient: They will do more intensive procedures on inappropriate patients, and fewer intensive procedures on patients who need them. Mullainathan and Obermeyer (2022) make the same observation in the context of heart attack treatment in the emergency department. Instead of the logit model used in some of the older studies, they use a machine learning model with gradient boosted trees and LASSO to predict which patients should be tested. They find that doctors make systematic errors matching procedures to patients, and that these decision errors have consequences for patient survival. Like Abaluck et al. (2016), they show that this is because physicians use the wrong weights on patient characteristics when deciding on treatments—they tend to overweight a few very salient features and underweight more subtle ones. As discussed further below, this finding is consistent with a large literature demonstrating that doctors use simple heuristics based on highly salient characteristics such as patient age to make decisions and that the use of these heuristics can lead to systematic error.

Notice that from equation (7), one can separately identify PI_j and PNI_j , and hence both τ_j and γ_j can be identified. The slope term is also affected by the physician’s beliefs about when invasive procedures are likely to be warranted via τ_j , and by any additional physician-specific factors that are included in λ_{ij} . For example, a doctor who believes that most women should have C-sections would have a very low decision threshold for C-section. Currie and MacLeod (2017) distinguish between τ_j and γ_j by noting that in a doctor-specific regression, the constant term in Equation (7) is affected only by τ_j , so given estimates of the constant term and the slope, it is possible to identify both τ_j and γ_j .

Finally, in patients with a high ex ante likelihood of having a C-section ($\rho(\vec{x}_i) \approx 1$), variation in patient outcomes is effectively independent of both diagnostic skill and the decision threshold. Hence, variation in outcomes for these patients mainly reflects procedural skill doing C-sections. A similar implication follows for patients with a very low likelihood of a C-section ($\rho(\vec{x}_i) \approx 0$). One of the findings in Currie and MacLeod (2017) is that there appears to be a positive correlation in procedural skill for both the intensive and non-intensive procedures, consistent with the hypothesis that some doctors are on average more skilled than others.⁶

Having estimated proxies for procedural skill and diagnostic skill, the estimated measures are then included in regressions of procedure choice and patient outcomes conditional on procedure prices, patient demographics, month, year, and zip code fixed effects. Two potential problems with this two-step method are that the skill measures are estimated and therefore measured with error, and that women may choose their physicians on the basis of their skills. In order to deal with these problems, Currie and MacLeod (2017) follow Kessler and McClellan (1996) and use leave-one-out, market-level averages of the skill measures as instruments for an individual doctor’s own skill measures. The doctor’s skill is empirically highly correlated with the skill of other doctors in the same market, and the average skill level of doctors in the obstetrics’s health care market is exogenously determined as long as patients did not choose their residential locations on the basis of these measures. The inclusion of zip code fixed effects makes this final condition more likely to be satisfied since it controls for fixed features of the location that might make it more or less attractive for people to live there. Having laid out this simple model of decision making, we use this model in the following sections to interpret the literature about factors that are thought to affect the quality of doctor decision making.

⁶In contrast, Chandra and Staiger (2007) hypothesize that physicians who are skilled in the intensive procedure will be less skilled in the non-intensive procedure and vice-versa.

3 Variation in Doctor Decisions and Health Equity

A vast literature shows that doctors treat patients with similar medical conditions differently depending on their income, education, gender, and race. Appendix Table 1 outlines a number of recent correspondence studies that provide further evidence about disparities in treatment. For example, Angerer, Waibel, and Stummer (2019) sent emails on behalf of mock patients who were trying to schedule doctor appointments in Austria. They found that doctors responded more quickly and offered lower wait times to patients whose signatures indicated that they had a PhD or MD degree. Button et al. (2020) conducted an innovative correspondence study in which fictive patients sought mental health appointments. The patients randomly signaled transgender or non-binary gender identities in the text of their requests. Race was also signaled using stereotypical Black and white names. They note that mental health professionals are more likely to work in solo practices than other providers, which might give them more scope for discrimination. The results suggest some complexity in physician responses across these groups: Transgender or non-binary (TNB) African Americans and Hispanics were 18.7% less likely to get a positive response than cisgender whites. There was no evidence of differential responses by TNB status for white patients.

As discussed below, some of these differences may be due to physician financial incentives, since higher income, or attributes correlated with higher income, could signal higher patient ability to pay. However, the evidence suggests that differences in average income are not a major part of the story. For example, Sommers et al. (2017) find that only a small fraction of reported racial differences in health care quality can be explained by the higher fraction of Black patients who lack of insurance coverage. Moreover, it is not clear that eliminating financial disparities would eliminate disparities in treatment. Brekke et al. (2018) study Norwegian data in which doctors were reimbursed similarly for all patients and found that patients with more education still got fewer and longer visits, while less educated patients got more visits and services (such as diabetes screenings) over the course of a year. The disparities might reflect physician affinity for spending time with more educated patients, or they might be an optimal response to differences in time costs and health needs. Chandra and Staiger (2010) replicate the well-known finding that female and minority patients receive fewer treatments than white male patients in a sample of Medicare patients. But they also find that the health benefit of treatment conditional on detailed patient observables is lower for these patients. This result is hard to reconcile with the view that useful treatments are being systematically withheld from female and minority patients who could benefit from them. As they point out, “the fact that providers may offer fewer treatments to women and minorities is not by itself evidence of prejudice.”

In order to try to get at the role of physician preferences and beliefs, it is necessary to go beyond purely observational data. Goyal et al. (2015), Hoffman et al. (2016), and Sabin and Greenwald (2012) focus on differences in the way Black and white patients are treated for pain. Goyal et al. (2015) consider children who arrive in the ED with appendicitis. The underlying assumption is that most children with acute appendicitis will be treated in hospital and that the clinician they get upon arrival at the ED will be approximately random. They find that Black children were less likely to receive any analgesia. Hoffman et al. (2016) explore the idea that racial disparities in treatment could be related to an erroneous belief that Black people have higher pain thresholds than white people. They find that doctors who endorse more erroneous beliefs about Black people’s biological responses to pain in a survey are also more likely to downrate Black patient’s pain when presented with patient vignettes. Similarly Sabin and Greenwald (2012) find that physicians with higher scores on an implicit bias test are less likely to say that they would give (clinically appropriate)

oxycodone to a Black child suffering pain after bone surgery, compared to how they say they would treat a white child.

Perhaps the most popular design for studying disparities is the concordance study. The focus in these studies is on whether patients who are more similar to doctors in terms of characteristics such as race and gender receive better treatment. In a compelling study, Cabral and Dillender (2024) obtained all Texas records for worker’s compensation and for the independent medical examinations that applicants received. Assignments to doctors were random conditional on geography and the doctor’s specialty. There were no effects of physician gender on the benefits received by male patients. However, female claimants seen by female doctors were 5.2 percent more likely to receive benefits and that the value of benefits received was 8.6 percent higher than for female claimants seen by male doctors. This finding is reminiscent of Eli, Logan, and Miloucheva (2019) who study U.S. civil war veterans and show that the same physician review boards were much less likely to recommend pensions for Black veterans than for white veterans with similar medical profiles. In turn, the lower pension benefits predicted lower life expectancy for these veterans.

Some studies suggest that discordance between physician and patient characteristics can have fatal consequences (Greenwood, Carnahan, and Huang (2018); Greenwood et al. (2020); Hill, Jones, and Woodworth (2023); McDevitt and Roberts (2014); Wallis et al. (2022)). As in Cabral and Dillender (2024), the effects are generally asymmetric: For example, Greenwood, Carnahan, and Huang (2018) find that in a matched sample, only female patients treated by male physicians are less likely to survive. Gender mismatch has no consequences for male patients treated by female physicians. In the case of racial discordance, Hill, Jones, and Woodworth (2023) focus on uninsured patients admitted to Florida hospitals through the ED and find that Black patients are 27 percent less likely to die when they have a Black physician. A nice feature of this study is that it takes the potential endogeneity of matching seriously and addresses it in two ways. First, their uninsured patient pool is unlikely to have a primary care physician who can help manage their stay in the hospital. And admission through the ED means that these are not scheduled admissions. Second, they develop an instrumental variables approach where the probability of concordance depends on the share of same-race physicians who are typically present during that shift (i.e. Friday nights) at the index hospital. Finally, they include hospital fixed effects to account for the fact that even Black and white patients who live in the same zip code may use different hospitals.

While these correspondence studies provide compelling evidence of disparate treatment, they generally shed little light on the reasons for it. Two possible channels are either explicit or implicit biases against some groups of patients, or, more subtly, difficulties communicating across groups. In turn, barriers to communication could degrade the quality of diagnosis and the efficacy of treatment. Figure 4 illustrates these two alternatives. The lower curve represents a doctor with a fixed level of diagnostic skill who has different views about patients A and B. These views are represented by the slopes of the lines tangent to the curve, which, as discussed above, capture differences in physicians beliefs about the efficacy of treatment to the two groups, as well as any differences in preferences for treating the two groups. As drawn, the physician is less likely to provide intensive treatment to patient B, whether it is appropriate or not. Hence, patient B will lose out on medically needed treatment when it is appropriate, but may also be shielded from inappropriate treatment. An example of the latter phenomena is that Black people were initially protected from the over-prescribing of prescription opioids at the start of the opioid epidemic by doctor’s lower propensity to prescribe painkillers to them, so that the opioid epidemic was initially concentrated

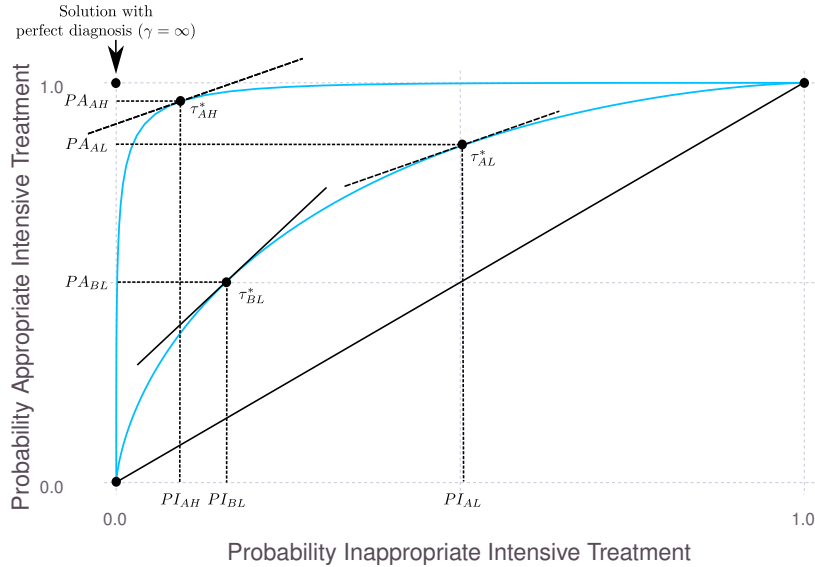


Figure 4: The Effects of Beliefs and Communication on Health Disparities

among white patients.

Alternatively, suppose that the physician treating A is unable to communicate well with A, and this barrier leads the physician to choose τ_{AL}^* . In this diagram, improvements in communication would move the physician's choice of a threshold for the aggressive procedure from τ_{AL}^* to τ_{AH}^* . This would reduce inappropriate procedure use and increase appropriate procedure use. If for example, female doctors listen more carefully to female patients or know better what questions to ask, then this could explain the better outcomes of female patients with female doctors. In this case the female doctor would be on the high diagnostic profile when treating female patients while the male doctor would be on the lower curve. It may also be the case that many Black patients have more trust in Black physicians which results in improved communication. Lack of trust in white physicians could result from many historical injustices inflicted on Black people, including the notorious Tuskegee experiment in which Black men with syphilis were not informed of their diagnosis and were left untreated so that researchers could study the untreated course of the disease. (Alsan and Wanamaker, 2018) shows that this specific incident generated a legacy of distrust that endures to the present day.

A handful of studies suggest that this trust \rightarrow communication \rightarrow diagnosis channel may be quite important. Greenwood, Carnahan, and Huang (2018) find that survival increases for female heart attack patients who are being treated by male doctors in the ED when there are more female physicians present and when the doctor has treated a larger number of female patients in the previous quarter. Alsan, Garrick, and Graziani (2019) conduct a concordance study in which Black male patients were recruited to a special clinic offering preventive care services. Patients initially received a signup sheet which included a picture of the doctor they were assigned to, who could be white or Black. They were asked to indicate on the sheet which services they wanted. At this stage, there was no difference in the number or type of services requested by the race of the doctor. The patients then actually saw the doctors, who tried to persuade them to take all of the recommended preventive services. Black doctors were much more successful at this step, increasing

take up of diabetes screening, cholesterol screening, and flu shots by 39, 53, and 27 percent, respectively.

Frakes and Gruber (2022) exploit data from the U.S. Military Health System and follow patients with severe but manageable chronic conditions whose provider changed race because of base relocation. They find that racial concordance leads to a 15 percent decline in Black mortality relative to white mortality. Over half of this decline is due to better patterns of medication use and adherence. In particular, Black patients are more likely to continue refilling prescriptions when they have a Black physician. Singh and Venkataramani (2022) show that racial disparities in in-hospital mortality increase when hospitals reach full capacity, suggesting that mistakes are more likely to be made in this kind of high-stress environment and that these mistakes have the greatest impact on the most vulnerable patients. Possibly, some mistakes involve lack of communication between patients and providers.

Tracking down the causes of disparate treatment is important because it may help to pinpoint possible solutions. As discussed above, differences in financial resources play a role, so equalizing access to insurance helps but is not likely to eliminate disparities. The pain studies, and studies directly investigating physician bias indicate that this is an important source of disparities in care, but eliminating bias has proven difficult. As Williams, Lawrence, and Davis (2019) point out, there is little evidence that interventions aimed at addressing bias have improved health outcomes. In a review of this literature in the *Annual Review of Public Health*, Vela et al. (2022) conclude that the effects of most anti-bias training interventions in medical settings are either null or extremely short-lived. They argue that this may be because the message in the anti-bias training is undermined and contradicted by other aspects of medical training. These findings may point to the importance of the larger social context both within and outside institutions in shaping physician behavior. They suggest that positive interactions with both providers and patients from historically marginalized groups could have a larger impact than formal anti-bias training in terms of resetting harmful provider beliefs.

The most obvious conclusion to be drawn from the concordance studies is that health equity would be improved by having larger numbers of female practitioners and practitioners of color. Women are still underrepresented in many medical specialties. For example, McDevitt and Roberts (2014) discuss urology and show that having even a single female urologist in a county is associated with fewer female deaths from bladder cancer. The situation is more extreme for Black physicians who make up only four percent of the workforce. It will take a long time to improve that number to the point where most Black patients could see a Black physician if they wanted to, or even to the point where most white physicians have experience working alongside Black doctors. Hence, if improving communication is a core issue, an important question for future work is whether there are additional ways to achieve this goal. Perhaps it is possible to leverage other medical professionals such as nurses or doulas as intermediaries so that important patient concerns are properly heard. However, empirical research on intermediaries such as doulas is limited Sobczak et al. (2023).

More generally, interventions that ensure that doctors correctly treat patients conditional upon their symptoms can be expected to reduce health disparities. We now turn to research that measures variation in doctor decisions that arise from variation in their skill and the conditions under which they are making choices.

4 Factors that Affect the Quality of Decision Making

4.1 Skill, Experience, and Training

An immediate implication of the theoretical framework is that doctors with lower skill levels should set different thresholds for using intensive procedures than doctors who are more skilled. For example, Doyle, Ewer, and Wagner (2010) have an elegant study in which hospital patients were randomly assigned to the “A team” or the “B team” of residents, where the A team was trained at a higher ranked medical school. Although the two groups of patients had similar outcomes on average, A-team patients had systematically shorter and cheaper hospital stays. The B team used more diagnostic and testing resources to arrive at the same outcomes, consistent with the idea that less skilled doctors set lower testing thresholds. In a related context, Chan, Gentzkow, and Yu (2022) suggest that since it is more costly to miss a pneumonia diagnosis than to erroneously admit a patient to hospital, less skilled radiologists will err on the side of caution by being more likely to admit a marginal patient. They find evidence consistent with this hypothesis. Currie and Zhang (2023) also find that more skilled physicians “do more with less” in the sense of achieving the same or better health outcomes with fewer inputs. Similarly, Gowrisankaran et al. (2022) find that in the Canadian province of Quebec, Emergency Department doctors with more intensive practice styles have worse patient outcomes on average. They rely on random assignment of patients to doctors within the ED, and they measure practice style and skill as doctor fixed effects in models of procedure choice and patient outcomes.

Several studies show that doctors with more training have better outcomes on average. For example, in models that control for hospital, quarter, and day of week effects as well as the number of doctors present, Doyle (2020) shows that EDs have better outcomes for heart failure patients when they have a cardiologist on staff. However, it is possible that cardiologists are positively selected in terms of doctor quality on average so it is difficult to distinguish between selection effects and the effects of additional training per se.

Schnell and Currie (2018) find that physicians from higher ranked schools prescribe fewer opioids, even within the same practice address. If physicians from higher ranked schools are more skilled, this could reflect either better training or the way that medical students are selected and sorted into schools of different ranks. But Schnell and Currie (2018) also show that in specialities that receive specific training in the use of opioids and other pain medicine, there is no difference in prescribing by medical school rank, as one might expect if doctors from higher ranked schools were just generally better. Hence, their results suggest that training is an important determinant of practice styles.

Chan and Chen (2022) expand beyond considering doctors as providers and compare outcomes for patients treated by nurse practitioners (NPs) or doctors in Veteran’s Administration Emergency Departments. They use the number of NPs who are on duty as an instrument for being treated by an NP. They find that, on average, being treated by an NP increases length of stay and health care costs, though being treated by an NP has relatively little effect on outcomes. These results echo Doyle, Ewer, and Wagner (2010)’s finding that the “B team” uses more resources to arrive at the same results. A more striking finding is that there is considerable variation in the skill levels of both groups - many NPs achieve better outcomes at lower cost than some doctors, even though NPs have much less lengthy and intensive training than doctors.

The evidence regarding the relationship between doctor experience and outcomes is mixed. van Parys (2016) finds that the least experienced ED physicians perform more procedures and spend more than ED physicians who have practiced for seven years or more. She also finds that these high spending physicians are

more likely to stop working in the ED so that the overall performance of ED doctors may rise slightly with experience due to positive selection in who stays. Epstein et al. (2016) focus on obstetricians and measure initial skill defined as a physician’s normalized, risk-adjusted maternal complication rate in the first year of practice. They find that even after 16 years, initial skill is most predictive of patient outcomes, and that years of experience have little impact. In contrast, Facchini (2022) estimate doctor fixed effects models and find that obstetricians have better infant health outcomes when they have done more C-sections in the last four weeks, suggesting that it may be very recent experience that matters. Finally, ? evaluate the extent to which primary care physicians promote medication adherence and positive health outcomes of patients on statins. Doctors whose patients do better on these measures are said to have better health management skills. Looking at patients who had to switch doctors, they find, however, that these skill measures appear to decay rather than to increase with a doctor’s age.

One way to operationalize the idea that experience matters in the context of the theoretical framework laid out above is to make diagnostic skill and procedural skill functions of experience. For example, Currie, MacLeod, and Van Parys (2016) compute γ_j as described above, but allow it to vary over time. Regressing this structural parameter on years of experience, they find that γ_j decreases sharply after 24 years of experience, consistent with the more negative views of the correlation between doctor experience and outcomes described above. It is possible for diagnostic skill and procedural skill to evolve in different directions with experience – a doctor might, for example, just decide that they were going to do C-sections for all patients. In this case their diagnostic skills might atrophy while, at the same time, they became very good at performing the procedure. However, the results of Epstein, Nicholson, and Asch (2016) suggest that procedural skill, s_{tj} , is fairly flat with respect to experience at least when it comes to doing C-sections. One difficulty with these comparisons is that we typically only observe doctors who have graduated from medical school so that we do not observe doctor skill levels during the period when returns to experience might be steepest.

On the whole, there has been little investigation of variation in procedural skill at the doctor level within the economics literature. Chandra and Staiger (2020) consider procedural skill at the hospital level. Arguably, while it is doctors who make decisions about how a given patient is to be treated, hospitals can influence that process. For example, a hospital can choose whether or not to have a heart catheterization facility, which will affect whether catheterizations can be performed. In terms of our framework, we can think of hospitals having a comparative advantage in either the intensive or the non-intensive procedure. Chandra and Staiger (2020) argue that hospital physicians have erroneous beliefs about their hospital’s comparative advantage, such that they overuse procedures which are not their comparative advantage. In a study of the treatment of heart attack patients in 45 states between February 1994 and July 1995, they conclude that eliminating such “allocative inefficiency,” i.e. having hospitals stick to their comparative advantage, would increase the benefits of treatment by 44%.

The papers discussed in this section are summarized in Appendix Table 2. Overall, the research suggests that training and experience affect doctor’s skill and practice styles. However, the effects of post-medical school training seem to be small. There is also less evidence that procedural skill improves with experience than one might expect, given the well known relationship between high surgical volumes and better surgical outcomes.⁷The evidence is also consistent with the hypothesis that selection matters, and that prospective

⁷For example, in a review of the literature, Chowdhury, Dagash, and Pierro (2007) find that 74% of studies find that higher volume surgeons have better outcomes and that specialist surgeons have better outcomes than general surgeons 91% of the time.

doctors vary in their innate ability to diagnose patients and execute procedures as well as in the extent to which they improve or keep up their skills. Overall it is unlikely that improvements in training or the accumulation of doctor experience alone will eliminate variations in the quality of doctor decision making and procedural skill.

4.2 Time Pressure and Fatigue

Doctors often work long hours in a fast-paced environment in which decisions must be made quickly and with little time for reflection. Time pressure could lead to mistakes if diagnostic skill, γ_j , falls with stress or fatigue. Figure 2 illustrates the idea that lowering diagnostic skill, γ_j , increases the probability of inappropriately choosing the intensive treatment and reduces the probability of appropriately choosing the intensive treatment. The more interesting point is that the increase in the use of inappropriate treatment is much greater for aggressive doctors (who move from τ_{AH}^* to τ_{AL}^*), while the decline in the probability that intensive treatments are appropriately rendered is greater for conservative doctors (who move from τ_{CH}^* to τ_{CL}^*). Hence, the same reduction in diagnostic skill has differing effects depending on the doctor’s baseline type, which in turn depends on their beliefs about the probability that an intensive treatment is likely to be appropriate and the relative efficacy of the procedures in their patient pools. This observation suggests that the effect of time pressures can be highly variable.

Studies focused on the impacts of time pressure and fatigue on doctor decision making are summarized in Appendix Table 3. They show a wide range of estimated effects. Tai-Seale and McGuire (2012) provided some early evidence about the importance of time pressures, showing that as the length of a visit increases, doctors are more likely to treat each new topic as the last to be covered during the visit. Subsequent authors focus on whether time pressures lead to more or less use of intensive procedures, with mixed results. For example, Freedman et al. (2021) find that unexpected increases in PCP patient waiting times result in fewer referrals, opioid prescriptions, and Pap tests, and increases in scheduled and unscheduled followup visits. Persson et al. (2019) find that within an orthopedic surgeon’s shift, each additional patient seen reduces the probability that a surgeon recommends surgery. On the other hand, Gruber, Hoe, and Stoye (2021) find that English ED doctors who were under pressure to reduce waiting times did so by admitting patients to the hospital, thereby increasing hospital costs by 4.9 percent without any effect on one year mortality, length of stay, or the number of in-patient procedures. Similarly, Chu et al. (2024) study ED doctors and find that when doctors are managing more cases simultaneously, they order more tests, perhaps substituting testing for their time and attention.

Chan (2018) studies ED doctors and finds that as they near the end of their shifts, they are increasingly likely to admit patients to the hospital, with a 21.19 percent increase in the last hour of the shift, resulting in 23.12 percent higher costs. There are no significant effects on 30-day mortality or “bounce back” of patients to the hospital. Chan (2018) also finds that these end-of-shift effects are not found when out-going doctors have sufficient time to hand off their patients to the incoming physician. He suggests that the changes in doctor behavior are not driven by fatigue or a higher probability of errors in judgment but by changes in doctors’ valuations of their leisure time over the course of a shift. In terms of the model, λ_{ij} , the payoff associated with the the intensive procedure rises leading to more bias in decision making.

The sign of the effect would depend on which course of action is most convenient for the doctor. In the ED, admitting the patient to the hospital may be the course of action that takes the least time, while in a

PCP office, skipping tests and referrals can save time. Costa-Ramón et al. (2018) report that in a Spanish hospital, the probability that an unscheduled delivery is via C-section rises between 11pm and 4am when, presumably the obstetrician on duty would like to quickly complete the delivery and go back to bed.

A related question is how the doctor’s emotional state impacts decision making. Chodick et al. (2023) look at the effect of a primary care doctor’s encounter with a patient who has been newly diagnosed with cancer. They find a short-lived (1 hour) but large effect on the doctor’s probability of ordering a wide variety of diagnostic tests, not just cancer screening tests. They discuss a number of reasons for this result, including a physician’s emotional response to the new diagnosis for their patient, or the need to test for comorbidities. Understanding the impact of a physician’s emotional state, broadly defined, could help to identify moments when doctors were particularly likely to make mistakes.

4.3 The Role of Peers and Teams

Research on the influence of peers and teams on doctor decision making has been motivated by the desire to explain geographical clusters in practice style. It is important to understand the size of these effects if we aim to improve doctor decision making. Proximity to peers and interactions with peers could have major effects on physician behavior through information channels, opportunities for matching patients with physicians (or physicians with physicians), and the creation or mitigation of moral hazard. Studies exploring these channels are reviewed in Appendix Table 4.

Several studies suggest that peers are an important source of information. For example, Agha and Molitor (2018) look at whether physical proximity to lead investigators in clinical trials for new cancer drugs leads to faster take up of those drugs and find that patients in the lead investigator’s hospital referral region are 36 percent more likely to get the new drug initially, with convergence across regions after four years. Theory predicts that a doctor’s threshold for using a drug or procedure is influenced by their beliefs about the proportion of patients in the population who are likely to benefit. In this case, doctors update their beliefs about whether the new drug will be beneficial for their patients more quickly when they have access to a lead investigator, or perhaps when they are more likely to see patients who have benefited from the new drug. The effects are largest in the areas with the slowest rate of new drug adoption.

Chen (2021) examines patients receiving heart procedures and finds that patients do better when the surgeon has worked longer with the other hospital physicians who are caring for the patient. The effects are large: A one standard deviation increase in shared work experience reduces 30-day mortality by 10 to 14 percent and reduces the utilization of medical resources and length of stay. The effect is greater for more complex cases. It is interesting to compare this example to Agha and Molitor (2018) in part because it does not involve information about new or more complex procedures. The effects presumably mainly reflect better communications among members of the team, which in turn improve patient outcomes.

Molitor (2018) explores another dimension of peer effects—the matching of like-minded physicians in the same geographic area. Using a “movers” design, he shows that when cardiologists move to a new hospital referral region, they quickly adapt their own treatment style to the predominant style in the new region: A one percentage point increase in cardiac catheterization in the new HRR raises the doctor’s own rate by 0.628 percentage points within one year. The effect is greater for doctors moving from low to high-intensity areas. Since physicians do not move randomly, it is possible that the cardiologists are choosing to move to areas in which others share their desired practice style. These moves would increase geographic dispersion

across regions and geographic concentration in practice styles within regions.

In other situations, doctors may have less choice about who their peers are or how much they must adapt to the practice styles of others. Chan (2021) studies teams of residents in a large teaching hospital where teams consist of junior residents who are led by a senior resident. He shows that almost all of the variation in decision making within a team is accounted for by the senior resident. The variation in the behavior of junior residents increases sharply after one year, when they become senior residents themselves. This finding is striking since medical residents presumably gain experience continuously over their first year of practice but only change their behavior discontinuously at the one year mark.

Silver (2021) focuses on teams of ED doctors and exploits variations in the composition of teams across shifts, arguing that these are essentially random. He finds that doctors work faster when they are placed with a fast-paced team and that, on average, the faster pace has no effect on the outcomes of discharged patients. However, the riskiest patients suffer increases in 30-day mortality. This result contrasts with Gruber, Hoe, and Stoye (2021) who, as discussed above, found that physicians working faster in response to a mandate to reduce ED wait times increased costs without having any negative effects on patient outcomes. Possibly, the American doctors were under greater pressure not to increase costs, but the contrasting results suggest that one should be cautious about extrapolating from any one study in this nascent peer effects literature.

While Silver (2021) and Gruber, Hoe, and Stoye (2021) suggest that doctors can choose to work faster or slower, Chan (2016) asks whether doctors who work more slowly are shirking, and thereby forcing other members of their team to work harder. His study focuses on two teams working in the same hospital. In the first team, doctors decided on how patients were allocated within their group. In the second team, patients were initially assigned to doctors by a nurse scheduler, and then later assigned by the doctors themselves. Chan (2016) shows that switching the nurse-managed team to being doctor-managed reduced wait times by 13.67 percent without any effects on costs, utilization, or outcomes. His interpretation is that the doctors had a better understanding of how long each patient should take so that they could detect shirking. Hence, the switch reduced moral hazard. The alternative explanation, that doctors are better able to match patients to the doctors who can treat them most efficiently, seems unlikely given the lack of any change in patient outcomes.

Currie, MacLeod, and Ouyang (2024) examine peer effects in physician prescribing to adolescents with mental health conditions. They point out that it can be difficult to identify peer effects if doctors with similar training and experience tend to have practice styles that evolve similarly over time and also cluster in the same locations. They look at correlations between the index physician's probability of prescribing inappropriately; the probability that physicians with similar training and experience from outside the area prescribe inappropriately; and the probability that physicians from the same area but with different training and experience prescribe inappropriately. They find that the "effect" of physicians from the same cohort but outside the area is about half the size of the effect of local physicians from different cohorts. Hence, some of what appears to be a peer effect actually reflects the co-evolution of practice styles among similar physicians. The size of the spillover effects are consistently larger for non-psychiatrists than for psychiatrists, indicating that specific training can mitigate the extent to which harmful practices spread through peer effects. However, physician fixed effects rather than peer effects appear to be the most important determinant of variations in practice style in their data.

These papers suggest that it is quite difficult to identify true peer effects outside of certain specialized

settings in which it is plausible to assume that doctors do not choose their peers. Hence, we are a long way from being able to use estimates of peer effects to think about influencing doctor behavior.

4.4 Financial Incentives

Health economists have long realized doctors can be influenced by financial incentives, an effect that the theory captures with the λ_{ij} parameter. Appendix Table 5 provides an overview of some post-2010 contributions to the large literature on financial incentives in health care markets. We focus on studies that either are particularly valuable in enhancing our understanding of physician decision making or that reflect concerns that have become focal in the literature post-2010. Handel and Ho (2021)’s chapter in the Handbook of Industrial Organization provides a review of some aspects of the healthcare market that impact financial incentives, including competition in hospital and insurance markets, negotiations between hospitals and insurers, and increasing vertical integration in hospital markets. In general, the IO literature they survey has focused on the larger players, such as hospitals and insurers which can be understood as “firms,” rather than on the decisions of individual physician providers. However, as more physicians work for large groups, and more practices become part of vertically integrated health care companies, this distinction may become less relevant. For example, Chernew et al. (2021) show that vertically-integrated physicians are much more likely to refer patients to expensive hospital-based MRI providers compared to non-vertically-integrated physicians (52 percent vs. 19 percent). Similarly, Currie, Karpova, and Zeltzer (2021) show that urgent care centers, which are increasingly owned by hospitals, increase inpatient hospital care for elderly patients rather than substituting for it.

Two overarching questions addressed in this section are whether (and how) governments and insurance plans can use financial incentives to reduce health care spending without worsening patient outcomes, and whether some types of patients are more or less vulnerable to the distortions in doctor decision making that are induced by financial incentives. Several studies look at changes in reimbursements from the U.S. Medicare program. Reducing spending in Medicare is of particular interest both to policy makers and economists as the population ages and advances in medical technology make Medicare spending an increasing part of the federal budget.⁸ Clemens and Gottlieb (2014) take advantage of a consolidation of Medicare reimbursement regions that raised reimbursements in some areas and lowered it in others. They show that higher reimbursement rates increased the use of elective procedures and the probability of hospitalization for heart attacks (acute myocardial infarction) within one year, without having any effect on four-year mortality rates. Note that if hospitalizations were driven primarily by consumer demand, higher prices would lead to lower quantities. Hence these results suggest that at the margin hospitalizations are driven by supply-side considerations.

A major complaint about Medicaid, the U.S. public health insurance program for low-income individuals, is that patients have difficulties getting an appointment. For example, Bisgaier and Rhodes (2011) report on an audit study in which patients on Medicaid were six times more likely to be denied an appointment and had to wait three weeks longer to see a provider if they did get one. Alexander and Schnell (2024) look at a Medicaid “fee bump” which resulted from the 2010 Affordable Care Act’s payments to states to equalize physician Medicaid reimbursements. The bump increased Medicaid payments by an average of 60 percent, with considerable variation across states. Their results suggest that closing the gap between the payments offered by Medicaid and those offered by private health insurance would eliminate disparities in

⁸Medicare accounted for 12 percent of the total federal budget in 2022. See <https://www.pgpf.org/budget-basics/medicare>.

access to primary care for children, and would also reduce access disparities by two thirds for adults. Dunn et al. (2024) consider another type of provider disincentive associated with Medicaid — an elevated risk of having a claim denied or otherwise unpaid. They find that 18 percent of Medicaid claims are denied, a much higher rate than under either Medicare or private insurance. They conclude that this high probability of non-payment is as great a barrier to accepting Medicaid patients as are lower fees.

Other authors focus on the effect of capitation – that is providing doctors with a fixed payment per patient. Most economists would predict that capitation would lower the intensity of service delivery relative to fee-for-service payment, which is exactly what has been found in empirical studies. For example Ding and Liu (2021) show that providers with capitated payments used 12.2 percent fewer resources (especially physical therapy and diagnostic testing) compared to non-capitated providers, with no change in outcomes. One issue with studies of capitation is that providers who are not being reimbursed for providing specific services may have little incentive to record them in claims data. Hence, some of the measured reduction in services rendered could be an artifact of changes in reporting practices.

Chorniy, Currie, and Sonchak (2018) show that diagnoses can be affected by incentives built into managed care contracts. In their South Carolina setting, providers who were switched to capitated payments got larger payments if patients had specific chronic conditions. Providers were also penalized if they screened children for chronic conditions at lower than average rates. In models that follow the same children over time as their providers were switched from fee-for-service to capitated contracts, they found an 11.6 percent increase in diagnoses of ADHD and an 8.2 percent increase in diagnoses of asthma without any effect on ED use or hospitalizations.

Several more tailored schemes for reducing health care costs without reducing quality have also been evaluated. Gupta (2021) studies the impact of the Hospital Readmissions Reduction Program (HRRP) which penalized hospitals with Medicare readmission rates that were higher than a given threshold. He deals with the potential for reversion to the mean by instrumenting the estimated probability of a penalty using the predicted penalty based on hospital characteristics from earlier years. He finds very large effects of the program—the HRRP was estimated to account for two-thirds of the observed reduction in readmission probabilities and to have reduced 1-year mortality by 8.87 percent. These positive effects were achieved by increasing the intensity of care during the initial hospital admission. One reason for the success of this scheme may have been that it applied to all hospitals. Alexander (2020) studies a New Jersey policy that allowed hospitals to select into a program that offered physicians incentives to lower the costs of care. Alexander (2020) finds that the program had no effects on costs or procedure use – instead, physicians were able to game the system by directing their lowest-cost patients to participating hospitals. This simple tactic lowered patient costs at these specific hospitals so that doctors could reap the incentive payments. However, it resulted in higher patient travel costs. Alexander and Currie (2017) show that doctors’ responses to incentives may also be affected by factors such as capacity constraints. They find that doctors are generally more likely to admit child respiratory patients when those patients have private insurance rather than lower-paying public insurance. This gap grows when beds are in high demand because of high flu caseloads.

Strong responses to physician financial incentives have also been found in European settings, where most countries have some form of universal health insurance coverage. For example, Wilding et al. (2022) focus on an English policy which imposed financial penalties on GPs when the fraction of hypertensive patients with blood pressure under control fell below a target. They show that stricter targets increased prescription of

anti-hypertensive medication. But doctors also showed evidence consistent with gaming: They did multiple tests on patients whose blood pressure initially exceeded the threshold (presumably trying to get a reading below the threshold), took actions to have patients declared exempt from testing requirements, and were more likely to report that patients exactly met the threshold suggesting greater use of rounding. In France, Coudin, Pla, and Samson (2015) show that the imposition of price controls increased the number of procedures by over 80 percent, suggesting that physicians increased quantities in order to make up for shortfalls in income due to the price controls.

As we pointed out at the start of this section, it might be more surprising to economists to find instances in which doctors did not respond to financial incentives. Some recent studies focus on factors that mute or mediate the expected relationship, including a variety of patient characteristics. For example, Johnson and Rehavi (2016) look at patients who are themselves physicians. They find that physician patients are about six percent less likely than other well educated patients to have unscheduled C-sections, and that financial incentives affect C-section rates only for non-physician patients. However, it is not entirely clear whether this null result reflects push back from informed consumers or physicians refraining from suggesting unnecessary C-sections to their peers.

Chen and Lakdawalla (2019) use the same change in Medicare billing areas as Clemens and Gottlieb (2014), and ask how physician responses to changes in Medicare reimbursements vary with the income of the patient. A key institutional detail is that fee-for-service Medicare patients have copays. Since richer patients are likely have greater willingness to pay than poorer ones, the authors predict that higher reimbursements will lead to larger increases in procedure use in richer patients because poorer patients are more likely to resist the higher copays. They show that increases in reimbursements increased the gap in services received between high and low-income patients.

Whether the physician has an ongoing relationship with a patient has also been shown to be an important mediator of the extent to which financial incentives affect patient care. Brekke et al. (2019) use Norwegian administrative data linking health, national insurance, and labor market participation to examine physician behavior with respect to the issuance of sick-leave certificates. In order for workers to claim sick-leave benefits, they must have a doctor sign a certificate. Physicians see patients both in their own practices and in EDs. They are likely to have on-going relationships with patients in their own practices but not with patients in the ED. Physicians may also be on fee-for-service or fixed salary contracts. The authors show that physicians are 34.63 percent more likely to issue sickness certificates for their own patients with fee-for-service contracts and 24.15 percent more likely with fixed salary contracts. However, for new GPs with fixed salaries, there is no gap in rates between own patients and ED patients, which may reflect the fact that new GPs do not yet have any on-going relationships with patients. The size of the gap in sick leave issuance between own patients and ED patients is greater in areas with larger numbers of GPs per capita and among GPs who have openings for new patients suggesting that competitive pressures also influence this behavior.

Currie, Li, and Schnell (2023) also examine the impact of competition on physicians, using state laws that allowed nurse practitioners to prescribe controlled substances independently as a source of exogenous variation in competition. They find that general practitioners responded by prescribing significantly more controlled anti-anxiety medications, more opioids, and more co-prescriptions of the two types of drugs. The impact of the change in laws was greater in areas with higher ratios of NPs per GP to begin with and was

concentrated in specialties that face the most competition from NPs. Their findings suggest that in some cases competition can have harmful effects on patients and lead to over-provision of services.

We will briefly touch on two other types of physician incentives here, those due to “detailing” and those due to malpractice. Detailing is the practice of marketing drugs and other medical equipment or products directly to physicians. In some cases, this may involve visits from company representatives providing information, but often detailing also involves a payment to the physician in cash or in kind (e.g. meals or travel expenses). U.S. sunshine laws passed as part of the 2010 Affordable Care Act require companies selling pharmaceuticals and medical devices to report most payments made to physicians to the federal government⁹. These disclosures have enabled researchers to learn more about these payments and their impacts on physician behavior. ? examine the impact of detailing on the use of generics and the efficacy of drugs prescribed. They find that even a small payment increases prescribing of the detailed drug by about 2 percent in the six months following receipt of a payment. However, doctors do not seem to be prescribing less effective drugs or delaying transitions to generics. Shapiro (2018) also suggests that the effects of detailing are relatively benign. He studies an antipsychotic drug, Seroquel. Two clinical trials showed that Seroquel had a more benign side-effect profile than leading competitors. Building on early work by Azoulay (2002) that suggested that the impact of drug research is amplified by marketing, Shapiro finds that these trials had little impact on prescribing unless they were accompanied by detailing visits. He interprets this as evidence that the new information from the trials was conveyed to doctors through detailing. Detailing visits after the trials resulted in small shifts in prescribing towards Seroquel, and more of these prescriptions were “on-label,” i.e. for indications approved by the U.S. Food and Drug Administration (FDA). In contrast to ? and Shapiro (2018), Newham and Valente (2024) find that payment to physicians increase prescribing of branded rather than generic diabetes drugs, raising costs. Carey, Daly, and Li (2024) also find that marketing payments increase expenditures on cancer drugs in Medicare with no subsequent improvement in patient mortality. As more years of CMS Open Payments data become available, further research will be possible to help clarify this issue, though the existence of these data may itself shape the course of pharmaceutical marketing in the years to come.

Agha and Zeltzer (2022) extend the peer effects literature discussed above to consider the impact of detailing on physicians who do not receive payments directly, but who share patients with doctors who received payments. Using the Medicare claims data, they find that such spillovers account for a quarter of the increased prescribing that results from detailing payments. The effects are larger for physicians who share more patients with the doctor who received drug company payments. This finding is particularly important in that it underscores the limitations of sunshine laws in tracking the influence of pharmaceutical companies on physicians.

Doctors themselves often cite fear of malpractice as a factor that influences them to practice defensive medicine - i.e. the practice of ordering unnecessary procedures and tests in order to protect against malpractice risk. In practice, the risk of financial loss is mitigated by malpractice insurance. And since malpractice insurance is not experience rated, doctors typically do not even face higher insurance premiums after a finding of malpractice. Hence, it may be the unpleasantness associated with being sued and the subsequent damage to their reputations that doctors wish to avoid rather than financial penalties per se.

⁹In response to the 2018 U.S. SUPPORT Act, CMS Open Payments started including payments to physician assistants, nurse practitioners, clinical nurse specialists, certified registered nurse anesthetists, anesthesiologist assistants, and certified nurse-midwives. Additional research is needed to study the effects of this expansion of reporting requirements.

A large literature leverages changes in state laws in order to assess the impact of malpractice on doctor behavior. Mello et al. (2020) offer a nice survey of this literature and conclude that while some authors find non-zero effects, the impact of changes in laws governing malpractice are typically quite small. Nevertheless, the National Academy of Sciences (Balogh, Miller, and Ball (2015)) notes that the malpractice system could have a negative systemic effect by inhibiting reporting of, and learning from, diagnostic errors.

Currie and MacLeod (2008) offer several possible reasons for the small estimated effects of malpractice reforms. First, most studies lump all changes in tort laws together, even though different types of laws are predicted to have different effects. For example, laws capping damages may encourage reckless behavior while reforms making physicians liable for the share of the damages that they caused (rather than allowing plaintiffs to sue the “deep pocket” in the case for 100 percent of damages)¹⁰ should have the opposite effect. Second, the impact of a law change is likely to depend on whether a physician is doing too many or too few intensive procedures to begin with. For example, if a doctor was causing harm by doing unnecessary C-sections, then stricter malpractice laws such as raising the cap on damages, might cause them to reduce the number of C-sections. On the other hand, if a doctor was doing too few C-sections, then the same law change might cause them to do more. A nice paper by Frakes (2013) captures this intuition. The key question in most malpractice cases is whether the doctor provided care consistent with accepted medical practice. As of the late 1970s, most states used a state standard to define accepted practice. But over time, many states moved to using national rather than state-level norms. Frakes (2013) shows that state C-section rates tended to converge to the national rate after this change, with no change in infant health outcomes.

In summary, there is a great deal of evidence that physicians respond to financial incentives, which should not surprise economists. Doctors adjust the services they provide, but they may also game the system by moving patients around or shading their reports about patient conditions. Hence, reliance on financial incentives alone to regulate health care markets is likely to have negative consequences for at least some patients. Research asking which types of patients are most affected by the unintended consequences of changes in financial incentives has provided some initial answers, but it is an interesting topic for further research. Research on other types of financial incentives such as those from detailing payments or threats of malpractice have so far suggested relatively mild effects on physician behavior, though higher non-payment risk appears to have large effects.

5 Improving the Quality of Doctor Decision Making

So far this survey has demonstrated that there is a great deal of variation in the quality of doctor decision making and that poor decisions can have a negative effect on patient health, increase health care costs, and widen health disparities. Not surprisingly then, there is a growing literature discussing possible ways to improve doctor decision making beyond adjusting payment systems. This section discusses research considering the effectiveness of providing information to doctors and/or patients, using heuristics or guidelines, or using new technologies, such as electronic medical records and decision support tools in an attempt to improve medical decision making. One way to think about these interventions is in terms of whether they target diagnosis (γ_j); whether they try to shift the doctor’s priors regarding the usefulness of a medical

¹⁰The default common law rule in the U.S. allows plaintiffs to sue any defendant for 100 percent of damages. Tort reforms introducing “joint and several liability” limit each defendants liability to share of the damages that they caused.

procedure for the two types of patients, $\Delta_{0NIj}/\Delta_{1Ij}$; or whether they affect the doctor’s beliefs about the relative proportions of patient types, p_{0j}/p_{1j} in the population. At the extreme (e.g. guidelines that specify or proscribe particular actions in specific cases) they might involve taking decision making out of the doctor’s hands.

5.1 Providing Information

A number of studies explore the consequences of providing information about practice style to either physicians, patients, or both. Appendix Table 6 summarizes several examples from this literature. The most straightforward studies are experiments in which letters were sent to treatment physicians while control physicians did not receive letters. For example, Sacarny et al. (2016) designed a randomized controlled trial targetting physicians who were high prescribers of Schedule II controlled substances (opioids, amphetamines, and barbituates) to Medicare patients. This intervention could be interpreted as an attempt to communicate to doctors that they were consistently over-estimating the share of patients in their practices who were likely to benefit from these drugs. Doctors in the treatment group received letters informing them that their prescribing patterns deviated significantly from those of their peers. These letters resembled comparative billing reports that Medicare routinely sends to providers comparing their billing practices to those of their peers and did not mention any sanctions. Regarding results, the title of the paper says it all: “Medicare Letters To Curb Overprescribing Of Controlled Substances Had No Detectable Effect On Providers.” There was no evidence of heterogeneous effects by prescriber specialty, region, or whether the prescriber had been investigated for fraud.

However, several subsequent studies have found significant effects of similar letters on physician prescribing. In a followup paper, Sacarny et al. (2018) targeted outlier prescribers of the antipsychotic drug quetiapine and sent them three letters highlighting their outlier status relative to peers. Over the nine months of the experiment, the number of days of quetiapine prescribed fell by 11.1 percent relative to the control mean, and the reduction lasted at least two years. This reduction was largest for patients with low-value indications, and there were no negative effects on patient outcomes. It is possible that receiving three letters over a short period made the intervention seem less like a routine “form letter” and more like an implied threat of some sort of sanction.

Ahomäki et al. (2020) report that a precautionary letter sent to Finnish physicians who were prescribing high numbers of paracetamol-codeine pills to new patients reduced the number of pills prescribed to new patients by 12.8 percent of the treatment group baseline, which is similar to the more recent Sacarny et al. (2018) paper. Again, the letter may have carried an implicit threat, since such letters are not routine in the Finnish context. Hence, the question raised by these papers is whether doctors are responding to the information contained in the letter, or whether they are afraid of being sanctioned for their outlier behavior. Possibly the important information being conveyed is not so much that they are outliers, but that an authority is watching their prescribing behavior.

In perhaps the most famous recent example of a letter-writing intervention, Doctor et al. (2018) started with vital statistics mortality data from California identifying people who died from overdoses of prescription drugs. Then, using the state’s prescription drug monitoring program (PDMP) records, they located the doctors who had prescribed the fatal drugs. The experimental intervention involved sending a letter to a treatment group drawn from these doctors informing them that their patient had died of a drug overdose. The

researchers could then monitor the treatment doctors' subsequent prescribing using the PDMP. They found a 9.7 percent reduction in morphine milligram equivalents in the three months following the intervention. Of the “letter experiments” discussed here, this one arguably comes closest to a pure information intervention. The researchers were not writing on behalf of any state or regulatory agency, so there was less of an implicit threat. And they were supplying information that doctors would not necessarily be able to acquire easily from other sources – when U.S. doctors treat a patient who does not return, they are not routinely informed about whether this is because the patient moved, switched physicians, stopped going to the doctor, or died.

A second group of “informational” studies seeks to measure the effect of new clinical knowledge on physician behavior. For example, in a meta analysis, Hammad, Laughren, and Racoosin (2006) suggested that selective serotonin reuptake inhibitors (SSRIs) increased suicidal thinking in children and young adults. A preliminary version of this study led the FDA to put a prominent warning label on SSRI drugs in 2004. Early studies such as Gibbons et al. (2007) indicated that these warnings led to sharp drops in prescribing to children and adolescents in the U.S. and Norway, as well as to declines in prescriptions of SSRIs generally. Building on this evidence, Dubois and Tunçel (2021) replicate the finding in French data and then build a random coefficient discrete choice logit model to examine changes in physician prescribing across several drug classes. They find reductions not only in SSRIs but in the prescribing of close substitutes, as well as an increase in the off-label use of other types of psychiatric drugs as treatments for depression. A quarter of physicians stopped prescribing SSRIs altogether, but considerable variation in physician prescribing remained both before and after the change. A limitation of their work is that their model must perforce rely on the strong assumption that the way physicians are matched to patients does not change following the announcement.

McKibbin (2023) presents another convincing study of the impact of new information. Since FDA approval is a lengthy process, many sick cancer patients do not have time to wait for the process to be completed but take promising new drugs “off label” (i.e. before they are FDA approved for that indication). McKibbin (2023) looks at what happens to off-label use of cancer drugs when new drug trial information becomes available. She finds that physician responses are sensitive to whether the p-value is less than 0.05. When the effect of the drug is deemed statistically significant, demand doubles in the year after the finding. If the drug is found not to have a significant effect, demand falls by a third over the next two years. Avdic et al. (2024) also find asymmetric responses to new information. Their study focuses on drug-eluting stents used in heart surgery. The new stents were first thought to be an improvement, and then shown to be inferior to older stents. Using Swedish data, Avdic et al. (2024) show that doctors were slow to take up the new stents but abandoned them quickly when the new information about their potentially harmful side effects came out.

Howard and Hockenberry (2019) ask how the uptake of new information from clinical studies is affected by physician age. The specific example is new information about episiotomies from clinical studies showing that they were ineffective in reducing complications of labor and delivery. They found that doctors with over 10 years of experience were much less likely to change their practice in response to the new information. However, they also found that the gap between new and old doctors was smaller in teaching hospitals, which are more likely to promote the adoption of evidence based medical practices.

Wu and David (2022) provide an example that fits nicely into the theoretical framework laid out above. They consider the choice of minimally invasive vs. “open” surgical procedures for hysterectomy. In 2014 the

FDA made an announcement that the minimally invasive procedure had a previously unappreciated risk of spreading a rare form of cancer. This announcement changed the expected benefit of the invasive relative to the non-invasive procedure ($\Delta_{0NIj}/\Delta_{1Ij}$). But the authors point out that this ratio also depends on the surgeon’s relative skill performing the two procedures. While overall use of the minimally invasive procedure fell, it actually rose among the subset of surgeons who were much better at performing the minimally invasive procedure than the open procedure.

Together with the “letter experiments” discussed above, these studies suggest that doctors pay more attention to some types of new information than others and that the impact of new information can vary with characteristics such as experience and skill. An important question going forward is what factors make information salient and whether these factors vary with physician characteristics in a predictable way.

Information provided to both physicians and consumers in forms such as “quality report cards” can also influence physicians. Kolstad (2013) considers two potentially important effects of the introduction of new report cards for coronary artery bypass graft (CABG) surgery. Report cards create an “extrinsic” incentive for surgeons to improve their scores in order not to lose business. But knowing how they are doing relative to other surgeons may also spur physicians to improve their practices for the “intrinsic” reason that they get utility from positive patient outcomes and realize that they could be doing better. Kolstad (2013) estimates a structural model of consumer demand in order to separate intrinsic from extrinsic motivations. Improvements made in response to predicted changes in consumer demand are thought to reflect extrinsic motivation, while the remaining change in doctor behavior after report cards are introduced is defined as change due to intrinsic motivation. He finds that intrinsic motivation is more important than extrinsic considerations, and that the response is greatest for physicians who are revealed to be worse than other surgeons in their own hospitals. This last finding opens the door for a third type of motivation – possibly surgeons who are worse than other surgeons in their own hospital fear loss of business or penalties due to interventions from hospital administrators. Alternatively, physicians may perceive other physicians who are more like themselves as a more relevant comparison group.

Finally, one can ask how extraneous information affects doctor decision making. Persson, Qiu, and Rossin-Slater (2021) ask how a sibling’s diagnosis of Attention Deficit Hyperactivity Disorder (ADHD) affects the probability that a child is diagnosed with the disorder. Since ADHD is heritable, sibling diagnoses could be correlated for legitimate reasons. Hence, the authors use variation in diagnoses induced by the interaction of the child’s birthdate with cutoff dates for school entry. It is well established that children who are “young-for-grade” are more likely to be diagnosed. These excess diagnoses are presumably spurious since small differences in children’s birthdates should not affect their underlying probability of having ADHD. Persson, Qiu, and Rossin-Slater (2021) show that the extra diagnoses induced by being young for grade spillover onto siblings since siblings are subsequently also more likely to be diagnosed with ADHD. This example shows that doctor’s decisions can be influenced by extraneous and perhaps erroneous information.

In sum, the research discussed in this section shows that information provision can have an impact on practice style. However, information provision does not eliminate undesirable variations in practice and does not always even lead to changes in the right direction. In view of the fact that a “helicopter drop” of information will not always have the desired effect, we next consider the role of various types of heuristics and guidelines.

5.2 Heuristics and Guidelines

Simon (1957) introduced the idea that because people are boundedly rational, they often take mental short-cuts and apply simple rules as aids in decision making. The properties of these rules, or heuristics, were further explored by Daniel Kahneman and Amos Tversky in many works (but see especially, Kahneman, Slovic, and Tversky (1982)). Heuristics are powerful because they often work well, though following them religiously can lead to systematic errors. We will use the term “guideline” to denote something more formal than a heuristic in that it is a set of rules that is laid down by an authority such as a professional association or a government agency. Guidelines usually do not have the force of law and there are typically few or no penalties for violating them, but they do provide clear expectations about appropriate (or inappropriate) behavior. Appendix Table 7 provides an overview of two types of studies. The first ask whether doctors tend to follow simple heuristic rules and if so, what effect they have on patient outcomes. The second group of studies asks whether patient outcomes would be improved by physician adherence to guidelines.

The use of simple decision rules is a ubiquitous human behavior, especially when complex decisions must be made under time pressure. So it would be surprising if doctors did not use them. What health economists have brought to the table is convincing evidence that these heuristics not only exist in medicine but can have important consequences for patient outcomes. In an ingenious paper, Almond et al. (2010) look at the treatment of newborns with birthweights on either side of a 1500 gram threshold that is used to define “very low birthweight.” They show that infants just below the threshold receive more medical care and are more likely to survive than infants just above the threshold. This result suggests that many infants above the threshold are erroneously denied the care that could save them because of a too literal adherence to the decision rule implied by the 1500 gram cutoff. Infants around the 1500 gram cutoff may be more or less sick depending on additional factors such as lung development. Closer attention to other indicators in addition to birth weight could improve the targeting of care. In a comment, Barreca et al. (2011) show that the regression discontinuity design employed by Almond, Doyle, Kowalski, and Williams (2010) is sensitive to measurement error (heaping) in birthweights at the threshold. However, Almond et al. (2011) show that their main results are robust to the use of a “donut” design that excludes observations that are very close to the threshold.

Geiger, Clapp, and Cohen (2021) use a similar donut regression discontinuity design to examine the effect of a designation of “advanced maternal age” (AMA) for pregnant women aged 35 or more on their expected delivery dates. They find that AMA mothers receive more screening and specialty visits and that this additional care has a large effect on perinatal mortality (infant death in the first month). As in Almond et al. (2010) this result suggests that rigid reliance on a simple heuristic based only on maternal age harms some patients who would have benefited from more care. The effects are greatest for pregnancies without obvious risk factors, suggesting that many apparently low-risk women would have to be more intensively screened and treated in order to prevent the marginal deaths.

Currie, MacLeod, and Van Parys (2016) find that doctors treating heart attack patients in Florida also appear to rely on age to ration treatment. They are less likely to treat older patients aggressively, even though all patients benefit from aggressive treatment in terms of a reduced risk of hospital readmission and mortality. Olenski et al. (2020) look more specifically at CABG surgery for heart patients using a regression discontinuity around a patient’s 80th birthday and find that patients admitted in the two weeks after their birthday are 28 percent less likely to receive this surgery than patients admitted in the two weeks

before. Coussens (2022) uses a regression discontinuity design to see whether the probability of being tested, diagnosed, or admitted for ischemic heart disease is higher when a patient is over 40. The results suggest that testing increases almost 10 percent at age 40 while diagnoses and admissions increase by 20 percent at this age threshold. Effects are larger for patients presenting without chest pain and for female patients, who are less likely to experience the stereotypical symptoms of heart disease. One might expect doctors to be more likely to use heuristics when they were busy but Coussens (2022) finds the reverse — the effect of the age threshold is larger when the ED is less busy and in the first half of the doctor’s shift. Geiger, Clapp, and Cohen (2021), Olenski et al. (2020), and Coussens (2022) all highlight that physicians have a tendency to “think discretely” about continuous patient characteristics such as age.

These articles provide strong evidence that doctors use simple heuristic cutoffs for providing care and that they do not necessarily assess each patient individually on the merits of their cases. Moreover, these decisions matter for patient health outcomes. However, this observation does not necessarily imply that heuristics are undesirable or inefficient. Only in a world with unlimited time and resources would we not want (or need) to use them. An important question then is whether these simple rules could be enriched in a way that meaningfully improves doctor’s choices and patient outcomes.

Guidelines tend to be more complex than simple heuristics and may be especially helpful for decisions that do not involve a simple zero-one choice. For example, Currie and MacLeod (2020) consider guidelines for drug treatment of adult depression. There are many possible treatment choices and it is not possible to know a priori which drug is best for a particular patient. In this case, there may be a trade-off between choosing the drug with the highest expected value and experimenting to find a drug that may be better for a particular patient. The downside of experimentation is that it can expose patients to the risk of poor outcomes because many drugs have side effects. A novel implication of their model is that experimentation is only useful if the doctor has enough diagnostic skill to learn from it and is willing to change their underlying beliefs about the efficacy of the treatment. Using claims data, they show that patients of more skillful doctors (psychiatrists) benefit from experimentation, while patients of less skillful doctors (GPs treating mental illness) derive little benefit from experimentation. The model predicts that higher diagnostic skill leads to greater diversity in drug choices across patients and better matching of drugs to patients even among doctors with the same initial beliefs regarding drug effectiveness. They also show that conditional on doctor skill, increasing the number of drug choices predicts poorer patient outcomes, by making it more likely that a drug that is a bad match will be chosen.

An important question is whether the use of guidelines could improve outcomes for complex and difficult to treat conditions? Medical guidelines vary from being very prescriptive (e.g., all heart failure patients should get beta blockers unless there are contraindications) to being rather loose and aimed not at mapping specific actions to specific conditions but at eliminating harmful choices. For example, a guideline might recommend that doctors do not use drug “cocktails”, without specifying which drugs they can use. Such guidelines may come from government agencies (as in the case of the English National Institute for Health and Care Excellence) or from professional associations such as the American Psychiatric Association. As in the case of heuristics, guidelines are usually not compulsory though physicians who violate guidelines could in some cases expose themselves to legal liability. Currie and MacLeod (2020) explore the rather loose guidelines that the American Psychiatric Association has drafted for adult depression treatment. These guidelines focus on changing drugs when an initial drug is found to be ineffective and on the use of “drug

cocktails.” They show that the patients of physicians who violate these guidelines have significantly worse outcomes than other patients.

Cuddy and Currie (2020) focus on guidelines for treatment of adolescent depression and anxiety patients. These guidelines are considerably more detailed and more prescriptive than those governing treatment of adults. Using claims data, they show that guideline violations are widespread. Cuddy and Currie (2024) build on this work by showing that these guideline violations are consequential. In order to deal with the possibility that patients are demanding treatment that violates a guideline, the treatment received is instrumented using measures of local practice style interacted with patient characteristics. The large number of possible instruments generated by this process is winnowed using the post-lasso TSLS procedure suggested by Belloni, Chen, Chernozhukov, and Hansen (2012). They find that individual patients who receive treatment that violates guidelines have higher health care costs over the next two years and are more likely to have been seen in the ED and to have been hospitalized. These results suggest that patients would indeed be better off if doctors followed professional guidelines.

Abaluck et al. (2021) asks several additional questions about the use of guidelines. First, when guidelines change, how quickly do doctors update their practice style? Second, if doctors fail to update, is this because they are unaware of the changes or is it for other reasons? Third, are some violations of guidelines justified by treatment effect heterogeneity? The context they study is the prescription of anticoagulants for patients with atrial fibrillation. Guidelines for treating these patients changed in 2006. Data from eight randomized controlled trials is available to try to explore treatment effect heterogeneity. They measure doctor awareness of the new procedures by using text mining of electronic medical records in order to find the first time the doctor mentioned them. After that date, the doctor is assumed to be aware of the new guidelines. The results suggest that doctors do move towards the new guidelines, but that adherence is highly imperfect. They estimate that stricter adherence to the new guidelines could have prevented 24 percent more strokes. They also find that departures from the guidelines do not seem to be justified by treatment effect heterogeneity, though the RCTs were not originally randomized on the dimensions of heterogeneity that they explore. Hence, even in cases where following guidelines has a clear health benefit, it appears to be difficult to achieve compliance.

Shurtz et al. (2024) have a similar finding with respect to colonoscopies. They find that when a doctor’s patient receives an unexpected colon cancer diagnosis, doctors are more likely to screen patients appropriately. But only for three months. The effects dies off by twelve months after the information shock.

Kowalski (2023) raises an additional issue – what if the guidelines are followed, but are flawed? She studies U.S. mammography screening guidelines which specify that women between ages 40 and 50 can make an individual decision in consultation with their doctors about whether mammography is warranted. Other countries, including Canada, recommend against the screening of asymptomatic women aged 40 to 50. The data come from a large Canadian RCT that assigned some women to a treatment group offered mammograms between 40 and 50. The control group were not offered mammograms at those ages. A novel feature of her analysis is that she differentiates between the rates of over diagnosis for women who always got mammography regardless of their assignment to treatment and control; women who are more likely to get mammograms if they are in the treatment group (the “compliers”); and those who never received mammograms regardless of their treatment status. She finds that under the voluntary screening regime, the women who are screened are disproportionately healthier and of higher socioeconomic status. Moreover, 14 percent of the cancers

uncovered in the complier group are “over diagnosed” in the sense that they were noninvasive cancers that would never have led to symptoms if they had remained undetected, while 36 percent of the cancers detected in the group that always took mammograms were over-diagnosed. The results imply that, if compliers in the U.S. are similar to those in Canada, bringing the U.S. guidelines into compliance with those of countries would be beneficial in the sense that it would eliminate over diagnosis that leads to harmful over treatment.

In sum, the limited economic research available suggests that guidelines have the potential to improve outcomes if doctors can be persuaded to follow them, and if they can be updated in a timely way when new knowledge becomes available. It is not known how current clinical practice is shaped by guidelines or what measures would be most effective in promoting adherence to guidelines. Finally, there has been little research on the optimal form of guidelines. Should they be very proscriptive (e.g. checklists), or should they be more in the nature of guardrails that make some treatments off limits but allow flexibility within relatively broad limits?

5.3 Can Technology Improve Medical Decision Making?

It may seem obvious that technology can improve medical decision making. For example, the invention of the mammogram meant that in many cases, doctors could tell whether a lump was likely to be cancerous or not. But as Kowalski’s study illustrates, a new tool can be overused or underused. Moreover, the use of the tool may expose patients to other dangers such as radiation and unnecessary surgery or chemotherapy in the case of mammograms. There is a large literature on the overuse of imaging technology more generally. As just one recent example, by comparing across state borders with and without Certificate of Need laws, Horwitz et al. (2024) show that such laws can reduce the probability of receiving low value magnetic resonance imaging without affecting high-value imaging. However, the same laws reduce the probability of getting even high-value computerized tomography scans, though low-value use falls even more.

In this section we focus on several technologies that have been touted as having the potential to revolutionize medicine. We focus on telemedicine (or telehealth), the adoption of electronic medical records (EMRs) and prescription drug monitoring programs (PDMPs); and on the use of algorithms to assist decision making. Some of the many studies in these areas are summarized in Appendix Table 8.

Telehealth is a technology with potentially widespread effects on medical decision making. Zeltzer et al. (2023) evaluate the introduction of a home health device that facilitated primary care visits by telemedicine by allowing patients to collect and upload some basic health data. The device reduced urgent care, ED, and inpatient visits and increased primary care visits, suggesting increases in the efficiency of medical care delivery, though it also increased the use of antibiotics substantially. Zeltzer et al. (2024) treat the COVID-19 pandemic as a shock that increased access to telemedicine in Israel in a long-lasting way. They find increases in primary care visits but a reduction in overall costs. There was no evidence of increases in missed diagnoses.

Dahlstrand (2024) suggests that telemedicine has the potential to improve patient outcomes by allowing sick patients to access skilled doctors regardless of their location. She estimates that matching patients at risk for avoidable hospitalization with the most skilled doctors would lead to an 8 percent reduction in such hospitalizations. It remains to be seen whether these kinds of hypothetical gains can be realized.

Goetz (2023) examines the impact of a change in the algorithm that provides patients with information about on-line talk therapists. Initially, the platform only displayed providers in the patient’s area. The

change occurred in areas with fewer than 20 providers. It allowed patients in these areas to see information about providers in other areas. He shows that the change caused the most skilled providers to stop offering sliding fees on-line, while less skilled providers were more likely to exit the platform. Presumably the former started receiving more requests for fee discounts, while the latter lost patients to out-of-area providers. These results suggest that the market for telehealth is sensitive to seemingly small differences in platform architecture. Both Dahlstrand (2024) and Goetz (2023) also highlight the potential for telehealth to change the boundaries of health care markets. Such a change could affect provider competition and, potentially, patient outcomes.

High quality information about a patient's condition is essential to patient care, whether it is provided in person or via telemedicine. The development of EMRs may enable and incentivize doctors to keep better records as well as facilitating the coordination of care across providers. In some cases, EMRs are combined with other types of decision support tools. In the U.S., the use of EMRs was incentivized by the 2009 HITECH Act, which was itself part of the federal government's response to the Great Recession. The Act set goals for the adoption of EMRs and gave providers financial incentives to encourage them to meet these goals. In retrospect, it is unfortunate that the Act did not set standards for the interoperability of different EMR systems. Today, while most providers use EMRs, there are many incompatible programs in use, limiting the extent to which EMR adoption can reduce the fragmentation of care. Other countries, such as England have also struggled to implement unified, inter-operable systems (Wilson and Khansa, 2018).

Most economic studies of EMRs have focused on whether adoption has improved the quality of care. Even in the absence of better care coordination, EMRs could improve the care provided by individual clinicians. By requiring doctors to fill in certain fields an EMR might prompt them to think about attributes of patients or care options that they would otherwise have neglected. An EMR might also lead to better care coordination within a practice or hospital, which could improve outcomes. A third possibility is that a more comprehensive track record encourages doctors to take more care lest they should be accused of malpractice. On the other hand, EMRs have proven unpopular with many clinicians who complain of information overload. One survey of primary care physicians in the U.S. Veterans Health Administration found that close to 90 percent of doctors found the number of alerts they received excessive, and over half of respondents said that the flood of information made it possible to overlook important information (Singh et al., 2013).

In one of the first papers on this topic, McCullough et al. (2010) examined the impact of EMR adoption on hospital-level (and hospital reported) measures of the quality of care. They find that only two of the measures they examine show any impact. Agha (2014) uses individual-level data from Medicare claims data and examines the impact of EMR adoption in models with hospital fixed effects. She finds that adoption increased health care spending by 1.3%, but had no impact on length of stay, intensity of care, care quality, readmissions, or 1-year mortality. In contrast to these two studies, Miller and Tucker (2011) use county-level data to examine the impact of EMR adoption over the 1995-2006 period. EMR adoption is instrumented using state medical privacy laws. They argue that by inhibiting the sharing of information, such laws make EMR adoption less attractive. They find that a 10% increase in EMR adoption reduces neonatal mortality by 3%. These reductions are due to prematurity and complications of labor and delivery and not to accidents, sudden infant death syndrome, or congenital defects. A caveat is that they cannot observe whether a particular baby was actually delivered in a hospital with EMRs, and there might be other changes in medical care in counties that happened to be rapid adopters of EMRs.

One interesting potential use of EMRs is to identify areas of particular concern so that they can be targeted for improvements. For example, in 2006, the State of California began an initiative to reduce maternal mortality. The first step was to identify hospitals with high rates, and to determine the most important cause of death for each hospital. For example, if a lot of people were dying of hemorrhage, there would be specific training of hospital staff aimed at that cause and a "crash cart" with everything necessary to intervene in one place (Main et al., 2020). This initiative reduced maternal mortality in California by 65 percent from 2006 to 2016, while rates continued to increase in the rest of the U.S.¹¹

PDMPs can be thought of as a specific and limited type of EMR. A PDMP is a state-level electronic registry of prescriptions for "scheduled" drugs (such as opioids). PDMPs can be searched by doctors, administrators, or law enforcement (depending on state rules) in order to identify patients or doctors who are using or prescribing drugs improperly. Because they are run at the state level, they come in many different flavors, but one of the most important distinctions is whether doctors are required to access the PDMP before prescribing. Several studies have found that the adoption of "must access" PDMPs reduced prescribing of opioids but had limited impacts on outcomes such as overdose deaths (Buchmueller and Carey (2018); Sacks et al. (2021); Neumark and Savych, 2023). One possible reason that PDMP adoption might have limited initial effects on overdoses is that it may take some time for a new opioid prescription to lead to addiction and death, so that the standard difference-in-differences framework may not be well suited to capturing these delayed effects.

Alpert, Dykstra, and Jacobson (2024) interpret a must access PDMP as something that imposes an additional "hassle cost" on prescribing compared to a PDMP that is not must access. They argue that if the PDMP operated mainly by providing information to prescribers about patients who were abusing opioids, then it should have no effect on opioid-naive patients. However, they show that the adoption of a must access PDMP affects both opioid naive and non-opioid naive patients, though it affects the latter more. They also note that patients who needed opioids the most, such as cancer patients, still received them so that adding an additional cost of prescribing improved targeting of treatment. They conclude that hassle costs, rather than increases in information available to providers, explain most of the observed decline in opioid prescribing. Another interpretation of these results is that the mere implementation of a must access PDMP may provide a signal to physicians about the risks associated with opioid prescribing.

In terms of other outcomes, Sacks et al. (2021) observe that PDMPs do not significantly affect "extreme use such as doctor shopping among new patients, because such behavior is very rare." This finding is ironic because the idea that addicted patients were "doctor shopping" to obtain multiple prescriptions of dangerous medications was one of the prime motivations for the creation of PDMPs.

Another technological approach to improving decision making is to use an algorithmic decision tool. Interest in using algorithms to assist physician decision making dates back at least to Meehl's 1954 book on the subject and the seminal article by Ledley and Lusted (1959) in *Science*. It is worthwhile to briefly discuss what an algorithm is, especially give the recent interest in large language models and their potential future impact on all areas of the labor market, including the market for doctors.

All algorithms are functions that take data in numerical form and produce a numerical output. For example, in the case of large language models, the text is mapped into a high dimension vector space (\mathfrak{R}^n , where n is a large number), and then transformed via a sequence of mathematical operations to produce

¹¹See <https://www.cmqcc.org/who-we-are>.

an output. In the context of a binary choice, as in the theory model introduced here, the output can be a probability that intensive treatment is best, say $\eta(x_i)$, where x_i is the vector representing all the information known about patient i . Then the algorithm will recommend intensive treatment if and only if $\eta(x_i) > 1/2$.¹²

Humans also make decisions based on data. Moreover, humans can quickly process vast quantities of information but much of that information comes through the visual field. Decades of research has shown, that in contrast to computers, humans cannot rapidly process large volumes of *numerical* information. When the numerical information provides a more accurate assessment of the benefits from a decision, then algorithms, even algorithms based on simple linear regressions, can perform better than a human decision maker.¹³

Thus, when good numerical data is available we should expect algorithms to provide high quality recommendations that can improve upon human decision makers. Ludwig, Mullainathan, and Rambachan (2024) point to the algorithm Mullainathan and Obermeyer (2022) developed to predict who should be tested for heart attacks and argue that the adoption of such an algorithm would amount to a “free lunch” in the sense that the social benefit would greatly outweigh the cost.

Yet, humans are capable of processing large volumes of visual data and making decisions in real time. A good doctor can tell at a glance that a wound is infected, or that a patient has hepatitis. The fact that humans are very good at processing visual information implies that in some cases the doctor is simply the most efficient agent to collect and act on information. For example, a patient coming into an ED may require immediate help, for example intravenous fluids. Getting the person’s weight and vital signs for the EMR takes time that might not be available. The attending doctor can estimate the patient’s weight and condition in less than a second, and then order or execute treatment. As Kahneman and Klein (2009) observe, there are many examples of experts with extraordinarily high levels of skill, and in principle, both algorithms and skilled experts can play a role in improving decision making. At the same time, as the evidence reviewed above illustrates, there is a great deal of variation in doctor skill. The question then is how best to incorporate the benefits of well designed algorithms, while also exploiting the knowledge of highly skilled doctors.

This problem turns out to be quite difficult. Agarwal et al. (2023) conducted a randomized experiment with radiologists who were asked to retrospectively diagnose patients in a laboratory setting that resembled their usual working environment. In some cases they received only an x-ray, while in other cases they were given either an AI prediction, additional contextual information about the patient’s history that was not considered by the AI tool, or both. The AI algorithm used has been shown to perform similarly to professional radiologists. The experimental subjects’ diagnoses were then compared to “ground truth” derived using the opinions of five expert radiologists. Agarwal et al. (2023) find that giving radiologists the AI prediction did not improve diagnostic accuracy, while giving them additional contextual information did. They estimate a model of belief updating which suggests that clinicians erroneously treat the AI prediction as independent of their own information, which causes it to bias their decision making. They argue that better results could have been achieved by using the AI prediction in cases where the tool had high confidence and allowing humans to make decisions without AI assistance in all other cases.

The problem of how to optimally combine algorithmic information and expert opinion arises in many other settings. For instance, Stevenson and Doleac (2022) find that judges given algorithmic assessments

¹²See Devroye, Györfi, and Lugosi (1996) on the mathematics of machine learning. Bengio, Lecun, and Hinton (2021) provide an up to date discussion of machine learning by three seminal contributors to the field.

¹³Kahneman (2003) noted in his Noble Prize lecture that he first recognized this point in the 1950s while work for the Israeli military. The seminal contribution by Dawes, Faust, and Meehl (1989) makes this point in the context of medical decision making.

of the probability of recidivism change their sentencing decisions but that use of the tool did not either reduce incarceration or improve public safety. Judges deviated from the algorithm in a way that increased incarceration but also reduced recidivism. Hoffman et al. (2018) look at manager hiring decisions before and after the introduction of formal job testing algorithms. In this case they find that managers who overrule the algorithmic recommendation hire worse people on average. Rambachan (2024) adds to the literature on bail decisions, arguing that well designed algorithms can improve upon judicial decisions.

The performance of AI models currently in clinical use is similarly mixed. (Obermeyer et al., 2019) point out that an algorithm that aims to minimize health care costs will tend to short-change Black patients if Black patients are under-treated in the training data. The issue is that minimizing health care costs does not necessarily maximize health outcomes so in a sense the algorithm was trained on the wrong objective. The authors point out however that it may be easier to correct such a problem in an algorithm than it is to get human decision makers to show less bias in the allocation of treatments.

Manz et al. (2023) conducted a large randomized trial to see whether a machine-learning generated nudge could encourage clinicians to engage in end-of-life conversations with terminally ill cancer patients. They found an increase in such conversations, and a reduction in systemic cancer therapy at the end of life, but no change in hospice, length of stay, or intensive-care admission at the end of life.

In another example, data from one of the largest purveyors of EMRs, EPIC, has been used to develop an AI tool for diagnosing sepsis. The tool has been adopted in hundreds of hospitals. However, Wong et al. (2021) found that the algorithm performed poorly in a large teaching hospital setting. It failed to identify 67 percent of patients with sepsis even though it generated an alert for 18 percent of all patients. Lyons et al. (2023) followed up on this finding by examining the performance of the tool in nine networked hospitals. They found that the tool did better in hospitals treating less sick patients with a lower average probability of sepsis.

As this example illustrates, even if an algorithm is trained on big data, it may not perform very well if the sample at hand is quite different than the one used to train the algorithm. Although economists have been aware of the selection problem since the famous work of Roy (1951) on wages and the self selection of workers to occupations, awareness of the selection problem in the machine learning literature is very recent (see Athey and Imbens (2019)). Many modern machine learning algorithms in medicine have access to large amounts of data, with patients who are allocated to different treatments. The problem is that if one does not incorporate the allocation (selection) mechanism in the machine learning model, then the predicted effects of treatment may be incorrect. For example, if clinicians only give an experimental treatment to the patients they believe are most likely to recover, then the effectiveness of the treatment is likely to be overstated. Moreover, Rambachan and Roth (2020) show that even if one knows the direction of the selection bias in the underlying data, the bias in the algorithm can be in any direction. This observation highlights the point that learning from large data sets requires more than simply choosing the right algorithm. It also entails understanding how the sample is selected, and testing the results in much the same way that new drugs are tested via rigorous randomized control trials.

In summary, these three new technologies, telemedicine, EMRs, and algorithmic decision tools, have considerable promise, but the available evidence suggests that the details of how they are implemented really matter. More research is required to understand how to use them to actually improve patient welfare.

6 Conclusions and Suggestions for Future Research

Modern medicine is a miracle. Human life span and well being has increased dramatically in the last century, in part due to rapid advances in science and medicine. But these advances have also increased the complexity of medical decision making. In a world where there was little that could be done for most ailments, there were few consequential decisions to be made. Today, medical decision making matters more than ever. This review highlights the variation in how doctors treat medically similar patients. The economics of doctor decision making is concerned with how to allocate resources to both increase the average quality of medical services provided, and to provide consistent quality to all individuals.

The model we have outlined has several moving parts. Doctors are assumed to care about patient welfare, but also about their own welfare which makes them imperfect agents. Doctors arrive at the bedside with a given training and experience, which results in a set of skills as well as prior beliefs about proper treatment. As humans, doctors are influenced by fatigue, time pressures, emotional states, prejudices, and peer effects. They may rely on simple decision rules in cases where it would be optimal to devote more focused attention. At present, no one has estimated a model that parses out the role of differences in patient population (α_i), doctor diagnostic skill (γ_j), procedural skill (s_{tj}), pecuniary factors (λ_{tj}) and doctor practice style (τ_j) in explaining why doctors vary in treatment rates, making this a possible goal for future research. As we have highlighted, existing models shut down one or more of these channels.

The fact that there are so many factors that affect medical decision making suggests that there is no one policy lever that will optimize care. In particular, the research reviewed here indicates that it can be difficult to tweak payment systems in a way that will have unambiguously positive effects on the allocation of medical care. Future work on the impacts of changes in payment systems (and other levers) should pay careful attention to heterogeneity in the effects on patients.

Similarly, it is unlikely that the quality of medical decision making can be greatly improved by additional training, at least in the short term. In part, this is because there is remarkably little research on the actual content of training, either in medical school or for doctors in practice, so it is hard to know what works. Economic studies tend to focus on crude measures such as years of training or type/rank of medical school. Moreover, chronic doctor shortages in many countries suggest that there will be continuing demand for the services of even the least skilled physicians, which may attenuate any incentives for continuous skill improvement. Short anti-bias trainings offer an interesting case in which the impact of a specific form of training has been evaluated and found to have little impact on physician behavior. Vela et al. (2022)'s hypothesis that the effect of the specific anti-bias training is counter-acted by the messages implicit in the rest of a doctor's training suggests that is necessary to better understand doctor training as a whole. Enhancing medical decision making by improving the concordance between the characteristics of doctors and patients will also take a long time. Research into other ways to enhance sympathy and communication between doctors and patients is sorely needed.

The fact that poor medical decision making is difficult to address with payment reforms or training (given what we now know about training effects) accounts for much of the excitement about guidelines, algorithms, and other emerging health care technologies among health economists. As researchers, we tend to have faith in the efficacy of providing information to economic agents, but the evidence reviewed here indicates that doctors pay more attention to some types of new information than others. Information provision alone does not eliminate undesirable variations in practice and does not always even lead to changes in the right

direction. A key question going forward is what factors make information salient and how these factors vary with physician characteristics.

Research suggests that adherence to clinical guidelines is helpful for patients, at least where the guidelines themselves represent best practice. But it is not known how current clinical practice is shaped by guidelines or what measures are most effective in promoting adherence to guidelines. There has also been little research on the optimal form of guidelines. Should they be very proscriptive (e.g. checklists), or should they be more in the nature of guardrails that make some treatments off limits but allow flexibility within relatively broad limits? Telemedicine, EMRs, and algorithmic decision tools, have considerable promise, but it is fair to say that we do not yet understand how to implement them to improve patient welfare. Like older medical technologies, these new tools can be over-used, under-used, and can lead to harmful consequences for patients when used inappropriately. Understanding how humans can interact with the tools to produce better outcomes is a first order question.

Many of the themes we highlight here are relevant to other labor markets with high-skilled workers. Health care data offer unique data opportunities to observe both physician decisions and consequences for patients. The literature we discuss speaks to questions about labor productivity, organizational economics, and the use of technology that are often impossible to analyze in other settings if only because it is usually so difficult to see the downstream consequences of an expert decision. That said, it can also be extremely difficult to see the consequences of expert decisions in medical settings particularly in the medium to longer-run. While most doctors are presumably aware that “correlation is not causation,” the principal is difficult to apply in practice. One of the most famous examples of the failure to distinguish correlation from causation in medicine is blood letting, a treatment that persisted for centuries even though it is now known to be more likely to harm than help patients.¹⁴

The empirical work we have reviewed wrestles with ubiquitous selection problems. Patients select doctors and may also choose procedures. Doctors may select patients. Medical schools and training programs select applicants. Doctors select peers. The most successful papers in this literature identify situations that approximate random assignment to one doctor or another, one treatment or another, or to a particular medical team in order to achieve causal identification.¹⁵

This work has shown both that different doctors treat medically similar patients differently, and that individual doctors often treat similar patients differently depending on patient characteristics such as age, race, and gender, or on time-varying doctor specific factors such as the time in their shift or the presence of peers. One caveat is that much of this work focuses on elderly Medicare patients for reasons of data availability, so extending these results to other populations and settings would be useful. A second caveat is that even when we can identify causal effects, it is difficult to understand the precise mechanisms and motivations underlying doctor decisions. Better understanding of mechanisms is necessary for the development of effective interventions.

Improving medical decision making is fundamentally an economic as well as a medical question. Training doctors, as well as collecting and analyzing information regarding patients to determine effective treatment is time consuming and hence costly. The efficient and equitable management of health services entails allocating

¹⁴See Parapia (2008) for an interesting history of attitudes towards blood-letting as a medical practice. In an era when many sick patients died, if a few patients survived after blood letting this might reinforce doctor beliefs in the benefits of the treatment.

¹⁵See Holland (1986) for discussion of the basic concepts, and Imbens and Rubin (2015) for a book length treatment of causal identification.

resources to these expensive activities. This review illustrates that a great deal has been learned regarding doctor decision making, but a great deal more research is needed if we are to fully realize the benefits of integrating skilled doctors into the complex health services industry.

References

- Abaluck, Jason, Leila Agha, Jr Chan, David C., Daniel Singer, and Diana Zhu.** 2021. “Fixing Misallocation with Guidelines: Awareness vs. Adherence.” Working Paper 27467, National Bureau of Economic Research.
- Abaluck, Jason, Leila Agha, Chris Kabrhel, Ali Raja, and Arjun Venkatesh.** 2016. “The Determinants of Productivity in Medical Testing: Intensity and Allocation of Care.” *American Economic Review* 106 (12): 3730–3764.
- Agarwal, Nikhil, Alex Moehring, Pranav Rajpurkar, and Tobias Salz.** 2023. “Combining Human Expertise with Artificial Intelligence: Experimental Evidence from Radiology.” Working Paper 31422, National Bureau of Economic Research.
- Agha, Leila.** 2014. “The Effects of Health Information Technology on the Costs and Quality of Medical Care.” *Journal of Health Economics* 34 19–30.
- Agha, Leila, and David Molitor.** 2018. “The Local Influence Of Pioneer Investigators On Technology Adoption: Evidence From New Cancer Drugs.” *The review of economics and statistics* 100 (1): 29–44.
- Agha, Leila, and Dan Zeltzer.** 2022. “Drug Diffusion through Peer Networks: The Influence of Industry Payments.” *American Economic Journal: Economic Policy* 14 (2): 1–33.
- Ahammer, Alexander, and Thomas Schober.** 2020. “Exploring Variations in Health-Care Expenditures—What Is the Role of Practice Styles?” *Health Economics* 29 (6): 683–699.
- Ahomäki, Iiro, Visa Pitkänen, Aarni Soppi, and Leena Saastamoinen.** 2020. “Impact of a Physician-Targeted Letter on Opioid Prescribing.” *Journal of Health Economics* 72 1–22.
- Alexander, Diane.** 2020. “How Do Doctors Respond to Incentives? Unintended Consequences of Paying Doctors to Reduce Costs.” *Journal of Political Economy* 128 (11): 4046–4096.
- Alexander, Diane, and Janet Currie.** 2017. “Are Publicly Insured Children Less Likely to Be Admitted to Hospital than the Privately Insured (and Does It Matter)?” *Economics & Human Biology* 25 33–51.
- Alexander, Diane, and Molly Schnell.** 2024. “The Impacts of Physician Payments on Patient Access, Use, and Health.” *American Economic Journal: Applied Economics*.
- Almond, Douglas, Joseph J. Doyle, Jr., Amanda E. Kowalski, and Heidi Williams.** 2010. “Estimating Marginal Returns to Medical Care: Evidence from At-risk Newborns.” *The Quarterly Journal of Economics* 125 (2): 591–634.
- Almond, Douglas, Joseph J. Doyle, Jr., Amanda E. Kowalski, and Heidi Williams.** 2011. “The Role of Hospital Heterogeneity in Measuring Marginal Returns to Medical Care: A Reply to Barreca, Guldi, Lindo, and Waddell*.” *The Quarterly Journal of Economics* 126 (4): 2125–2131.

- Alpert, Abby, Sarah Dykstra, and Mireille Jacobson.** 2024. “Hassle Costs versus Information: How Do Prescription Drug Monitoring Programs Reduce Opioid Prescribing?” *American Economic Journal: Economic Policy* 16 (1): 87–123.
- Alsan, Marcella, Owen Garrick, and Grant Graziani.** 2019. “Does Diversity Matter for Health? Experimental Evidence from Oakland.” *American Economic Review* 109 (12): 4071–4111.
- Alsan, Marcella, and Marianne Wanamaker.** 2018. “Tuskegee and the Health of Black Men*.” *The Quarterly Journal of Economics* 133 (1): 407–455.
- Angerer, Silvia, Christian Waibel, and Harald Stummer.** 2019. “Discrimination in Health Care: A Field Experiment on the Impact of Patients’ Socioeconomic Status on Access to Care.” *American Journal of Health Economics* 5 (4): 407–427.
- Arnold, David, Will Dobbie, and Peter Hull.** 2022. “Measuring Racial Discrimination in Bail Decisions.” *American Economic Review* 112 (9): 2992–3038.
- Athey, Susan, and Guido W. Imbens.** 2019. “Machine Learning Methods That Economists Should Know About.” *Annual Review of Economics* 11 (1): 685–725.
- Avdic, Daniel, Stephanie von Hinke, Bo Lagerqvist, Carol Propper, and Johan Vikström.** 2024. “Do Responses to News Matter? Evidence from Interventional Cardiology.” *Journal of Health Economics* 94.
- Azoulay, Pierre.** 2002. “Do Pharmaceutical Sales Respond to Scientific Evidence?” *Journal of Economics & Management Strategy* 11 (4): 551–594.
- Badinski, Ivan, Amy Finkelstein, Matthew Gentzkow, and Peter Hull.** 2023. “Geographic Variation in Healthcare Utilization: The Role of Physicians.” Working Paper 31749, National Bureau of Economic Research.
- Balogh, Erin P., Bryan T. Miller, and John R. Ball.** eds. 2015. *Improving Diagnosis in Health Care*. Washington, D.C.: National Academies Press.
- Barreca, Alan I., Melanie Guldi, Jason M. Lindo, and Glen R. Waddell.** 2011. “Saving Babies? Revisiting the Effect of Very Low Birth Weight Classification.” *The Quarterly Journal of Economics* 126 (4): 2117–2123.
- Belloni, A., D. Chen, V. Chernozhukov, and C. Hansen.** 2012. “Sparse Models and Methods for Optimal Instruments With an Application to Eminent Domain.” *Econometrica* 80 (6): 2369–2429.
- Bengio, Yoshua, Yann Lecun, and Geoffrey Hinton.** 2021. “Deep Learning for AI.” *Communications of the ACM* 64 (7): 58–65.
- Berndt, Ernst R., Robert S. Gibbons, Anton Kolotilin, and Anna Levine Taub.** 2015. “The Heterogeneity of Concentrated Prescribing Behavior: Theory and Evidence from Antipsychotics.” *Journal of Health Economics* 40 26–39.
- Brekke, Kurt R., Tor Helge Holmås, Karin Monstad, and Odd Rune Straume.** 2018. “Socio-Economic Status and Physicians’ Treatment Decisions.” *Health Economics* 27 (3): e77–e89.
- Brekke, Kurt R., Tor Helge Holmås, Karin Monstad, and Odd Rune Straume.** 2019. “Competition and Physician Behaviour: Does the Competitive Environment Affect the Propensity to Issue Sickness Certificates?” *Journal of Health Economics* 66 117–135.

- Buchmueller, Thomas C., and Colleen Carey.** 2018. “The Effect of Prescription Drug Monitoring Programs on Opioid Utilization in Medicare.” *American Economic Journal: Economic Policy* 10 (1): 77–112.
- Button, Patrick, Eva Dils, Benjamin Harrell, Luca Fumarco, and David Schwegman.** 2020. “Gender Identity, Race, and Ethnicity Discrimination in Access to Mental Health Care: Preliminary Evidence from a Multi-Wave Audit Field Experiment.” Working Paper 28164, National Bureau of Economic Research.
- Cabral, Marika, and Marcus Dillender.** 2024. “Gender Differences in Medical Evaluations: Evidence from Randomly Assigned Doctors.” *American Economic Review* 114 (2): 462–499.
- Carey, Colleen, Michael Daly, and Jing Li.** 2024. “Nothing for Something: Marketing Cancer Drugs to Physicians Increases Prescribing Without Improving Mortality.” Working Paper 32336, National Bureau of Economic Research.
- Chan, David C.** 2016. “Teamwork and Moral Hazard: Evidence from the Emergency Department.” *Journal of Political Economy* 124 (3): 734–770.
- Chan, David C.** 2018. “The Efficiency of Slacking off: Evidence from the Emergency Department.” *Econometrica* 86 (3): 997–1030.
- Chan, David C.** 2021. “Influence and Information in Team Decisions: Evidence from Medical Residency.” *American Economic Journal: Economic Policy* 13 (1): 106–137.
- Chan, David C, Matthew Gentzkow, and Chuan Yu.** 2022. “Selection with Variation in Diagnostic Skill: Evidence from Radiologists.” *The Quarterly Journal of Economics* 137 (2): 729–783.
- Chan, Jr, David C., and Yiqun Chen.** 2022. “The Productivity of Professions: Evidence from the Emergency Department.” Working Paper 30608, National Bureau of Economic Research.
- Chandra, Amitabh, and Jonathan Skinner.** 2012. “Technology Growth and Expenditure Growth in Health Care.” *Journal of Economic Literature* 50 (3): 645–680.
- Chandra, Amitabh, and Douglas O. Staiger.** 2007. “Productivity Spillovers in Health Care: Evidence from the Treatment of Heart Attacks.” *Journal of Political Economy* 115 (1): pp.103–140.
- Chandra, Amitabh, and Douglas O. Staiger.** 2010. “Identifying Provider Prejudice in Healthcare.” Working Paper 16382, National Bureau of Economic Research.
- Chandra, Amitabh, and Douglas O Staiger.** 2020. “Identifying Sources of Inefficiency in Healthcare.” *The Quarterly Journal of Economics* 135 (2): 785–843.
- Chen, Alice, and Darius N. Lakdawalla.** 2019. “Healing the Poor: The Influence of Patient Socioeconomic Status on Physician Supply Responses.” *Journal of Health Economics* 64 43–54.
- Chen, Yiqun.** 2021. “Team-Specific Human Capital and Team Performance: Evidence from Doctors.” *American Economic Review* 111 (12): 3923–3962.
- Chernew, Michael, Zack Cooper, Eugene Larsen Hallock, and Fiona Scott Morton.** 2021. “Physician Agency, Consumerism, and the Consumption of Lower-Limb MRI Scans.” *Journal of Health Economics* 76 102427.
- Chodick, Gabriel, Yoav Goldstein, Ity Shurtz, and Dan Zeltzer.** 2023. “Challenging Encounters and Within-Physician Practice Variability.” *The Review of Economics and Statistics* 1–27.

- Chorniy, Anna, Janet Currie, and Lyudmyla Sonchak.** 2018. “Exploding Asthma and ADHD Caseloads: The Role of Medicaid Managed Care.” *Journal of Health Economics* 60 1–15.
- Chowdhury, M M, H Dagash, and A Pierro.** 2007. “A Systematic Review of the Impact of Volume of Surgery and Specialization on Patient Outcome.” *British Journal of Surgery* 94 (2): 145–161.
- Chu, Bryan, Ben Handel, Jonathan Kolstad, Jonas Knecht, Ulrike Malmendier, and Filip Matejka.** 2024. “Cognitive Capacity, Fatigue and Decision Making: Evidence from the Practice of Medicine.” Technical report, UC Berkeley.
- Clemens, Jeffrey, and Joshua D. Gottlieb.** 2014. “Do Physicians’ Financial Incentives Affect Medical Treatment and Patient Health?” *American Economic Review* 104 (4): 1320–1349.
- Costa-Ramón, Ana María, Ana Rodríguez-González, Miquel Serra-Burriel, and Carlos Campillo-Artero.** 2018. “It’s about Time: Cesarean Sections and Neonatal Health.” *Journal of Health Economics* 59 46–59.
- Coudin, Elise, Anne Pla, and Anne-Laure Samson.** 2015. “GP Responses to Price Regulation: Evidence from a French Nationwide Reform.” *Health Economics* 24 (9): 1118–1130.
- Coussens, Stephen.** 2022. “Behaving Discretely: Heuristic Thinking in the Emergency Department.”
- Cuddy, Emily, and Janet Currie.** 2020. “Treatment of Mental Illness in American Adolescents Varies Widely within and across Areas.” *Proceedings of the National Academy of Sciences* 117 (39): 24039–24046.
- Cuddy, Emily, and Janet Currie.** 2024. “Rules vs. Discretion: Treatment of Mental Illness in U.S. Adolescents.” *Journal of Political Economy*.
- Currie, Janet, Anastasia Karpova, and Dan Zeltzer.** 2021. “Do Urgent Care Centers Reduce Medicare Spending?.” July.
- Currie, Janet, Anran Li, and Molly Schnell.** 2023. “The Effects of Competition on Physician Prescribing.” Working Paper 30889, National Bureau of Economic Research.
- Currie, Janet M., and W. Bentley MacLeod.** 2017. “Diagnosis and Unnecessary Procedure Use: Evidence from C-section.” *Journal of Labor Economics* 35 (1): 1–42.
- Currie, Janet M, and W Bentley MacLeod.** 2020. “Understanding Doctor Decision Making: The Case of Depression Treatment.” *Econometrica : journal of the Econometric Society* 88 (3): 847–878.
- Currie, Janet, and W. Bentley MacLeod.** 2008. “First Do No Harm? Tort Reform and Birth Outcomes.” *Quarterly Journal of Economics* 123 (2): 795–830.
- Currie, Janet, W. Bentley MacLeod, and Jessica Van Parys.** 2016. “Provider Practice Style and Patient Health Outcomes: The Case of Heart Attacks.” *Journal of Health Economics* 47 64–80.
- Currie, Janet, and Jonathan Zhang.** 2023. “Doing More with Less: Predicting Primary Care Provider Effectiveness.” *The Review of Economics and Statistics* 1–45.
- Cutler, David M.** 2014. *The Quality Cure: How Focusing on Health Care Quality Can Save Your Life and Lower Spending Too.* Volume 9. of The Aaron Wildavsky Forum for Public Policy, Berkeley: University of California Press.

- Cutler, David, Jonathan S. Skinner, Ariel Dora Stern, and David Wennberg.** 2019. “Physician Beliefs and Patient Preferences: A New Look at Regional Variation in Health Care Spending.” *American Economic Journal: Economic Policy* 11 (1): 192–221.
- Dahlstrand, Amanda.** 2024. “Defying Distance? The Provision of Services in the Digital Age.”
- Dawes, Robyn M., David Faust, and Paul E. Meehl.** 1989. “Clinical Versus Actuarial Judgment.” *Science* 243 (4899): 1668–1674.
- Devroye, Luc, László Györfi, and Gábor Lugosi.** 1996. *A Probabilistic Theory of Pattern Recognition*. New York, NY: Springer-Verlag.
- Ding, Yu, and Chenyuan Liu.** 2021. “Alternative Payment Models and Physician Treatment Decisions: Evidence from Lower Back Pain.” *Journal of Health Economics* 80.
- Doctor, Jason N., Andy Nguyen, Roneet Lev, Jonathan Lucas, Tara Knight, Henu Zhao, and Michael Menchine.** 2018. “Opioid Prescribing Decreases after Learning of a Patient’s Fatal Overdose.” *Science* 361 (6402): 588–590.
- Doyle, Joseph J., Steven M. Ewer, and Todd H. Wagner.** 2010. “Returns to Physician Human Capital: Evidence from Patients Randomized to Physician Teams.” *Journal of Health Economics* 29 (6): 866–882.
- Doyle, Jr., Joseph J.** 2020. “Physician Characteristics and Patient Survival: Evidence from Physician Availability.” Working Paper 27458, National Bureau of Economic Research.
- Dubois, Pierre, and Tuba Tunçel.** 2021. “Identifying the Effects of Scientific Information and Recommendations on Physicians’ Prescribing Behavior.” *Journal of Health Economics* 78 102461.
- Dunn, Abe, Joshua D Gottlieb, Adam Hale Shapiro, Daniel J Sonnenstuhl, and Pietro Tebaldi.** 2024. “A Denial a Day Keeps the Doctor Away.” *The Quarterly Journal of Economics* 139 (1): 187–233.
- Eli, Shari, Trevon D. Logan, and Boriana Miloucheva.** 2019. “Physician Bias and Racial Disparities in Health: Evidence from Veterans’ Pensions.” Working Paper 25846, National Bureau of Economic Research.
- Epstein, Andrew J., Sean Nicholson, and David A. Asch.** 2016. “The Production of and Market for New Physicians’ Skill.” *American Journal of Health Economics* 2 (1): 41–65.
- Facchini, Gabriel.** 2022. “Forgetting-by-Not-Doing: The Case of Surgeons and Cesarean Sections.” *Health Economics* 31 (3): 481–495.
- Fadlon, Itzik, and Jessica van Parys.** 2020. “Primary Care Physician Practice Styles and Patient Care: Evidence from Physician Exits in Medicare.” *Journal of Health Economics* 71 102304.
- Fawcett, Tom.** 2006. “An Introduction of ROC Analysis.” *Pattern Recognition Letters* 27 861–874.
- Frakes, Michael.** 2013. “The Impact of Medical Liability Standards on Regional Variations in Physician Behavior: Evidence from the Adoption of National-Standard Rules.” *American Economic Review* 103 (1): 257–276.
- Frakes, Michael D., and Jonathan Gruber.** 2022. “Racial Concordance and the Quality of Medical Care: Evidence from the Military.” Working Paper 30767, National Bureau of Economic Research.

- Freedman, Seth, Ezra Golberstein, Tsan-Yao Huang, David J. Satin, and Laura Barrie Smith.** 2021. “Docs with Their Eyes on the Clock? The Effect of Time Pressures on Primary Care Productivity.” *Journal of Health Economics* 77 102442.
- Geiger, Caroline K., Mark A. Clapp, and Jessica L. Cohen.** 2021. “Association of Prenatal Care Services, Maternal Morbidity, and Perinatal Mortality With the Advanced Maternal Age Cutoff of 35 Years.” *JAMA Health Forum* 2 (12): e214044.
- Gibbons, Robert D., C. Hendricks Brown, Kwan Hur, Sue M. Marcus, Dulal K. Bhaumik, Joëlle A. Erkens, Ron M.C. Herings, and J. John Mann.** 2007. “Early Evidence on the Effects of Regulators’ Suicidality Warnings on SSRI Prescriptions and Suicide in Children and Adolescents.” *American Journal of Psychiatry* 164 (9): 1356–1363.
- Goetz, Daniel.** 2023. “Telemedicine Competition, Pricing, and Technology Adoption: Evidence from Talk Therapists.” *International Journal of Industrial Organization* 89 102956.
- Gowrisankaran, Gautam, Keith Joiner, and Pierre Thomas Léger.** 2022. “Physician Practice Style and Healthcare Costs: Evidence from Emergency Departments.” *Management Science*.
- Goyal, Monika K., Nathan Kuppermann, Sean D. Cleary, Stephen J. Teach, and James M. Chamberlain.** 2015. “Racial Disparities in Pain Management of Children With Appendicitis in Emergency Departments.” *JAMA Pediatrics* 169 (11): 996–1002.
- Greenwood, Brad N., Seth Carnahan, and Laura Huang.** 2018. “Patient–Physician Gender Concordance and Increased Mortality among Female Heart Attack Patients.” *Proceedings of the National Academy of Sciences* 115 (34): 8569–8574.
- Greenwood, Brad N., Rachel R. Hardeman, Laura Huang, and Aaron Sojourner.** 2020. “Physician–Patient Racial Concordance and Disparities in Birthing Mortality for Newborns.” *Proceedings of the National Academy of Sciences* 117 (35): 21194–21200.
- Gruber, Jonathan, Thomas P. Hoe, and George Stoye.** 2021. “Saving Lives by Tying Hands: The Unexpected Effects of Constraining Health Care Providers.” *The Review of Economics and Statistics* 1–45.
- Gupta, Atul.** 2021. “Impacts of Performance Pay for Hospitals: The Readmissions Reduction Program.” *American Economic Review* 111 (4): 1241–1283.
- Hammad, Tarek A., Thomas Laughren, and Judith Racoosin.** 2006. “Suicidality in Pediatric Patients Treated with Antidepressant Drugs.” *Archives of General Psychiatry* 63 (3): 332–339.
- Handel, Benjamin R., and Kate Ho.** 2021. “Industrial Organization of Health Care Markets.” August.
- Hill, Andrew J., Daniel B. Jones, and Lindsey Woodworth.** 2023. “Physician–Patient Race-Match Reduces Patient Mortality.” *Journal of Health Economics* 92.
- Hoffman, Kelly M., Sophie Trawalter, Jordan R. Axt, and M. Norman Oliver.** 2016. “Racial Bias in Pain Assessment and Treatment Recommendations, and False Beliefs about Biological Differences between Blacks and Whites.” *Proceedings of the National Academy of Sciences* 113 (16): 4296–4301.
- Hoffman, Mitchell, Lisa B Kahn, and Danielle Li.** 2018. “Discretion in Hiring.” *The Quarterly Journal of Economics* 133 (2): 765–800.

- Holland, Paul W.** 1986. "Statistics and Causal Inference." *Journal of the American Statistical Association* 81 (396): 945–960.
- Horwitz, Jill, Austin Nichols, Carrie H. Colla, and David M. Cutler.** 2024. "Technology Regulation Reconsidered: The Effects of Certificate of Need Policies on the Quantity and Quality of Diagnostic Imaging." February.
- Howard, David H., and Jason Hockenberry.** 2019. "Physician Age and the Abandonment of Episiotomy." *Health Services Research* 54 (3): 650–657.
- Imbens, Guido W., and Donald B. Rubin.** 2015. *Causal Inference for Statistics, Social, and Biomedical Sciences*. Cambridge University Press.
- Johnson, Erin M., and M. Marit Rehavi.** 2016. "Physicians Treating Physicians: Information and Incentives in Childbirth." *American Economic Journal: Economic Policy* 8 (1): 115–141.
- Kahneman, Daniel.** 2003. "Les Prix Nobel 2002." Chap. Autobiography, Stockholm, Sweden: Almqvist & Wiksell International.
- Kahneman, Daniel, and Gary Klein.** 2009. "Conditions for Intuitive Expertise: A Failure to Disagree." *American Psychologist* 64 (6): 515–526.
- Kahneman, Daniel, Paul Slovic, and Amos Tversky.** 1982. *Judgment Under Uncertainty: Heuristics and Biases*. Cambridge University Press.
- Kessler, Daniel, and Mark McClellan.** 1996. "Do Doctors Practice Defensive Medicine?" *Quarterly Journal of Economics* 111 (2): 353–90.
- Kling, Jeffrey R.** 2006. "Incarceration Length, Employment, and Earnings." *American Economic Review* 96 (3): 863–876.
- Kolstad, Jonathan T.** 2013. "Information and Quality When Motivation Is Intrinsic: Evidence from Surgeon Report Cards." *American Economic Review* 103 (7): 2875–2910.
- Kowalski, Amanda E.** 2023. "Behaviour within a Clinical Trial and Implications for Mammography Guidelines." *The Review of Economic Studies* 90 (1): 432–462.
- Ledley, Robert S., and Lee B. Lusted.** 1959. "Reasoning Foundations of Medical Diagnosis." *Science* 130 (3366): 9–21.
- Ludwig, Jens, Sendhil Mullainathan, and Ashesh Rambachan.** 2024. "The Unreasonable Effectiveness of Algorithms." February.
- Lyons, Patrick G., Mackenzie R. Hofford, Sean C. Yu, Andrew P. Michelson, Philip R. O. Payne, Catherine L. Hough, and Karandeep Singh.** 2023. "Factors Associated With Variability in the Performance of a Proprietary Sepsis Prediction Model Across 9 Networked Hospitals in the US." *JAMA Internal Medicine* 183 (6): 611–612.
- Main, Elliott K., Shen-Chih Chang, Ravi Dhurjati, Valerie Cape, Jochen Profit, and Jeffrey B. Gould.** 2020. "Reduction in Racial Disparities in Severe Maternal Morbidity from Hemorrhage in a Large-Scale Quality Improvement Collaborative." *American journal of obstetrics and gynecology* 223 (1): 123.e1–123.e14.
- Manz, Christopher R., Yichen Zhang, Kan Chen et al.** 2023. "Long-Term Effect of Machine Learning–Triggered Behavioral Nudges on Serious Illness Conversations and End-of-Life Outcomes Among Patients With Cancer: A Randomized Clinical Trial." *JAMA Oncology* 9 (3): 414–418.

- McCullough, Jeffrey S., Michelle Casey, Ira Moscovice, and Shailendra Prasad.** 2010. “The Effect Of Health Information Technology On Quality In U.S. Hospitals.” *Health Affairs* 29 (4): 647–654.
- McDevitt, Ryan C., and James W. Roberts.** 2014. “Market Structure and Gender Disparity in Health Care: Preferences, Competition, and Quality of Care.” *The RAND Journal of Economics* 45 (1): 116–139.
- McKibbin, Rebecca.** 2023. “The Effect of RCTs on Drug Demand: Evidence from off-Label Cancer Drugs.” *Journal of Health Economics* 90 102779.
- Meehl, Paul E.** 1954. *Clinical versus Statistical Prediction: A Theoretical Analysis and a Review of the Evidence*. Clinical versus Statistical Prediction: A Theoretical Analysis and a Review of the Evidence, Minneapolis, MN, US: University of Minnesota Press, x, 149.
- Mello, Michelle M., Michael D. Frakes, Erik Blumenkranz, and David M. Studdert.** 2020. “Malpractice Liability and Health Care Quality: A Review.” *JAMA* 323 (4): 352–366.
- Miller, Amalia R., and Catherine E. Tucker.** 2011. “Can Health Care Information Technology Save Babies?” *Journal of Political Economy* 119 (2): 289–324.
- Molitor, David.** 2018. “The Evolution of Physician Practice Styles: Evidence from Cardiologist Migration.” *American Economic Journal: Economic Policy* 10 (1): 326–356.
- Mullainathan, Sendhil, and Ziad Obermeyer.** 2022. “Diagnosing Physician Error: A Machine Learning Approach to Low-Value Health Care.” *The Quarterly Journal of Economics* 137 (2): 679–727.
- Neumark, David, and Bogdan Savych.** 2023. “Effects of Opioid-Related Policies on Opioid Utilization, Nature of Medical Care, and Duration of Disability.” *American Journal of Health Economics* 9 (3): 331–373.
- Newham, Melissa, and Marica Valente.** 2024. “The Cost of Influence: How Gifts to Physicians Shape Prescriptions and Drug Costs.” *Journal of Health Economics* 95 102887.
- Obermeyer, Ziad, Brian Powers, Christine Vogeli, and Sendhil Mullainathan.** 2019. “Dissecting Racial Bias in an Algorithm Used to Manage the Health of Populations.” *Science* 366 (6464): 447–453.
- Olenski, Andrew R., André Zimmerman, Stephen Coussens, and Anupam B. Jena.** 2020. “Behavioral Heuristics in Coronary-Artery Bypass Graft Surgery.” *New England Journal of Medicine* 382 (8): 778–779.
- Parapia, Liakat Ali.** 2008. “History of Bloodletting by Phlebotomy.” *British Journal of Haematology* 143 (4): 490–495.
- Persson, Emil, Kinga Barrafreem, Andreas Meunier, and Gustav Tinghög.** 2019. “The effect of decision fatigue on surgeons’ clinical decision making.” *Health Economics* 28 (10): 1194–1203.
- Persson, Petra, Xinyao Qiu, and Maya Rossin-Slater.** 2021. “Family Spillover Effects of Marginal Diagnoses: The Case of ADHD.” Working Paper 28334, National Bureau of Economic Research.
- Rambachan, Ashesh.** 2024. “Identifying Prediction Mistakes in Observational Data*.” *The Quarterly Journal of Economics* qjae013.

- Rambachan, Ashesh, and Jonathan Roth.** 2020. “Bias In, Bias Out? Evaluating the Folk Wisdom.” December.
- Roy, A. D.** 1951. “Some Thoughts on the Distribution of Earnings.” *Oxford Economic Papers* 3 (2): pp.135–146.
- Sabin, Janice A., and Anthony G. Greenwald.** 2012. “The Influence of Implicit Bias on Treatment Recommendations for 4 Common Pediatric Conditions: Pain, Urinary Tract Infection, Attention Deficit Hyperactivity Disorder, and Asthma.” *American Journal of Public Health* 102 (5): 988–995.
- Sacarny, Adam, Michael L. Barnett, Jackson Le, Frank Tetkoski, David Yokum, and Shantanu Agrawal.** 2018. “Effect of Peer Comparison Letters for High-Volume Primary Care Prescribers of Quetiapine in Older and Disabled Adults: A Randomized Clinical Trial.” *JAMA Psychiatry* 75 (10): 1003–1011.
- Sacarny, Adam, David Yokum, Amy Finkelstein, and Shantanu Agrawal.** 2016. “Medicare Letters To Curb Overprescribing Of Controlled Substances Had No Detectable Effect On Providers.” *Health Affairs* 35 (3): 471–479.
- Sacks, Daniel W., Alex Hollingsworth, Thuy Nguyen, and Kosali Simon.** 2021. “Can Policy Affect Initiation of Addictive Substance Use? Evidence from Opioid Prescribing.” *Journal of Health Economics* 76 102397.
- Schnell, Molly, and Janet Currie.** 2018. “Addressing the Opioid Epidemic: Is There a Role for Physician Education?” *American Journal of Health Economics* 4 (3): 383–410.
- Shapiro, Bradley T.** 2018. “Informational Shocks, Off-Label Prescribing, and the Effects of Physician Detailing.” *Management Science* 64 (12): 5925–5945.
- Shurtz, Ity, Yoav Goldstein, and Gabriel Chodick.** 2024. “Realization of Low-Probability Clinical Risks and Physician Behavior: Evidence from Primary Care Physicians.” *American Journal of Health Economics* 10 (1): 132–157.
- Silver, David.** 2021. “Haste or Waste? Peer Pressure and Productivity in the Emergency Department.” *The Review of Economic Studies* 88 (3): 1385–1417.
- Simon, Herbert Alexander.** 1957. *Models of Man: Social and Rational; Mathematical Essays on Rational Human Behavior in Society Setting.* Wiley.
- Singh, Hardeep, Christiane Spitzmueller, Nancy J. Petersen, Mona K. Sawhney, and Dean F. Sittig.** 2013. “Information Overload and Missed Test Results in Electronic Health Record–Based Settings.” *JAMA Internal Medicine* 173 (8): 702–704.
- Singh, Manasvini, and Atheendar Venkataramani.** 2022. “Rationing by Race.” Working Paper 30380, National Bureau of Economic Research.
- Sobczak, Alexandria, Lauren Taylor, Sydney Solomon et al.** 2023. “The Effect of Doulas on Maternal and Birth Outcomes: A Scoping Review.” *Cureus* 15 (5): .
- Sommers, Benjamin D., Caitlin L. McMurtry, Robert J. Blendon, John M. Benson, and Justin M. Sayde.** 2017. “Beyond Health Insurance: Remaining Disparities in US Health Care in the Post-ACA Era.” *The Milbank Quarterly* 95 (1): 43–69.
- Stevenson, Megan T., and Jennifer L. Doleac.** 2022. “Algorithmic Risk Assessment in the Hands of Humans.” September.

- Tai-Seale, Ming, and Thomas McGuire.** 2012. “Time Is up: Increasing Shadow Price of Time in Primary-Care Office Visits.” *Health Economics* 21 (4): 457–476.
- van Parys, Jessica.** 2016. “Variation in Physician Practice Styles within and across Emergency Departments.” *PLOS ONE* 11 (8): .
- Vela, Monica B., Amarachi I. Erundu, Nichole A. Smith, Monica E. Peek, James N. Woodruff, and Marshall H. Chin.** 2022. “Eliminating Explicit and Implicit Biases in Health Care: Evidence and Research Needs.” *Annual Review of Public Health* 43 477–501.
- Wallis, Christopher J. D., Angela Jerath, Natalie Coburn et al.** 2022. “Association of Surgeon-Patient Sex Concordance With Postoperative Outcomes.” *JAMA Surgery* 157 (2): 146–156.
- Wilding, Anna, Luke Munford, Bruce Guthrie, Evangelos Kontopantelis, and Matt Sutton.** 2022. “Family Doctor Responses to Changes in Target Stringency under Financial Incentives.” *Journal of Health Economics* 85.
- Williams, David R., Jourdyn A. Lawrence, and Brigette A. Davis.** 2019. “Racism and Health: Evidence and Needed Research.” *Annual Review of Public Health* 40 (Volume 40, 2019): 105–125.
- Wilson, Karen, and Lara Khansa.** 2018. “Migrating to Electronic Health Record Systems: A Comparative Study between the United States and the United Kingdom.” *Health Policy* 122 (11): 1232–1239.
- Wong, Andrew, Erkin Otles, John P. Donnelly et al.** 2021. “External Validation of a Widely Implemented Proprietary Sepsis Prediction Model in Hospitalized Patients.” *JAMA Internal Medicine* 181 (8): 1065–1070.
- Zeltzer, Dan, Liran Einav, Joseph Rashba, and Ran D Balicer.** 2024. “The Impact of Increased Access to Telemedicine.” *Journal of the European Economic Association* 22 (2): 712–750.
- Zeltzer, Dan, Liran Einav, Joseph Rashba, Yehezkel Waisman, Motti Haimi, and Ran D. Balicer.** 2023. “Adoption and Utilization of Device-Assisted Telemedicine.” *Journal of Health Economics* 90 102780.

Table 1: Variation in Physician Practice Style

Paper	Research Question	Data	Empirical Methods	Results	Heterogeneous Effects?
Abaluck et al. (AER 2016)	Variation in physician propensity to test for pulmonary embolism (PE) and effect of test misallocation on health outcomes.	20% sample Part B Medicare Claims 2000-2009; Part A claims with PE diagnosis; patient chart and billing data from two academic medical centers.	See text.	The average doctor tests if she believes the likelihood of a positive test is higher than 5.6 percent (SD = 5.4). Doctors react strongly to clinical symptoms but not to known PE risk factors from the patient's medical history.	
Ahammer and Schober (Health Economics, 2020)	How much of the variation in Austrian health expenditures is explained by GP practice style?	Upper Austrian Health Insurance Fund data (2005–2012); Medical Chamber data on doctor demographics; inpatient records.	AKM decomposition with patient and GP FEs, exploiting patients who change GPs over time. Card et al. (2013) decomposition of variance.	Accounting for patient demand, patients of high-usage GPs have 20 to 148.5% higher expenditures than patients seeing an average GP.	Older doctors, female doctors, and doctors practicing in areas in higher GP density have higher expenditures.
Badinski et al. (NBER Working Paper, 2023)	How does geographic variation in physician practice intensity affect healthcare utilization?	20% random sample of Medicare fee-for-service claims 1998–2013.	Movers design exploiting patients and physician moves between HHRs and differences in utilization within HHRs estimated using patient and physician FE models.	A 1 SD increase in an HHR's average physician practice intensity increases utilization per visit 13%. 3/5 of the variation in an HHR's average physician practice intensity comes from variation within specialties and the rest from differences in physician specialty mix across HHRs.	Variation in PCP intensity across HHRs explains 19% of variation primary care utilization. Variation in cardiologist intensity explains only 3% of variation in cardiology utilization.
Berndt et al. (JHE 2015)	How concentrated are antipsychotic prescribing practices? (Do doctors have favorite drugs?)	10% sample from IMS retail prescriptions data, with refreshment each year; linked to the AMA Masterfile.	Descriptive.	Two thirds of a physician's prescriptions are for the same drug. The Herfindahl in prescribing concentration is decreasing in the log of total yearly antipsychotic prescriptions suggesting learning by doing.	The relationship between the volume of prescribing and the Herfindahl is larger for primary care physicians than for psychiatrists.

Chan, Gentzkow, and Yu (QJE, 2022)	Does radiologists' diagnostic skill affect diagnosis and outcomes for suspected pneumonia patients?	Veteran's Health Administration Emergency Department data Oct. 1999 to Sept. 2015.	See text.	Variation in skill explains 39% of the variation in diagnostic decisions and 78% of the variation in outcomes for suspected pneumonia patients. Diagnostic thresholds increase with skill.	
Currie, MacLeod and Van Parys (JHE, 2016)	Characterize practice style and describe how variation in practice style affects outcomes of heart attack patients?	Florida hospital discharge data for AMI patients admitted through the ED, 1992-2014; Data on providers from Florida medical license database.	Define appropriateness for invasive procedure using teaching programs. Regress use of invasive procedures on appropriateness and examine intercept (aggressiveness) and slope (responsiveness).	Within hospitals and years, patients with more aggressive providers have higher costs and better outcomes. Providers who follow "best practices" do too few procedures on healthy elderly suggesting over-reliance on age as a criterion.	Young, male providers from top schools are more aggressive.
Currie and MacLeod (JOLE, 2017)	How do variations in physician diagnostic and surgical skill affect outcomes of pregnancy?	~1 million NJ electronic birth records for 1997-2006.	See text.	Better diagnosis would reduce C-sections for low-risk mothers and increase C-sections for high-risk births, which would prevent infant death. Better surgical skills increase C-section rates and improve outcomes across the board.	Reducing C-section rates across the board would harm infants in high-risk pregnancies.
Cutler et al. (AEJ:EP 2019)	How does the percentage of "cowboys" and "comforters" in an area relate affect end-of-life spending.	Random sample of 598 cardiologists, 967 PCPs and 2,882 Medicare patients; Medicare expenditures from Dartmouth Atlas; Measures from the "Hospital Care" database.	Categorization of physicians based on survey results. Cowboys are physicians who recommend intensive care beyond current guidelines. Comforters recommend palliative care for the severely ill. Categories not mutually exclusive.	A 1 SD increase in the share of cowboys leads to 10.66-13.12% higher spending in last 2 years and a 2.15-3.56% higher 1-year spending for AMI patients. A 1 SD increase in the share of comforters leads to a 2.68-5.51% fall in spending in last 2 years, and a 0.82-1.2% fall in 1-year spending for AMI patients. Shares not significantly associated with survival.	
Fadlon and Van Parys (JHE 2020)	How does PCP practice style affect patient health care utilization?	20% sample of Medicare enrollees with at >=one month of traditional	Event study/d-in-d exploiting PCP changes when a patient's PCP relocates or retires.	Switching to a PCP whose patients spend \$10 more on primary care (PC) increases per capita spending 4.07%. Switching to a PCP whose patients have 1 SD more	Distinguish PCP switches within and between practices. Results similar

		Medicare enrollment in the year.		PC visits increases visits 38.20%. Similar effects for #diagnoses, flu vaccines, and diabetes care.	indicating variation is associated with individual PCPs.
Marquardt (R&R JPE 2021)	How does physician practice style affect diagnosis of ADHD? What doctor characteristics predict practice style.	Electronic medical records from 129 doctors (12,311 pediatric patients) in a large healthcare system, Jan. 2014-Sept. 2017. Physician characteristics from the web.	Use natural language processing to measure child's suitability for ADHD diagnosis. Regress diagnosis on suitability. Examine intercept (intensiveness) and slope (compliance with guidelines). Regress doctor-specific estimates on doctor characteristics.	A physician with the median intensity (intercept) and median compliance (slope) diagnoses patients with the median symptom level 3.46% of the time. Increasing physician intensity by 1 SD increases diagnosis probability to 22.45%. Increasing physician compliance 1 SD increases diagnosis probability to 20.0%.	Less experienced male physicians have lower intercepts. Less experienced female physicians have higher slopes. Physicians who see patients with higher average severity have lower intensity and higher compliance.

Notes: See glossary for abbreviations.

Appendix for “First Do No Harm? Doctor Decision Making and Patient Outcomes”

Janet Currie*, W Bentley MacLeod†and Kate Musen‡

June 14, 2024

*Princeton University and NBER
†Princeton University and NBER
‡Columbia University

Appendix for Theory in Section 2.

This appendix lays out the detailed proofs of the model discussed in the text. We first recap the framework.

Consider a population of patients where patient $i \in \mathcal{N}$ seeks treatment from doctor $j \in J$. It is assumed that some treatment is optimal, but neither patient or physician is sure which is the best choice. The doctor chooses between a non-intensive or a more intensive treatment, denoted by $t \in \{NI, I\}$. It is assumed that there is a best choice for the patient given by their *unobserved* state $\alpha_i \in \{0, 1\}$. If $\alpha_i = 0$, then the non-intensive treatment is optimal, while $\alpha_i = 1$ implies that the intensive treatment is more appropriate. This modeling strategy is based on Savage (1972 (first published 1954)’s model of Bayesian choice in which the goal of the model is not to provide a complete representation of the patient’s condition, but to highlight only those aspects of a patient’s state that are relevant for the decision at hand.¹

Let the fraction of patients in \mathcal{N} for which $\alpha_i = 0$ be given by $p_0 \in (0, 1)$, while a fraction $p_1 = 1 - p_0$ are in state $\alpha_i = 1$. Doctor j cannot perfectly observe the patient’s state, but observes a signal:

$$T_{ij} = \alpha_i + \epsilon/\gamma_j, \quad (1)$$

where $\epsilon \sim N(0, 1)$ and γ_j is the diagnostic skill of the doctor. An increase in diagnostic skill implies a more precise assessment of a person’s state. The doctor is never perfectly sure of the patient’s condition since it is observed with error.

T_{ij} is increasing with α_i so it follows that the optimal diagnostic rule for the treatment $t_{ij} \in \{NI, I\}$ takes the form:

$$t_{ij} = \begin{cases} I, & T_{ij} \geq \tau_j, \\ NI, & T_{ij} < \tau_j, \end{cases}$$

where τ_j is the doctor’s *decision threshold* for deciding when to implement the intensive treatment. Increasing the threshold reduces the probability that the intensive treatment is chosen.

The quality of diagnosis is measured by the likelihood that a patient is assigned to the correct treatment. Here there are two measures of performance corresponding to whether patients correctly or incorrectly receive the intensive treatment. Suppose a patient is in state $\alpha_i = 1$ and hence should be assigned to intensive treatment. The probability that the patient correctly receives the intensive treatment is given the doctors decision threshold, τ_j , and diagnostic skill. γ_j :

$$\begin{aligned} PI(\tau_j, \gamma_j) &\equiv \Pr[T_{ij} \geq \tau_j | \alpha_i = 1], \\ &= \Pr[1 + \epsilon/\gamma_j \geq \tau_j], \\ &= F(\gamma_j(1 - \tau_j)), \end{aligned} \quad (2)$$

¹See the discussion in Chapter 2 of MacLeod (2022).

where $F(\cdot)$ is the Normal cumulative probability distribution.

The probability that a patient who needs non-intensive treatment ($\alpha_i = 0$) receives intensive treatment, given by:

$$\begin{aligned} PNI(\tau_j, \gamma_j) &\equiv \Pr [T_{ij} \geq \tau_j | \alpha = 0] \\ &= \Pr [\epsilon/\gamma_j \geq \tau_j] \\ &= F(-\gamma_j \tau_j). \end{aligned} \quad (3)$$

The Optimal Decision Threshold (τ_j^*)

This section derives the optimal decision threshold, τ_j^* , given a doctor's diagnostic skill, γ_j , and potential patient outcomes. The optimal rule takes the form of a threshold value such that if the patient's severity is over the threshold, the doctor will perform the intensive procedure.

Given patient type $\alpha_i \in \{0, 1\}$, doctor j 's utility from administrating treatment $t \in \{NI, I\}$ is given by:

$$U_{\alpha t j} = u_{\alpha t j} + \lambda_{t j}, \quad (4)$$

where $u_{\alpha t j}$ is the expected medical benefit to a patient of type $\alpha \in \{0, 1\}$, getting treatment $t \in \{NI, I\}$ from doctor j . For the same patient type, the outcome $u_{\alpha t j}$ can differ by doctor, a variation that we associate with a doctor's *procedural skill*. Additional factors that affect treatment, such as a payment that the doctor receives from administering the treatment, are captured by $\lambda_{t j}$. We begin by assuming that the doctor cares only about the medical benefit to the patient and ignore the $\lambda_{t j}$ term.

Assume that if a patient is of type $\alpha_i = 0$, then non-intensive treatment is preferred ($u_{0NIj} > u_{0Ij}$), while for type $\alpha_i = 1$ intensive treatment is preferred ($u_{1Ij} > u_{1NIj}$). If this were not the case, then they would be no diagnostic decision to make - all patients would be assigned to either intensive or non-intensive treatment. Let $\Delta_{1Ij} = \{u_{1Ij} - u_{1NIj}\} > 0$ and $\Delta_{0NIj} = \{u_{0NIj} - u_{0Ij}\}$ be the increase in utility for patients who receive the appropriate treatment. The doctor's *ex ante* belief regarding the appropriate treatment for a patient in this pool of potential patients is given by:

$$p_{1j} = \Pr [\alpha = 1]$$

while the belief that the probability that $\alpha_i = 0$ is $p_{0j} = 1 - p_{1j}$.

The expected utility of doctor j with who chooses decision threshold τ for patient i is given by:

$$\begin{aligned} u_{ij}(\tau) &= (u_{1Ij} \Pr [T_{ij} \geq \tau | \alpha = 1] + u_{1NIj} \Pr [T_{ij} < \tau | \alpha = 1]) \Pr [\alpha = 1] \\ &\quad + (u_{0Ij} \Pr [T_{ij} \geq \tau | \alpha = 0] + u_{0NIj} \Pr [T_{ij} < \tau | \alpha = 0]) \Pr [\alpha = 0] \\ &= (u_{1NIj} + \Delta_{1Ij} \Pr [T_{ij} \geq \tau | \alpha = 1]) p_{1j} \\ &\quad + (u_{0Ij} - \Delta_{0NIj} \Pr [u \geq \tau | \alpha = 0]) p_{0j}, \\ &= u_j^0 + \Delta_{1Ij} PNI(\tau_j, \gamma_j) \times p_{1j} - \Delta_{0NIj} PNI(\tau_j, \gamma_j) \times p_{0j}, \end{aligned} \quad (5)$$

where:

$$\begin{aligned} u_j^0 &= u_{1NIj} \Pr[\alpha = 1] + u_{0Ij} \Pr[\alpha = 0], \\ &= u_{1NIj} \times p_{1j} + u_{0Ij} \times p_{0j}. \end{aligned}$$

The quantity u_j^0 is the *worst* possible payoff for doctor j . It occurs when all individuals with type $\alpha = 1$ are given the non-intensive treatment, and all type $\alpha = 0$ individuals are given the intensive treatment. The payoff to a doctor can now be written in terms of the expected gains, beliefs and expected patient outcomes. The optimal decision threshold for each physician is $\tau_{ij}^* = \arg \max_{\tau \in \mathfrak{R}} u_{ij}(\tau, \gamma)$. The solution is given by the following proposition.

Proposition 1. *The medically optimal decision threshold solves $\tau_{ij}^* = \arg \max_{\tau \in \mathfrak{R}} u_{ij}(\tau, \gamma)$, and satisfies the likelihood ratio condition:*

$$L(\tau_{ij}^*, \gamma_j) = \frac{\Delta_{0NIj}}{\Delta_{1Ij}} \times \frac{p_{0j}}{p_{1j}},$$

where the likelihood ratio is given by:

$$L(\tau_{ij}^*, \gamma_j) = \frac{f(\gamma_j(1 - \tau_{ij}^*))}{f(-\gamma_j\tau_{ij}^*)},$$

and $f(\cdot)$ is the Normal density function.

Proof. The optimal solution satisfies the first order condition:

$$\begin{aligned} 0 &= \partial u_{ij}(\tau, \gamma_j) / \partial \tau, \\ &= \Delta_{1Ij} \partial PI(\tau, \gamma_j) / \partial \tau \times p_{1j} - \Delta_{0NIj} \partial PNI(\tau, \gamma_j) / \partial \tau \times p_{0j}, \\ &= \Delta_{1Ij} f(\gamma_j(1 - \tau))(-\gamma_j) \times p_{1j} - \Delta_{0NIj} f(-\gamma_j\tau)(-\gamma_j) \times p_{0j}. \end{aligned}$$

The first order condition follows from the last line. \square

The first order conditions imply a unique optimal decision threshold, τ_{ij}^* :

$$L(\tau_{ij}^*, \gamma_j) = \frac{f(\gamma_j(1 - \tau_{ij}^*))}{f(-\gamma_j\tau_{ij}^*)} = \frac{\Delta_{0NIj}}{\Delta_{1Ij}} \times \frac{p_{0j}}{p_{1j}}.$$

The first order condition characterizes the global optimum, which follows from the Neyman-Pearson lemma showing that likelihood ratios are the most powerful form of hypothesis test (Neyman and Pearson (1933)).² The model yields a closed form solution for the optimal diagnostic rule τ_{ij}^* , which is given by the following proposition:

²We thank Han Hong for pointing out the link between optimal choice and the Neyman-Pearson lemm.

Proposition 2. *The medically optimal threshold satisfies:*

$$\tau_{ij}^* = \frac{1}{2} + b_{ij}^*/\gamma_j^2, \quad (6)$$

where $b_{ij}^* \equiv (\ln(\Delta_{0NIj}/\Delta_{1Ij}) + \ln(p_{0j}/p_{1j}))$ is the *optimal threshold shifter*.

Proof. Observe:

$$\begin{aligned} \frac{f(\gamma_j(1-\tau_{ij}^*))}{f(-\gamma_j\tau_{ij}^*)} &= \frac{\exp - \{\gamma_j(1-\tau_{ij}^*)\}^2 / 2}{\exp - \{-\gamma_j\tau_{ij}^*\}^2 / 2} \\ &= \exp \left(- \{\gamma_j(1-\tau_{ij}^*)\}^2 + \{\gamma_j\tau_{ij}^*\}^2 \right) / 2 \end{aligned}$$

Taking the logarithm of the first order condition gives us:

$$\begin{aligned} \left(- \{\gamma_j(1-\tau_{ij}^*)\}^2 + \{\gamma_j\tau_{ij}^*\}^2 \right) / 2 &= b_{ij}, \\ \gamma_{ij}^2(-1/2 + \tau_{ij}) &= b_{ij}, \end{aligned}$$

giving the desired result:

$$\tau_{ij}^* = \frac{1}{2} + b_{ij}/\gamma_j^2, \quad (7)$$

where $b_{ij} \equiv (\ln(\Delta_{0NIj}/\Delta_{1Ij}) + \ln(p_{0j}/p_{1j}))$. \square

Equation (6) shows that the optimal decision threshold depends on diagnostic skill, γ_j , the relative effectiveness of non-intensive and intensive treatments for the two types of patients, $\Delta_{0NIj}/\Delta_{1Ij}$, and the doctor's beliefs about the relative proportions of patient types, p_{0j}/p_{1j} , in the population. When the doctor believes that there is a higher probability that the patient needs non-intensive treatment, she adopts a higher threshold resulting in less use of the intensive treatment. Similarly, if the relative benefit from intensive treatment is higher, then this results in a lower threshold.

Corollary 1. *Let $\delta_\alpha^u = u_{\alpha I} - u_{\alpha NI}$ be the relative gain when a patient in state α gets the intensive treatment. An increase in the relative gain from intensive treatment increases the likelihood that the patient gets the intensive treatment, $\frac{\partial b_{ij}}{\partial \delta_\alpha^u} < 0, \alpha \in \{0, 1\}$. The size of these effects increase with a decrease in doctor diagnostic skill.*

Proof. The proof follows immediately from differentiating b_{ij} , and observing that the size of this effect increases as γ_j becomes smaller. \square

As diagnostic skill increases, both patient types are more likely to be allocated to the appropriate treatment. The optimal decision rule entails patients getting the appropriate treatment with probability close to one as diagnostic skill increases. Conversely, as diagnostic skill falls, the b_{ij} term dominates. When $b_{ij} > 0$, treatment is biased in favor of the non-intensive treatment and the probability that patients are treated with the non-intensive procedure rises

as diagnostic skill falls. When $b_{ij} < 0$, treatment is biased in favor of intensive treatment and the probability of intensive treatment rises as diagnostic skill falls. In effect, as diagnostic skill falls, physicians choose the treatment that they believe is optimal for most patients, and more patients receive the same treatment. These observations are formally summarized in the following proposition:

Proposition 3. *As diagnostic skill increases each patient receives treatment that close to optimal for their type. More precisely:*

$$\lim_{\gamma_j \rightarrow \infty} \tau_{ij}^* = 1/2,$$

$$\lim_{\gamma_j \rightarrow \infty} u_{ij}^* = \begin{cases} u_{1Ij}, & \text{if } \alpha_i = 1, \\ u_{0NIj}, & \text{if } \alpha_i = 0. \end{cases}$$

As diagnostic skill falls all patients get the same treatment depending upon the sign of the decision shifter, b_{ij} :

$$\lim_{\gamma_j \rightarrow 0} \tau_{ij}^* = \begin{cases} \infty, & \text{if } b_{ij} > 0, \\ 1/2, & \text{if } b_{ij} = 0 \\ -\infty, & \text{if } b_{ij} < 0. \end{cases}$$

$$\lim_{\gamma_j \rightarrow \infty} u_{ij}^* = \begin{cases} u_{1NIj}, & \text{if } \alpha_i = 1, b_{ij} > 0, \\ u_{0NIj}, & \text{if } \alpha_i = 0, b_{ij} > 0, \\ (u_{1NIj} + u_{1Ij})/2, & \text{if } \alpha_i = 1, b_{ij} = 0, \\ (u_{0NIj} + u_{0Ij})/2, & \text{if } \alpha_i = 0, b_{ij} = 0, \\ u_{1Ij}, & \text{if } \alpha_i = 1, b_{ij} < 0, \\ u_{0Ij}, & \text{if } \alpha_i = 0, b_{ij} < 0. \end{cases}$$

Proof. The proof of this proposition follows from equation (7). \square

We now return to consideration of the λ_{tj} term. The role of this term is quantified in the following corollary:

Corollary 2. *Let $\delta_j = \lambda_{Ij} - \lambda_{NIj}$ denote the doctor's pecuniary preference for intensive treatment. If $\delta_j \in (\Delta_{0NIj}, -\Delta_{1Ij})$ then Doctor j chooses a decision threshold for treatment choice to satisfy:*

$$\tau_{ij}^0 = \frac{1}{2} + b_{ij}^0/\gamma_j^2, \quad (8)$$

where $b_{ij}^0 \equiv (\ln(\Delta_{0NIj} - \delta_j) - \ln(\Delta_{1Ij} + \delta_j) + \ln(p_{0j}/p_{1j}))$ is the optimal threshold shifter. If $\delta_j > \Delta_{0NIj}$ then the doctor chooses the intensive procedure for all patients, while the non-intensive procedure is chosen for all patients if $\delta_j < \Delta_{1Ij}$.

Proof. When $\delta_j \in (\Delta_{0NIj}, -\Delta_{1Ij})$ the result follows from the proofs of propositions (1) and (2). Notice that as $\delta_j \rightarrow \Delta_{0NIj}$ then $b_{ij} \rightarrow \infty$, $\tau_{ij}^0 \rightarrow \infty$,

and hence the doctor always chooses the intensive treatment. It follows that for $\delta_j > \Delta_{0NI_j}$ only the intensive treatment is chosen. A similar argument holds for $\delta_j < \Delta_{1I_j}$, in which case the doctor only chooses the non-intensive treatment. \square

Identifying the Doctor Diagnostic Rule, Diagnostic Skill, and Procedural Skill From Data

From observational data one observes the doctor's treatment choice ($t_{ij} \in \{NI, I\}$), and some measure of patient outcomes following treatment, as well as some information on patient type that may be available in medical records. Let \vec{x}_i be patient characteristics that are observable in the data. One can use the vector of observed patient characteristics, \vec{x}_i , to estimate the patient's appropriateness for the intensive procedure and treat this estimated propensity as an index of appropriateness, i.e. the medical benefit of the procedure. Call this index $\rho(\vec{x}_i)$. It is assumed that there is a monotonic ordering of patients from the least appropriate for the intensive procedure to the most appropriate, and that all providers are in general agreement about this ordering. The following proposition shows that a doctor with better diagnostic skill will be more responsive to the measure of patient appropriateness for the procedure, $\rho(\vec{x}_i)$:

Proposition 4. *The doctor's estimated likelihood of performing an intensive procedure is:*

$$\Pr [t_{ij} = I | j, \vec{x}_i] = (PI_j - PNI_j) \rho(\vec{x}_i) + PNI_j, \quad (9)$$

and the slope term, $\theta_j = (PI_j - PNI_j)$ is a doctor-specific measure that increases with doctor diagnostic skill:

$$\frac{d\theta_j}{d\gamma_j} > 0,$$

where $\Pr [t_i = I] = \rho(\beta\vec{x}_i)$ is the estimated probability of intensive treatment.

Proof. The probability of a C-section is:

$$\begin{aligned} \Pr [T_{ij} \geq \tau_j] &= PI_j \times \sigma(\vec{x}_i) + PNI_j \times (1 - \sigma(\vec{x}_i)), \\ &= (PI_j - PNI_j) \sigma(\vec{x}_i) + PNI_j. \end{aligned}$$

Then we have using the optimal decision rule from proposition (1):

$$\begin{aligned}
\frac{d\theta_j}{d\gamma_j} &= \frac{dF(\gamma_j(1 - \tau_{ij}^*))}{d\gamma_j} - \frac{dF(-\gamma_j\tau_{ij}^*)}{d\gamma_j} \\
&= \frac{dF(\gamma_j/2 - b_{ij}/\gamma_j)}{d\gamma_j} - \frac{dF(-\gamma_j/2 - b_{ij}/\gamma_j)}{d\gamma_j} \\
&= (1/2 + b_{ij}/\gamma_j^2) f(\gamma_j/2 - b_{ij}/\gamma_j) + (1/2 - b_{ij}/\gamma_j^2) f(-\gamma_j/2 - b_{ij}/\gamma_j) \\
&= f(-\gamma_j\tau_{ij}^*) \left\{ (1/2 + b_{ij}/\gamma_j^2) \exp(b_{ij}) + (1/2 - b_{ij}/\gamma_j^2) \right\} \\
&= f(-\gamma_j\tau_{ij}^*) \left\{ (\exp(b_{ij}) + 1)/2 + \frac{b_{ij}}{\gamma_j^2} (\exp(b_{ij}) - 1) \right\}
\end{aligned}$$

When $b_{ij} \geq 0$ then $\exp(b_{ij}) > 1$ and hence the right hand side is positive. When $b_{ij} < 0$, then $\exp(b_{ij}) - 1 < 0$ and hence $\frac{b_{ij}}{\gamma_j^2} (\exp(b_{ij}) - 1) > 0$. Thus the derivative is positive. \square

Notice that from equation (9) one can separately identify PI_j and PNI_j . Hence we can identify both τ_j and γ_j .

The slope term is also affected by the physician's beliefs about when invasive procedures are likely to be warranted via τ_j , and by any additional physician-specific factors that are included in λ_{ij} . For example, a doctor who believes that most women should have C-sections (and therefore has a very low decision threshold for C-section). Currie and MacLeod (2017) attempt to distinguish between τ_j as well as γ_j by noting that in a doctor-specific regression, the constant term in Equation (9) is affected only by τ_j so given two estimated parameters and two unknowns, it is possible to identify both.

Finally, notice that patients with high *ex ante* likelihood of having a C-section ($\rho(\vec{x}_i) \approx 1$) then variation in patient outcomes is effectively independent of both diagnostic skill and the decision threshold. Hence, we can associate variation in outcomes with procedural skill. A similar implication follows for patients with a low likelihood of a C-section ($\rho(\vec{x}_i) \approx 0$).

Appendix Describing Research Papers Organized by Topic

Appendix Table 1: Health Disparities

Paper	Research Question	Data	Empirical Methods	Results	Heterogeneous Effects?
Alsan, Garrick, and Graziani (AER 2019)	How does physician race affect Black men’s take up of preventative care services?	Experimental data with 1,374 recruited Black male participants, with 637 completing the study.	Field experiment with random assignment to either a Black or non-Black physician in a special clinic offering preventive care. Doctor race was signaled to patients by a headshot.	Viewing the headshot did not significantly affect intended take-up of services. But patients who saw a Black patient increased demand for services exposte by 38.79% for diabetes screening, 52.77% for cholesterol screening and 26.54% for flu shots.	No differences by income, education, or age. Effects greater for patients without a recent medical screening, with more ER visits, and with higher levels of measured medical mistrust.
Angerer, Waibel, and Stummer (AJHE 2019)	What is the effect of socioeconomic status, signaled by education level, on the probability of receiving a medical appointment and on response times?	Experimental data for April 26-June 2, 2017, with email requests for appointments sent to 1,249 Austrian specialists.	Correspondence study via email with varying email signatures to signal no degree, a doctoral degree, or a medical degree	Patients with degrees are more likely to receive an appointment and have lower response times and lower waiting times. Whether patients are offered an appointment depends on the assistant, while response and waiting times depend on the doctor.	The effects are driven by practices that do not contract with social insurance.
Button et al. (NBER WP 2020)	How does being nonbinary or transgender interact with patient race to affect the probability of getting an appointment with a mental health care provider (MHP)?	Experimental correspondence data from 1,000 emails sent to MHPs between Jan. 28, 2020-May 15, 2020, with number of emails per zip code proportional to population.	Emails sent through an MHP appointment request website with randomly assigned content disclosing trans or nonbinary status. Names signal gender and race. Randomize whether condition is depression, anxiety, or “stress.”	Transgender or non-binary African Americans and Hispanics are 18.7% less likely to get a positive response than cisgender whites. No evidence of differential responses by TNB status for whites.	N/A
Brekke et al. (HE 2018)	What is the relationship between SES of Type II diabetes patients and GP treatment decisions?	Norwegian administrative health data 2008- 2012; patient and GP characteristics from Statistics Norway.	GP FE models of service provision conditional on patient characteristics. Additional results using GP quits, retirements, and moves.	High ed. patients get fewer, longer visits. Less ed. patients get more medical tests and services over a year. E.g. high ed. 14.79% more likely to get a visit over 20 minutes. Less ed. 3.94% more 2+ HbA1C tests.	Results are similar when disaggregated by patient age and GP sex, age, specialty, number of patients, and fixed payment vs. fee-for-service.
Cabral and Dillender (AER 2024)	How does gender concordance between claimants and doctors	Open records request for Texas worker’s compensation claims	Assignment to doctors is random conditional on doctor’s credential and the	Female claimants seen by a female doctor are 5.2% more likely to receive benefits compared to when female	Differences are not statistically significant but suggest larger effects for

	performing independent medical evaluations for workers compensation affect disability determinations?	2013-17, and independent medical evaluations 2005-2017; NPI registry; novel survey of 1,519 adults 30-64, 2021.	claimants' county. Estimate OLS with an interaction between female doctor and female claimant controlling for main effects, credential, and county.	claimants are seen by male doctors. Physician gender does not affect likelihood of receiving benefits for male claimants. Female claimants seen by a female doctor receive 8.6% higher benefits than female claimants seen by male doctors.	those with lower earnings, in less dangerous industries, but with worse injuries.
Chandra and Staiger (NBER WP 2010)	Are differences in the treatment of Black and female AMI patients due to physician preferences or statistical discrimination?	Clinical records for 200,000+ patients admitted for AMI in 1994 & 1995 from the Cooperative Cardiovascular Project.	Propensity score estimation; taste based discrimination implies that similar patients who receive fewer services will suffer worse outcomes.	Black and female patients receive less treatment but also receive slightly lower benefits from treatment suggesting that they are not being denied beneficial treatment due to discrimination.	N/A.
Eli, Logan, and Miloucheva (NBER WP 2019)	Use union army pension awards to examine the effect of income on mortality. Investigate differences in a board's disability evaluations by race of applicant.	Union Army and United States Colored Troops (USCT) sample from the Early Indicators Project; Rosters of Examining Surgeons from the National Archives.	Instrument pension income using leave-one-out mean of pension board's decisions. Include board FEs. First stage shows the same boards were less generous to Black veterans.	Pension income significantly increased life expectancy. Bias against Black veterans in determining pension eligibility is substantial and accounts for much of the racial mortality gap in this population.	Bias against Black veterans is strongest for conditions where valuations may be more subjective, such as digestive diseases.
Frakes and Gruber (NBER WP 2022)	How does the availability of Black physicians on a military base affect Black Tricare patients' outcomes?	Military Health System Data Repository fiscal years 2003–2013	Mover-based ITT design exploiting differences in racial shares of physicians across bases.	1 SD increase in share of Black physicians reduces Black patients' mortality from diabetes, hypertension, high cholesterol, and cardiovascular disease by 15%. 55–69% of the effect attributed to medication adherence.	N/A.
Goyal et al. (JAMA Pediatrics 2015)	How does treatment of pain in the ED vary by race for child appendicitis patients?	National Hospital Ambulatory Medical Care Survey 2003-2010.	Multivariate logistic regression.	Black patients less likely to receive any analgesia, adjusted OR=0.1 for moderate pain and 0.2 for severe pain. Black patients were less likely to receive opioids, adjusted OR= 0.2.	The authors test for interactions between race and sex but do not find any.
Greenwood, Carnahan, and Huang (PNAS 2018)	How does patient-attending gender concordance affect mortality from heart attacks among patients admitted to the ED? Do	Census of patients admitted to hospitals in Florida 1991- 2010 from Florida's Agency for Healthcare Administration.	Assume patient assignment to physicians is conditionally random in the ED and either include physician FEs or hospital-quarter FEs. They also	In the full sample with hospital-quarter FEs, relative to male or female patients treated by female physicians, female patients treated by male doctors are 1.80% less likely to survive and male patients treated by	Female survival increases with more female physicians in ED, even in patients treated by male physicians. Survival of female patients treated by

	male doctors with more female colleagues or AMI patients have better female survival?		estimate additional specifications using matching.	male doctors are 0.90% less likely to survive. In the matched sample, only female patients treated by male doctors have lower survival rates.	male doctors increases with the number of female patients seen by their doctor in the prior quarter.
Greenwood et al. (PNAS 2020)	How does infant and maternal mortality vary as a function of patient-doctor racial concordance?	Census of patients admitted to hospitals in Florida 1992- 2015 from Florida’s Agency for Healthcare Administration.	OLS with controls including physician FEs in some models.	Racial concordance between infant and physician corresponds to about a 40% reduction in gap in mortality between Black and white infants. No significant racial concordance effects are found for mothers.	Effects more precisely estimated for infants with ≥ 1 comorbidity and in hospitals with more Black patients. Effects similar in % terms for pediatricians and non-pediatricians.
Hill, Jones, and Woodworth (JHE 2023)	Does physician-patient race concordance affect within-hospital mortality of uninsured non-Hispanic, patients admitted through the ED?	Florida Hospital Discharge Data File from October 2011 to December 2014; Florida Physician Workforce Survey from 2008-2016.	IV measures “the lagged share of same-race physicians typically present at the indexed hospital on the weekday and shift” when patient admitted.	Physician-patient race concordance reduces mortality by 27%.	The largest effects are for subgroups of patients with high variance in number of procedures and in total charges.
Hoffman et al. (PNAS 2016)	How do false beliefs about biological racial differences among white doctors mediate racial differences in recommended for hypothetical patients?	Experimental and survey data from U.S. medical students and residents (N=222 after restricting to white, US-born, native English-speaking).	Surveys and experimental vignettes.	Participants one SD above the mean in terms of false beliefs rated the Black patient as having 0.45 less pain than the white patient on a scale of 1-10 and were less accurate in recommendations for the Black patients.	Some statistics are disaggregated by medical school year or resident status, but sample sizes are too small to draw inferences.
McDevitt and Roberts (RAND 2014)	How does the availability of female urologists relate to rates of bladder cancer death among female patients?	AMA data on urologists 2006-2009; Florida hospital discharge data Jan. 2006 -June 2008; Florida Licensure Data; NCI’s State Cancer Profiles; Census, BEA, ARF for each market.	Descriptive statistics and a structural model to explain the distribution of female urologists across counties and the lack of entry.	Counties that have one more female urologist per 100,000 residents have 29.08% fewer female bladder cancer deaths per 100,000 residents. No significant associations between female urologists and male bladder cancer deaths or overall cancer deaths.	N/A.
Sabin and Greenwald (AJPH 2012)	Association between pediatricians’ scores on an implicit bias test (IAT) and racial	Survey data from 86 academic pediatricians conducted during Oct. and Sept. 2005.	Online survey with IAT tests plus patient vignettes describing children with pain following femur	Pro-white bias in the IAT predicts not giving oxycodone to the Black vignette patient in pain after bone surgery ($p < 0.05$).	N/A.

	differences in treatment?		fracture, UTIs, ADHD, asthma.		
Singh and Venkataramani (NBER WP 2022)	How do racial disparities in in-hospital mortality vary with strain on hospital capacity?	EHR with time stamps from 2 “highly regarded” academic hospitals serving predominantly Black patients.	OLS with rich controls; Assume hospital capacity strain at patient arrival is conditionally independent of mortality risk.	No racial difference in conditional patient mortality in quintiles 1-4 of hospital capacity strain. In 5th quintile, Black patients are 0.4 pp more likely to die on a baseline of 2%.	Effects are larger for Black women and Black patients without insurance. Effects driven by high-risk patients.
Wallis et al. (JAMA Surgery 2022)	How does surgeon-patient sex concordance affect post-operative outcomes?	Ontario Health Insurance Plan data; CIHI Discharge Abstracts and Ambulatory Care Reporting Services System; Registered Persons Data; Corporate Provider Database.	Population-based, retrospective cohort study.	Sex discordance was associated with increased likelihood of death (adjusted odds ratio 1.07) and complications (adjusted odds ratio 1.09), but not readmission.	They disaggregate by patient sex and find that effects are driven by male surgeons treating female patients. They also find stronger effects for cardiothoracic surgery.

Note: See Glossary for abbreviations.

Appendix Table 2: Effect of Experience and Training on Doctor Skills

Paper	Research Question	Data	Empirical Methods	Results	Heterogenous Effects?
Chan and Chen (NBER WP 2023)	How do NPs compare to doctors with respect to patient outcomes and resource use in the ED? How does variation in provider skill vary across and within professions?	Administrative health records from the VHA for ED visits between 01/2017 and 01/2020 (1.1 million cases, 44 EDs) linked to death records.	Use number of NPs on duty as IV for assignment to an NP vs. a doctor on arrival at the ED.	Assignment to an NP increases patient length of stay by 11%, increases cost of care by 7%, and increases 30-day preventable hospitalizations by 20%. Productivity variation greater within than between each profession.	NP-physician performance gap smaller for experienced providers and larger for patients with complex or severe conditions. Many NPs more skilled than some doctors.
Currie and Zhang (ReStat 2023)	Are some physicians more effective in promoting patient health? Correlation in effectiveness across domains of patient care? Do effective providers have lower/higher costs?	EHR data from the Veterans Health Administration's Corporate Data Warehouse for 2004 to Feb. 2020, VHA Vital Status files, CDC National Death Index Plus files.	Quasi-random assignment of veterans to PCP teams in the VHA system; value-added measure of provider effectiveness.	PCPs with 1 SD higher effectiveness re: mental health, circulatory conditions, or ACSC have a 27-44% reduction in adverse outcomes. Effectiveness measures positively correlated. PCPs with a 1 SD higher effectiveness 3.6-4.2% lower mortality and 2.5-5.4% lower costs over the next three years.	Provider effectiveness increases with provider age and number of patients seen.
Doyle, Ewer, and Wagner (JHE 2010)	Do residents from highly ranked programs do better than residents from lower ranked programs re: costs and health outcomes?	Veteran's Administration inpatient data 1993-2006; 2000 Census zip code level data.	Residency teams randomly assigned to patients based on the last digit of the SSN.	Patients assigned residents from lower ranked program had 11.96% longer stays and 13.31% higher costs. No differences in health outcomes.	Differences in costs were higher for more serious conditions.
Doyle (NBER WP 2020)	Does having cardiologists in the ER affect treatment and outcomes for patients with heart failure? Does additional experience with heart failure patients affect outcomes?	Medicare claims data (1998-2002) linked to mortality data; AMA's Masterfile for physician characteristics.	Estimate the effect of the share of physicians of different types in the ER, conditional on hospital*quarter *day-of-week FE.	Given number of physicians available, 1-year mortality falls by 1.10% with each additional cardiologist. More cardiologists increase intensity of care. Seeing 10 more heart failure patients yearly reduces mortality 1.2%.	Mortality point estimates larger for patients with higher predicted mortality, in high-volume hospitals, and for patients seen on slow days but differences imprecisely estimated.
Epstein, Nicholson, and Asch (AJHE 2016)	Compare effect of initial skill to the effect of experience in predicting obstetrician performance?	Florida and New York all-payer discharge databases (1992 to 2012); AMA Physician Masterfile; AMA FREIDA identifiers of hospitals with OB residency training.	Initial skill defined as physician's normalized, risk-adjusted maternal complication rate in the 1 st year.	Without hospital FE, initial skill explains much of the variance in performance. After 16 years, it explains 39-75% of performance. With hospital FEs initial skill explains only 1-9%, suggesting better doctors go to better hospitals.	Privately insured patients respond to recent measures of physician skill. Robustness checks with physician "stayers" only show similar results.

				Experience explains little.	
Facchini (HE 2022)	Does the recent volume of C-sections performed affect the outcomes of a surgeon performing a nonelective C-section?	Birth certificates from a large public hospital in Tuscany, Italy (2011 to 2014)	Patients cannot select their surgeon though more skilled surgeons may get harder cases. Include surgeon FEs.	Recent experience defined as #C-sections in the last 4 weeks. A one SD increase in experience reduces NICU admission 13.86% and reduces low APGAR 13.19%.	N/A.
Gowrisankaran, Joiner, and Léger (Management Science 2023)	How are measures of physician practice style and of physician skill correlated in the context of patients visiting the ED?	La Régie de l'assurance maladie du Québec (RAMQ) data on Montreal patients who visited an ED between April and Dec. 2006.	Identification via conditional random assignment of patients within an ED. Physician practice style and skill estimated from physician FEs.	Physicians with more intensive practice style have worse outcomes on average. Practice intensity correlated across conditions, as is skill.	Negative correlation intensive practice style and patient outcomes strongest for appendicitis, weakest for transient ischemic attacks.
Schnell and Currie (AJHE 2018)	How does a doctor's medical school rank affect their propensity to prescribe opioids? How does this relationship vary over time and between specialties with different levels of training in pain relief?	QuintilesIMS opioid prescription data 2006-2014; US News and World Reports; CMS provider utilization and payment data; ACS data; Mortality data.	FE models (specialty, county of practice, practice address).	Physicians from the lowest ranked medical school are 121% more likely to prescribe any opioids and prescribe 160% more than physicians trained at the top school.	Rank doesn't matter for specialties with pain medicine training. Rank matters less for more recent cohorts. Foreign physicians from low prescribing areas have low prescription rates.
Simeonova, Skipper, and Thingholm (JHR 2024)	Do health management skills (HMS) of primary care physicians affect medication adherence and hospitalizations for cardiovascular (CV) disease, and CV hospital costs of patients on statins? Do skills change with age?	Danish registry data on population of statin users and their PCPs (01/2004-06/2008). However, cannot observe PCP for 54% of clinics.	Leave-one-out adherence rates for each physician adjusted for patient and physician observables. Event studies after changes in PCP induced by clinic closures or patient moves.	A one SD increase in PCP HMS is associated with a 1.10% increase in medication adherence and 1.47% fall in CV hospitalization. CV hospital expenditures fall by 0.298%. Skill declines with physician age.	N/A.
Van Parys (PLOS One 2016)	How are variations in ED physicians' treatment of minor injuries related to physician characteristics including experience? Does practice style explain persistence as an ED physician?	All Florida ED visits for minor injuries 2005-2011 matched to Florida Healthcare Practitioner Database; HCUP databases.	OLS assuming little systematic matching of physicians and patients, conditional on observables.	Physicians with <2 years of experience spend 4.60% more and perform 3.46% more procedures than physicians with 7+ years. High-cost physicians are 3% less likely to work in a Florida ED 2 years after start.	Differences in care intensity fall with experience after 2-7 years of experience.

Note: See Glossary for abbreviations.

Appendix Table 3: Time Pressure and Fatigue

Paper	Research Question	Data	Empirical Methods	Results	Heterogeneous Effects?
Chan (Econometrica 2018)	How does ER physician decision-making change over the course of a shift?	Data on physician shifts from the ER in a large, U.S. academic, tertiary-care center 06/2005-12/2012.	Exploits randomness and pre-determination of shifts and overlap in shifts. Counter-factual simulations of patient assignments.	8.70% shorter visits in the 4th to last hour before shift ends, 44.40% shorter in last hour. Patients arriving in last hour have 10.44% more tests/treatments, a 5.7 pp (21.19%) higher likelihood of admission, and 23.12% higher total costs. No significant effects beyond the last hour. No effects found with respect to 30-day mortality or 14-day bounce back.	The effects on workload-adjusted length-of-stay are greater in the daytime and disappear if the index physician has enough time to offload cases to the incoming physician.
Chu et al. (Working Paper 2024)	How does cognitive load affect how a physician takes notes, orders tests, and treats patients?	High frequency “click stream” data from EHRs, for patients over 18 at the UCSF ED (2017-2019)	Cognitive load proxied by complexity of caseloads. Predict physician orders from past orders; measure deviations in actual orders as a function of load.	When load is high, physicians reduce note editing by 7-14% and increase diagnostic orders by 2-5%, with higher entropy in diagnostic tests. For every 1 SD from expected orders induced by cognitive load, probability of admission increases 3.4 p.p. (14%).	N/A.
Costa-Ramón et al. (JHE 2018)	How does time of delivery affect unscheduled C-sections, and infant health.	6163 births in 4 Spanish public hospitals 2014-2016. Scheduled and breech excluded.	IV estimation using an indicator for births between 11 p.m. and 4 a.m.	Unplanned C-sections increase by 53.21% between 11 p.m. and 4 a.m. There is a negative effect on 1-minute and 5-minute APGAR (-0.992 and -0.936).	N/A
Freedman et al. (JHE 2021)	Unexpected scheduling changes and decisions of PCPs.	EMR data on all visits to 31 primary care centers in a health system 2005- 2015.	Physician FE models with unexpected schedule changes in minutes as the independent variable.	10-minute increase in waiting time reduces total/new (0.19%/0.14%), referrals (0.32%), opioid Rx (0.33%), pap tests (0.39%). Increases scheduled/unscheduled follow ups (0.80%/0.50%), inpatient visits within 14/30 days (1.15%/1.85%), and hospital care within 30 days (0.17%). No effect on ER visits, imaging, antibiotic Rx, diabetes management.	Effects with respect to PT referrals and opioid Rx among opioid-naïve patients are not significant in the baseline specification.
Gruber, Hoe, and Stoye (ReStat 2021)	Studies an English policy limiting ER wait times to 4	Records of all visits to public hospitals at the visit level	Bunching estimator using the four-hour target. Assumes that	Wait times fell 8% in patients with wait times of 180-400 minutes, and by 59 minutes for patients moved from the post-threshold period to the pre-	Larger wait time effects and mortality for sicker patients. No significant difference in

	hours for 95% of patients at public hospitals.	linked to vital statistics mortality records for 4/2011-03/2013.	only patients around the four-hour mark are affected.	threshold period. Increased 30-day total costs (4.9%); hospital admissions (12.2%); tests in the ER (4.6%); Decreased 30/90-day mortality (13.8%/7.9%); discharge probability (7%); referrals (8.9%). No effect on 1-year mortality, length of stay or number of inpatient procedures.	probability of hospital admission. Most mortality reduction driven by circulatory, respiratory, and digestive problem deaths.
Linder et al. (JAMA IM 2014)	How does time in shift relate to the decision to prescribe antibiotics?	Billing and EMRs for 23 Partners HealthCare-affiliated PCPs 05/2011-09/2012.	Logistic regression.	Relative to the first hour of a shift, adjusted odds ratios of antibiotic prescribing in the 2nd, 3rd, and 4th hours were 1.01, 1.14, and 1.26. 44.46% of the sample was prescribed antibiotics.	N/A.
Neprash et al. (JAMA HF 2023)	What is the association between primary care visit length and inappropriate prescribing?	Claims and EHR data from AthenaHealth Inc., 2017.	Descriptive; linear probability models with physician FEs and patient covariates.	An additional minute of visit duration decreases inappropriate antibiotic prescribing 0.11 pp (0.2%), opioid and benzodiazepine co-prescribing for pain 0.01 pp (0.3%), and a prescribing of medications from the Beers List to older adults 0.004 pp (0.4%).	For patients with an anxiety and pain, each additional minute of visit duration decreased dangerous opioid and benzodiazepine co-prescribing 0.05 pp.
Shurtz et al. (RAND 2022)	Do PCPs increase treatment intensity and screening in response to time pressure caused by absent colleagues?	Administrative data from the largest HMO in Israel covering all primary care visits in Jerusalem 2011-2014.	Event studies at physician-day level. IV for visit length is %caseload missing physicians. (Alt. IV= any doctors missing). Nonparametric methods bound ATE.	A 1 minute longer visit increases use of any diagnostic input 4.50% and referrals 7.93%. No significant effects on imaging, pain killer Rx, antibiotic Rx, additional visits.	Effects on use of diagnostic tools bigger for older patients (>60 years) and patients with higher predicted utilization of primary care.
Persson et al. (HE 2019)	How are orthopedic surgeons' decisions affected by the number of patients already seen in a shift?	848 Swedish orthopedic clinic visits spanning 133 work shifts by eight surgeons between 10/2015-12/2015.	Logits with surgeon FEs, assuming patient allocation to time slots is exogenous conditional on observables.	Every additional patient already seen decreases the odds an operation is scheduled by 10.5% (OR = 0.895, CI 0.842 to 0.951). Patients seen in the afternoon are 1.955x more likely to be scheduled for surgery (CI 1.110 to 3.486). Surgery prescribed in 32% of cases.	N/A.
Tai-Seale and McGuire (HE 2012)	Do physicians have a target time per patient?	385 video-taped visits 1998-2000 with 35 PCPs; patient surveys.	Logits on the probability of a topic being the last of the visit.	Topics in the 1 st 5 minutes=reference group. Probability of a topic being last increases by 16.8 pp, 26.8 pp, and 35.7 pp for topics raised at 5-10, 10-15, 15+ minutes.	Academic medical centers demonstrated sharpest increase in the shadow price of time.

Note: See Glossary for abbreviations.

Appendix Table 4: Peer Effects and Team Dynamics

Paper	Research Question	Data	Empirical Methods	Results	Heterogeneous Effects?
Agha and Molitor (ReStat 2018)	Does proximity to lead investigators in new cancer drug trials increase the propensity to prescribe new drugs?	Medicare Part B claims 1998-2008; Dartmouth Atlas data; FDA drug application data.	DiD, patient location IV (secondary analysis).	Cancer patients in lead investigator's HHR 4.04 pp (36%) more likely to get new cancer drug, with convergence after 4 years. No effect in other authors' HHRs. IV estimates smaller.	Effects bigger in areas with slower drug adoption. Convergence suggests lead investigators are not in areas with higher latent demand for the cancer drug.
Chan (JPE 2016)	Is doctor shirking reduced when doctors vs. nurse schedulers do patient assignments?	6 years of ED data from an academic medical center. ED had 2 pods of doctors.	A nurse-managed pod became doctor-managed, as the other pod was.	The doctor-managed system reduced patient wait times by 13.67% with no significant effects on quality, cost, or utilization.	Patient assignment is more negatively correlated with a physician's number of patients in doctor-managed systems.
Chan (AEJ: EP 2021)	How much influence do senior residents have on team decisions? How do junior resident's decisions vary with experience?	Five years of data from the internal medicine residency program of a large teaching hospital.	RE model exploiting discontinuity caused by promotion of junior residents to senior.	There is a jump in the SD of log costs after promotion. Senior residents are responsible for almost all of the variance in decision making within a team of residents.	The jump in practice variation is highest for diagnostic spending (vs. medication, blood work, or nursing). No differences by patient characteristics.
Chen (AER 2021)	How does the length of time that PCI/CABG surgeons and other hospital physicians have worked together affect patient outcomes?	20% of Medicare claims 2008-2016 linked to Vital Statistics, MD-PPAS 2008-2016, Physician Compare 2014-2017.	1. Use admissions through ED. Include FE for proceduralists. 2. TWFE model with FE for proceduralists and PCPs.	1 SD increase in shared work experience reduces 30-day mortality by 10 to 14%. Shared work experience decreases use of medical resources and length of stay.	Effect of shared work experience declines with individual physicians' experience, but this decline is small. The effect is larger for more complex cases.
Molitor (AEJ: EP 2018)	How are cardiologists affected when they move to areas with different practice styles?	Medicare fee-for-service claims 1998-2012; AMA Masterfile;	"Movers" design follows cardiologists moves across HRRs; event study and difference-in-differences.	A 1pp increase in cardiac catheterization in the new HRR increases physician's own rate 1.36%. A 1pp increase in the rate at the physician's hospital leads to a 1.72% increase in the physician's own rate.	Effects larger for moves from low to high-intensity areas. Effects similar for moving earlier vs. later in their careers. Effects of moving are larger for more marginally appropriate patients.
Silver (ReStud 2021)	How do peer-groups affect speed and outcomes in the ED?	ED visits from New York (2005-2013). Linked to state physician licenses, public physician profiles, VS mortality data.	Peers vary across shifts. Decompose variation in outcomes due to physicians and physician-peer matches. Peer group is IV for outcomes.	First-Stage: A 10% increase in the speed of peers increases own speed 1.47% with controls. 2SLS: A peer group that increases physician speed by 10% decreases charges by 2.17% with no significant effect on the 30-day mortality of discharged patients.	Physicians work faster in smaller groups and when all peers are male. 2SLS: In at-risk patients, peer groups that increase speed by 10% decrease charges 2.55% and increase 30-day mortality in discharged patients by 5.65%.

Appendix Table 5: U.S. Financial Incentives

Paper	Research Question	Data	Empirical Methods	Results	Heterogeneous Effects?
Alexander (JPE 2020)	When hospitals offer incentives to physicians to lower costs, does it affect (1) who is admitted (2) which hospital they are admitted to, and (3) how intensely patients are treated?	New Jersey Uniform Billing Records (2006-2013); American Hospital Association annual survey; Medicare cost-to-charge ratio series.	DiD with doctor FEs using the New Jersey Gainsharing Demonstration as a policy experiment.	The policy doesn't reduce costs or change procedure choice. But lower predicted cost patients are sorted towards participating hospitals.	Effects are less precisely estimated for surgical patients, where there is less opportunity for gaming.
Alexander and Currie (EHB 2017)	What is the effect of private vs. public insurance on propensity to be admitted to hospital from ED? How are effects moderated by capacity constraints?	New Jersey Uniform Billing Records 2006-2012.	Exogenous variation in hospital bed supply due to local flu conditions; hospital FEs.	In high flu weeks, publicly insured children are .3 p.p. (6.4%) less likely to be admitted for non-flu conditions compared to privately insured children. Outcomes are no worse for marginal children.	Effects are larger when restricting to diagnoses with mid-range admissions rates.
Alexander and Schnell (AEJ:AE forthcoming)	What was the impact of increasing Medicaid PCP payments in 2013 and 2014 to comply with the ACA?	State-level Medicaid reimbursement rates; NHIS (2009–2015); NAEP (2009, 2011, 2013).	DiD and event studies exploiting variation in effect of ACA rule given pre-ACA reimbursement rates.	A \$10 rise in payments decreases prob. doctors decline new Medicaid patients 10%, decreases prob. that parents have trouble finding a doctor for child by 25%. Increased payments increase doctor visits, improve reported health, and reduce school absences.	Effects on school absences are larger and more precisely estimated for younger students than older students.
Bisgaier and Rhodes (NEJM 2011)	How does public vs. private insurance affect the probability that specialists will accept new pediatric patients, and wait times?	Experimental with 546 paired calls to 273 specialty clinics.	Audit study. One call with public insurance and one with a month later with private insurance.	Medicaid-CHIP callers were 6.2 times more likely to be denied an appointment. Conditional on getting an appointment, Medicaid-CHIP callers waited 22 days longer.	N/A.
Chen and Lakdawalla (JHE 2019)	Do physician responses to changes in Medicare reimbursement vary with patient income?	Medicare Current Beneficiary Survey (MCBS) 1993- to 2002; Federal Registers from 1993 to 2002.	2SLS: Instruments= changes in fees from 1997 consolidation of Medicare areas and 1999 changes in expense estimation.	A 10% increase in patient income corresponds to a 0.0508 increase in the price elasticity for services (53% of the mean).	Diff. in physician responses to reforms with respect to patient income explain 53% of the increase in the gap in services received by high-income vs. low-income patients.

Chorniy, Currie, and Sonchak (JHE 2018)	How does switching from FFS to MMC affect children's treatment of asthma and ADHD?	60% random sample of all South Carolina (SC) Medicaid enrollees < 17, 2005-2015; Vital Statistics	Staggered roll out of MMC contracts with higher capitated payments for children with chronic conditions; child FEs.	Switching to MMC increased ADHD caseloads by 11.6% and asthma caseloads by 8.2%. No significant effects on hospitalization and increases in ER use.	
Clemens and Gottlieb (AER 2014)	How do changes to Medicare physician payment rates affect provision of care, technology adoption, and patient health?	Medicare Part B claims 1993-2005.	Natural experiment: 1997 consolidation of Medicare geographic areas. Event study with nearest-neighbor matching on counties.	Higher fees increase elective procedures and RVUs per physician. Imprecise effects on MRIs by non-radiologists. No effect on 4-year mortality for cardiac patients. Higher hospitalization for AMI within 1 year.	Point estimates suggest higher care elasticities of care for older patients and patients from states with more intense care.
Dickstein (WP 2017)	Are there differences in how physicians in capitated plans prescribe compare to physicians in non-capitated plans?	MarketScan: 2003-2005 Commercial Claims & Benefit Plan Design Data; County-level IRS Income; National Ambulatory Medical Care Survey.	Structural model, instrumenting drug price with sum of price changes within an insurer's plan for all other drugs.	Prescribers in capitated plans more likely to choose generic Rx. Patients have higher adherence and less medication switching but also higher relapse rates.	Lower drug switching may promote adherence but have negative effects on patients at highest risk of relapse.
Ding and Liu (JHE 2021)	How does capitation affect treatment of lower back pain?	MarketScan Commercial Claims 2003- 2006.	Plan history FEs and physician FEs	Providers with capitation use 12.2% fewer medical resources to evaluate and treat lower back pain with no effect on relapse probabilities.	Effects are biggest for physical therapy and diagnostic testing. But do capitated providers report all procedures?
Gupta (AER 2021)	Effects of the Hospital Readmissions Reduction Program (HRRP) on care quality and admissions for patients with heart attacks, heart failure, and pneumonia?	Medicare fee-for-service claims 07/2006-07/2006; 20% sample of all Medicare beneficiaries.	DiD, IV using baseline predicted readmission rate.	HRRP reduced 30-day readmissions by 10.5% and 30-day returns to the hospital by 6.92%. Little effect on admission decisions or upcoding. Increases in procedures for AMI patients and 8.87% fall in 1-year mortality.	Readmission rates only lower for patients initially admitted to index hospital, not transfer patients. Government hospitals responded less. Hospitals in at-risk systems and higher volume hospitals responded more.
Johnson et al. (NBER 2016)	Are OBs more/less likely to do unscheduled C-sections on own patients? Effects recent patients' laceration rates?	EMR and billing databases for three practice groups.	They use rotating call schedules of OB groups as a plausibly exogenous source of OB assignments.	OBs are 4 pp (25.97%) more likely to perform a C-section and 2.5 pp (25.0%) less likely to use vacuum or forceps on their own patients vs. another OB's.	Higher rates of recent lacerations increase the probability of C-section for an OB's own patients but not for other patients.

Johnson and Rehavi (AEJ: EP 2016)	How is the probability of C-section affected if the patient is a physician? Is there an interaction with financial incentives?	Confidential CA Vital Statistics data, 1996-2005; CA physician licensure data; TX birth data 1996-2003 and 2005-2007.	Comparison group is educated mothers. Nearest neighbor matching regressions for CA. Hospital fixed effects.	In California physicians are 1.17 pp (6.13%) less likely to have an unscheduled C-section at non-HMO hospitals. In Texas physicians are 2.09 pp (6.39%) less likely to receive any C-section.	Effects greater for physician parents in specialties related to childbirth. Comparing HMO-owned and non-HMO-owned hospitals in CA, financial incentives affect C-section rates only among non-physicians.
<i>Financial Incentives in Other countries</i>					
Allen, Fichera, and Sutton (HE 2016)	Examined a policy that increased payments 24% for outpatient vs. inpatient surgeries for cholecystectomies in English hospitals?	Hospital Episode Statistics from the NHS Information Centre for Health and Social Care from 12/2007-03/2011.	DiD using control procedures with similar recommended outpatient rates that were not affected by the policy.	Planned outpatient increased by 27% of mean in year before the change. Actual outpatient rate increased 29%. Reversion from laparoscopic to open surgery decreased. No effect on deaths or readmissions.	N/A.
Brekke et al. (JHE 2019)	How does GP compensation and relationship with patients affect their propensities to issue sick-leave certificates patients need to claim benefits?	Norwegian administrative data 2006-2014 linking health, national insurance, and labor market data.	Physicians see patients both in their own practices and in EDs where they do not face reputational effects. Models with physician and patient FEs.	GPs with a FFS contract are 34.63% more likely to issue sickness certificates for own patients vs. ED patients. For GPs with fixed salaries the gap is 24.15%.	For GPs with new practices, effects similar with FFS but disappear for fixed salary. The 24.15% effect for fixed salary driven by relationships with patients. Effects larger in areas with more GPs per capita and GPs with more openings.
Coudin, Pla, and Samson (HE 2015)	How did a French reform that increased the proportion of GPs subject to price regulation, affect the provision of health services?	Administrative INSEE-CNAMTS-DGFIP File on physicians for 2005-2008.	Fuzzy RD using increase in the requirements for GPs to “bill freely” in their contracts with public health insurance.	Price regulation increased the supply of medical care by 66.53% and the number of procedures by 84.23%.	Male GPs increase their labor supply more in response to the reform than female GPs. Males also increased home visits and prescriptions while female GPs did not.
Fortin et al. (JAE 2021)	Compare FFS contracts vs. contracts that pay a per diem plus a smaller amount per service. Effects on care rendered by pediatricians?	Doctor time-use survey linked to records from Health Insurance Organization of Quebec (19962002).	Structural discrete choice model with variation from a reform introducing an optional per diem plus payment contract.	Small changes in time spent with patients, but large changes in services by 5-12%.	Female doctors and younger doctors are more likely to switch to the per diem contract, but model does not consider the interaction of age and gender.
Wilding et al. (JHE 2022)	How did increased stringency of blood pressure	EHRs from Clinical Practice Research	Temporary increase in the stringency of	Stricter targets did not increase diagnoses of hypertension in new	Consider practice-level heterogeneity in meeting the

	targets for patients <80 affect English GPs' treatment and testing decisions for hypertensive patients?	Datalink (04/2010-03/2017); Health Survey for England.	targets that doctor's patients had to meet in order for them to be paid. DiD comparing patients over and under 80; bunching estimators.	patients but increased antihypertensive prescriptions 1.2 pp. Doctors did multiple tests when patients failed, reported more patients as exempt from reporting, and increased reports of patients exactly meeting targets, suggesting rounding.	target in the pre period. Lower-performing practices increased reporting of patients as exempt more than higher-performing practices, but other effects were similar. No data on health outcomes.
<i>Other types of incentives</i>					
Agha and Zeltzer (AEJ: EP 2022)	How do pharma payments affect the prescribing of physicians who only share patients with physicians who receive payments?	Medicare Part D (2014–2016); Open Payments database (2013–2016); CMS Referral Patterns; Physician Compare.	Event studies; DiD-style regressions with doctor-drug and drug-quarter-specialty FEs	Peers of physicians who receive payments for speaking, consulting, etc., increase prescribing of the promoted drug 1.8%. Spillovers account for ¼ of increased prescribing from payments.	Effects are larger for peer physicians with more shared patients with the physician receiving payments.
Carey, Daly, and Li (NBER WP 2024)	How do pharma payments affect the prescribing of physician-administered cancer drugs in Medicare?	Open Payments database; claims from 20% sample of Medicare FFS (2014–2018).	DiD and event study models with physician-drug and time-drug FEs.	Payments increase Rx for the marketed drug by 4% in the year after payment is received. No improvement in patient mortality following payments.	Targeted doctors increase treatment of patients with lower expected mortality.
Carey, Lieber, and Miller (JPubE 2021)	How does detailing affect physician prescribing behavior in terms of drug efficacy, and use of generics?	20% sample of Medicare Part D 2013-2015; Open Payments database; Drug efficacy data.	Event studies with physician by drug FEs	Prescribing of the detailed drug increases by 2.2% in the 6 months following payment. No significant effects on efficacy or transitions to generics.	Results are similar when restricting sample to physicians who receive small payments.
Chernew et al. (JHE 2021)	How much of the variation in prices for lower-limb MRIs is explained by physician referral patterns vs. patient characteristics?	Insurance claims from a large national insurer for 2013; data from the company's online price comparison tool; SK&A physician-level dataset.	Restrict to lower-limb MRIs without contrast since these are "shoppable, homogeneous MRI scans." Estimate models with referrer FEs.	Referrer FEs explain 52% of the variance in patient spending on lower-limb MRIs. Patient cost-sharing and characteristics explain less than 1%. Patient HHR FEs explain 2%. Going to the cheapest provider within the same driving distance would reduce spending 35.83%.	The mean vertically- integrated physician refers 52% of patients to a hospital-based MRI provider compared to 19% for non-vertically integrated physicians.
Frakes (AER 2013)	Does physician treatment of AMI and c-section converge towards national averages when states change	National Hospital Discharge Survey (1977-2005), Natality Data (1978-2004);	Event study exploiting variation in states adoption of national-standard rules;	After adoption of a national-standard rule, the deviations between state and national C-section rates fall by 4.87 pp (48.31%). Estimates for AMI are	Disaggregates by whether states have rates that are initially higher or lower rates than the national rate.

	malpractice laws to consider national rather than local norms of behavior?	Mortality Data (1977-2004).		noisier. No convergence in outcomes.	Convergence occurs in subsamples.
Howard and McCarthy (JHE 2021)	Did a DOJ investigation of Medicare fraud re: implantable cardiac defibrillators (ICDs) change practice?	All-payer data from Florida; ED data from Florida's Agency for Healthcare Administration.	DiD using ICD procedures not subject to the investigation as a control.	The investigation plus new checklists that were part of the settlement caused a 22% decline in unnecessary ICD implantations.	Decline in ICDs stronger for hospitals involved in the lawsuit. Decline for Medicare patients smaller in percent but larger in absolute terms compared to patients with other insurance.
Newham and Valente (JHE 2024)	How do gifts to doctors from pharmaceutical companies affect antidiabetic drug prescribing patterns and costs?	Open Payments database; Medicare Part D data (2014–2017); demographic and health data from ACS and CDC.	Compare physicians with similar propensities to receive payments. Leverage randomness in timing. Regress residuals from outcome models on residuals from payment models.	A \$10 payment increases Rx of branded antidiabetic drugs by 2.3%, increasing costs of Rx drugs.	Effects are higher for doctors in areas with a higher proportion of patients receiving subsidies for out-of-pocket drug costs for low-income individuals.
Shapiro (MS 2018)	Compare effect of new information from clinical trials and detailing on PCP prescribing behavior for Seroquel.	AlphaImpactRx monthly panel of 1,762 PCPs 2002–2006 (linked self-reported detailing and patient treatment information).	Two clinical trials over sample period, plus record of detailing. Examine effects in models with physician and month FEs.	No effect of the clinical trial information. Detailing increased after both trials. Detailing increased Seroquel Rx 26% in the month of the visit.	One third of the increase in prescribing occurred in off-label uses.

Note: See Glossary for abbreviations.

Appendix Table 6: Doctor Responses to New Information

Paper	Research Question	Data	Methods	Results	Heterogeneous Effects?
Avdic et al. (JHE 2024)	New stents were first thought to reduce complications and then to increase them. How did cardiologists respond to new information and guidelines?	Swedish Coronary Angiography and Angioplasty Registry 2002-2011.	Separate models for periods after positive info, after negative info, and after guidelines allow physician-specific intercepts and trends.	Doctors responded more quickly to negative information than to the initial positive information.	Doctors slow to take up new stents were more likely to use the appropriate stent and had better patient outcomes. No heterogeneity within hospitals. Slow responders more likely to practice in teaching hospitals.
Ahomaki, Pitkanen, Soppi, and Saastamoinen (JHE 2020)	Experiment with letters sent to Finnish doctors who prescribed 100+ paracetamol-codeine pills to a new patient.	National Prescription Register including all purchases, merged to Nordic Product Number and physician characteristics.	DiD using new patients where non-targeted physicians are the control. "Treatment" is intent-to-treat.	Significant 6.13 tablet decrease in number of pills purchased by new patients of treated doctors relative to patients of untreated doctors (12.8% of treatment group baseline).	Treatment effects larger for high prescribers. Top 5 specialties have similar effect size. The decrease in large purchases was greatest in urban areas and not significant in rural areas.
Bradford & Kleit (HE 2015)	The effect of the 2005 Blackbox warning on NSAID prescriptions, and how it was mediated by advertising, media coverage, and patient characteristics.	EMRs from the Primary Care Practices Research Network; media data from Competitive Media Reporting, Inc. and Lexis/Nexis; NSAID sample dispensation data from IMS health.	Probit models on having active prescription for non-COX-2 inhibitor NSAIDs, COX-2 inhibitor NSAIDs, opioids, and other analgesics.	Blackbox warnings resulted in a 2.8pp (54.90%) decrease in prescriptions for COX-2-inhibitors and 2.8pp (23.14%) increase in prescriptions for a non-COX-2-inhibitor (p<.001).	Patients with cardiovascular disease had a similar decrease in prescription of COX-2-inhibitors, but no significant increase in non-COX-2-inhibitors. These patients substituted toward opioids and other analgesics.
Doctor et al. (Science 2018)	Effect of notification of patient death by overdose on future opioid prescribing.	Opioid dispensing from California's Prescription Drug Monitoring Program database.	RCT with intent-to-treat analysis. Letters from CA's Chief Medical Examiner.	Milligram morphine equivalents prescribed down 9.7% in treatment vs. control 3 months after intervention.	N/A
Dubois and Tuncel (JHE 2021)	How did French physicians respond to the 2004 information that SSRIs increase suicidal thinking in children?	Cegedim proprietary longitudinal patient data covering all prescriptions by 386 GPs. Includes doctor and patient demographics, and visit-level information.	DiD estimation, older patients are control. Random coefficient discrete choice logit examines choice across drug categories.	Child SSRI prescriptions fell 9.9 pp (19.8%). The baseline effect for adults was -2.8 pp (5.6%). Many physicians decreased prescription of other classes of anti-depressants but substituted to off-label use of other drugs.	25% of the physicians prescribe an SSRI for depression <20% of the time before the warning, and 25% prescribe an SSRI >73% of the time. Over 25% of physicians never prescribe SSRIs to children after the warning.
Howard, David, and Hockenberry	Variation in surgeon responses to the	Outpatient claims data from Florida's State	Triple DiD, alternative model using	Preferred specification: if free-standing centers responded like	Disaggregating by procedure type, the differential decline

(JEMS 2016)	information that arthroscopic knee surgery is ineffective by whether it is a hospital or a free-standing surgery.	Ambulatory Surgery Database, 1998-2000. Surgeons cannot be linked over time. Analysis at facility level.	differential trends in the ratio of knee to shoulder surgeries (preferred specification).	hospitals the number of surgeries would be reduced 6.27-11.37% on a baseline of 34,000 each year.	between free-standing centers and hospital centers is driven by meniscectomies, which have received more insurance company scrutiny.
Howard and Hockenberry (HSR 2019)	How is physician age related to the response to new information that episiotomies are ineffective?	Pennsylvania Inpatient Hospital Discharge Data (1994–2010)	Descriptive. LPM with hospital FEs.	Physicians who started delivering babies 10 years earlier are 6 pp (19.5%) more likely to perform an episiotomy.	Relationship between physician age and episiotomy rate has decreased over time and is weaker in teaching hospitals.
Kolstad (AER 2013)	Effects of quality “report cards” for Coronary Artery Bypass Graft (CABG) surgeries. Is provider response profit motivated?	Pennsylvania Health Care Cost Containment Council data for 89,406 CABG surgeries 1994-1995, 2000, and 2002-2003 merged with surgeon tenure. Focus is on the surgeons’ mortality rate before report cards less the report card risk-adjusted rate.	Reduced form responses to differences between own mortality rates and other doctors’. Structural model of consumer demand separates “intrinsic” and “extrinsic” motivations.	Counterfactuals indicate that “extrinsic” incentives induced a 3.5% decline in predicted risk-adjusted mortality whereas “intrinsic” incentives induced a 13% decline in predicted risk-adjusted mortality.	The response is larger for surgeons who are worse than other surgeons in their own hospital compared to surgeons who are just worse than expected.
McKibbin (JHE 2023)	How do physicians change prescribing of off-label cancer drugs in response to new information from RCTs?	Data on FDA approvals and RCT results, 100% Outpatient and 20% Carrier Claims files for Medicare part B, 1999-2013.	Event studies comparing drug-cancer pairs with and without newly presented RCT evidence from academic conferences.	8 quarters after a conference, prescriptions of drugs with confirmed efficacy up 192%. Prescribing falls by 33% over 8 quarters with negative information.	Responses discontinuous around p-value 0.05. When the abstract has no mention of improvements in quality of life or side effects, adoption and de-adoption rates are less asymmetric.
Olson and Yin (HE 2021)	Physician responses to changes in drug labeling from the FDA's 1997 Pediatric Exclusivity provision (provides 6 months of exclusivity in return for conducting Pediatric trials).	NAMC prescription data; Label changes and exclusivity from FDA; journal publication data from Benjamin et al. (2006) and PubMed; IMS health data on drug promotions; disease prevalence from MEPS.	DiD with treatment group defined as children <18 years old and controls as adults >35 (using a zero-inflated negative binomial model).	In their preferred specification, the marginal effect of a pediatric label change is 2.09 fewer prescriptions (12.67 %) for children.	Negative information added to the label reduces prescribing more than positive information. Magnitudes are larger for physicians in solo practice. No clear pattern by child age group. Estimates somewhat sensitive to included controls.
Persson et al. (NBER WP 2021)	Do doctors consider the diagnosis of an older	Swedish population register 1990-2018, (2016	Birthday cut-off RD using older sib or	An older sibling born after the school entry cutoff decreases the	Effects on younger siblings are greater before older siblings

	sibling when evaluating children for ADHD?	for HS records); prescription drug claims July 2005-Dec. 2017; birth records data from NHBW, 1996-2016.	cousin's birth date and school eligibility cutoffs to use "young for grade" sib's higher prob. of ADHD diagnosis.	probability of ADHD diagnosis by 0.59 pp (12.04%) and decreases the probability of ADHD drug claims by 0.55 pp (9.82%). Smaller results for cousins.	graduate from HS. Spillovers greater in cities with more funding for special needs children. Cousin spillover effects are greater when cousins are in the same municipality.
Sacarny, Yokum, Finkelstein, and Agrawal (HA 2016)	Effect of letters from Medicare to outlier prescribers of controlled substances on future opioid prescriptions.	CMS Integrated Data Repository-- records for prescription drugs covered by Medicare Part D with prescriber ID.	RCT with analysis of intent-to-treat.	Statistically insignificant increase of 0.8% relative to the control mean after 90 days, 95% CI (-1.38%, 2.91%).	No evidence of heterogeneity by prescriber specialty, geographic region, prescribing pre-treatment, and whether the physician had been investigated for fraud.
Sacarny, Barnett, Le, Tetkoski, Yokum, and Agrawal (JAMA Psych 2018).	Effect of three letters sent by Medicare to outlier prescribers of quetiapine on future quetiapine prescriptions.	100% Medicare claims data 2013-2017; enrollment data 2015-2017; risk-adjustment data 2013-2014.	RCT with analysis of intent-to-treat.	11.1% fewer days over 9 months vs. control mean (11.99% of the sample mean). Effects lasted 2+ years. No negative effects on patients.	The reduction in prescribing was larger for patients with low-value indications and smaller for guideline-concordant patients.
Wu and David (JHE 2022)	How did relative procedural skill affect the prob. that doctors abandoned laparoscopic hysterectomy after a negative info shock about the safety of the procedure?	All hospital inpatient and outpatient visit data for patients receiving hysterectomies in Florida (January 2012 – Sept. 2015)	Leave-one-out IV for physician skill at laparotomy/ laparoscopic hysterectomy; DiD event study estimates before/after 2014 FDA announcement.	A 1 SD increased in relative skill in laparoscopic hysterectomy decreased prob. of abandoning the procedure by 4.6–4.9 p.p. (6.2–6.5% reduction from pre-period mean). Only top laparotomy doctors increased laparotomies.	Patients with characteristics that indicate less appropriateness for the laparoscopic procedure had greater reductions in likelihood of receiving a laparoscopic procedure after the announcement.

Note: See Glossary for abbreviations.

Appendix Table 7: Heuristics and Guidelines

Paper	Research Question	Data	Empirical Methods	Results	Heterogeneous Effects?
Abaluck et al. (NBER WP 2021)	How does the proportion of physicians following guidelines for anticoagulants for atrial fibrillation patients change after 2006 guidelines? Is lack of implementation due to awareness or nonadherence?	Text mining of EMRs from the VA for patients newly diagnosed with atrial fibrillation between Oct. 2002-Dec. 2013; Patient-level data for 8 clinical trials of anticoagulants.	Causal-forest model to estimate heterogenous treatment effects using data from eight RCTs; Chernozhukov et al. (2018) approach to calculating best linear predictions of conditional average treatment effects.	After 1 st mention of guidelines, physicians become more compliant. Stricter adherence could prevent 24% more strokes.	Most departures from guidelines are not justified by measurable treatment effect heterogeneity (though RCTs were not originally randomized on the observables analyzed).
Almond et al. (QJE 2010)	Does the care of newborns change discretely at the threshold for being classified “very low birthweight” and does this affect mortality?	NCHS linked birth/infant death files (1983-1991, 1995-2002); linked birth, death, hospital discharge data from California (1991-2002); HCUP for AZ, NJ, MD, NY.	RD centered around threshold of 1,500 grams.	Relative to the means just above the threshold, VLBW classification has an 11.11% effect on spending and a 5.93% effect on length of hospital stay.	Effects are greater for non-NICU and Level 0/1/2 NICU hospitals than for Level 3A-3D NICU hospitals.
Coussens (Working Paper 2022)	Do doctors use simple heuristics in patient age to make treatment decisions for ischemic heart disease (IHD)?	Truven Commercial Claims and Encounters database 2005-2013; ED records from a large Boston-area hospital 01/2010-05/2015.	Regression discontinuity centered at age 40	Turning 40 increases the probability of being tested, diagnosed, or admitted for IHD by 0.887pp, 0.131pp, and 0.068pp, respectively. Changes relative to intercepts are 9.51%, 19.29%, and 17.80%.	Effects are larger for women and patients presenting without chest pain. Effects are also stronger when the ED is less busy and in the 1 st half of a physician’s shift.
Cuddy and Currie (PNAS 2020)	What is the probability that adolescents with private insurance receive appropriate care following an initial diagnosis of mental illness? What factors are related to the type of care received?	Claims data for a large national insurer. Children covered for at least a year between 2012 and 2018 who were ever diagnosed with a mental health condition.	Observational study using linear probability models. Define “red-flag” treatment as prescribing that falls outside accepted guidelines.	Only 75% of adolescents receive follow-up care within 3 months. 44.85% with any drug receive “red flag” drugs. Composition of clinicians affects treatment: More psychiatrists \Rightarrow more drug use vs. more therapists \Rightarrow more therapy.	Any treatment, drug treatment, red-flag drugs increase with age. Girls more likely to be treated, to get therapy, and to get red-flag drugs. Variation across zip codes explains less than half of overall treatment variation.
Cuddy and Currie (JPE forthcoming)	Would adherence to guidelines improve outcomes? Is there a	Claims data for a large national insurer. Children diagnosed with depression	Instrument individual prescriptions with area-level practice style	Outcomes for red-flag vs. grey-area vs. FDA approved drug treatment after 24 months:	P(drug treatment) is higher for girls, older children, and children whose 1 st visit

	difference between “grey-area” prescribing sanctioned by professional societies but not by FDA, and “red-flag” prescribing not sanctioned by either?	or anxiety for the first time 2012-2018. Measures of local practice style computed from IQVIA and from the claims data.	measures interacted with patient characteristics (use Lasso to choose instrument set).	P(self-harm): 5.8%; 4.9%; 3.8%. P(ED or hosp.): 33.6%; 18.6%; 26.8%. Total costs: \$9557; \$1745; \$9658. Red-flag has highest costs and worst outcomes.	resulted in hospitalization.
Currie and MacLeod (Econometrica 2020)	Would adherence to professional guidelines improve outcomes? Does the answer to this question vary with the physician’s skill?	Claims data for a large national insurer. Adults ever diagnosed with depression 2013-2016; NPPES; Propensity to experiment measured using dispersion of prescriptions across drugs in IQVIA Xponent data.	Patient FE models of effects of having more experimental doctors and of violations of guidelines. Simulations measure benefits of experimentation for different skill groups. (Psychiatrists assumed more skilled than GPs).	Violations of professional guidelines are associated with worse subsequent outcomes (spending, hospitalizations, ED visits) for all patients.	Among patients seeing psychiatrists, switching to a more experimental doctor improves outcomes (a 0.25 increase reduces P(ED visit or hospitalization) by 10.2%). No effect of experimentation with less skilled doctors.
Geiger et al. (JAMA HF 2021)	What is the effect of a designation of “advanced maternal age” (AMA) on prenatal care and birth outcomes?	Claims and monthly enrollment data from a large, nationwide commercial insurer 2008-2009; zip-code level public ACS data.	Focus on discontinuities in care for mothers 35+ on expected delivery date. Donut RD excluding women with due dates within 7 days of their 35 th birthday.	AMA increases screening, specialty visits; decreases perinatal mortality by 0.39pp or 42.39% of sample mean. No effects on severe maternal morbidity, preterm birth, or low birth weight.	As a percentage of baseline, the effects on prenatal care services and perinatal mortality are much greater for low-risk pregnancies than for the full sample.
Kowalski (ReStud 2023)	Are women who are more likely to receive mammograms different from women who are less likely? How does the probability of being “over-diagnosed” vary with the propensity to receive mammograms?	RCT data from the Canadian National Breast Cancer Screening Study (CNBSS) linked to cancer registries and the mortality data. Allows long-term follow up to see cancers that are detected but would not have caused symptoms.	Extension of Imbens and Angrist (1994) framework in the context of an RCT (which provides identifying variation).	In treated compliers w.r.t. screening guidelines, 14% of breast cancers are “over-diagnosed”. For always takers, over 36% of breast cancers are over-diagnosed. Results suggest current guidelines should be revised to reduce screening in women 40-50.	Women who are more likely to receive mammograms are healthier and of higher socioeconomic status on average.
Olenski et al. (NEJM 2020)	Do doctors use simple heuristics in patient age to make treatment decisions for Coronary Artery Bypass Graft Surgery (CABG)?	Medicare data from 2006 to 2012.	Regression discontinuity at age 80.	Patients admitted in 2 weeks after their 80 th birthday were 28.05% less likely to get CABG than patients admitted 2 weeks before their birthday.	N/A.

Notes: See Glossary for abbreviations.

Appendix Table 8: Technology

Paper	Research Question	Data	Empirical Methods	Results	Heterogeneous Effects?
Agarwal et al. (NBER WP 2024)	How do radiologists use AI predictions and clinical histories in diagnosis? What is optimal use of AI?	Patient cases from Stanford University healthcare; data from an experiment on radiologist decisions and decision time.	2x2 experiment with radiologists. Add AI prediction, clinical history from referring doctor, or both; random forest regression.	AI does not improve performance. Access to clinical history reduces deviation from diagnostic standards by 4%. Optimal to have AI decide cases when confident and radiologists decide all other cases w/o AI.	When the AI tool has high confidence, AI improves radiologist diagnosis. When the tool has low confidence, AI worsens radiologist diagnostic accuracy.
Agha (JHE 2014)	Impact of EMRs plus clinical decision supports on quality and cost of care.	20% sample of Medicare claims, 1998- 2005; Health Information and Management System Survey.	Exploits differential timing of Health Information Technology (HIT) adoption at hospital level w FE.	HIT adoption increases spending 1.3%. No effect on 1-year patient mortality, length of stay, #physicians seen within a year of admission, intensity of care, 30-day readmissions, complications, or an index of care quality.	No evidence of higher returns to more comprehensive HIT systems. Do not see larger effects in larger hospitals.
Alpert, Dystra, and Jacobson (AEJ: EP 2024)	How much does information versus hassle costs from MA-PDMPs affect opioid prescribing?	Claims data from Optum’s Clinformatics Data Mart (2006–2016).	DiD and event studies using policy change in Kentucky. Triple differences comparing opioid naïve and non-naïve patients.	Hassle and information explain 69% and 31% of fall in opioid Rx respectively. MA-PDMPs reduce opioid Rx 6.8% for opioid naïve patients, 10.6% for non-naïve patients, and 16% for patients with opioid-inappropriate conditions.	Declines in prescribing to opioid non-naïve patients occur for patients with history of doctor shopping or high dose/quantity of opioid use.
Arrow, Bilir, and Sorenson (AEJ: AE 2020)	Does access to an electronic database for pharmaceuticals affect doctors’ prescribing of cholesterol drugs?	IMS Health Xponent database 2000-2010; data from the firm that owns the studied electronic reference database.	Models with zip-code-month FEs, physician FEs, and physician-specific time trend; IV doctor’s access using share of area doctors using database.	Database increases prescribing of generic Rx in its 1st year by 1.3 pp (3.7%). No effect on new branded Rx. New and old generic Rx increase; Old branded Rx decrease. Providers prescribe 0.7% more unique Rx.	In zip codes with more pharmaceutical patenting, database has less effect on drug adoption. Effects stronger for providers who access the database more frequently upon adoption.
Buchmueller and Carey (AEJ: EP 2018)	How do MA-PDMPs versus PDMPs without must-access provisions affect opioid use in Medicare?	PDMP info from Prescription Drug Abuse Policy System; 5% Medicare beneficiaries in Part D and FFS in any year 2007–2013.	DiD and event study models using variation in state-level policy.	Without must-access provisions PDMPs have no effect on opioid utilization. MA-PDMPs reduce doctor shopping by 8% and pharmacy shopping by 15%. Neither PDMP significantly affects opioid poisoning rates.	Effect sizes are larger must access provisions are broader.
Buchmueller,	Effect of Kentucky’s	Kentucky (2006-2016)	DiD comparing	Quarterly morphine equivalents per	Providers who initially

Carey, and Meille (HE 2020)	must-access PDMP program on opioid prescribing.	and Indiana (2012-2016) PDMPs; CDC data on opioid prescriptions; ARCOS 2006-2016.	Kentucky (treated) to Indiana (control).	capita fell 11–13% in KY vs. IN. Providers prescribing any opioids fell by 3.8 pp (5%). The number of patients prescribed fell 16% among providers prescribing any opioids.	prescribed fewer opioids were more likely to stop prescribing. Greater falls in prescribing for patients with multiple opioid Rx and doctor-shoppers.
Dahlstrand (Working Paper, 2021 updated 2024)	How much could patient outcomes be improved by using an algorithm to match patients and GPs?	Data from Sweden’s largest digital healthcare platform (2016–2018) matched to Swedish registry health data.	Physician skill estimated using leave-one-out measures with shrinkage. Match effects use platform’s conditional random assignment of patients.	Using an algorithm with positive assortative matching could reduce avoidable hospitalizations by 8%, all hospitalizations by 3%, and counter-guideline antibiotic Rx by 3%.	Effects are smaller for patients seeing a doctor within the day/hour. In urban areas, similar improvements are possible by restricting matches to doctors who patients can travel to see in person.
Ellyson, Grooms, and Ortega (HE 2022)	Do the effects of must-access PDMPs vary by specialty?	CMS Part D public use files 2010–2017; AMA Physician Masterfile; PDMP start dates from Prescription Drug Abuse Policy System.	DiD and event study.	Primary care doctors decrease opioid prescribing by 4% after MA-PDMP implementation. No significant effect for providers in IM, EM, surgery, palliative care, oncology, and pain medicine.	Primary care and IM providers with initially low prescribing stop prescribing opioids after MA-PDMP.
Goetz (International Journal of Industrial Organization 2023)	How does an increase in competition on a telehealth platform affect providers’ pricing and exit decisions?	Therapist data from Psychology Today in 2020; controls from Canadian government sources and Facebook’s Movement Range maps.	Propensity score matched DiD exploiting change in how platform shows providers to patients. For areas with <20 providers, platform made providers outside area visible.	Increased competition caused by the platform displaying more providers decreases the likelihood that affected providers provide sliding scale discounts by 8.9%.	Providers with more training respond to competition by stopping sliding scale offers; providers with less training exit the platform. Bigger effects on late adopters of teletherapy.
Horwitz et al. (NBER Working Paper 2024)	How do Certificate of Need (CON) laws affect imaging? How does this vary by the value of imaging?	Hand-coded laws; AHA’s Annual Survey 2018; accreditor data on free-standing CT/MRIs; 20% sample Medicare FFS claims 2009–2014.	RDD at state borders where one state has a CON law and the other does not.	The prob. of receiving an MRI is 2% lower on the CON side of the state border, compared to the mean on the non-CON side. Overall, no effect on prob. of a CT.	The prob. of receiving a high-value MRI does not change at border, the prob. of receiving a high-value CT on the CON side falls by 6% of non-CON mean. Low-value imaging falls 20–26%.
McCullough et al. (HA 2010)	How is quality of care related to EMR adoption 2004-2007?	AHA’s annual survey; Health Information and Management Systems Society	OLS with hospital and year fixed effects. Estimated effect on the one-year lag of EMR	Pneumococcal vaccination rates up 2.1pp (3.2%); use most appropriate antibiotic for pneumonia up 1.3pp (1.6%). No effect on other quality of	The relationship between quality measures and EMR adoption is stronger in academic vs. non-academic hospitals.

		Analytics database.	adoption.	care measures studied.	
Miller and Tucker (JPE 2011)	Does EMR adoption lower neonatal mortality.	Linked birth and infant death data 1995–2006; AHA surveys; BEA Regional Accounts; CBP; HIMS Analytics Data; Georgetown Health Privacy Project; Lexis-Nexis for laws.	Construct balanced county-level panel over 12 years. OLS w county and year FEs; IV for EMR adoption using state medical privacy laws.	A 10% increase in EMR adoption reduces neonatal mortality by 3%. Reductions are due to prematurity and complications not to accidents, SIDS, or congenital defects.	Larger effects when EMRs combined with digital storage, and obstetric-specific/decision support technologies. Larger gains for mothers who are Black, Hispanic, unmarried, or have < high school education.
Neumark and Savych (AJHE 2023)	How do MA-PDMPs and laws that limit initial opioid Rx length for patients with work-related injuries?	Workers Compensation Research Institute claims for workers injured Oct. 2009 – March 2018.	DiD using state-level variation in laws.	Laws that limit opioid Rx length have no effect on opioid Rx (w/pre-trend w/o state trends). MA-PDMPs reduce opioid Rx on intensive but not extensive margin. For neuro spine pain, non-opioid pain Rx increase 14%.	Effects of MA-PDMPs are larger for neurologic spine pain, spine sprains and strains, and other sprains and strains cases.
Obermeyer et al. (Science 2019)	Is there racial bias in algorithms used to target care for high-risk patients? Do doctors correct for algorithmic biases?	Data from all primary care patients enrolled in risk-based contracts at a large academic medical center, 2013–2015.	Descriptive statistics and simulations.	Conditional on chronic condition, Black patients get less recommended care. Black patients have 26% more chronic conditions at the 97 th percentile of the risk score. Simulations suggest that physicians do not counteract bias in the algorithms.	Algorithm trained on spending. Conditional on diagnosis, Black patients have lower spending. Changing algorithm to target health outcomes could potentially resolve the problem.
Mullainathan and Obermeyer (QJE 2022)	Ask how the actual decision to test for heart attacks differs from algorithmically predicted risks and explore health implications.	“Large urban hospital’s” HER from Jan. 2010 to May 2015 linked to Social Security Death Index; 20% sample Medicare FFS claims Jan. 2009 to June 2013.	Descriptive comparisons of output from risk model and actual physician decisions; shift-to-shift variation in average testing rates associated with triage team.	Physicians over test low-risk patients and under test high-risk patients because they focus on salient and representative symptoms, ignoring more complicated predictors of risk. High risk patients who arrive at the ED during high-testing shifts have 32% lower 1-year mortality.	Stress testing is more overused than catheterization. More experienced physicians test less but more accurately target tests toward high-risk patients.
Sacks et al. (JHE 2021)	What are the effects of MA-PDMPs and laws that limit initial opioid Rx length on opioid-naïve patients?	Commercial claims from “large, national insurer” (20% sample and 100% sample for patients w/opioid Rx) Jan. 2007–Apr. 2018.	DiD using state-level variation in laws.	MA-PDMPs decrease hazard of a new opioid Rx by 4.7%. Laws that limit initial Rx length increase hazard of new opioid Rx by 8.7%—reductions in Rx for >7 days are more than offset by increase in Rx for <7 days.	Rise in new opioid Rx in response to laws that limit initial opioid Rx length is stronger for PCPs. Laws may inadvertently signal that short prescriptions are safe.

Van Parys and Brown (NBER WP 2023)	Did broadband access improve the outcome of joint replacement surgeries?	FCC data on broadband roll-out; Medicare Current Beneficiary Survey for internet use; Medicare FFS Claims 1999–2014.	DiD exploiting staggered rollout of broadband; discrete choice model.	Broadband access explains 16% of the improvement in joint replacement outcomes between 1999-2008. 10% stems from patients seeking better providers and 6% stems from improvements in care conditional on patient demand.	Improvements in outcomes due to hospital access to broadband are driven by hospitals in markets with less competition.
Zeltzer et al. (JHE 2023)	How does the adoption of a digital device to assist with telehealth visits affect health care?	EHR data from Israeli Clait Health Services (covering ~1/2 the Israeli population) 2018–2022.	Matched DiD and event study.	Device-assisted telemedicine increases primary care visits 12%, increases antibiotic use 15.6%, and decreases urgent care/ED/inpatient visits 11–24% compared to baseline mean.	Adults have a smaller increase in primary care use and a larger decrease in urgent care/ER/inpatient visits than pediatric patients.
Zeltzer et al. (JEEA 2024)	Impact of increased access to telemedicine during COVID-19 after lockdowns lifted were in May–June 2020.	EHR data from Israeli Clait Health Services from January 2019 to June 2020.	DiD at the patient level. Treatment is a patients’ propensity to use telemedicine during the initial March–May 2020 lockdown.	Having a PCP who was a high user of telemedicine increased the prob. of a primary care visit by 3.6% but reduced visit costs by 5.7% (of the pre-lockdown mean). Visits had fewer Rx and referrals. No evidence missed diagnoses for patients of high adopters.	Effects measured in % changes with respect to baseline are similar across patient age, gender, and SES. Reduction in Rx larger for providers who prescribed more in the pre-period.

Note: See Glossary for abbreviations.

Glossary of Table Abbreviations

AHA – American Hospital Association
AKM– Abowd, Kramarz, and Margolis (1999)
AMA – American Medical Association
AMI/MI –Acute myocardial infarction
ATE—Average Treatment Effect
CCI—Charlson Comorbidity Index
CDC – Center for Disease Control and Prevention
CMS –Centers for Medicare and Medicaid Services
CPOE – Computerized provider order entry
DEA – Drug Enforcement Authority
DiD – Difference in differences
DO – Doctor of Osteopathic Medicine
ED/ER – Emergency department
EMR/EHR – Electronic medical/health record
FCC—Federal Communications Commission
FDA— Food and Drug Administration (United States)
FE – Fixed effects
FFS—Fee-for-service
GP—General Practitioner
HCUP – Health care utilization project
HIT – Health information technology
HRR – Hospital referral regions (from the Dartmouth Atlas)
IV –Instrumental variable
MA-PDMP – Must-Access Prescription Drug Monitoring Program
MD – Medical Doctor
MMC—Medicaid managed care
NCHS -- National Center for Health Statistics
NHS—National Health Service (U.K., Norway)
NPI – National Provider Identifier
OR – Odds ratio
PCP –Primary care provider
PDMP – Prescription drug monitoring program
pp – Percentage point
PSI – Patient safety indicator
RCT – Randomized controlled trial

RD—Regression discontinuity
Rx—Prescription
SES – Socioeconomic status
SSRI—Selective Serotonin Reuptake Inhibitor
VHA—Veterans Health Administration (United States)
VS—Vital Statistics

Bibliography for Table Citations

- Abaluck, Jason, Leila Agha, Jr Chan David C., Daniel Singer, and Diana Zhu, “Fixing Misallocation with Guidelines: Awareness vs. Adherence,” Working Paper, 2021 (National Bureau of Economic Research).
- Abowd, John M., Francis Kramarz, and David N. Margolis. “High Wage Workers and High Wage Firms.” *Econometrica* 67, no. 2 (1999): 251–333.
- Agarwal, Nikhil, Alex Moehring, Pranav Rajpurkar, and Tobias Salz, “Combining Human Expertise with Artificial Intelligence: Experimental Evidence from Radiology,” Working Paper, 2023 (National Bureau of Economic Research).
- Agha, Leila, “The effects of health information technology on the costs and quality of medical care,” *Journal of Health Economics*, 34 (2014), 19–30.
- Agha, Leila, and David Molitor, “The Local Influence Of Pioneer Investigators On Technology Adoption: Evidence From New Cancer Drugs,” *The review of economics and statistics*, 100 (2018), 29–44.
- Agha, Leila, and Dan Zeltzer, “Drug Diffusion through Peer Networks: The Influence of Industry Payments,” *American Economic Journal: Economic Policy*, 14 (2022), 1–33.
- Ahomäki, Iiro, Visa Pitkänen, Aarni Soppi, and Leena Saastamoinen, “Impact of a physician-targeted letter on opioid prescribing,” *Journal of Health Economics*, 72 (2020), 1–22.
- Alexander, Diane, “How Do Doctors Respond to Incentives? Unintended Consequences of Paying Doctors to Reduce Costs,” *Journal of Political Economy*, 128 (2020), 4046–4096 (The University of Chicago Press).
- Alexander, Diane, and Janet Currie, “Are publicly insured children less likely to be admitted to hospital than the privately insured (and does it matter)?,” *Economics & Human Biology*, In Honor of Nobel Laureate Angus Deaton: Health Economics in Developed and Developing Countries, 25 (2017), 33–51.
- Alexander, Diane, and Molly Schnell, “The Impacts of Physician Payments on Patient Access, Use, and Health,” *American Economic Journal: Applied Economics*, (2024).
- Allen, Thomas, Eleonora Fichera, and Matt Sutton, “Can Payers Use Prices to Improve Quality? Evidence from English Hospitals,” *Health Economics*, 25 (2016), 56–70.
- Almond, Douglas, Joseph J. Doyle Jr., Amanda E. Kowalski, and Heidi Williams, “Estimating Marginal Returns to Medical Care: Evidence from At-risk Newborns,” *The Quarterly Journal of Economics*, 125 (2010), 591–634.
- Alpert, Abby, Sarah Dykstra, and Mireille Jacobson, “Hassle Costs versus Information: How Do Prescription Drug Monitoring Programs Reduce Opioid Prescribing?,” *American Economic Journal: Economic Policy*, 16 (2024), 87–123.

- Alsan, Marcella, Owen Garrick, and Grant Graziani, “Does Diversity Matter for Health? Experimental Evidence from Oakland,” *American Economic Review*, 109 (2019), 4071–4111.
- Angerer, Silvia, Christian Waibel, and Harald Stummer, “Discrimination in Health Care: A Field Experiment on the Impact of Patients’ Socioeconomic Status on Access to Care,” *American Journal of Health Economics*, 5 (2019), 407–427 (The University of Chicago Press).
- Arrow, Kenneth J., L. Kamran Bilir, and Alan Sorensen, “The Impact of Information Technology on the Diffusion of New Pharmaceuticals,” *American Economic Journal: Applied Economics*, 12 (2020), 1–39.
- Avdic, Daniel, Stephanie von Hinke, Bo Lagerqvist, Carol Propper, and Johan Vikström, “Do responses to news matter? Evidence from interventional cardiology,” *Journal of Health Economics*, 94 (2024).
- Bisgaier, Joanna, and Karin V. Rhodes, “Auditing Access to Specialty Care for Children with Public Insurance,” *New England Journal of Medicine*, 364 (2011), 2324–2333 (Massachusetts Medical Society).
- Bradford, W. David, and Andrew N. Kleit, “Impact of FDA Actions, DTCA, and Public Information on the Market for Pain Medication,” *Health Economics*, 24 (2015), 859–875.
- Brekke, Kurt R., Tor Helge Holmås, Karin Monstad, and Odd Rune Straume, “Socio-economic status and physicians’ treatment decisions,” *Health Economics*, 27 (2018), e77–e89.
- , “Competition and physician behaviour: Does the competitive environment affect the propensity to issue sickness certificates?,” *Journal of Health Economics*, 66 (2019), 117–135.
- Buchmueller, Thomas C., and Colleen Carey, “The Effect of Prescription Drug Monitoring Programs on Opioid Utilization in Medicare,” *American Economic Journal: Economic Policy*, 10 (2018), 77–112.
- Buchmueller, Thomas C., Colleen M. Carey, and Giacomo Meille, “How well do doctors know their patients? Evidence from a mandatory access prescription drug monitoring program,” *Health Economics*, 29 (2020), 957–974.
- Button, Patrick, Eva Dils, Benjamin Harrell, Luca Fumarco, and David Schwegman, “Gender Identity, Race, and Ethnicity Discrimination in Access to Mental Health Care: Preliminary Evidence from a Multi-Wave Audit Field Experiment,” Working Paper, 2020 (National Bureau of Economic Research).
- Cabral, Marika, and Marcus Dillender, “Gender Differences in Medical Evaluations: Evidence from Randomly Assigned Doctors,” *American Economic Review*, 114 (2024), 462–499.
- Carey, Colleen, Michael Daly, and Jing Li, “Nothing for Something: Marketing Cancer Drugs to Physicians Increases Prescribing Without Improving Mortality,” Working Paper, 2024 (National Bureau of Economic Research).
- Carey, Colleen, Ethan M. J. Lieber, and Sarah Miller, “Drug firms’ payments and physicians’ prescribing behavior in Medicare Part D,” *Journal of Public Economics*, 197 (2021), 104402.
- Chan, David C., “Teamwork and Moral Hazard: Evidence from the Emergency Department,” *Journal of Political Economy*, 124 (2016), 734–770 (University of Chicago).
- , “The efficiency of slacking off: Evidence from the emergency department,” *Econometrica*, 86 (2018), 997–1030.
- , “Influence and Information in Team Decisions: Evidence from Medical Residency,” *American Economic Journal: Economic Policy*, 13 (2021), 106–137.
- Chan, Jr, David C., and Yiqun Chen, “The Productivity of Professions: Evidence from the Emergency Department,” Working Paper, 2022 (National Bureau of Economic Research).
- Chandra, Amitabh, and Douglas O. Staiger, “Identifying Provider Prejudice in Healthcare,” Working Paper, 2010 (National Bureau of Economic Research).

- Chen, Alice, and Darius N. Lakdawalla, “Healing the poor: The influence of patient socioeconomic status on physician supply responses,” *Journal of Health Economics*, 64 (2019), 43–54.
- Chen, Yiqun, “Team-Specific Human Capital and Team Performance: Evidence from Doctors,” *American Economic Review*, 111 (2021), 3923–3962.
- Chernew, Michael, Zack Cooper, Eugene Larsen Hallock, and Fiona Scott Morton, “Physician agency, consumerism, and the consumption of lower-limb MRI scans,” *Journal of Health Economics*, 76 (2021), 102427.
- Chorniy, Anna, Janet Currie, and Lyudmyla Sonchak, “Exploding asthma and ADHD caseloads: The role of medicaid managed care,” *Journal of Health Economics*, 60 (2018), 1–15.
- Chu, Bryan, Ben Handel, Jonathan Kolstad, Jonas Knecht, Ulrike Malmendier, and Filip Matejka, “Cognitive Capacity, Fatigue and Decision Making: Evidence from the Practice of Medicine,” 2024 (UC Berkeley).
- Clemens, Jeffrey, and Joshua D. Gottlieb, “Do Physicians’ Financial Incentives Affect Medical Treatment and Patient Health?,” *American Economic Review*, 104 (2014), 1320–1349.
- Costa-Ramón, Ana María, Ana Rodríguez-González, Miquel Serra-Burriel, and Carlos Campillo-Artero, “It’s about time: Cesarean sections and neonatal health,” *Journal of Health Economics*, 59 (2018), 46–59.
- Coudin, Elise, Anne Pla, and Anne-Laure Samson, “GP responses to price regulation: evidence from a French nationwide reform,” *Health Economics*, 24 (2015), 1118–1130.
- Coussens, Stephen, “Behaving Discretely: Heuristic Thinking in the Emergency Department,” 2022 (Columbia University).
- Cuddy, Emily, and Janet Currie, “Treatment of mental illness in American adolescents varies widely within and across areas,” *Proceedings of the National Academy of Sciences*, 117 (2020), 24039–24046 (Proceedings of the National Academy of Sciences).
- , “Rules vs. Discretion: Treatment of Mental Illness in U.S. Adolescents,” *Journal of Political Economy*, (2024).
- Currie, Janet M., and W. Bentley MacLeod. 2017. “Diagnosis and Unnecessary Procedure Use: Evidence from C-section.” *Journal of Labor Economics* 35 (1): 1–42.
- Currie, Janet M, and W Bentley MacLeod, “Understanding doctor decision making: The case of depression treatment,” *Econometrica : journal of the Econometric Society*, 88 (2020), 847–878.
- Currie, Janet, and Jonathan Zhang, “Doing More with Less: Predicting Primary Care Provider Effectiveness,” *The Review of Economics and Statistics*, (2023), 1–45.
- Dahlstrand, Amanda, “Defying Distance? The Provision of Services in the Digital Age,” (2024).
- Dickstein, “Physician vs. Patient Incentives in Prescription Drug Choice,” 2017.
- Ding, Yu, and Chenyuan Liu, “Alternative payment models and physician treatment decisions: Evidence from lower back pain,” *Journal of Health Economics*, 80 (2021).
- Doctor, Jason N., Andy Nguyen, Roneet Lev, Jonathan Lucas, Tara Knight, Henu Zhao, and Michael Menchine, “Opioid prescribing decreases after learning of a patient’s fatal overdose,” *Science*, 361 (2018), 588–590 (American Association for the Advancement of Science).
- Doyle, Joseph J., Steven M. Ewer, and Todd H. Wagner, “Returns to physician human capital: Evidence from patients randomized to physician teams,” *Journal of Health Economics*, 29 (2010), 866–882.
- Doyle, Jr., Joseph J., “Physician Characteristics and Patient Survival: Evidence from Physician Availability,” Working Paper, 2020 (National Bureau of Economic Research).
- Dubois, Pierre, and Tuba Tunçel, “Identifying the effects of scientific information and recommendations on physicians’ prescribing behavior,” *Journal of Health Economics*, 78 (2021), 102461.

- Eli, Shari, Trevon D. Logan, and Boriana Miloucheva, “Physician Bias and Racial Disparities in Health: Evidence from Veterans’ Pensions,” Working Paper, 2019 (National Bureau of Economic Research).
- Ellyson, Alice M., Jevay Grooms, and Alberto Ortega, “Flipping the script: The effects of opioid prescription monitoring on specialty-specific provider behavior,” *Health Economics*, 31 (2022), 297–341.
- Epstein, Andrew J., Sean Nicholson, and David A. Asch, “The Production of and Market for New Physicians’ Skill,” *American Journal of Health Economics*, 2 (2016), 41–65 (The University of Chicago Press).
- Facchini, Gabriel, “Forgetting-by-not-doing: The case of surgeons and cesarean sections,” *Health Economics*, 31 (2022), 481–495.
- Fortin, Bernard, Nicolas Jacquemet, and Bruce Shearer, “Labour supply, service intensity, and contracts: Theory and evidence on physicians,” *Journal of Applied Econometrics*, 36 (2021), 686–702.
- Frakes, Michael, “The Impact of Medical Liability Standards on Regional Variations in Physician Behavior: Evidence from the Adoption of National-Standard Rules,” *American Economic Review*, 103 (2013), 257–276.
- Frakes, Michael D., and Jonathan Gruber, “Racial Concordance and the Quality of Medical Care: Evidence from the Military,” Working Paper, 2022 (National Bureau of Economic Research).
- Freedman, Seth, Ezra Golberstein, Tsan-Yao Huang, David J. Satin, and Laura Barrie Smith, “Docs with their eyes on the clock? The effect of time pressures on primary care productivity,” *Journal of Health Economics*, 77 (2021), 102442.
- Geiger, Caroline K., Mark A. Clapp, and Jessica L. Cohen, “Association of Prenatal Care Services, Maternal Morbidity, and Perinatal Mortality With the Advanced Maternal Age Cutoff of 35 Years,” *JAMA Health Forum*, 2 (2021), e214044.
- Goetz, Daniel, “Telemedicine competition, pricing, and technology adoption: Evidence from talk therapists,” *International Journal of Industrial Organization*, 89 (2023), 102956.
- Gowrisankaran, Gautam, Keith Joiner, and Pierre Thomas Léger, “Physician Practice Style and Healthcare Costs: Evidence from Emergency Departments,” *Management Science*, (2022) (INFORMS).
- Goyal, Monika K., Nathan Kuppermann, Sean D. Cleary, Stephen J. Teach, and James M. Chamberlain, “Racial Disparities in Pain Management of Children With Appendicitis in Emergency Departments,” *JAMA Pediatrics*, 169 (2015), 996–1002.
- Greenwood, Brad N., Seth Carnahan, and Laura Huang, “Patient–physician gender concordance and increased mortality among female heart attack patients,” *Proceedings of the National Academy of Sciences*, 115 (2018), 8569–8574 (Proceedings of the National Academy of Sciences).
- Greenwood, Brad N., Rachel R. Hardeman, Laura Huang, and Aaron Sojourner, “Physician–patient racial concordance and disparities in birthing mortality for newborns,” *Proceedings of the National Academy of Sciences*, 117 (2020), 21194–21200 (Proceedings of the National Academy of Sciences).
- Gruber, Jonathan, Thomas P. Hoe, and George Stoye, “Saving Lives by Tying Hands: The Unexpected Effects of Constraining Health Care Providers,” *The Review of Economics and Statistics*, (2021), 1–45.
- Gupta, Atul, “Impacts of Performance Pay for Hospitals: The Readmissions Reduction Program,” *American Economic Review*, 111 (2021), 1241–1283.
- Hill, Andrew J., Daniel B. Jones, and Lindsey Woodworth, “Physician-patient race-match reduces patient mortality,” *Journal of Health Economics*, 92 (2023).
- Hoffman, Kelly M., Sophie Trawalter, Jordan R. Axt, and M. Norman Oliver, “Racial bias in pain assessment and treatment recommendations, and false beliefs about biological differences between blacks and whites,” *Proceedings of the National Academy of Sciences*, 113 (2016), 4296–4301 (Proceedings of the National Academy of Sciences).

- Horwitz, Jill, Austin Nichols, Carrie H. Colla, and David M. Cutler, “Technology Regulation Reconsidered: The Effects of Certificate of Need Policies on the Quantity and Quality of Diagnostic Imaging,” Working Paper Series, Working Paper, 2024 (National Bureau of Economic Research).
- Howard, David H., Guy David, and Jason Hockenberry, “Selective Hearing: Physician-Ownership and Physicians’ Response to New Evidence,” *Journal of Economics & Management Strategy*, 26 (2017), 152–168.
- Howard, David H., and Jason Hockenberry, “Physician age and the abandonment of episiotomy,” *Health Services Research*, 54 (2019), 650–657.
- Howard, David H., and Ian McCarthy, “Deterrence effects of antifraud and abuse enforcement in health care,” *Journal of Health Economics*, 75 (2021), 102405.
- Johnson, Erin M., and M. Marit ReHAVI, “Physicians Treating Physicians: Information and Incentives in Childbirth,” *American Economic Journal: Economic Policy*, 8 (2016), 115–141.
- Johnson, Erin, M. Marit ReHAVI, Jr Chan David C., and Daniela Carusi, “A Doctor Will See You Now: Physician-Patient Relationships and Clinical Decisions,” Working Paper, 2016 (National Bureau of Economic Research).
- Kolstad, Jonathan T., “Information and Quality When Motivation Is Intrinsic: Evidence from Surgeon Report Cards,” *American Economic Review*, 103 (2013), 2875–2910.
- Kowalski, Amanda E, “Behaviour within a Clinical Trial and Implications for Mammography Guidelines,” *The Review of Economic Studies*, 90 (2023), 432–462.
- Linder, Jeffrey A., Jason N. Doctor, Mark W. Friedberg, Harry Reyes Nieva, Caroline Birks, Daniella Meeker, and Craig R. Fox, “Time of Day and the Decision to Prescribe Antibiotics,” *JAMA Internal Medicine*, 174 (2014), 2029–2031.
- MacLeod, W Bentley. 2022. *Advanced Microeconomics for Contract, Institutional and Organizational Economics*. Boston, MA: MIT Press.
- McCullough, Jeffrey S., Michelle Casey, Ira Moscovice, and Shailendra Prasad, “The Effect Of Health Information Technology On Quality In U.S. Hospitals,” *Health Affairs*, 29 (2010), 647–654 (Health Affairs).
- McDevitt, Ryan C., and James W. Roberts, “Market structure and gender disparity in health care: preferences, competition, and quality of care,” *The RAND Journal of Economics*, 45 (2014), 116–139.
- McKibbin, Rebecca, “The effect of RCTs on drug demand: Evidence from off-label cancer drugs,” *Journal of Health Economics*, 90 (2023), 102779.
- Miller, Amalia R., and Catherine E. Tucker, “Can Health Care Information Technology Save Babies?,” *Journal of Political Economy*, 119 (2011), 289–324 (The University of Chicago Press).
- Molitor, David, “The Evolution of Physician Practice Styles: Evidence from Cardiologist Migration,” *American Economic Journal: Economic Policy*, 10 (2018), 326–356.
- Mullainathan, Sendhil, and Ziad Obermeyer, “Diagnosing Physician Error: A Machine Learning Approach to Low-Value Health Care,” *The Quarterly Journal of Economics*, 137 (2022), 679–727.
- Neprash, Hannah T., John F. Mulcahy, Dori A. Cross, Joseph E. Gaugler, Ezra Golberstein, and Ishani Ganguli, “Association of Primary Care Visit Length With Potentially Inappropriate Prescribing,” *JAMA Health Forum*, 4 (2023), e230052.
- Neumark, David, and Bogdan Savych, “Effects of Opioid-Related Policies on Opioid Utilization, Nature of Medical Care, and Duration of Disability,” *American Journal of Health Economics*, 9 (2023), 331–373 (The University of Chicago Press).
- Newham, Melissa, and Marica Valente, “The cost of influence: How gifts to physicians shape prescriptions and drug costs,” *Journal of Health Economics*, 95 (2024), 102887.
- Neyman, Jerzy, and Egon Sharpe Pearson. 1933. “IX. On the Problem of the Most Efficient Tests of Statistical Hypotheses.” *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character* 231 (694-706): 289–

- Obermeyer, Ziad, Brian Powers, Christine Vogeli, and Sendhil Mullainathan, “Dissecting racial bias in an algorithm used to manage the health of populations,” *Science*, 366 (2019), 447–453 (American Association for the Advancement of Science).
- Olenski, Andrew R., André Zimerman, Stephen Coussens, and Anupam B. Jena, “Behavioral Heuristics in Coronary-Artery Bypass Graft Surgery,” *New England Journal of Medicine*, 382 (2020), 778–779 (Massachusetts Medical Society).
- Olson, Mary K., and Nina Yin, “New clinical information and physician prescribing: How do pediatric labeling changes affect prescribing to children?,” *Health Economics*, 30 (2021), 144–164.
- van Parys, Jessica, “Variation in Physician Practice Styles within and across Emergency Departments,” *PLOS ONE*, 11 (2016) (Public Library of Science).
- van Parys, Jessica, and Zach Y. Brown, “Broadband Internet Access and Health Outcomes: Patient and Provider Responses in Medicare,” Working Paper Series, Working Paper, 2023 (National Bureau of Economic Research).
- Persson, Emil, Kinga Barrafreem, Andreas Meunier, and Gustav Tinghög, “The effect of decision fatigue on surgeons’ clinical decision making,” *Health Economics*, 28 (2019), 1194–1203.
- Persson, Petra, Xinyao Qiu, and Maya Rossin-Slater, “Family Spillover Effects of Marginal Diagnoses: The Case of ADHD,” Working Paper, 2021 (National Bureau of Economic Research).
- Sabin, Janice A., and Anthony G. Greenwald, “The Influence of Implicit Bias on Treatment Recommendations for 4 Common Pediatric Conditions: Pain, Urinary Tract Infection, Attention Deficit Hyperactivity Disorder, and Asthma,” *American Journal of Public Health*, 102 (2012), 988–995 (American Public Health Association).
- Sacarny, Adam, Michael L. Barnett, Jackson Le, Frank Tetkoski, David Yokum, and Shantanu Agrawal, “Effect of Peer Comparison Letters for High-Volume Primary Care Prescribers of Quetiapine in Older and Disabled Adults: A Randomized Clinical Trial,” *JAMA Psychiatry*, 75 (2018), 1003–1011.
- Sacarny, Adam, David Yokum, Amy Finkelstein, and Shantanu Agrawal, “Medicare Letters To Curb Overprescribing Of Controlled Substances Had No Detectable Effect On Providers,” *Health Affairs*, 35 (2016), 471–479 (Health Affairs).
- Sacks, Daniel W., Alex Hollingsworth, Thuy Nguyen, and Kosali Simon, “Can policy affect initiation of addictive substance use? Evidence from opioid prescribing,” *Journal of Health Economics*, 76 (2021), 102397.
- Savage, Leonard J. 1972 (first published 1954). *The Foundations of Statistics*. New York, N.Y.: Dover Publications.
- Schnell, Molly, and Janet Currie, “Addressing the Opioid Epidemic: Is There a Role for Physician Education?,” *American Journal of Health Economics*, 4 (2018), 383–410 (The University of Chicago Press).
- Shapiro, Bradley T., “Informational Shocks, Off-Label Prescribing, and the Effects of Physician Detailing,” *Management Science*, 64 (2018), 5925–5945 (INFORMS).
- Shurtz, Ity, Alon Eizenberg, Adi Alkalay, and Amnon Lahad, “Physician workload and treatment choice: the case of primary care,” *The RAND Journal of Economics*, 53 (2022), 763–791.
- Silver, David, “Haste or Waste? Peer Pressure and Productivity in the Emergency Department,” *The Review of Economic Studies*, 88 (2021), 1385–1417.
- Simeonova, Emilia, Niels Skipper, and Peter Rønø Thingholm, “Physician Health Management Skills and Patient Outcomes,” *Journal of Human Resources*, 59 (2024), 777–809 (University of Wisconsin Press).
- Singh, Manasvini, and Atheendar Venkataramani, “Rationing by Race,” Working Paper, 2022 (National Bureau of Economic Research).
- Tai-Seale, Ming, and Thomas McGuire, “Time is up: increasing shadow price of time in primary-care office visits,” *Health Economics*, 21 (2012), 457–476.

Wallis, Christopher J. D., Angela Jerath, Natalie Coburn, Zachary Klaassen, Amy N. Luckenbaugh, Diana E. Magee, Amanda E. Hird, Kathleen Armstrong, Bheeshma Ravi, Nestor F. Esnaola, Jonathan C. A. Guzman, Barbara Bass, Allan S. Detsky, and Raj Satkunasivam, "Association of Surgeon-Patient Sex Concordance With Postoperative Outcomes," *JAMA Surgery*, 157 (2022), 146–156.

Wilding, Anna, Luke Munford, Bruce Guthrie, Evangelos Kontopantelis, and Matt Sutton, "Family doctor responses to changes in target stringency under financial incentives," *Journal of Health Economics*, 85 (2022).

Wu, Bingxiao, and Guy David, "Information, relative skill, and technology abandonment," *Journal of Health Economics*, 83 (2022), 102596.

Zeltzer, Dan, Liran Einav, Joseph Rashba, and Ran D Balicer, "The Impact of Increased Access to Telemedicine," *Journal of the European Economic Association*, 22 (2024), 712–750.

Zeltzer, Dan, Liran Einav, Joseph Rashba, Yehezkel Waisman, Motti Haimi, and Ran D. Balicer, "Adoption and utilization of device-assisted telemedicine," *Journal of Health Economics*, 90 (2023), 102780.