

NBER WORKING PAPER SERIES

NEGATIVE TESTS AND THE EFFICIENCY OF MEDICAL CARE:  
WHAT DETERMINES HETEROGENEITY IN IMAGING BEHAVIOR?

Jason Abaluck  
Leila Agha  
Christopher Kabrhel  
Ali Raja  
Arjun Venkatesh

Working Paper 19956  
<http://www.nber.org/papers/w19956>

NATIONAL BUREAU OF ECONOMIC RESEARCH  
1050 Massachusetts Avenue  
Cambridge, MA 02138  
March 2014

Thanks to Brian Abaluck, Joe Altonji, Joshua Aronson, Judy Chevalier, Michael Dickstein, David Dranove, Amy Finkelstein, Howard Forman, Jonathan Gruber, Nathan Hendren, Vivian Ho, Mitch Hoffman, Lisa Kahn, Jon Kolstad, Amanda Kowalski, Danielle Li, Costas Meghir, David Molitor, Blair Parry, Michael Powell, Constana Esteves-Sorenson, Ashley Swanson, Bob Town, and Heidi Williams as well as seminar participants at AHEC 2012, AEA meeting 2013, Boston University, Cornell, HEC Montreal, IHEA 2013, the National Bureau of Economic Research, NIA Dartmouth research meeting, the National Tax Association annual meeting, Northwestern, Stanford, University of Houston, and Yale. Funding for this work was provided by NIA Grant Number T32-AG0000186 to the NBER. The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

At least one co-author has disclosed a financial relationship of potential relevance for this research. Further information is available online at <http://www.nber.org/papers/w19956.ack>

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 They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2014 by Jason Abaluck, Leila Agha, Christopher Kabrhel, Ali Raja, and Arjun Venkatesh. 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.

Negative Tests and the Efficiency of Medical Care: What Determines Heterogeneity in Imaging Behavior?

Jason Abaluck, Leila Agha, Christopher Kabrhel, Ali Raja, and Arjun Venkatesh

NBER Working Paper No. 19956

March 2014, Revised August 2015

JEL No. I0,I12

**ABSTRACT**

We develop a model of the efficiency of medical testing based on rates of negative CT scans for pulmonary embolism. The model is estimated using a 20% sample of Medicare claims from 2000- 2009. We document enormous across-doctor heterogeneity in testing decisions conditional on patient risk and show it explains the negative relationship between physicians' testing frequencies and test yields. Physicians in high spending regions test more low-risk patients. Under calibration assumptions, 84% of doctors test even when costs exceed expected benefits. Furthermore, doctors do not apply observables to target testing to the highest risk patients, substantially reducing simulated test yields.

Jason Abaluck  
Yale School of Management  
Box 208200  
New Haven, CT 06520-8200  
and NBER  
jason.abaluck@yale.edu

Ali Raja  
Department of Emergency Medicine  
Brigham and Women's Hospital  
75 Francis St  
Boston, MA 02115  
asraja@partners.org

Leila Agha  
School of Management  
Boston University  
595 Commonwealth Avenue  
Boston, MA 02215  
and NBER  
lagha@bu.edu

Arjun Venkatesh  
Yale School of Medicine  
Department of Internal Medicine  
arjun.venkatesh@yale.edu

Christopher Kabrhel  
Department of Emergency Medicine  
Massachusetts General Hospital  
55 Fruit St, Clinics Building 115  
Boston, MA 02114  
ckabrhel@partners.org

# 1 Introduction

Many have argued that current medical practice involves large amounts of wasteful spending, with little cross-sectional correlation between regional health spending and health outcomes (Wennberg et al. 1996). But determining the best approach to lower costs and maintain quality depends critically on the nature of the inefficiency: is the problem that physicians are spending to the “flat of the curve” where marginal returns to treatment are low, or are physicians treating the wrong patients (Garber and Skinner 2008)? And more fundamentally, does heterogeneity in medical spending imply care patterns are inefficient or is this heterogeneity fully explicable by unobserved differences in returns to treatment across patient populations?

In this paper, we develop an econometric framework for evaluating whether the wide variation in the use of diagnostic testing across doctors is due to heterogeneity in patients’ benefits from testing or heterogeneity in physicians’ practice styles (i.e. differences in behavior when treating identical patients). The model also identifies whether physicians are weighting patient observable risk factors to maximize the incidence of positive tests. Our model builds on classical econometric selection models originally developed by Heckman (1979) and refined by Chandra and Staiger (2011). We apply a novel instrumental variable identification strategy to study the behavior of physicians who make repeated decisions selecting which patients will receive a diagnostic test. A similar modeling approach and identification strategy could be applied to any setting meeting the criteria of a standard selection model (where outcomes are only observed among the treated) *and* where we observe repeated selection decisions by the same decisionmaker, including decisions of loan officers, judges, hiring directors, and many others.

We apply our model to analyze CT scans that test for pulmonary embolism (PE). Estimation of the model requires that we can observe test outcomes among patients selected for testing, as well as the structural assumption that doctors will order a CT scan to test for PE if the patient’s *ex ante* risk of PE exceeds a doctor-specific testing threshold. This threshold is our patient invariant measure of physician practice style and we seek to recover it for each doctor in our sample. Identifying differences in physicians’ practice styles separately from patient heterogeneity typically requires either random assignment of patients to physicians or estimates of potentially heterogeneous causal effects of medical treatment for each patient. Prior research, including Chandra and Staiger (2011) and Currie and MacLeod (2013), has argued that reliable estimates of causal treatment effects can be obtained using detailed chart data to control for all patient characteristics observable to doctors, but such data is typically only available in limited samples. This stumbling block makes it difficult to investigate both the extent and the determinants of healthcare overuse or misuse.

A key insight of this paper is that the *ex post* value of a diagnostic test, in this case chest CT scans, is partially observable in insurance claims records based on whether the test results in the relevant diagnosis. A doctor who performs many negative CT scans, which have little *ex post* value for improving patient health, is likely to have a low testing threshold. Physicians with identical testing thresholds may have different rates of negative tests if their patient population varies in the *ex ante* risk of PE. Our model accounts for heterogeneity in patient PE risk and shows how to recover physicians’ testing thresholds. Using these estimated testing thresholds, we investigate the

role of medical training, malpractice environment, hospital characteristics and regional factors in shaping practice styles.

The model also allows investigation of whether doctors are misweighting observable patient risk factors in selecting which patients to test for PE. By comparing how observable risk factors predict physicians' testing decisions to how those same variables predict rates of positive tests amongst tested patients, we can identify whether physicians are targeting CT scans to the patients with the highest risk of PE based on demographics and comorbid conditions.

Previous research has identified important differences in practice style and skill across physicians. Chandra and Staiger (2011) conclude that overuse of care explains a large amount of variation in treatment for heart attacks across hospitals. Currie and MacLeod (2013) uncover substantial heterogeneity in diagnostic skill across obstetricians. Finkelstein et al. (2014) find that roughly half of the variation in medical spending across regions is driven by provider behavior (rather than patient preferences or health risks), and Molitor (2012) reports that environmental factors explain much of the variation in physician's rates of cardiac catheterization.

We extend this prior literature by not only estimating heterogeneity in physician practice styles, but also explicitly demonstrating that differences in practice style explain why physicians who use more medical resources have lower average medical returns to utilization. We then estimate the resulting welfare loss from the measured variation in practice styles. Further, to our knowledge, this is the first paper to test for physicians' systematic underweighting and overweighting of patient risk factors and to assess how failure to target medical resources to the patients with the highest expected returns may impact realized health benefits and total welfare. It introduces a new dimension to the existing literature on practice styles by highlighting an additional mechanism by which physician decisions may influence returns to testing or treatment.

PE is the third most common cause of death from cardiovascular disease, behind heart attack and stroke (Goldhaber and Bounameaux 2012), and CT scans are the primary tool for diagnosis of PE. Yet given the financial costs and medical risks of testing, PE CT scans are commonly thought to be overused in emergency care. The American College of Radiology targeted PE CT as a key part of the *Choosing Wisely* campaign aimed to reduce overuse of medical services. Despite the concern about overuse, the Office of the Surgeon General (2008) estimates that approximately half of PE cases are undiagnosed, based on analysis of autopsy reports. The simultaneous concern in the medical community about overuse and missed diagnoses raises the question of whether diagnostic testing for PE is currently being targeted to maximize PE detection.

We analyze 1.9 million emergency department visits drawn from a 20% sample of Medicare claims data, 2000-2009. We observe whether each patient is tested with a chest CT, and whether this test leads to diagnosis of PE. We present reduced form evidence of a sharply negative relationship between physician testing rates and test yields: those physicians who test most have the lowest rate of positive tests. We apply a structural model to show that this pattern is explained by enormous heterogeneity in doctors' testing thresholds; doctors who test more move further down the net benefit curve and test patients who are less likely to test positive. Less experienced doctors and doctors in higher spending regions tend to have lower risk thresholds at which they deem CT imaging worthwhile.

Further, physicians fail to target the test to the highest risk patients. Recognized risk factors, some of which are included in popular PE risk scores, continue to receive too little weight in physicians' testing decisions. Black patients are tested less often than other patients despite their higher risk of PE. Finally, physicians overttest patients who have been previously diagnosed with one of several conditions which have similar clinical symptoms to PE: rather than infer the patient is having a recurrent episode of their existing condition, the physician may order a PE CT despite the low predicted risk.

Applying calibration assumptions about the cost of testing, the benefits of treating PE and the likelihood of false positives, we compare our estimated distribution of physician testing thresholds to the calibrated socially optimal threshold. This comparison tells us the degree of allocative inefficiency: whether doctors are overttesting or undertesting from a social standpoint.<sup>1</sup> Under our preferred calibration assumptions, 84% of doctors are overttesting in the sense that for their marginal patients, the costs of testing exceed the benefits. In a simulation where no doctors overttested, the total social benefits from chest CTs would increase by 60% and the number of chest CT scans would fall by 50%. The calibration also allows us to assess the degree of productive inefficiency from physician misweighting of patient risk factors. Weighting observable comorbidities to maximize test yields would increase the net benefits of testing by more than 300%, primarily by leading to additional testing and appropriate diagnosis of affected patients.

The paper is organized as follows. Section 2 provides some background on chest CT scans for PE. Section 3 describes the data and uses reduced form evidence to motivate the structural model. Section 4 lays out our structural model of testing behavior and describes our estimation strategy. Section 5 reports results from estimating our structural model. Section 6 probes the robustness of these results to alternative modeling approaches that relax or vary key identifying assumptions. Section 7 conducts simulations to uncover the welfare implications of our findings, and Section 8 concludes.

## 2 Background on PE CTs

We study testing behavior in the context of chest CT scans performed in the emergency department (ED) to detect PE. A PE occurs when a substance, most commonly a blood clot that originates in a vein, travels through the bloodstream into an artery of the lung and blocks blood flow through the lung. It is a serious and relatively common condition, with an estimated 350,000 diagnosed cases of PE per year in the United States (Office of the Surgeon General 2008). Left untreated, the mortality rate from a PE depends on the severity and has been estimated to be 2.5% within three months for a small PE (Lessler et al. 2010), with most of the risk concentrated within the first hours after onset of symptoms (Rahimtoola and Bergin 2005). Accurate diagnosis of PE is necessary for appropriate follow-up treatment; even high risk patients are unlikely to be treated presumptively.

---

<sup>1</sup>We are defining allocative and productive inefficiency from the standpoint of production functions with spending on CT scans as the input and patient health status as the output. With "overttesting", one is testing to the flat of the curve, where the marginal health returns to additional spending are very low; this is an allocative inefficiency. With misweighting, a higher production function—producing greater health gains for a given level of spending—is achievable if doctors would correctly weight observable risk factors; misweighting thus creates productive inefficiency.

CT scans to test for PE have a number of attractive features for our purposes: they are a frequently performed test; they introduce significant health risks and financial costs; a positive test is almost always followed up with immediate treatment, observable in Medicare claims records; and a negative test provides little information to the physician about alternative diagnoses or potential treatments. We discuss each of these features in more detail in Appendix B, explaining how the clinical context supports our modeling assumptions.

The symptoms of PE are both common and nonspecific: shortness of breath, chest pain, or bloody cough. Hence, there is a broad population of patients who may be considered for a PE evaluation. Practice guidelines recommend that physicians also consider several additional risk factors before determining whether to pursue a workup for PE.<sup>2</sup> Because PE is an acute event with a sudden onset, the workup must be completed urgently and knowing the results of previous CT scans is not a critical part of the evaluation of PE.

Many argue that PE CT scans are widely overused (Coco and O’Gurek 2012, Mamlouk et al. 2010 and Costantino et al. 2008). Recent estimates by Venkatesh et al. (2012) suggest that one third of CT scans in a sample of 11 US emergency departments would have been avoidable if physicians had followed National Quality Forum guidelines on CT usage. The nonspecific symptoms of PE and significant mortality risk likely both contribute to overuse, particularly in the emergency care setting.

A CT angiogram is the standard diagnostic tool for PE. The average allowed charge in the Medicare data is around \$320 per PE CT when the bill is not covered by a capitation payment. Payment goes to the radiologist for interpreting the scan and to the hospital for the technician and capital investment required to perform the scan. The emergency department doctor responsible for ordering the test has, at most, a diffuse incentive to ensure the hospital’s financial health and reduce his malpractice risk, but he receives no direct payments from Medicare or the hospital for ordering a scan.

PE CT scans also come with small but important medical risks. The most significant risk arises from false positive CT scans which lead to additional unnecessary treatment with anticoagulants, incurring financial costs and creating significant risk of bleeding. In addition, there is an estimated 0.02% chance of a severe reaction to the contrast, which then carries a 10.5% risk of death (Lessler et al. 2010), although this cost is small relative to the billed financial costs of a CT scan. Finally radiation exposure may increase downstream cancer risk, although the additional lifetime cancer risk is minimal for the elderly Medicare population in this study.

The key simplifying assumption we make to evaluate the net benefits of testing is that a negative test has no value. This assumption is not true in general for all tests: a negative test may rule out one treatment thus justifying treatment for an alternative, or a negative test might prevent an otherwise costly treatment. However, in our setting—CT scans for PE—a positive test is followed by an inpatient admission and treatment with blood thinners while a negative test does not suggest any further interventions or testing for related problems.

---

<sup>2</sup>Popular practice guidelines use the following factors to calculate a risk score: age, elevated heart rate, recent immobilization or surgery, history of deep vein thrombosis or PE, recent treatment for cancer, coughing up blood, lower limb pain or swelling, and chances of an alternative diagnosis.

### 3 Data

We combine data from five sources: Medicare claims records, the American Hospital Association annual survey, the American Medical Association Masterfile, the Medicare Physician Identification and the Eligibility Registry, and the Avraham Database of State Tort Law Reforms. Using a 20% sample of Medicare Part B claims from 2000 through 2009, we identify patients evaluated in an emergency department and observe whether they were tested for PE, as well as whether any such test succeeded in detecting PE.

#### 3.1 Medicare claims data

We begin by identifying all patients evaluated in the emergency department (ED), using physician-submitted Medicare Part B claims for evaluation and management.<sup>3</sup> The physician submitting this claim for evaluation and management is responsible for the patient’s emergency care; it is his decision whether or not to order testing for PE. Using physician identifiers, we track the behavior of all doctors who routinely evaluate Medicare patients in the ED.

We identify which ED patients are tested for a PE using bills submitted by radiologists for the interpretation of chest CTs with contrast, when the CT is performed within 1 day of the ED visit.<sup>4</sup> We restrict our sample to physicians who order at least seven CT scans between 2000-2009, since very low-volume doctors provide too little information to accurately estimate physicians’ testing thresholds.<sup>5</sup>

While diagnosis of PE is the most common purpose of a chest CT performed in the emergency care setting, there are a small handful of other, less common indications, including pleural effusion, chest and lung cancers, traumas, and aortic dissection. For this reason, we exclude patients from the sample who are coded with a diagnosis related to trauma, pleural effusion, chest or lung cancer, or patients with a history of aortic aneurysm, aortic dissection, or other arterial dissection. We also exclude patients with a history of renal failure, since these patients are likely ineligible for a CT scan with contrast, due to risks of the contrast agent. These sample restrictions are designed to limit the sample to patients who may be eligible for a chest CT scan and for whom the scan is highly likely to have been ordered to detect PE; these assumptions are discussed in more detail in Appendix C.

Once we have identified relevant CT scans in billing data, we then need to code the test outcome, i.e. whether or not the scan detected a PE. Patients with acute PE are typically admitted to the hospital for monitoring and to begin a course of blood thinners or place a venous filter to reduce clotting risk. From the sample of patients tested in the emergency department with a chest CT, we identify positive tests on the basis of Medicare Part A hospital claims that include a diagnosis code for PE among any of the diagnoses associated with the hospital stay.

We have validated this approach to identifying positive tests by using cross-referenced patient chart and hospital billing data from two large academic medical centers. The evidence from these

---

<sup>3</sup>In particular, we identify patients based on CPT codes for emergency department evaluation and management: 99281, 99282, 99283, 99284, 99285, and place of service 23 (i.e. hospital emergency department).

<sup>4</sup>We begin by identifying all bills for chest CTs on the basis of CPT codes 71260, 71270, and 71275.

<sup>5</sup>Note this restriction suggests the sample may be selected towards physicians who have lower testing thresholds. Unfortunately, some such restriction is required for the model to be estimable and to ensure that all in-sample physicians have routine access to a CT scanner.



centers suggests that we are unlikely to understate physicians’ testing thresholds due to undercounting of positive test results. More detail on this data validation exercise is presented in Appendix D.

In addition to measuring whether patients were tested and the testing outcome, we also document a number of characteristics that allow us to predict the patient’s propensity to be diagnosed with a PE, including age, race, sex, and medical comorbidities. We code comorbidities from both Medicare’s Chronic Condition Warehouse and from the Elixhauser et al. (1998) definitions; while these sets of conditions overlap, the Chronic Condition Warehouse utilizes outpatient claims to code comorbidities whereas the Elixhauser comorbidities are based only on inpatient medical history, so they typically encode different levels of disease severity. We augment these standard sets of medical comorbidities to include several measures that are specific to PE risk: whether the patient was previously admitted to the hospital with a diagnosis of PE, thoracic aortic dissection, abdominal aortic dissection, or deep vein thrombosis, and any cause admission to the hospital or surgical hospital admission within 7 days or 30 days.

### 3.2 Physician, hospital, and regional data

After using the Medicare claims data to estimate the testing threshold applied by each doctor, we explore predictors of physicians’ practice styles by linking testing thresholds to physician, hospital, and regional characteristics.

We draw physician data from two sources, the Medicare Physician Identification and Eligibility Registry (MPIER) and the American Medical Association Masterfile (AMA data). The MPIER and AMA both identify the medical school and graduation year for each physician, which we have linked to the US News & World Report medical school rankings. We bin schools according to whether they are typically ranked in the top 50 for either primary care or research rankings.

Hospital characteristics are drawn from the American Hospital Association annual survey. We use these data to observe whether the physician typically practices at a for profit hospital or an academic hospital, defined as a hospital with a board certified residency program.

Using provider zip codes, we identify the hospital referral region (HRR) in which each patient is treated. HRRs are regional health care markets defined by the Dartmouth Atlas to reflect areas within which patients commonly travel to receive tertiary care. There are 306 HRRs in total. Using data from the Dartmouth Atlas, we link each HRR to the average spending per Medicare beneficiary to capture a broad measure of regional care intensity.

Finally, data on state malpractice environment is from Avraham (2011) Database of State Tort Law Reforms. Following prior work by Currie and MacLeod (2006) and Avraham et al. (2012), we focus on two key measures of malpractice law: whether a state has enacted malpractice damage caps on award amounts, and joint and several liability reform.

### 3.3 Summary statistics

There are 1.9 million emergency department visit evaluations in our dataset, after making the sample exclusions noted above. Of these patients evaluated in the ED, 3.8% of them are tested with a chest



CT scan with contrast. Amongst tested patients, 6.9% of them receive a positive test, i.e. are admitted to the hospital within 24 hours with a diagnosis of PE.

Summary statistics are reported in Table 1, with results reported separately for patients who do not receive a CT scan (column A), patients who receive a negative test (column B), and patients with a positive test (column C). We observe the testing behavior of over 6600 physicians, with an average of 284 ED patients per physician.

Patient demographics are similar across the untested and tested patient groups. The average age is 78 years in the untested sample and slightly lower (77 years) in the sample of patients with negative or positive tests. Patients who test negative are more than twice as likely to have a history of PE as untested patients; patients with positive tests are five times more likely to have a history of PE than untested patients.

We note a few modest differences in physician background and practice environment across patient groups. Patients with negative tests are evaluated by doctors with five months less experience on average than patients with positive tests, and were treated in regions with 1% higher Medicare spending per beneficiary, compared to patients with positive tests. Among tested patients, those with positive tests were 1 percentage point more likely to have been evaluated by a doctor trained at a top tier medical school. In the structural model, we will decompose to what extent these differences may be driven by differential sorting of high risk patients and to what extent they reflect differences in physician practice styles.

### 3.4 Reduced form evidence of heterogeneity in doctor testing behavior

Before describing our model, we consider reduced form evidence of heterogeneity in doctors' testing behavior. We first divide doctors in our sample into 10 deciles according to the average fraction of patients tested. We observe average testing rates that range from 1.7% of ED patients in the lowest physician decile to 8.2% of ED patients in the highest physician decile. We want to know whether this variation reflects differences in doctor behavior for patients with similar PE risk, or differences in patient PE risk for physicians with similar testing practices.

We can separate these hypotheses by comparing rates of positive tests conditional on testing behavior. If doctors who test more do so because their patients are at higher risk of PE, we should expect that doctors with higher testing rates will also have a higher fraction of positive tests among tested patients.<sup>6</sup> Alternatively, if doctors who test more do so because they are the type that tests more for any given level of patient risk, then we expect to find that physicians who test more also have a lower fraction of positive tests among tested patients. In the latter case, physicians could differ in the threshold probability at which they think testing is worthwhile, and physicians who test more are moving further down the expected benefits curve.

To illustrate this point, we have sketched a stylized picture of the testing decision in Figure 1. Patients are sorted along the x-axis according to their risk of PE,  $q_{id}$ , from highest risk to lowest risk. The x-axis corresponds to the cumulative fraction of patients, and the y-axis corresponds

---

<sup>6</sup>In particular, both doctors would have similar test yields among marginal tested patients, but the doctor who tests more would have a higher test yield among the higher risk inframarginal patients. We formalize the points in this section in the context of our structural model in section 4.

to the marginal patient’s PE risk  $q_{id}$ , so that each point  $(x, y)$  along the plotted curve shows the fraction of patients  $x$  for whom  $q_{id} \geq y$ . For example, at point  $(T^A = 2/3, \tau^A = 1/2)$  in Panel A, the graph indicates that 2/3 of patients have a risk of PE that equals or exceeds 1/2. (We use this unrealistically high risk for illustrative purposes.)

In Panel A, we consider two doctors with the same patient distribution of PE risk, but with different testing thresholds. Doctor A tests every patient whose personal PE risk  $q_{id}$  exceeds Doctor A’s testing threshold  $\tau^A$ , and likewise Doctor B tests all patients for whom  $q_{id} > \tau^B$ . Because Doctor B’s threshold is lower than Doctor A’s, i.e.  $\tau_B < \tau_A$ , Doctor B tests a greater fraction of patients,  $T^B > T^A$ . Doctor B’s tested patients have a lower average PE risk than Doctor A’s tested patients, so Doctor B’s test yield  $Z^B$ —i.e. the fraction of positive tests among tested patients—is lower than Doctor A’s test yield  $Z^A$ , as can be seen in the graph. In this panel, there is a downward sloping relationship between the fraction of patients each doctor tests and his average test yield.

In Panel B, we consider an alternate scenario which could also explain why Doctor B continues to test a greater fraction of his patients than Doctor A, i.e. why  $T^B > T^A$ . In this example, doctor A and Doctor B have the same testing threshold, so  $\tau'_B = \tau'_A$ . Given the same expected patient PE risk, Doctors A and B would arrive at the same testing decision. However, the two doctors now face different distributions of patient PE risk. For any given probability of a positive test, Doctor B sees (weakly) more patients with  $q_{id}$  exceeding the common threshold for testing. In other words, Doctor B’s patient population is higher risk than Doctor A’s. As can be seen in the graph, Doctor B’s test yield  $Z^{B'}$  will be higher than Doctor A’s test yield  $Z^{A'}$ , even though both doctors have the same testing threshold, since more of the mass in Doctor B’s distribution of patient risk is concentrated at higher risk levels. In contrast with Panel A, there is now an upward sloping relationship between the fraction of patients each doctor has tested and his average test yield.

Now turning to our observed Medicare data, we use a simple binned scatterplot to explore whether variation in risk for PE or variation in testing behavior can explain the differences in physicians’ testing propensities. We begin by binning physicians into deciles according to the fraction of patients they test; next we calculate the fraction of tested patients for whom PE was detected within each decile. This relationship between fraction tested and average test yield is plotted in Figure 2. The graph displays a generally downward sloping relationship between average testing probability along the x-axis and fraction of tested patients with detected PE along the y-axis. Doctors who test a greater fraction of their patients are less likely to find positive test outcomes among tested patients. This figure suggests that differences in testing thresholds across doctors may be an important determinant of observed heterogeneity in testing behavior. It appears that doctors who are more likely to test their patients compared to their peers are also testing more low-risk patients.

Our structural model formalizes the intuition described above. It is designed to disentangle (observable and unobservable) differences in patient PE risk from differences in physician testing thresholds and evaluate the contribution of each to observed variation in testing behavior, following the intuition of this simple empirical exercise. We discuss the structural model in more detail in Section 4 below.

### 3.5 Reduced form evidence of misweighting patient PE risk factors

In addition to considering heterogeneity in physicians' testing thresholds, we also investigate whether physicians are successfully identifying observable risk factors associated with the highest probability of positive tests and testing patients with those characteristics. Determining which patients should be tested requires complex, subtle judgments about clinical risk on the basis of many factors. In our data, we capture some of the most common and relevant comorbidities by analyzing patients' claims histories. Guided by the structural analysis that follows, we motivate our exploration of misweighting PE risk with a few simple examples.

Consider a comparison of patients with a history of prostate cancer to those with no such history. Patients with a history of prostate cancer are no more likely to be tested for PE than patients without that condition; in fact, testing rates are slightly lower among prostate cancer patients (3.7%) compared to the rest of the population (3.8%). However, it turns out that among tested individuals, prostate cancer patients are over 50% more likely to be diagnosed with PE than patients with no such history.

A PE risk score popularly used to guide physicians on whether to order diagnostic testing includes treatment for cancer malignancy among its 7 risk criteria (Wells et al. 1995; Wells et al. 1998; Wells et al. 2000). And yet, although cancer is a recognized clinical risk factor for PE, a relationship supported by our data, it appears that patients with a history of prostate cancer are no more likely to be tested than the average ED patient. This provides the first suggestive evidence that physicians may not be properly accounting for the increased PE risk associated with prostate cancer, and thus may be under-testing prostate cancer patients relative to the rest of the population.

In Table 2, we highlight the basic summary statistics for eight of the clinical factors that show significant evidence of misweighting in the structural model that follows. Similar to the case of prostate cancer, we find that black patients are less likely to be tested than non-black patients, even though among tested patients, the rate of positive tests is much higher for black patients. A reverse pattern holds for patients with ischemic heart disease, atrial fibrillation or chronic obstructive pulmonary disease (COPD); they are tested at similar or higher rates than patients without those conditions, despite the fact that tested patients with these conditions are approximately 30% *less* likely to have a PE detected.

For other conditions, physicians respond in the right direction but overweight or underweight that condition relative to what would maximize the incidence of positive tests. The model implies that, everything else held equal (including other patient characteristics and physician thresholds), two comorbidities which have the same marginal impact on testing behavior should also have the same marginal impact on the conditional likelihood of a positive test. Our model identifies a few factors which appear to have a disproportionate impact on the likelihood of a positive test given their impact on testing behavior: a past history of PE, deep vein thrombosis, or a recent hospital admission are associated with 20 to 90 percent higher rates of testing but are 140 to 200 percent more likely to have a PE detected, a disproportionate increase relative to other factors in our model with a similar impact on testing behavior.

This simple exploration of misweighting relies on the presumption that patients with and without a particular risk factor don't differ in their other comorbidities and are sorting to ED physicians with

similar testing thresholds. In the structural model, we formalize this analysis, explicitly modeling differences in testing rates that may be driven by physician’s testing thresholds or other PE risk factors.

## 4 Model of testing behavior

### 4.1 Theory

Our structural model allows us to decompose the observed variation in physician testing rates into variation due to patient PE risk and variation due to doctor preferences, even if patient PE risk varies across doctors and is not fully captured by observable comorbidities. This approach to modeling physician preferences is based upon classical selection models developed by Heckman (1979) and Heckman and MaCurdy (1980) and subsequently refined and applied to a healthcare setting by Chandra and Staiger (2011).

Direct estimation of the selection model developed by Chandra and Staiger (2011) requires observing the individual-specific return to treatment for all treated individuals, a difficult object to recover in most empirical settings; we adapt the model to cover diagnostic testing, where test results (positive or negative) can proxy for the impact of treatment on the treated.<sup>7</sup> Further, we extend the Chandra and Staiger (2011) model to allow for the possibility that physicians are not appropriately weighting observable risk factors to select patients for testing with the highest expected PE risk.

To understand variation in physician testing decisions, we begin by studying the link between a physician’s decision to test a patient and the outcome of the test among tested patients. Assume that the suitability of a patient for testing is determined entirely by the *ex ante* likelihood of a positive test. We define  $q_{id}$  to be the conditional probability of a positive test for patient  $i$  evaluated by doctor  $d$ , given all the information available to the doctor:

$$q_{id} = x_{id}\beta + \alpha_d + \eta_{id} \tag{1}$$

where  $x_{id}$  are observed patient characteristics,  $\alpha_d$  are doctor fixed effects, and  $\eta_{id}$  are factors observable to the doctor but unobservable to the econometrician which impact the likelihood that a test is positive. Note that the inclusion of physician fixed effects  $\alpha_d$  allows the population risk of PE to vary across doctors in ways that are not captured by the included patient covariates. Following the typical structure of Heckman selection models, we begin by assuming that  $\eta_{id}$  is independently and identically distributed across doctors; we refer to this as the “ignorability assumption” following the prior literature. (We explore relaxing the ignorability assumption in Section 6.) Further assume that  $\eta_{id}$  is bounded with full support.

Following Chandra and Staiger (2011), we make the structural modeling assumption that physicians apply the same decision rule to each patient. Suppose that physicians test if and only if the probability of a positive test  $q_{id}$  exceeds a physician-specific threshold  $\tau_d$ . That is, they test if and

---

<sup>7</sup>Given our assumption that negative test results do not improve patient health *ex post*, the testing outcome can proxy for the impact of treatment on the treated, as long as the benefits of treating a detected PE are constant across patients. The clinical basis for this assumption is discussed at greater length in Appendix B.

only if:

$$Test_{id} = 1 \leftrightarrow q_{id} = x_{id}\beta + \alpha_d + \eta_{id} > \tau_d \quad (2)$$

which implies that:

$$Pr(Test_{id} = 1) = f(x_{id}\beta + \alpha_d - \tau_d) \quad (3)$$

where the functional form of  $f(x_{id}\beta + \alpha_d - \tau_d) = Pr(\eta_{id} > -(x_{id}\beta + \alpha_d - \tau_d))$  depends on the distribution of  $\eta_{id}$ . Our goal is to recover the parameters  $\beta$ ,  $\alpha_d$  and  $\tau_d$ . Separately identifying the two sets of physician fixed effects,  $\alpha_d$  and  $\tau_d$ , will allow us to decompose variation in observed testing rates into variation due to differences in patient PE risk ( $x_{id}\beta + \alpha_d$ ) and differences due to physician thresholds for testing ( $\tau_d$ ). From equation 3 alone,  $\alpha_d$  and  $\tau_d$  are not separately identified; to separate them, we will need to use data on test outcomes.

By estimating equation 3, we can calculate the predicted conditional probability that a patient is tested, which will be a nonlinear function of the testing propensity index  $I_{id} = x_{id}\beta + \alpha_d - \tau_d$ . Let  $Z_{id}$  denote a binary variable for tested patients indicating whether the test is positive or negative. If every patient were tested, we would observe  $Z_{id}$  for the entire sample and could recover  $\beta$  and  $\alpha_d$  by estimating the linear probability model implied by equation 1 using OLS. (Of course, if every patient were tested, there would be no variation in doctor testing thresholds.) In practice, we only observe whether a test is positive or negative for those patients whom doctors choose to test, so there is a selection problem; this is the standard selection problem originally studied by Heckman (1979).

Formally, we model testing outcomes as follows:

$$\begin{aligned} E(q_{id}|Test_{id} = 1) = E(Z_{id}|q_{id} > \tau_d) &= x_{id}\beta + \alpha_d + E(\eta_{id}|q_{id} > \tau_d) \\ &= x_{id}\beta + \alpha_d + h(x_{id}\beta + \alpha_d - \tau_d) \\ &= \tau_d + \lambda(I_{id}) \end{aligned} \quad (4)$$

where  $h(x_{id}\beta + \alpha_d - \tau_d) \equiv E(\eta_{id}|q_{id} > \tau_d) = E(\eta_{id}|\eta_{id} > -I_{id})$  and  $\lambda(I_{id}) \equiv I_{id} + h(I_{id})$ . Because we only observe whether a test is positive conditional on patients being tested, a regression of the indicator for a positive test  $Z_{id}$  on  $x_{id}$  and doctor fixed effects would produce biased estimates of  $\beta$  and  $\alpha_d$  unless we properly control for the selection correction. Equation 4 is the primary equation of interest, and equation 2 governs selection into that sample.

The binned scatterplot of testing rates and test yields described in section 3.4 can provide some intuition for understanding this model. Variation in testing propensities  $I_{id}$  could be driven by differences in patient PE risk, either through differences in observed comorbidities  $x_{id}$  or unobserved population risk  $\alpha_d$ . Alternatively, differences in testing propensities could be explained by differences in physician testing thresholds  $\tau_d$ . If all variation across doctors in testing behavior were driven by patient PE risk, then we would typically find that physicians with higher average testing propensities have higher test yields.<sup>8</sup> On the other hand, variation in physician testing thresholds  $\tau_d$  will lead to a downward sloping relationship between testing propensities  $I_{id}$  and test yields  $E(Z_{id}|q_{id} > \tau_d)$ ,

---

<sup>8</sup>This is satisfied as long as  $E(\eta_{id} + I_{id}|\eta_{id} + I_{id} > 0)$  is upward sloping in the function  $I_{id}$ . This restriction holds for many general distributions of  $\eta_{id}$ , including, for example, under distributions meeting the restriction that  $\eta_{id}$  is symmetric and mean 0, or if the density of  $\eta_{id}$  is non-decreasing.

as can be seen from equation 4 above.<sup>9</sup> This derivation formalizes the intuitive argument made in section 3.4, which interpreted the observed downward sloping relationship between doctors' average fraction of patients tested and test yield as evidence of variation in testing thresholds.

## 4.2 Misweighting of patient risk

A key difference between our model and Chandra and Staiger (2011) is that we extend the model laid out above to allow for the possibility that doctors misweight observable characteristics in deciding which patients to test. We previously assumed that the coefficients  $\beta$  attached to patient observables when doctors decide which patients to test reflect the true relationship between those characteristics and the likelihood of a positive test. This need not be the case. Doctors may under- or over-weight the importance of different risk factors, so that testing is not necessarily targeted at the highest risk patients. Assume that each doctor's belief about the probability of a positive test is given by:

$$q'_{id} = x_{id}\beta' + \alpha'_d + \eta_{id} \quad (5)$$

while the actual probability remains:

$$q_{id} = x_{id}\beta + \alpha_d + \eta_{id} \quad (6)$$

We define the new testing propensity  $I'_{id} = x_{id}\beta' + \alpha'_d - \tau_d$  to reflect the observed propensity given physician beliefs about  $\beta'$  and  $\alpha'_d$ .

With this change, we can rewrite the test outcomes equation:

$$\begin{aligned} E(Z_{id}|Test_{id} = 1) &= E(q_{id}|q'_{id} > \tau_d) \\ &= E(q'_{id}|q'_{id} > \tau_d) + x_{id}(\beta - \beta') + \alpha_d - \alpha'_d \\ &= \tau_d + x_{id}(\beta - \beta') + \alpha_d - \alpha'_d + \lambda(I'_{id}) \end{aligned} \quad (7)$$

The model with misweighting now has three different doctor fixed effects,  $\tau_d$ ,  $\alpha'_d$ , and  $\alpha_d$ . As a result, to separately identify  $\tau_d$  based on equation 7 we will require one additional assumption: we assume that doctors may misweight observable risk factors but are correct on average about the probability that their tested patients will have PEs. This assumption can be written as:

$$E_d(q'_{id}|Test_{id} = 1) = E_d(q_{id}|Test_{id} = 1) \quad (8)$$

where  $E_d$  denotes the conditional expectation for doctor  $d$ . Note that this implies that doctors' beliefs about  $\alpha'_d$  must offset misweighting in each doctor's patient population so that doctors have correct expectations about the overall rate of positive tests among their tested patients. Then an analogous derivation to equation 4 gives:

$$E(Z_{id}|Test_{id} = 1) = \tau_d + (x_{id} - E_d(x_{id}|Test_{id} = 1))(\beta - \beta') + \lambda(I'_{id}) \quad (9)$$

---

<sup>9</sup>Note that  $I_{id}$  will increase as  $\tau_d$  falls. Increases in  $I_{id}$  will cause  $E(\eta_{id}|\eta_{id} > -I_{id})$  to fall and the test yield among tested patients  $E(Z_{id}|q_{id} > \tau_d)$  will decrease in turn.



This is identical to equation 4, except now the (demeaned) observables  $x_{id}$  directly enter the test outcomes equation, even after conditioning on the propensity to test. In other words, the model implies that if observables  $x_{id}$  continue to have explanatory power after conditioning on the propensity  $I_{id}$ , then physicians are not weighting those observables in the manner that would maximize the incidence of positive tests.

### 4.3 Identification

As is typical for Heckman selection models,  $\lambda(\cdot)$  can in principal be identified using functional form restrictions, but more desirable identification is only feasible with exclusion restrictions.

In the simpler form of the model without misweighting, as presented by Chandra and Staiger (2011) and outlined in Section 4.1, identification may come from the fact that  $x_{id}$  only enters the test outcome equation, equation 4, via  $\lambda(I_{id})$ . In that model,  $x_{id}$  are excluded from directly entering the test outcomes equation and we can think of them as instrumental variables which aid in the estimation of  $\lambda(\cdot)$ , parallel to the standard instrumental variables identification in Heckman selection models (e.g. Mulligan and Rubinstein 2008; Chandra and Staiger 2011). This restriction is no longer valid if physicians incorrectly assess the PE risk associated with some observable comorbidities and demographics  $x_{id}$ . In the model with misweighting, equation 9 above shows that  $x_{id}$  directly enters the test outcomes equation with coefficients that are not known from estimating the equation governing selection into testing.

In order to generalize the model to the case where doctors fail to appropriately weight observable risk factors in deciding whom to test, we consider a new set of exclusion restrictions. We exploit the fact that  $\tau_d$  can be directly estimated for physicians testing patients we can identify as marginal.<sup>10</sup> Marginal tested patients are those with the lowest observed values of the testing propensity  $I'_{id}$  who are still tested. We estimate the average probability of a positive test among these marginal tested patients. For these patients who are “just barely worth testing,” the observed probability of a positive test reveals the threshold at which doctors are willing to test.

Formally, since  $\eta_{id}$  is bounded with full support, there exists some value of the propensity in the testing equation  $\underline{I}$  such that patients are only tested for  $I'_{id} > \underline{I}$ . For those marginal tested patients with  $I'_{id} \rightarrow \underline{I}$ , we know the realization of  $\eta_{id}$  is just barely sufficient to tip these patients across the testing threshold, so that  $h(\underline{I}) = E(\eta_{id}|q'_{id} = \tau_d) = -\underline{I}$ . Since  $\lambda(I_{id}) = I_{id} + h(I_{id})$ , it follows that  $\lambda(\underline{I}) = 0$  for these marginal tested patients.

Let  $QQ_d$  denote the average rate of positive tests  $Z_{id}$  among tested marginal patients for doctor  $d$ ; taking the expectation of equation 9 yields:

$$QQ_d = \tau_d + (E_{m,d}(x_{id}|Test_{id} = 1) - E_d(x_{id}|Test_{id} = 1))(\beta - \beta') \quad (10)$$

In the equation above,  $E_{m,d}(x_{id}|Test_{id} = 1)$  denotes the expectation of  $x_{id}$  only among doctor  $d$ 's tested marginal patients  $m$ . The likelihood of a positive test for those tested patients with the lowest testing propensities is given by the physician's threshold  $\tau_d$  plus an adjustment for the fact that the actual likelihood of a positive test for these patients differs from physician's beliefs because  $\beta \neq \beta'$ .

<sup>10</sup>More precisely,  $\tau_d$  is known modulo a misweighting adjustment we spell out below.



This provides an exclusion restriction—after subtracting the average yield among marginal patients from both sides, doctor fixed effects are excluded for physicians with tested marginal patients in equation 9. A more detailed derivation of this result is in Appendix E.

Intuitively, suppose that by studying marginally tested patients, we uncover multiple physicians with identical thresholds  $\tau_d$ ; these doctors may still differ in their patients’ risk of PE  $\alpha_d$  and thus in the propensity to test for identical observables. This variation in  $\alpha'_d$  across doctors known to have identical testing thresholds can then identify the function  $\lambda(\cdot)$  when we estimate the testing outcomes equation (cf. equation 9); in other words, we can estimate the slope of the selection term by asking to what extent doctors with higher testing propensities also tend to have more positive tests, when comparing doctors with similar thresholds who treat patients with similar observables.

This identification argument naturally raises the question: why not simply estimate  $\tau_d$  directly for all doctors using only tested marginal patients, without estimating equation 9? If we observed a large number of patients for each doctor in our sample, this approach would be feasible (although it still would not identify misweighting). Unfortunately, since many doctors only test a small number of in-sample patients, we cannot recover  $\tau_d$  for most doctors using this method. Instead, we can think of the doctor fixed effects for physicians who test marginal patients as excluded from estimation of equation 9 (since their coefficient  $\tau_d$  is known); the variation in those known fixed effects aids in identifying  $\lambda(\cdot)$ , parallel to the role of instrumental variables in standard Heckman selection models, and in turn allows us to recover  $\tau_d$  for doctors with non-marginal patients.

In addition to the validity of the exclusion restrictions, the other crucial identifying restriction under this estimation approach is the ignorability assumption:  $\eta_{id}$  is i.i.d. across doctors and patients. The ignorability assumption implies that the function  $\lambda(\cdot)$  is the same for different doctors and patients. If this assumption were violated and  $\eta_{id}$  were distributed differently across doctors, the function  $\lambda(\cdot)$  could be doctor-specific. In Section 6.2, we consider one such model and show that it does not materially impact our results.

The identification of misweighting also relies on the ignorability assumption, i.e. the assumption that the error term  $\eta_{id}$  is i.i.d. The ignorability assumption implies that if doctors were optimally assessing PE risk, any two conditions with the same  $\beta'$  weight in the testing equation should induce the same change in the fraction of positive tests amongst tested patients, holding all other comorbidities and testing thresholds constant. If two conditions with the same  $\beta'$  weight in the testing equation lead to different changes in the fraction of positive tests, then we identify misweighting; we conclude the risk factor that induces the larger increase in positive tests is underweighted relative to the other factor. The slope of the function  $\lambda(\cdot)$  with respect to  $\alpha'_d$  pins down how  $x_{id}$  should impact test outcomes  $Z_{id}$  given  $\beta'$ —so we can in principle identify misweighting even with just a single  $x$  variable. This strategy echoes the logic of the reduced form evidence on misweighting presented in section 3.5, but the additional structure allows us to make more detailed comparisons of weighting and risk across conditions, after accounting for differences in patient risk and testing thresholds across doctors.

Empirically, the ignorability assumption may be undermined if the distribution of unobserved patient PE risk differs across conditions. For example, if fewer patients with the risk factor that appears to be under-weighted present to the ED with the relevant PE symptoms (e.g. chest pain,

shortness of breath, elevated heart rate), then it may be that physicians are already testing every patient in the relevant at-risk population. This assumption is directly analogous to the standard exogeneity assumption used in virtually all structural models; e.g. just as discrete choice models assume that observed product characteristics are independent of the error term, our misweighting model is identified by assuming that observed characteristics are not systematically related to unobserved determinants of PE risk.

#### 4.4 Estimation of the structural model

Let us now specify precisely how we estimate the structural model outlined in the previous sections. Define  $\theta'_d = \alpha'_d - \tau_d$ . Plugging our specification for the probability of a positive test from equation 5 into the testing equation 2 yields the final form of the testing equation:

$$Test_{id} = 1 \leftrightarrow x_{id}\beta' + \theta'_d + \eta_{id} \geq 0 \quad (11)$$

These assumptions yield a binary choice model of testing. In our baseline specification, we assume that  $\eta_{id}$  is i.i.d. across doctors and patients with a parametric distribution we describe below. Thus, patients' ex ante risk distributions may have different means ( $x_{id}\beta + \alpha_d$ ) but are assumed to be otherwise identically distributed.

Specifically, we assume that each  $\eta_{id}$  is drawn from a two parameter distribution which is a mixture of a Bernoulli and a uniform distribution. With probability  $1 - p$ ,  $\eta_{id} \sim U[-\eta, \eta]$  and with probability  $p$ ,  $\eta_{id} \sim U[v - \eta, v + \eta]$ . Intuitively, this distribution captures the idea that most patients are not candidates for a CT scan. A small fraction of patients  $p$  present with symptoms of PE such as chest pain and given those symptoms, there is a range of ex ante risks parameterized by  $\eta$ .<sup>11</sup> We assume that patients are never tested unless they receive the shock  $v$  (i.e. unless they present with PE symptoms). In Appendix E, we show that this implies:

$$Pr(Test_{id} = 1) = \max \left\{ 0, \frac{p}{2} + \frac{p(I'_{id} + v)}{2\eta} \right\} \quad (12)$$

where  $I'_{id} = x_{id}\beta' + \theta'_d$ . Estimation of this equation by non-linear least squares allows us to recover  $\hat{\beta}' = \beta' \frac{p}{2\eta}$  and  $\hat{\theta}' = \frac{p}{2} + \frac{p(\theta'_d + v)}{2\eta}$  which we use to construct an estimate of the testing propensity  $\tilde{I}'_{id} = \frac{p}{2} + \frac{p(I'_{id} + v)}{2\eta}$ .

Following the steps outlined in the previous section, the testing threshold parameters  $\tau_d$  can be recovered from a regression of test outcomes (i.e. positive or negative for detecting PE) on doctor fixed effects, controlling for the propensity  $I'_{id}$  estimated from the testing equation. Note that under the parametric assumptions we have made so far,  $E(\eta_{id} | \eta_{id} > -I'_{id}) = \frac{\eta - I'_{id} + v}{2}$ . As shown in more

---

<sup>11</sup>Methodologically, we use this mixture distribution rather than simply assuming a uniform for two reasons: firstly, because if we assume that  $p = 1$  (the uniform case), the estimated variance of  $\eta$  is so large that it implies  $q_{id} < 0$  in some cases, which is inconsistent since  $q_{id}$  is a probability. Secondly, since testing is a low probability event, a uniform distribution would imply that more precise information (a higher variance of  $\eta_{id}$ ) leads doctors to test more everything else held equal; the mixture distribution allows for the possibility that more precise information leads to less testing. This second point is especially relevant in the heteroskedastic model considered in the robustness section where the variance of  $\eta_{id}$  is allowed to vary across doctors.

detail in Appendix E, this implies that:

$$E(Z_{id}|Test_{id} = 1) = \tau_d + (x_{id} - E_d(x_{id}|Test_{id} = 1))(\beta - \beta') + \frac{\eta \tilde{I}'_{id}}{p} \quad (13)$$

As discussed in section 4.3, we avoid relying on functional form identification for the coefficient on  $\tilde{I}'_{id}$  by imposing the additional restriction that  $\tau_d$  can be estimated directly for doctors with tested marginal patients based on the observed average rate of positive tests among those marginal patients,  $\widehat{QQ}_d$ . We define marginal patients as patients in the first decile of  $\tilde{I}'_{id}$  among tested patients; this definition is conservative from the standpoint of detecting overtesting since more restrictive definitions (e.g. the first percentile) will tend to lead to lower estimated thresholds.

Imposing the formally correct version of this constraint yields our estimating equation:

$$Y_{id} = (1 - M_d)\tau_d + \frac{\eta \tilde{I}'_{id}}{p} + X_{id}(\beta - \beta') + \epsilon_{id} \quad (14)$$

where  $Y_{id} = Z_{id}$  for doctors with no tested marginal patients and  $Y_{id} = Z_{id} - \widehat{QQ}_d$  for doctors with marginal patients,  $M_d$  is an indicator for whether a doctor has marginal patients,  $X_{id} = (x_{id} - E_{m,d}(x_{id}))$  for doctors with marginal patients and  $(x_{id} - E_d(x_{id}))$  for doctors with no marginal patients.

Least squares estimation of equation 14 will allow us to recover the constant  $\frac{\eta}{p}$  and doctor fixed effects  $\tau_d$  for non-marginal patients which, when combined with our estimates for marginal patients from  $\widehat{QQ}_d$ , can be used to recover the full distribution of estimated  $\hat{\tau}_d$ .

The distribution of  $\hat{\tau}_d$  combines both the true underlying variation in  $\tau_d$  and estimation error from the fact that each  $\tau_d$  is imprecisely estimated. To correct for estimation error, we apply an “empirical Bayes” technique to recover moments of the true underlying distribution of  $\tau_d$ . Our approach is described in detail in Appendix F.<sup>12</sup> Unlike more standard estimators (such as Kane and Staiger 2008), this technique is robust to the fact that we observe only a small number of observations per doctor and makes no distributional assumptions about either the true distribution of  $\tau_d$  or the estimation error. The true distribution cannot be nonparametrically identified, but we can recover moments of that distribution; we report the mean and standard deviation. Simulation results do require us to recover a posterior estimate of  $\tau_d$  for each doctor, and for these exercises we impose a further assumption that  $\tau_d$  is log-normally distributed as described in Appendix F.

## 5 Results

In this section, we report results of the estimation strategy described in section 4.4 above. First, we describe the recovered distribution of physician testing thresholds and explore how differences in test yields across physicians depend on differences in testing thresholds. Next, we test how physicians’ training and practice environment are related to practice styles. Then, we report results on which risk factors are under- and over-weighted in physicians’ risk assessments relative to the weighting

<sup>12</sup>We use quotation marks since our procedure is not a traditional empirical Bayes approach: we do not derive our estimator as the posterior of any specific distribution.

that would maximize detection of positive tests and consider possible clinical explanations for these patterns. Finally, we simulate how misweighting affects physicians’ test yields.

### 5.1 Distribution of physician testing thresholds and relationship to test yields

After estimating the model laid out in Section 4 and applying the empirical Bayes adjustment, we find the mean value of  $\tau_d$  is 0.056 and the standard deviation is 0.054.<sup>13</sup> In other words, the average doctor is willing to test a patient provided the doctor’s estimate of the probability of a positive test exceeds 5.6%. Note that this positive test rate includes tests which detect actual PEs and false positives. The standard deviation of 0.054 suggests that there is a large amount of heterogeneity across doctors in their testing thresholds, with some doctors testing almost all patients displaying the relevant symptoms, and other doctors testing only patients with very substantial PE risk. Considering that the overall test yield in our sample is only 6.9%, it is likely that this variation in testing thresholds may affect testing decisions for many patients.

To quantify the role that heterogeneity in testing thresholds plays in the observed patterns of testing behavior and test yields, we return to the graph of physician testing rates and test yields. Now, rather than binning physicians by the average fraction of patients tested as we did in Figure 2, we bin physicians by the structural analogue: the average estimated testing propensity  $\hat{I}'_{id}$  across their patients. Recall the observation from the reduced form analysis in section 3.4 that physicians with the highest average testing rates also had the lowest test yields. This downward sloping relationship is what we would expect to find if heterogeneity in  $\tau_d$  were the primary driver of observed variation in testing rates across doctors.

We can explore this hypothesis more formally by using our model to simulate what the relationship between average physician testing propensities and positive test rates would have been if all doctors had the same testing threshold. We simulate testing decisions and test outcomes under a counterfactual where  $\tau_d$  is held constant across doctors, at the estimated average value  $E(\tau_d) = 0.056$ . Details of this simulation are provided in Appendix G.

Results of this exercise are pictured in Figure 3. The solid black line depicts the downward sloping relationship between physicians’ average testing propensities and their test yields in our observed data. As we suggested earlier, if all doctors had the same testing threshold, the remaining variation in doctors’ average testing propensities would be driven by differences in patient risk of PE. As a result, the relationship between doctors’ average testing propensities and their test yields would become upward sloping over most of the domain. The dashed line with the “x” markers display the results of this simulation in Figure 3. Now the doctors with higher testing rates are those with the highest risk patients; these doctors test the greatest fraction of their patients and experience the highest test yields, as evidenced by the upward slope in the simulated plot.<sup>14</sup>

<sup>13</sup>Note that of course this would not be consistent with a normal distribution since in this case  $\tau_d > 0$  for all doctors or they would test every patient. In our welfare exercises we assume a log-normal distribution.

<sup>14</sup>If we graphed testing propensities vs. simulated rates of positive tests at the individual patient level, fixing  $\tau_d = E(\tau_d)$ , our model implies that the resulting relationship would be monotonic. Because we are aggregating to the physician level in the figure, this relationship also depends on the variance in testing propensities for a given physician; the slight non-monotonicity at the lowest deciles arises because doctors with the lowest average testing propensities have more heterogeneous patients (driven by variation in observed comorbidities  $x_{id}$ ) than those in adjacent deciles. At these low average testing propensities, higher variance in  $I_{id}$  is associated with more positive tests amongst tested

## 5.2 Determinants of physician testing thresholds

We next consider regressions of the estimated testing thresholds  $\hat{\tau}_d$  on doctor, hospital and regional characteristics to explore the determinants of practice style. Specifically, we regress  $\hat{\tau}_d$  on variables capturing doctor experience (the number of years since the doctor graduated from medical school), whether the medical school the doctor attended is ranked in the top 50 for research or primary care by US News & World Report, whether the hospital where the physician practices is a for profit hospital or an academic hospital, regional medical spending, the state tort environment, and average income in the region.

We consider OLS estimates as well as FGLS estimates which take into account the estimation error in the dependent variable  $\tau_d$ .<sup>15</sup> For each specification, we consider models with and without hospital fixed effects. Including hospital fixed effects to identify the impact of within-hospital variation in physician characteristics obviates the concern that our model omits unobserved differences in the cost of testing at the hospital level. For example, there may be variation in the opportunity cost of testing, depending on whether the CT scan is used to capacity. This heterogeneity will be absorbed into the hospital fixed effect.

Table 3 reports the results. We find that doctors in higher spending regions have lower testing thresholds, i.e. they are more likely to test low risk patients. A 10% increase in regional spending, as reported by the Dartmouth Atlas, is associated with a 0.4 percentage point decline in testing thresholds, significant at the 1% level. This finding provides empirical support for the hypothesis that high spending regions are providing lower marginal value, “flat of the curve” medical care.

We also find evidence that more experienced doctors have higher testing thresholds: a 10-year increase in doctor experience is associated with 0.7 percentage point higher testing thresholds, significant at the 1% level. This relationship persists after controlling for hospital fixed effects, suggesting that even within the same institution, more experienced doctors are less likely to test low-risk patients. Unfortunately, we do not observe enough testing decisions per physician to estimate the model with year-specific testing thresholds for each physician, and as a result we cannot disentangle cohort and experience effects.

Many factors predicted to influence care quality, such as the quality of the physician’s training, the financial structure of the hospital (for profit or otherwise), its status as an academic institution, and the income of the patients served have no significant relationship to testing thresholds. Estimates relating physician’s medical school rank to testing thresholds are imprecisely estimated, with the upper bound of the 95% confidence interval at a 1.2 percentage point higher threshold for those attending a top 50 research institution. Point estimates suggest slightly higher thresholds for academic hospitals and lower thresholds among for-profit hospitals, but the 95% confidence intervals bound the differences in average thresholds to less than one percentage point.

Finally, exploiting cross-sectional variation in enactment of tort reform, including joint and several liability and malpractice damage caps, we find no consistent relationship between the malpractice environment and testing thresholds. The FGLS estimates point to a significant, negative

---

patients due to the convexity of the relationship between  $I_{id}$  and positive testing rates at the individual level.

<sup>15</sup>The FGLS estimates are based on Lewis and Linzer (2005), where the error term consists of both a homoskedastic  $\epsilon_{id}$  with unknown variance and a heteroskedastic component with known variance. The heteroskedastic component arises from the estimation error in  $\hat{\tau}_d$  which is in turn recovered from estimation of equation 14.

relationship between testing thresholds and malpractice damage caps, which would be the opposite prediction of theory suggesting physicians are more likely to test low-risk patients in states with damage caps. The coefficient is much smaller in magnitude and no longer statistically significant in the OLS specification. Our lack of power to estimate year-specific testing thresholds precludes us from undertaking a difference-in-differences analysis of malpractice law.

Given the large estimated variation in  $\tau_d$ , with a standard deviation of 0.054 after adjusting for statistical noise, observed factors can explain only a small fraction of the estimated variation in physician practice style. This observation implies that policy responses targeted at reducing testing rates in specific hospital types (e.g. for profit hospitals) or policies aimed at raising the qualifications of emergency department doctors are unlikely to lead to substantial reductions in testing variation. Instead, focusing on policies which target the decision-making process rather than physician credentials or practice environment may have greater scope for reducing heterogeneity in practice style. This parallels the finding in the teacher fixed effects literature that there is substantial variation in teacher productivity not explained by teacher credentials or other observable factors (Jackson et al. 2014).

### 5.3 Identifying misweighted comorbidities

Next, we explore physicians' misweighting of observable PE risk factors. As outlined in section 4.2, we focus on measuring aggregate misweighting: factors which appear to be systemically under- or over-weighted in physicians' assessments of patient PE risk. The model implies that physicians are overweighting a given risk factor if they are substantially more likely to test a patient with that factor (holding constant other observable patient characteristics), but this variable does not yield a commensurate increase in the rate of positive tests among tested patients. The evidence of both under- and over-weighting suggests that physicians could perform the same total number of tests, but detect more PE cases, if they improved targeting of the tests by applying different weights to many important risk factors.

Results are reported in Table 4 and Appendix Table A.1. For each risk factor in our model, column 1 reports the marginal effect of this variable on testing probability based on the coefficient  $\beta'$  from the testing equation (cf. equation 5). Column 2 reports the estimated error in physicians' assessment of the PE risk associated with each comorbidity, implied by how the weights attached to each comorbidity in their testing decisions compare to the conditional influence of each comorbidity on test outcomes (cf. equation 13). Finally, columns 3 and 4 report the standard error and t-statistic on estimated misweighting, respectively. Variables are sorted by their t-statistic in this table.

Given our nonlinear model, the reported marginal effects in column 1 hold for all patients for whom  $\tilde{I}'_{id} > 0$ , which is true for the average patient in our data. (Marginal effects are zero for patients with negative values of  $\tilde{I}'_{id}$ .) All included risk factors are binary variables; variables with the most misweighting will have the largest absolute value of misweighting reported in column 2. We report robust standard errors that don't account for estimation error in the testing propensity index  $\tilde{I}'_{id}$ , although this adjustment would be very small given the large sample of patients identifying  $\tilde{I}'_{id}$ .

We find evidence of substantial under- and over-weighting of key risk factors, relative to the weights that would maximize test yields. Comparing physician's implied prediction of PE risk for



each patient with the estimated actual risk, we find that physicians appear to be misestimating a patient’s probability of a positive test by 2.3 percentage points on average, accounting for all comorbidities and averaging the absolute value of each patient’s aggregate misweighting to include both under- and over-estimates. This degree of misestimation has the potential to affect testing decisions for many patients.

The strongest evidence of underweighting comes from physicians’ implicit estimate of the PE risk associated with a recent inpatient admission history. While immobilization is a commonly known risk factor for PE, popular risk scores highlight the role of recent surgery but do not broadly include other types of hospitalization. Perhaps as a result, we see evidence that physicians have adequately increased testing rates for patients with a recent surgical history, but do not place sufficient weight on recent hospital admissions that did not include a surgical procedure. The marginal effect reports that physicians are 0.9 percentage points *less* likely to test a patient with a prior inpatient admission within the past 30 days, implying that doctors have underestimated these patients’ PE risk by 11 percentage points after account for the role of other observed comorbidities.

In addition, several specific cancer diagnoses, and a history of PE or the related condition deep vein thrombosis also show evidence of substantial underweighting, suggesting that physicians are failing to adequately consider these risks when assessing a patient for PE.<sup>16</sup> For all but one of these conditions (metastatic cancer), physicians are indeed more likely to test patients with the observed condition, holding constant other patient risk factors, but the response is not adequate given the large influence of this preexisting condition on the current risk for PE. This pattern is occurring despite the fact that both cancer treatment and history of PE or deep vein thrombosis are two of the seven risk factors in a popular PE risk-scoring algorithm known as the Wells score. This suggests that physicians are continuing to under-respond to these critical risk factors despite their recognized role in PE risk.<sup>17</sup>

A few other risk factors also show evidence of significant underweighting, including rheumatoid arthritis, obesity and paralysis, all of which are known risk factors for PE documented in the medical literature, although not explicitly included in popular risk scoring algorithms. A complete list of underweighted risk factors is reported in the top panel of Table 4.

Turning to demographic variables, we find evidence that black patients are under-tested. They are less likely to be tested for PE than non-black patients, despite the fact that they are at higher risk of PE. Given the structure of our model, these differences in testing patterns of black and white patients cannot be explained by differential sorting to physicians, since we have controlled for differences in physicians’ testing thresholds. This finding provides new empirical support for the concern about racial disparities and possible provider prejudice in medical treatment (cf. Nelson 2002). The result stands in contrast to results from Chandra and Staiger (2010) that applied a related analytic framework to a different clinical setting and found that while blacks receive less treatment for heart attacks, differences were fully explained by their lower benefits from treatment.

---

<sup>16</sup>Prostate cancer, metastatic cancer, endometrial cancer and colorectal cancer all have significant underweighting.

<sup>17</sup>Whether the underweighting of these risk factors is driven by failure to adhere to Wells’ score criteria or whether the Wells score inadequately weights these risks is not something we can directly assess in our data. Complete calculation of the Wells’ score would require information that is difficult to observe in claims data or even retrospective study of patient charts. For example, the most highly weighted factor in the score is the physician’s clinical opinion that PE is the most likely diagnosis, or equally likely to the other possible diagnosis.



In the setting of testing for PE, differences in test yields do not explain disparities in testing rates. Notably, these disparities are arising among patients who all have Medicare insurance coverage, although they may differ in their subscription to wrap-around private insurance, and all of whom have arrived at the emergency room for evaluation by a physician with access to a CT scanner.

A number of different factors show evidence of overweighting: these are conditions where test yields are predicted to improve if physicians became less likely to test patients with these particular conditions. Several of these overweighted conditions, including the three with the most significant evidence of overweighting (i.e. atrial fibrillation, chronic obstructive pulmonary disease, and ischemic heart disease), have chest pain and difficulty breathing as hallmark symptoms; these are also key clinical symptoms of PE. Patients who visit the ED with an exacerbation of another previously diagnosed condition could be suspected of having PE due to similar symptoms and thus may be tested at a higher rate even though our data suggests they are not at higher risk of PE, holding constant their other risk factors. Given that these other conditions must have been diagnosed prior to the ED visit in order to be included on our comorbidity list, physicians should be aware of them at the time they are evaluating the patient for PE. Of course, failure to take an appropriate medical history or limited access to patients' prior health records could hinder evaluation and contribute to the observed overweighting of these conditions.

Taken together, these results suggest that misassessments of the clinical risk associated with pre-existing comorbidities may lead to substantially diminished test yields. It is possible that physicians could detect more PE cases while performing a similar number of tests, by adjusting the targeting.

An alternative explanation for these patterns of apparent misweighting would be that the value of detecting PE differs for patients with these varying risk factors. For example, if the value of detecting PE were substantially lower in patients with a recent hospital admission or a cancer diagnosis, that could explain the apparent underweighting. Conversely, if the value of detecting PE were higher for patients with ischemic heart disease, COPD or atrial fibrillation, then that could also help rationalize the observed testing behavior. We find no obvious link between these conditions and the value of PE detection. In fact, our results on age-related risk suggests that physicians are undertesting younger patients, for whom the value of PE detection should be particularly high, since they have a longer life expectancy and accordingly higher value of statistical life.

#### **5.4 The impact of misweighting on test yields**

We now return to the graph that displays the relationship between physicians' average testing propensities and test yields to see how misweighting impacts this relationship. Recall that the graph is downward sloping in our observed data: much of the variation in average testing propensities is driven by differences in physician testing thresholds, and doctors with lower testing thresholds have lower test yields among tested patients. In section 5.1, we found that if there were no variation in physician testing thresholds, then the relationship between average testing propensities and test yields would become upward sloping, since variation in testing propensities would now be solely driven by differences in patient PE risk.

In this section, we consider the role of misweighting in determining the relationship between testing propensities and yield. We simulate the counterfactual relationship between physicians'

average testing propensities and test yields that would be observed if there were no heterogeneity in testing thresholds *and* no misweighting of observable risk factors. Eliminating misweighting should increase the test yield for all values of the testing propensity index by improving the targeting of PE CT tests. Details of the simulation exercise are described in Appendix G.

Results of this simulation are pictured in Figure 3 and plotted with the dashed line with triangle markers. We see that for every decile of physicians’ average testing propensity, the predicted test yield is higher in the simulation with no misweighting than was observed in both our actual data or the simulation that only eliminated threshold variation. We predict more detected positive tests if physicians attached appropriate weights to observable risk factors, and the increase is largest at lower testing propensities. (We quantify the precise increase in test yields and their welfare consequences in section 7.3.) Inframarginal patients are likely to be tested even with misweighting, but the set of marginal patients changes—some patients who are less likely to test positive are no longer tested and others who were previously not tested but have a higher likelihood of testing positive are now tested. This exercise suggests that misweighting is a substantial contributor to low test yields, and attention to better targeting of testing resources is warranted, rather than focusing solely on reducing variation in testing rates.

## 6 Robustness

The results discussed in the previous sections depend on a number of modeling assumptions. The critical identifying assumption can be framed in terms of the specification  $\eta_{id}$  term, the factors influencing testing choices that are observable to the doctor but unobservable to the econometrician. In our baseline specification, we assume that  $\eta_{id}$  is i.i.d. across patients and doctors and follows a specific parametric distribution. In the robustness checks described below, we test the sensitivity of our results to these assumptions. Specifically, we consider the robustness of our results to the set of included covariates (which essentially tests robustness with respect to a particular form of heteroskedasticity); we estimate a version of our model where the variance of  $\eta_{id}$  is allowed to vary flexibly across doctors; and we estimate a semiparametric model where  $\eta_{id}$  is once again assumed to be homoskedastic but now with an arbitrary distribution.

### 6.1 Stability of results to inclusion of alternate patient controls

In the spirit of Altonji et al. (2008), we explore the sensitivity of our results to the set of included variables to assess potential bias from unobservable risk factors. The rationale for this exercise is that omitting the variables  $x_{id}^{omit}$  from the baseline specification could generate heteroskedasticity, if the resulting error term  $\eta'_{id} = \eta_{id} + x_{id}^{omit}\beta$  is not i.i.d. across doctors and patients. If this heteroskedasticity substantially changes our estimates of the distribution of  $\tau_d$  or the degree of misweighting for the remaining variables, this might suggest that including additional unobserved variables would change our estimates further.

The model outlined above included four main classes of patient level risk factors: PE specific risk factors, chronic condition warehouse comorbidities, Elixhauser comorbidities, and patient demographic variables. Because some variation in comorbidities is required to appropriately identify

this model, we retain the PE specific risk factors and the chronic condition warehouse comorbidities throughout, and test the stability of our findings to excluding the Elixhauser comorbidity set and the vector of demographic variables.<sup>18</sup> Results from this exercise are reported in Table 5; the empirical Bayes correction has been applied before reporting the mean and standard deviation of physician’s testing thresholds.

The mean estimated value of physician’s testing thresholds ranges between 5.6% and 6.6%, and shows evidence of substantial dispersion in all models. The standard deviation of  $\tau_d$  ranges between 3.9% and 5.4%, depending on the set of included patient risk factors.

Dropping covariates does appear to increase the value of the estimated mean  $\tau_d$  although the range of values across specifications is only around 1/4 of the estimated across-doctor standard deviation. If including additional covariates would cause estimates of  $\tau_d$  to decrease, this suggests that our results may be conservative with respect to the amount of overtesting. Controlling for the full set of risk factors also appears to increase the variance in estimated testing thresholds, providing suggestive evidence that the observed variation in thresholds is not driven by the exclusion of unobserved risk factors from the model. In all of these cases, variation in testing thresholds is sufficient to imply large differences in testing probabilities for identical patients depending on which doctor they visit.

All specifications also predict substantial misweighting of included risk factors. The average absolute value of misweighting in physicians’ assessment of PE risk ranges from 0.020 to 0.023 percentage points. Perhaps unsurprisingly, the full model which includes all available risk factors as candidate sources of misweighting recovers the largest predicted amount of misweighting. In all cases, misweighting is sufficiently large that it has the potential to change testing decisions for many marginal patients.

In results reported in Appendix Table A.2, we find that the specific misweighted factors identified in Table 4 and discussed in section 5.3 continue to show evidence of misweighting of similar direction and magnitude, even as we vary the set of other included comorbidities. For example, the PE risk associated with recent hospital admissions and history of PE or deep vein thrombosis appears significantly underweighted in all specifications; black patients also show evidence of being under-tested in both specifications that include demographic variables. Similarly, a consistent set of conditions shows evidence of overweighting across specifications, including ischemic heart disease, chronic obstructive pulmonary disease and atrial fibrillation. These findings are not sensitive to the choice of other included covariates.

## 6.2 Estimation with physician-specific heteroskedasticity

Even if our results are not sensitive to dropping some covariates, we might worry that PE risk factors we cannot observe from insurance claims vary systematically across doctors. Differences across doctors in the variance of  $\eta_{id}$  could arise for at least three reasons. First, doctors may differ in their skill at assessing risk factors unobservable to the econometrician. A doctor with more diagnostic skill may have a higher variance in  $\eta_{id}$  across his patients, since he is more discerning in

---

<sup>18</sup>Recall that we rely on comorbidities to identify the marginal tested patients, and then calculate test outcomes among that group for high-volume doctors to implement our instrumental variables strategy.

his judgement of which patients should be tested on the basis of clinical presentation and symptoms. Second, doctors may differ in the variance of latent PE risk present in their patient population. A doctor with a more heterogeneous patient population may have a higher variance in  $\eta_{id}$  across his patients. Finally, doctors may simply make “errors” that lead them to deviate from typical practice patterns; a doctor who frequently deviates from his peers’ practice patterns in assessing PE risk may have a higher variance in  $\eta_{id}$ . The model we develop in this section allows us to isolate differences in physician testing thresholds that are unrelated to possible differences in the variance of  $\eta_{id}$  across physicians.

Recall the assumption we made in Section 4.4 that  $\eta_{id}$  followed a mixture of a Bernoulli and uniform distribution. We maintain the basic shape of the distribution but now allow both the Bernoulli probability and the variance of the uniform distribution to vary across doctors, so that  $\eta_{id} \sim U(-\eta_d, \eta_d)$  with probability  $1 - p_d$  and  $\eta_{id} \sim U[v - \eta_d, v + \eta_d]$  with probability  $p_d$ .

Following the derivation in Appendix E, the more flexible distributional assumption implies the testing equation takes this form:

$$Pr(Test_{id} = 1) = \max \left\{ 0, \frac{p_d}{2} + \frac{p_d(I'_{id} + v)}{2\eta_d} \right\} \quad (15)$$

From the testing equation above, we can see that heteroskedasticity in  $\eta_{id}$  is identified by the fact that observables are less predictive of testing behavior for doctors with a high variance in  $\eta_{id}$ , i.e. a smaller value of  $\frac{p_d}{\eta_d}$ . As described in the appendix, the testing equation can be used to estimate  $C \frac{p_d}{2\eta_d}$ , where  $C$  is an unknown scaling constant. For computational tractability given the demands of this more flexible estimation strategy, we randomly exclude half of the physicians from our sample to reduce sample size, and drop the Elixhauser comorbidities and demographic risk factors from our list of included covariates.

With the introduction of heteroskedasticity, the conditional probability of a positive test is given by:

$$E(q_{id}|Test_{id} = 1) = \tau_d + \frac{C \tilde{I}'_{id}}{2 \hat{\eta}_d} + (x_{id} - E_d(x_{id}))(\beta - \beta') \quad (16)$$

where  $\hat{\eta}_d = C \frac{p_d}{2\eta_d}$  are the variances estimated in the testing equation. Further details of the estimation strategy are provided in Appendix E.

Table 5 reports the results of this analysis in panel 4, which can be compared to results from the baseline model with the same excluded comorbidity set, as reported in panel 3. The mean value of  $\tau_d$  is 7.0% in the model allowing for heteroskedasticity compared to 6.6% in the baseline model with the same covariates; allowing for heteroskedasticity slightly raises our estimate of the average testing threshold. Estimates of the standard deviation of  $\tau_d$  are also higher at 5.1 percentage points in the heteroskedastic model compared to 3.9 percentage points in the homoskedastic model. Thus, the cross-physician variation in testing behavior is not explained by differences in the variance of  $\eta_{id}$  across doctors. This provides reassuring evidence that the assumption of homoskedasticity in the baseline model was not leading us to overstate differences across physicians in testing thresholds. Finally, the degree of misweighting remains very similar to the original estimates, with the average

absolute value of misweighting estimated at 0.021 in the heteroskedastic model compared to 0.020 in the baseline model.

As described earlier, one potential driver of heteroskedasticity across physicians could be “mistakes” physicians make that lead them to deviate from typical practice patterns in ways that do not improve patient outcomes. This idea of physician diagnostic errors is central to the model of cesarean section births studied by Currie and MacLeod (2013). Currie and MacLeod (2013) begin with a normative model of the returns to performing a cesarean section and argue that physicians who deviate more from the predicted optimal treatment choice according to that model have worse diagnostic skill—i.e. the best doctors respond only to the index of observables. They corroborate this interpretation by demonstrating that outcomes are indeed worse for physicians who deviate more frequently from the model’s recommendations.

In contrast, our model suggests that physicians do have private information about PE risk. The scaling factor  $C$  is positive, which suggests that at least in this context, physicians with higher variance  $\eta_{id}$  are not making random mistakes when they deviate from the predicted testing behavior—rather, they are selecting patients on unobservables in order to increase their test yields.

The other two potential drivers of heteroskedasticity identified previously—differences in physician ability to assess PE risk based on clinical symptoms or differences in the latent distribution of population PE risk—remain as potential explanations for the observed heteroskedasticity. We cannot directly distinguish these hypotheses.

The role of physician diagnostic judgment in driving testing behavior and outcomes was previously explored by Doyle, Ewer, and Wagner (2010). In a natural experiment, they find that physicians from more prestigious residency programs achieve similar patient outcomes at 10-25% lower cost compared to their less skilled peers. One potential explanation for this phenomenon is that physicians from less prestigious schools prefer to administer more low-value care and could achieve the same outcomes at lower cost if they cut back some services. In the language of our model, these less skilled physicians might have lower testing thresholds, i.e. smaller  $\tau_d$ . A second explanation is that these less skilled physicians just need to use more medical resources to achieve the same quality of care, because they are less accurate in their assessments of ex ante patient risk. In the language of our model, this decreased diagnostic accuracy would correspond to a lower variance of  $\eta_{id}$ , since these less skilled physicians would be failing to incorporate clinical information about patient risk to improve test targeting. Our results suggest that the heterogeneity in measured  $\tau_d$  across physicians persists even after allowing for heterogeneous variance of  $\eta_{id}$  across doctors. This finding raises the possibility that cost variance across physicians is driven in part by lower marginal value services provided by doctors with lower expected benefit thresholds.

### 6.3 Estimation of a semiparametric selection model

Next we test whether our results are sensitive to the shape of the distribution assumed for the unobserved component of their PE risk,  $\eta_{id}$ . We previously imposed a strict distributional assumption, requiring  $\eta_{id}$  to be distributed according to a mixture of Bernoulli and Uniform distributions. Now, we relax this assumption by estimating Equation 11 as a semiparametric binary choice model, using the Klein and Spady (1993) binary choice estimator. This robustness exercise will ensure

that differences in testing thresholds observed in the previous sections are not driven solely by the strong distributional assumptions which restricted the functional form of the testing equation and the shape of the selection correction function  $\lambda(\cdot)$ . To implement the semiparametric model, we return to our original, strong version of the ignorability assumption that  $\eta_{id}$  is i.i.d. across physicians and patients.

Estimation of the semiparametric model proceeds as follows. Let  $g$  denote the probability that patient  $i$  is tested given index  $I'_{id} = x_{id}\beta' + \theta'_d$ . The log likelihood is given by:

$$L(\beta, g) = \sum_i [Test_{id} \ln g(x_{id}\beta' + \theta'_d) + (1 - Test_{id})(1 - \ln g(x_{id}\beta' + \theta'_d))] \quad (17)$$

The idea of the Klein-Spady estimator is to approximate  $g$  using a “leave-one-out” estimator which predicts the probability of testing for a particular patient, giving more weight to patients with nearby indices  $I'_{id}$ . Specifically, we substitute for  $g$  using the following function:

$$\hat{g}_{-i,d} = \frac{\sum_{j \neq i} k\left(\frac{I'_{jd} - I'_{id}}{h}\right) Test_j}{\sum_{j \neq i} k\left(\frac{I'_{jd} - I'_{id}}{h}\right)} \quad (18)$$

We use a 4th-order Gaussian Kernel,  $k(\cdot)$ , and empirically select for the smallest bandwidth  $h$  such that  $\hat{g}$  is a monotonic function of the index  $I'_{id}$ .

Given the propensity to test index  $I'_{id}$  from estimating equation 11 by the Klein-Spady procedure, the next step is to estimate the testing outcome equation. Echoing the derivation in Section 4.2, the probability of a positive test among tested patients is given by:

$$E(Z_{id}|Test_{id} = 1) = \tau_d + (x_{id} - E_d(x_{id}|Test_{id} = 1))(\beta - \beta') + \lambda(I'_{id}) \quad (19)$$

where  $\lambda(I'_{id}) = I'_{id} + h(I'_{id})$ . Because we no longer assume a particular distribution of  $\eta_{id}$ , we now fit the function  $\lambda(\cdot)$  flexibly, reporting results with  $\lambda(\cdot)$  as a linear function and as a cubic polynomial, and estimate the net benefit equation by OLS.

Note that the Klein-spady estimator only recovers  $I'_{id}$  up to a location and scale normalization. The scale normalization is embedded in the function  $\lambda(\cdot)$ . We impose the appropriate location normalization so that at the smallest value of  $I'_{id}$  among tested patients,  $\underline{I}$ , we have  $\lambda(\underline{I}) = 0$  as shown in Section 4.3.<sup>19</sup>

Results of the semiparametric estimation are reported in Table 5, panels 5 and 6. This semiparametric estimation approach estimates the mean value of  $\tau_d$  at 6.7% (linear) or 6.6% (cubic), similar to the parametric model estimate of 6.6% in the sample with identical comorbidities. We continue to find a large amount of cross-doctor dispersion in estimated testing thresholds. The standard deviation of  $\tau_d$  is 5.4% across doctors, compared to 3.9% in the parametric model with the same covariates (but interestingly nearly identical to the parametric model with the full set of covariates included). Our assessment of misweighting continues to be highly consistent across models, with

---

<sup>19</sup>This normalization can be implemented by omitting the constant term from the polynomial  $\lambda(\cdot)$  and subtracting a constant  $\underline{I}$  from  $\hat{I}'_{id}$ ; thus the resulting polynomial  $\lambda(I'_{id} - \underline{I})$  will equal 0 for  $I'_{id} = \underline{I}$ . To avoid sensitivity to outliers, we normalize  $I'_{id}$  so that  $\lambda(\underline{I}) = 0$  for  $I'_{id}$  in the 10th percentile amongst tested patients, which agrees with our definition of marginal patients in Section 4.3.



an average absolute value of the error due to misweighting at 2.1% in the semiparametric model, compared to 2.0% in the parametric model.

Taken together, these robustness checks, including varying the set of included covariates, allowing for physician-specific heteroskedasticity, and estimating a semiparametric selection model, all suggest that our findings on the dispersion in testing thresholds and amount of misweighting are very stable across alternative modeling assumptions. We find substantial variance in testing thresholds of similar magnitude in all specifications, suggesting that much of the observed variation in testing behavior may be driven by differences in practice styles. Further, doctors are mis-assessing patient PE risk by similar amounts in percentage point terms across all models.

## 7 Welfare cost of overtesting and misweighting

We now turn to the welfare implications of the models estimated in the previous sections. In order to assess the welfare cost of overtesting and misweighting, we will need to make additional assumptions about the costs of testing and the dollar-equivalent benefits of detecting and treating a PE. Given these assumptions, we can evaluate whether the observed variation in testing thresholds reflects overuse and compare the welfare cost of overuse to the welfare cost of misweighting. Applying the structure and estimates of our baseline estimation procedure, we perform simulations to determine how welfare would change if doctors behaved optimally from a social standpoint. We begin by simulating worlds with no overtesting but maintaining the observed patterns of misweighting; next, we simulate a world with no misweighting but maintain the observed distribution of testing thresholds. In each case, we decompose the sources of estimated welfare gains into financial costs, medical costs and medical benefits.

This section proceeds first by describing the calibration of the optimal testing threshold  $\tau^*$ , then exploring the welfare implications of the measured variation in physician testing thresholds, and finally estimating the welfare costs of misweighting the PE risk associated with patient comorbidities. All of the calibrations in this section are implemented in our baseline model as outlined and reported in Sections 4 and 5.

### 7.1 Calibration of parameters

In order to proceed with welfare calculations, we make several additional assumptions about the costs of testing and the benefits of a positive test. We assess these costs and benefits from a social standpoint; e.g. if some physicians test more due to reimbursement incentives, this would appear in our model as measured heterogeneity in  $\tau_d$  that deviates from the social optimum we compute below.

If physicians are behaving optimally, they should test a patient if and only if:  $NUq_{id} - c > 0$  where  $NU$  represents the net utility of detecting a positive test,  $c$  represents the cost of the test and as above,  $q_{id}$  denotes the likelihood of a positive test. This yields a socially optimal testing threshold  $\tau^* = \frac{c}{NU}$  such that physicians should test only if  $q_{id} > \tau^*$ .

If there were no false positive or false negative tests, the net utility would correspond to the net medical benefits of treating PE minus any financial costs of treatment. However, CT scans, like



many other medical tests, can generate both false positive and false negative results (Stein et al. 2006). It turns out that an important cost of overtesting is a consequence of type I and type II errors: overtesting leads to unneeded treatment which can have adverse consequences. Patients with false positive test results receive medical treatment as if they truly had a PE; this treatment will incur medical risks and financial costs without conferring any medical benefit on the patient, since they do not truly have the condition being treated.

Let  $fp$  denote the likelihood of a false positive,  $s$  the sensitivity of the test (one minus the probability of a false negative),  $MB$  the medical benefits of treating a PE,  $MC$  the medical costs and  $CT$  the financial costs of treatment. In Appendix H, we show that allowing for false positives and false negatives results in a model which is isomorphic to the one above with  $NU$  replaced by  $\hat{NU} = \frac{s}{s-fp}MB - MC - CT$  and  $c$  replaced by  $\hat{c} = c + \frac{s \cdot fp}{s-fp}MB$ .

Table 6 reports the values of the parameters that we use to compute  $\tau^* = \frac{\hat{c}}{\hat{NU}}$ . Parameters specifying test sensitivity and specificity, the medical benefits of testing, and the medical costs of testing are drawn from the existing medical literature. Note that our calibration of both the medical benefits and the medical cost of treatment depend on an estimate of the value of a statistical life (VSL); following Murphy and Topel (2006) we assume a VSL of \$1 million.<sup>20</sup> We estimate the financial cost of testing and the financial cost of PE treatment directly from our Medicare claims data. Appendix Table A.4, which we discuss below, explores the sensitivity of our welfare findings to these calibration parameters.

One parameter of this calibration turns out to be of particular importance and remains a source of uncertainty in the medical literature: the rate of false positive tests. To our knowledge, the single piece of medical evidence on chest CT scans' false positive rate derives from a comparison of CT imaging results to older diagnostic methods, VQ scanning and ultrasonography; the authors estimate the false positive rate at 4% (Stein et al. 2006). We report results with a false positive rate of 4% as our preferred welfare calibration, but also show the welfare implications of assuming a 3% or 0% false positive rate. Lower false positive rates boost the net utility associated with treating a positive test, and thus provide more conservative estimates of the costs of overtesting.

Table 7 reports the optimal testing threshold  $\tau^*$  under these calibration assumptions. With a false positive rate of 4%, we find physicians should optimally test all patients with an ex ante likelihood of a positive test greater than or equal to 6.2%. The optimal threshold decreases to 5.0% at a false positive rate of 3%; at the (unlikely) extreme of no false positive test results, the optimal threshold falls to 1.5%.

## 7.2 Welfare impact of eliminating overtesting

The model implies welfare loss whenever a physician's testing threshold  $\tau_d$  does not equal the optimal value  $\tau^*$ . We focus on the welfare consequences of overtesting, where  $\tau_d$  is below this calibrated optimum, for two reasons. First, overtesting is empirically the larger problem in our sample, with an estimated 84% of doctors overtesting under our preferred calibration assumptions. Second, unlike the overtesting case, we find that the welfare loss due to under-testing is highly dependent on the

<sup>20</sup>The choice of a lower VSL estimate in this context is driven by the fact that we are studying an elderly population, with an average age of around 77.

distribution we assume for  $\tau_d$  when applying an empirical Bayes technique to recover the posterior distribution of  $\tau_d$ . Previously, we were agnostic about the distribution of  $\tau_d$  and recovered only the posterior mean and variance, but for welfare calculations, a specific distributional assumption is required. For some distributions of  $\tau_d$ , even a small number of doctors under-testing can lead to large welfare losses if the right tail of the  $\tau_d$  distribution is sufficiently thick.

To determine the percentage of doctors overtesting we need to extend our empirical Bayes analysis to recover a posterior estimate of  $\tau_d$  for each physician; proceeding requires an assumption about the shape of the underlying  $\tau_d$  distribution. First, note that  $\tau_d$  is bounded below at the false positive rate. We assume that  $\tau_d$  minus the false positive rate is log-normally distributed with the posterior mean and variance of the  $\tau_d$  distribution as previously calculated. Table 7 reports the percentage of doctors overtesting at each false positive rate, given this distributional assumption.

Our initial estimates of  $\tau_d$  are in units of the probability of a positive test. For example, in our baseline specification, we find that the average doctor tests a patient if the probability of a positive test exceeds 5.6%. We want to know: how would testing behavior change for each physician if all physicians with testing thresholds below  $\tau^* = 6.2\%$  instead adopted a threshold of 6.2%? If we observed  $q_{id}$  for each patient, this would be a simple matter of counting the number of inframarginal patients. But  $q_{id}$  is not observed—instead, we know the probability of a positive test as a function of the propensity to test. Our model allows us to determine how changes in  $\tau_d$  impact the propensity to test using the scaling factor  $\frac{q}{p}$ , the estimated coefficient on the selection term in equation 14. Equation 14 also allows us to compute how the probability of a positive test conditional on testing changes for each observation. More details are provided in Appendix H.

Combined with our assumptions about costs and net utility, we compute separately the realized medical benefits of testing, the medical costs of testing, the financial costs of testing and the net benefits of testing given the estimated  $\hat{\tau}_d$  as well as a counterfactual where  $\tau_d = \tau^*$  for all doctors with  $\hat{\tau}_d < \tau^*$ . These results are shown in Table 7, under a series of different assumptions about the false positive rate.

At a false positive rate of 4% (the estimate in the medical literature), we estimate that 84% of the physicians in our sample are overtesting on the margin, i.e. they apply a testing threshold that is lower than the 6.2% threshold probability of a positive test the calibration suggests is optimal. At a false positive rate of 3%, the proportion of doctors overtesting falls to 67.2%. To illustrate the importance of the false positive rate in assessing welfare, note that if there were *no* false positive tests, the optimal testing threshold  $\tau^*$  drops substantially to 1.5% and only 10% of physicians are overtesting on the margin, i.e. have a testing threshold lower than 1.5%.

At a false positive rate of 3% or 4%, eliminating overtesting would decrease the total number of patients tested by more than 30% or 50%, respectively. Why such large effects? Recall that with a false positive rate of 4%, the minimum possible perceived probability ( $q_{id}$ ) of a positive test is 4%. The median physician in our sample has a  $\tau_d$  which is less than .05 (much less than the mean, since the distribution is bounded from below by .04). Increasing  $\tau_d$  to 0.062 thus greatly increases the range of probabilities  $q_{id}$  which would not be tested for many physicians.

In these scenarios, the financial and medical costs of testing would fall by an amount proportional to the decline in tested patients. There would be a small offsetting decline in the medical benefits

of testing because the patients not tested in the counterfactual world have a very low probability of truly having a PE. Eliminating overttesting leads to a 12.5% increase in net benefits at a false positive rate of 3% and a more than 60% increase in net benefits at a false positive rate of 4%; the increase in net benefits per test is of course much larger. This exercise illustrates both the large welfare implications of overuse of medical testing and the sensitivity of this result to the false positive rate. As detailed in Table 7, most of the net benefit increase comes from eliminating the financial costs associated with testing low-probability patients for PE and unneeded treatment of patients with false positive test results.

Given the widespread incidence of overttesting under our preferred calibration, it is worth considering a few possible explanations. As we illustrate in Table 7, the estimated overttesting behavior of a majority of doctors in our sample could be explained if they were behaving as if there were no false positive test results. Similarly, if physicians ignored the financial costs associated with testing and treating PE, this could also explain much of the overttesting behavior. However, the only way to rationalize the entire estimated posterior distribution of physician testing patterns would be to allow physicians to vary substantially in their assessment of financial costs or the false positive rate.

One could also interpret variation in  $\tau_d$  as variation in the patients’ “value of knowing” that they do not have a PE. In contrast to the case of Huntington’s disease (Oster, Shoulson, and Dorsey 2011), the value of knowing seems an unlikely driver of testing decisions in this context, since in most cases a PE has a very low ex ante probability and the rate of false negatives is sufficiently high that even after testing one has only somewhat reduced that probability. Further, Finkelstein et al. (2014) find that variation in patient demand (i.e. both patient preferences and medical needs) explains only 14% of the regional variation in spending on imaging, suggesting a very limited role for patient preferences in explaining variation in imaging decisions.

Finally, the socially optimal testing threshold depends on the cost of scanning a patient, which we estimate directly from the Medicare claims data. The \$300 financial cost of testing is calculated based on the allowed charges which compensate for the technician’s time to run the scan, the radiologist’s time to interpret the scan and capital depreciation. If some of this reimbursement is intended as compensation for the high fixed costs of owning a CT scanner, then we may be overstating the social cost of testing. We believe this concern is mitigated by calculating costs directly from the Medicare data, where reimbursement for CT scans remains much below the estimated fees paid by insured consumers (cf. Healthcare Blue Book which estimates the typical fee at \$517 to \$577 depending on the precise billing code). In addition, there may be opportunity costs of scanning a patient not accounted for in our calibration if the hospital is capacity constrained in its allocation of time in the CT scanner or time spent awaiting a scan in an ED bed. If present, opportunity costs would lead us to understate the true costs of performing a scan, and thus understate the amount of overttesting in our data.

Panel A of Table A.4 explores how our results on the net welfare cost of overttesting vary with the calibrated parameters. The results do not vary much with the calibration of test sensitivity. Changing either the VSL or the cost of the test shifts the optimal testing threshold  $\tau^*$  and thus the welfare benefits. For example, with a VSL of \$500,000 rather than \$1 million, the optimal threshold increases from 6.2% to 14.3%. Due to this dramatic increase in  $\tau^*$ , simulations with no physicians

overtesting involve more dramatic declines in the fraction of patients tested, and the net benefits of eliminating overtesting almost double vis-a-vis the baseline calibration results. If the VSL is \$1.5 million rather than \$1 million, the number of patients tested in a world with no overtesting increases by 50%, and the net benefits of eliminating overtesting likewise fall. Similarly, if the cost of the test is \$0 (i.e. if there is zero marginal social cost of running a CT scan), the optimal threshold  $\tau^*$  falls to 4.8%, there is substantially less overtesting and the overtesting that does occur has much lower social cost (only the costs from overtreatment of false positive tests). If the cost of the test is \$500 (comparable to the fees paid to private insurers per CT scan) rather than \$300, the net benefits of eliminating overtesting almost double.

### 7.3 Welfare impact of eliminating misweighting of patient risk factors

Table 8 reports results from a simulation in which doctors select patients for testing by weighting observable comorbidities in the manner the model suggests would maximize detection of positive tests. In other words, we simulate physician behavior if they were to use the true weights  $\beta$  rather than the observed weights  $\beta'$  to assess PE risk. In this simulation, we maintain the distribution of physician testing thresholds at their baseline values, so we allow for the observed patterns of under- and overtesting. We report results at our preferred calibration of the false positive rate, 4%; the welfare consequences of eliminating misweighting would be even larger at lower false positive rates.

Structurally, this exercise is very similar to the exercise where we simulate alternative values of  $\tau_d$ . Our initial estimates tell us the degree of misweighting in units of the probability of a positive test. We want to determine how the propensity to test would differ if physicians did not misweight; the scaling factor  $\frac{\eta}{p}$  allows us to translate the estimated degree of misweighting into the same units as the testing propensity and calculate the testing propensity and expected test outcomes if there were no misweighting. We demonstrate this explicitly in Appendix H.

We find that properly weighting observables to improve PE detection would lead the fraction of patients tested to increase from 3.8% to 4.3%, by moving some patients just over their estimated physician’s testing threshold. But by far the predominant welfare impact comes from the predicted increase in the rate of PE detection. The medical benefits due to treatment of PE nearly double and the net benefits of testing more than triple. The total welfare loss from misweighting (\$35.9 million in our sample) is more than 4 times as large as the welfare loss from overtesting (\$8.1 million) even in the model with the highest rate of false positives.

To investigate whether a small number of risk factors account for most of the observed costs of misweighting, we conduct an exercise where we correct the weights applied to each variable, one at a time. Results from this exercise with more detailed notes are reported in Appendix Table A.3. First, it is worth noting that in this simulated second-best world where physicians do not all share the optimal testing threshold  $\tau^*$  and where other factors are misweighted, correcting misweighting of a single risk factor in isolation can sometimes worsen total welfare; certain misweighting errors offset some of the costs associated with overtesting. However, in most cases, correcting a single variable’s weight weakly improves estimated welfare.

Correcting the weighting on 30-day inpatient admissions accounts for approximately 20% of the total potential gains from eliminating misweighting. Expanding the list to include the 5 highest-

impact covariates (30-day admission history, 1-week admission history, 1-year surgical history, chronic obstructive pulmonary disease, and ischemic heart disease) accounts for roughly 60% of the total potential gains. These covariates are both substantially misweighted and common enough to induce large welfare consequences.

Intuitively, given our estimates of misweighting in Section 5.3, it is not surprising that the welfare loss from misweighting substantially exceeds the welfare losses from overtesting. Several factors combine to make misweighting a more serious problem. Physicians behave as if they are misestimating a patient’s PE risk by 2.3 percentage points on average by failing to weight observable characteristics to maximize detection of positive tests. By comparison, the average difference between  $\tau_d$  and  $\tau^*$  for physicians who are overtesting is only 1.7 percentage points in the calibration with a false positive rate of 4%. The welfare cost of misweighting errors or suboptimal values of  $\tau_d$  increases with the square of the deviation—as the bias grows, both the number of patients impacted and the average severity of the error among those patients increases. Further, the welfare costs of overtesting are bounded. The worst outcome of overtesting is that a patient is tested with no chance of having a PE and incurs the cost of the test (a few hundred dollars) plus the potential financial costs and medical risk of treatment if they receive a false positive test result. The potential costs of misweighting are substantially greater since you might fail to treat a patient with a substantial risk of death.

Panel B of Table A.4 explores how our results on the net welfare cost of misweighting vary with the calibrated parameters. The positive impact of misweighting on testing behavior does not depend on the calibration (unlike the case of overtesting, since the calibration determines which physicians overtest). The welfare cost of misweighting is not too sensitive to the false positive rate, the sensitivity of the test or the cost of the test, but it is sensitive to the VSL. Misweighting creates more welfare loss from undertesting than overtesting: the welfare costs of overtesting are bounded by the financial costs of the test plus the costs of treating false positive test results, while the costs of undertesting in the worst case is the 2.5% chance of mortality from a missed PE. These latter costs are roughly proportional to the VSL.

Undiagnosed PE is thought to be a major public health problem, with the Office of the Surgeon General (2008) estimating that approximately half of PE cases are never diagnosed; analysis of autopsy reports have found it to be a frequently missed mortality risk. By improving physician assessment of patient PE risk, our model suggests that the rate of undiagnosed PE could fall substantially. Although there is policy attention in the medical community on the risks associated with the perceived overuse of PE CT, this evidence suggests that there may be even larger gains possible from improving the targeting of CT scans.

## 7.4 National scale of welfare estimates

Our welfare calculations are based on a 20% sample of patients enrolled in Medicare Parts A and B over a 10-year period, and the numbers reported in Tables 7 and 8 reflect potential gains to this sample only. To understand the annual welfare loss for Medicare patients associated with the inefficiencies we identify in this sample, we do an informal scaling exercise. We first scale the estimates up by a factor of 5 to account for the entire population of Medicare fee for service enrollees, then adjust to account for the 28% of Medicare patients who enroll in a Medicare Advantage plan,

and finally divide by 10 to calculate annual estimates. We recover a \$5.5 million annual welfare loss from overuse of PE CT due to low testing thresholds, and a \$25 million annual loss from misweighting observable patient risk factors, for emergency department CT scans among elderly patients. More speculatively, if we further scale the number of diagnosed PEs in our sample to represent total national incidence of 350,000 PE cases per year (Office of the Surgeon General 2008), we estimate \$560 million in annual welfare loss from overtesting and \$2.5 billion in annual welfare loss from misweighting. This final scaling requires extrapolating our results to the many PEs diagnosed in settings other than the emergency department and among the non-elderly population. Yet even these final scaled welfare gains from the efficient application of PE CT may represent only a small fraction of the total welfare benefit available from more efficient diagnostic testing and treatment decisions across a variety of medical conditions.

## 8 Conclusion

While it is commonly believed that the US health care system spends significant resources on services that have low medical returns and high costs, there is little consensus on how this waste could be reduced. Wasteful spending is characterized both by overuse of medical care (allocative inefficiency) and mistargeting of medical resources (productive inefficiency). This paper investigates both forms of inefficiency, analyzing whether doctors efficiently select patients for medical testing and how physicians vary in the risk thresholds at which they test patients. We study these inefficiencies in the context of emergency department CT scans to diagnose pulmonary embolism (PE). We document both widespread variation in physician use of CT scans for PE unexplained by differences in patient risk, and also systemic failure to target medical testing to the highest risk patients.

Estimating the model to study physicians' CT scanning decisions in a national sample of Medicare claims, we find substantial variation in physician's use of diagnostic scans on low-risk patients. This variation generates a negative relationship between testing propensities and test yield across physicians, since physicians who test more also test lower risk patients on average. Investigating the role of training and practice environment in explaining practice styles, we find that physicians practicing in high-spending Dartmouth Atlas regions and those with less experience are more likely to scan low-risk patients. Other factors, such as hospital ownership or quality of medical school training are not significantly related to testing behavior. Taken as a whole, observable characteristics can explain only a small fraction of the total variation in testing thresholds. Applying further calibration assumptions suggests that 84% of physicians in our sample are overtesting on the margin by applying a risk threshold that is lower than the calibrated optimum.

We also find that doctors do not weight observable patient risk factors in a way that would maximize test yields. Physicians systematically underweight certain important predictors of PE risk, including recent prior hospitalizations and metastatic cancer. These apparent errors occur despite the fact that physicians are widely encouraged to use diagnostic scoring systems such as the Wells or Geneva score to assess the risk of PE before deciding whether to order a CT scan. The continued prevalence of risk assessment mistakes despite the popularity of these PE risk scoring systems may reflect shortcomings in the scoring systems themselves or failures to make adequate



use of these scores. (The data used in this project cannot disentangle these possibilities.) Other preexisting conditions that have similar clinical symptoms to PE are over-weighted in the testing decision. Together, these mistakes in assessing patient PE risk lead to significant welfare losses from failing to target the test to the highest risk patients according to our welfare simulations.

The model developed in this paper could be applied to a variety of empirical contexts—it is applicable whenever economic actors make repeated decisions about whom to “treat, ”as long as we observe outcomes for “treated” cases and can assume the actor is applying the same decision threshold in each case. For example, the model could be used to evaluate the decisions of loan officers to extend credit, hiring directors who select among potential job applicants, admissions officers attempting to predict which students will perform most highly, or juvenile court judges who must assess which children will benefit from detention. Positively, one could investigate the degree to which observed heterogeneity in treatment rates is due to decision-maker discretion. Normatively, many of these organizations have specific objectives they seek to optimize (e.g. reducing default on loans or recidivism among parolees) and one could use the model developed here to investigate whether observed selection patterns are successfully optimizing these outcomes.

Our findings suggest that both overuse and misuse of medical resources are important drivers of high spending and low medical returns to care. By measuring physician-level preferences for testing, we are able to explore the training and environmental factors that contribute to overuse. Future work could pair this framework for estimating the overuse of diagnostic testing with experimental or quasi-experimental variation in physician’s training or practice environment; together, these estimates could more directly inform policy by causally identifying how these changes to a physician’s education or training affect the efficiency of the medical care delivered. Given more detailed patient-level data, our model could be used to formulate optimal guidelines and risk scores, overcoming the selection problems that may lead to biased estimates of risk under popular existing methodologies. Our findings underscore the fact that purely cost-focused health reform may be insufficient to achieve efficiency in healthcare delivery—there are potentially large benefits to patients from physicians making better use of the available information to target medical resources to those patients with the highest returns.

## References

- Altonji, J. G., T. E. Elder, and C. R. Taber (2008). Using selection on observed variables to assess bias from unobservables when evaluating swan-ganz catheterization. *The American Economic Review* 98(2), pp. 345–350.
- Avraham, R. (2011). Database of state tort law reforms (dstlr 4th). *U of Texas Law, Law and Econ Research Paper* (184).
- Avraham, R., L. S. Dafny, and M. M. Schanzenbach (2012). The impact of tort reform on employer-sponsored health insurance premiums. *Journal of Law, Economics, and Organization* 28(4), 657–686.
- Chandra, A. and D. Staiger (2011). Expertise, Overuse and Underuse in Healthcare. *Working Paper*.



- Chandra, A. and D. O. Staiger (2010, September). Identifying provider prejudice in healthcare. Working Paper 16382, National Bureau of Economic Research.
- Coco, A. S. and D. T. O’Gurek (2012, January-February). Increased emergency department computed tomography use for common chest symptoms without clear patient benefits. *Journal of the American Board of Family Medicine* 25(1), 33–41.
- Costantino, M. M., G. Randall, M. Gosselin, M. Brandt, K. Spinning, and C. D. Vegas (2008, August). Ct angiography in the evaluation of acute pulmonary embolus. *American Journal of Roentgenology* 191(2), 471–474.
- Currie, J. and W. B. MacLeod (2006). First do no harm?: Tort reform and birth outcomes. Technical report, National Bureau of Economic Research.
- Currie, J. and W. B. MacLeod (2013). Diagnosis and unnecessary procedure use: Evidence from c-section. Technical report, National Bureau of Economic Research.
- David, S., P. Beddy, J. Babar, and A. Devaraj (2012, Feb). Evolution of ct pulmonary angiography: referral patterns and diagnostic yield in 2009 compared with 2006. *Acta Radiologica* 53(1), 36–43.
- Doyle, J. J., S. M. Ewer, and T. H. Wagner (2010). Returns to physician human capital: Evidence from patients randomized to physician teams. *Journal of health economics* 29(6), 866–882.
- Elixhauser, A., C. Steiner, D. Harris, and R. Coffey (1998). Comorbidity measures for use with administrative data. *Medical Care* 36(1), 8–27.
- Finkelstein, A., M. Gentzkow, and H. Williams (2014). Sources of geographic variation in health care: Evidence from patient migration. Technical report, National Bureau of Economic Research.
- Garber, A. M. and J. Skinner (2008). Is american health care uniquely inefficient? Technical report, National Bureau of Economic Research.
- Goldhaber, S. Z. and H. Bounameaux (2012). Pulmonary embolism and deep vein thrombosis. *The Lancet* 379(9828), 1835–1846.
- Heckman, J. and T. MaCurdy (1980). A life cycle model of female labour supply. *The Review of Economic Studies* 47(1), 47–74.
- Heckman, J. J. (1979). Sample selection bias as a specification error. *Econometrica* 47(1), pp. 153–161.
- Jackson, C. K., J. E. Rockoff, and D. O. Staiger (2014). Teacher effects and teacher-related policies. *Annual Review of Economics* 6(1), 801–825.
- Kane, T. J. and D. O. Staiger (2008). Estimating teacher impacts on student achievement: An experimental evaluation. Technical report, National Bureau of Economic Research.
- Klein, R. and R. Spady (1993). An efficient semiparametric estimator for binary response models. *Econometrica: Journal of the Econometric Society*, 387–421.

- Lessler, A. L., J. A. Isserman, R. Agarwal, H. I. Palevsky, and J. M. Pines (2010, April). Testing low-risk patients for suspected pulmonary embolism: A decision analysis. *Annals of Emergency Medicine* 55(4), 316–326.
- Lewis, J. B. and D. A. Linzer (2005). Estimating regression models in which the dependent variable is based on estimates. *Political Analysis* 13(4), 345–364.
- Mamlouk, M. D., E. vanSonnenberg, R. Gosalia, D. Drachman, D. Gridley, J. G. Zamora, G. Casola, and S. Ornstein (2010, August). Pulmonary embolism at ct angiography: Implications for appropriateness, cost, and radiation exposure in 2003 patients. *Radiology* 256, 625–632.
- Meszaros, I., J. Morocz, J. Szlavi, J. Schmidt, L. Tornoci, L. Nagy, and L. Szep (2000, May). Epidemiology and clinicopathology of aortic dissection. *Chest* 117(5), 1271–1278.
- Molitor, D. (2012). The evolution of physician practice styles evidence from cardiologist migration. Technical report, MIT working paper.
- Mulligan, C. B. and Y. Rubinstein (2008). Selection, investment, and women’s relative wages over time. *The Quarterly Journal of Economics* 123(3), 1061–1110.
- Murphy, K. M. and R. H. Topel (2006). The value of health and longevity. *Journal of Political Economy* 114(5), 871–904.
- Nelson, A. (2002). Unequal treatment: confronting racial and ethnic disparities in health care. *Journal of the National Medical Association* 94(8), 666.
- Office of the Surgeon General (2008). The surgeon general’s call to action to prevent deep vein thrombosis and pulmonary embolism.
- Oster, E., I. Shoulson, and E. Dorsey (2011). Optimal expectations and limited medical testing: evidence from huntington disease. Technical report, National Bureau of Economic Research.
- Rahimtoola, A. and J. D. Bergin (2005, February). Acute pulmonary embolism: An update on diagnosis and management. *Current Problems in Cardiology* 30, 61–114.
- Stein, P. D., S. E. Fowler, L. R. Goodman, A. Gottschalk, C. A. Hales, R. D. Hull, J. Kenneth V. Leeper, J. John Popovich, D. A. Quinn, T. A. Sos, H. D. Sostman, V. F. Tapson, T. W. Wakefield, J. G. Weg, and P. K. Woodard (2006, June 1). Multidetector computed tomography for acute pulmonary embolism. *New England Journal of Medicine* 354(22), 2317–27.
- Venkatesh, A., J. A. Kline, and C. Kabrhel (2013, Jan. 28). Computed tomography in the emergency department setting—reply. *Journal of the American Medical Association Internal Medicine* 173(2), 167–168.
- Venkatesh, A. K., J. A. Kline, D. M. Courtney, C. A. C. Jr, M. C. Plewa, K. E. Nordenholz, C. L. Moore, P. B. Richman, H. A. Smithline, D. M. Beam, and C. Kabrhel (2012, July 9). Evaluation of pulmonary embolism in the emergency department and consistency with a national quality measure: Quantifying the opportunity for improvement. *Archives of Internal Medicine* 172(13), 1028–1032.
- Wells, P. S., D. R. Anderson, M. Rodger, J. S. Ginsberg, C. Kearon, M. Gent, A. Turpie, J. Bormanis, J. Weitz, M. Chamberlain, D. Bowie, D. Barnes, and J. Hirsh (2000). Derivation of a

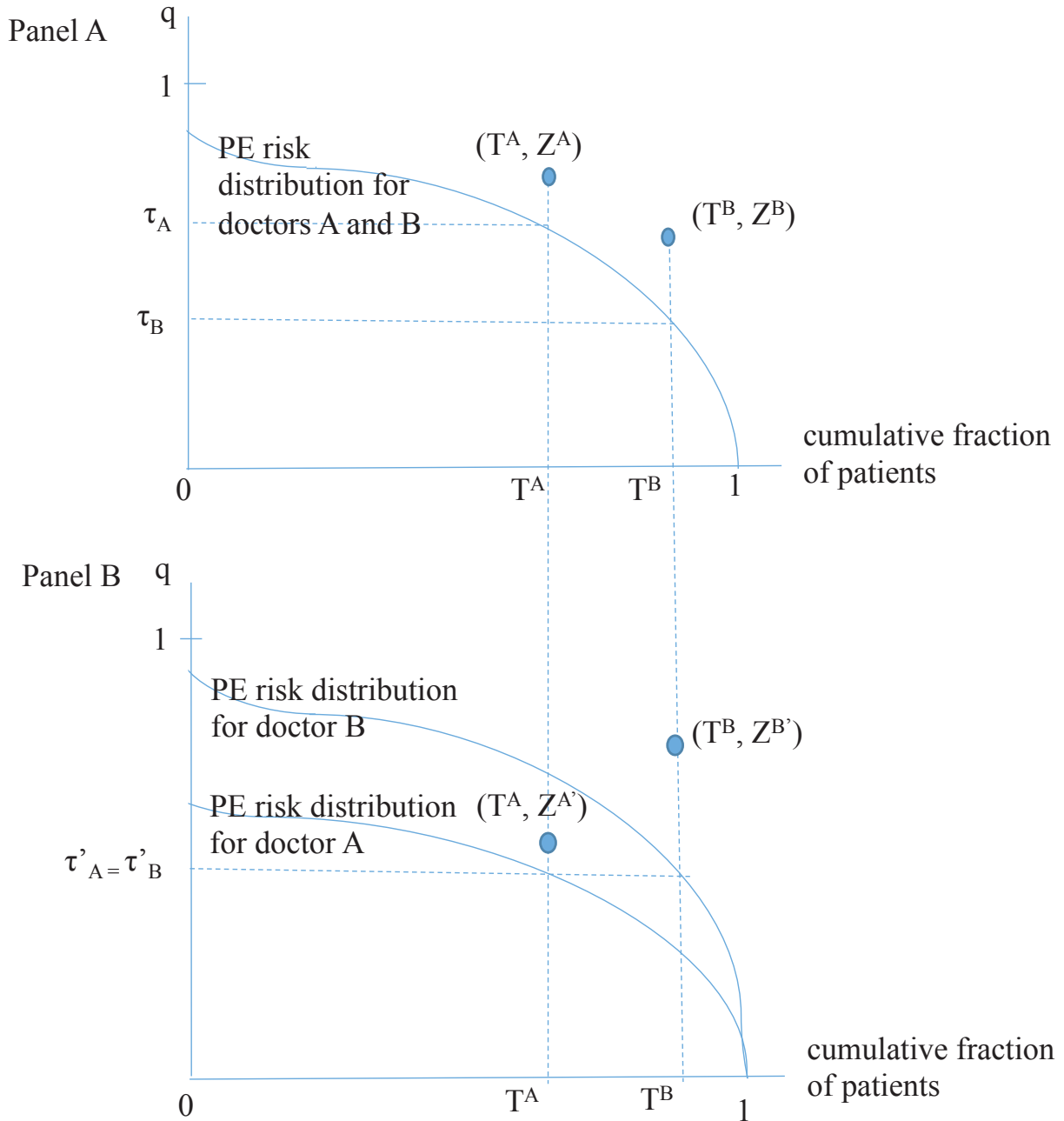
simple clinical model to categorize patients probability of pulmonary embolism-increasing the models utility with the simplified d-dimer. *Thrombosis and Haemostasis* 83(3), 416–420.

Wells, P. S., J. S. Ginsberg, D. R. Anderson, C. Kearon, M. Gent, A. G. Turpie, J. Bormanis, J. Weitz, M. Chamberlain, D. Bowie, D. Barnes, and J. Hirsh (1998). Use of a clinical model for safe management of patients with suspected pulmonary embolism. *Annals of internal medicine* 129(12), 997–1005.

Wells, P. S., J. Hirsh, D. R. Anderson, A. W. A. Lensing, G. Foster, C. Kearon, J. Weitz, R. D’Ovidio, A. Cogo, P. Prandoni, A. Girolami, and J. S. Ginsberg (1995). Accuracy of clinical assessment of deep-vein thrombosis. *The Lancet* 345(8961), 1326–1330.

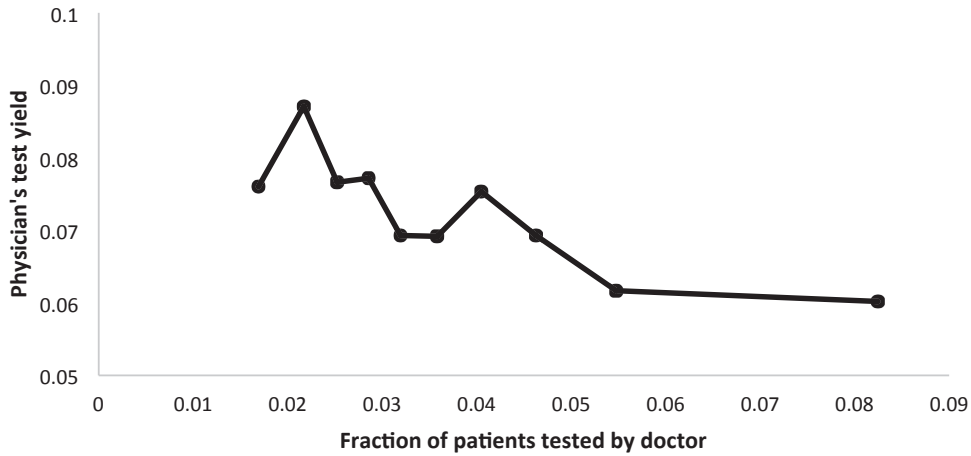
Wennberg, J., M. Cooper, et al. (1996). The Dartmouth atlas of health care in the United States. *Chicago, IL: American Hospital Association.*

Figure 1: Stylized relationship between testing thresholds, testing rates, and test yields



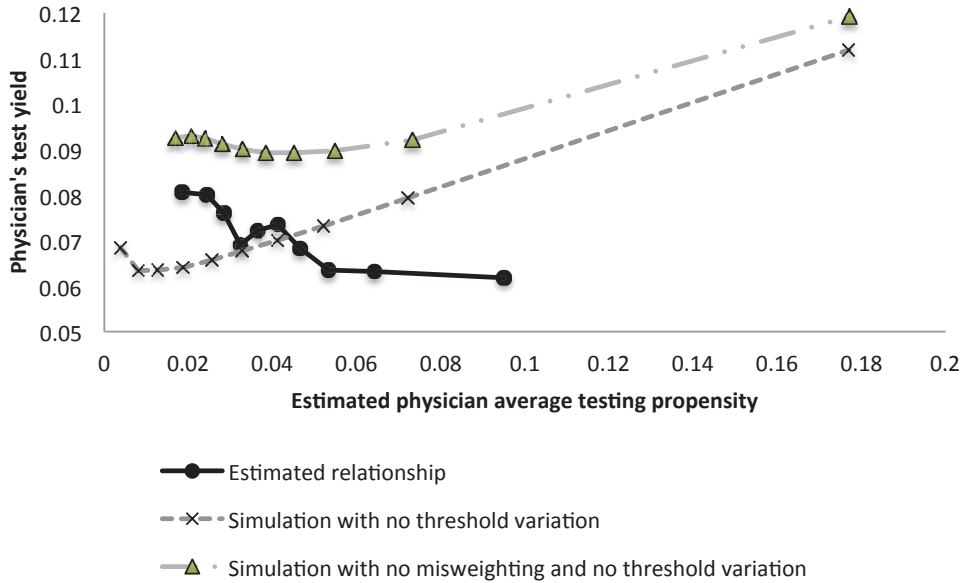
Notes: Figure illustrates the theoretic relationship between testing thresholds, test yields and fraction of patients tested for two hypothetical doctors, A and B. Patients are sorted along the x-axis according to their risk of PE,  $q_{id}$ , from highest risk to lowest risk. Each point  $(x, y)$  along the plotted curve shows the fraction of patients  $x$  for whom  $q_{id} \geq y$ . For example, at point  $(T^A = 2/3, \tau^A = 1/2)$  in Panel A, the graph indicates that 2/3 of patients have a risk of PE that equals or exceeds 1/2.  $\tau_A$  denotes doctor A's testing threshold,  $T^A$  denotes the fraction of patients tested by doctor A,  $Z^A$  denotes doctor A's test yield (among tested patients), and likewise for doctor B. In Panel A, both doctors face patient populations with the same distribution of PE risk. In Panel B, Doctor B's patients are higher risk, i.e. for any given probability of a positive test  $q$ , a greater fraction of doctor B's patients meet or exceed that threshold compared to doctor A.

Figure 2: Binned scatterplot of physician test yield by fraction of patients tested



Notes: Figure displays a binned scatterplot based on our sample of Medicare claims data. Physicians are binned into deciles according to the fraction of patients they test (along the x-axis). The y-axis indicates rate of positive test results among tested patients within each physician decile.

Figure 3: Binned scatterplot of physician test yield by testing propensity index: Estimation results and simulations



Notes: Figure displays a binned scatterplot based on our estimation and simulation results; physicians are binned into deciles based on the average estimated value of the testing propensity index  $\tilde{I}'_{id}$ . The solid black line with circle markers plots the relationship between physicians' actual test yields and physicians' average  $\tilde{I}'_{id}$ . The dashed line with X markers displays the simulated relationship between testing propensities and test yields under a counterfactual with no variation in physician testing thresholds, and instead all physicians assigned the average testing threshold  $E(\tau_d)$ . The line with triangle markers displays the simulated relationship between testing propensities and test yields if there were no variation in physician testing thresholds *and* there were no misweighting of observable risk factors.

Table 1: Summary statistics

	<i>A. Untested patients</i>	<i>B. Patients with negative tests</i>	<i>C. Patients with positive tests</i>
<i>Patient characteristics</i>			
Age	77.6	76.8	76.9
Female	0.586	0.602	0.600
Black	0.082	0.066	0.083
History of PE	0.003	0.006	0.017
<i>Doctor, hospital and region characteristics</i>			
Doctor experience	16.5 (8.3)	16.4 (8.4)	16.8 (8.5)
Top 50 research med. school	0.28	0.29	0.30
Top 50 primary med. school	0.26	0.27	0.28
Academic hospital	0.33	0.34	0.356
For profit hospital	0.12	0.13	0.120
HRR avg spending (in \$)	8,198 (959)	8,173 (972)	8,089 (936)
Average income in region	22,771 (5521)	23,005 (5490)	23,039 (5710)
Joint and several liability	0.69	0.70	0.692
Malpractice damage caps	0.70	0.76	0.747
Number of observations	1,819,015	66,677	4,968

Notes: Table reports means and standard deviations (in parentheses). Data is from the Medicare claims 2000-2009, the American Hospital Association annual survey, the American Medical Association Masterfile, the Dartmouth Atlas, and the Avraham Database of State Tort Law Reform.

Table 2: Summary statistics illustrating potential misweighting of risk factors

	<i>A. Fraction tested</i>	<i>B. Test yield</i>
<i>Selected candidates for under-weighting</i>		
Prostate cancer (CCW)	0.0370	0.1019
No prostate cancer (CCW)	0.0380	0.0677
Black	0.0313	0.0851
Non-black	0.0385	0.0682
History of PE	0.0726	0.1881
No history of PE	0.0378	0.0686
History of deep vein thrombosis	0.0507	0.1656
No history of deep vein thrombosis	0.0378	0.0685
Prior hospital visit within 30 days	0.0465	0.1976
No prior hospital visit within 30 days	0.0377	0.0656
<i>Selected candidates for over-weighting</i>		
Chronic obstructive pulmonary disease (CCW)	0.0466	0.0524
No chronic obstructive pulmonary disease (CCW)	0.0360	0.0742
Atrial fibrillation	0.0742	0.0520
No atrial fibrillation	0.0388	0.0713
Ischemic heart disease	0.0376	0.0566
No ischemic heart disease	0.0382	0.0786

Notes: Table reports summary statistics for selected comorbidities to motivate the examination of misweighting. Variables are selected on the Column A reports average rates of testing for patients with and without the listed conditions. Column B reports average rate of positive tests among tested patients with and without the listed conditions. Data is from the Medicare claims 2000-2009.

Table 3: Regressions of testing threshold on physician characteristics and practice environment

<i>Independent variables:</i>	<i>Dependent variable: Physician testing threshold <math>\tau_d</math></i>			
	OLS (1)	FGLS (2)	OLS (3)	FGLS (4)
Doctor experience	0.0007*** (0.0001)	0.0007*** (0.0001)	0.0007*** (0.0002)	0.0008*** (0.0001)
Top 50 research medical school	0.0047 (0.0038)	0.0050 (0.0031)	0.0053 (0.0047)	0.0032 (0.0037)
Top 50 primary care medical school	-0.0062 (0.0039)	-0.0042 (0.0032)	-0.0077 (0.0048)	-0.0030 (0.0037)
Academic hospital	0.0006 (0.0026)	0.0007 (0.0022)		
For profit hospital	-0.0004 (0.0041)	-0.0018 (0.0032)		
Log(HRR average Medicare spending)	-0.0391*** (0.0109)	-0.0474*** (0.0093)		
Average income in region (in \$10k)	0.0000 (0.0025)	0.0000 (0.0019)		
Joint and several liability	0.0001 (0.0027)	0.0003 (0.0023)		
Malpractice damage caps	-0.0029 (0.0028)	-0.0053** (0.0023)		
Hospital Fixed Effects	No	No	Yes	Yes

Notes: Each column reports results from a regression of estimated physician testing thresholds  $\tau_d$  on characteristics of the physician’s training and practice environment. Even numbered columns report FGLS estimates which account for estimation error in  $\tau_d$ . Columns 3 and 4 include hospital fixed effects. An observation is an individual doctor; there are 6636 observations. \* significant at the 10% level \*\*significance at the 5% level; \*\*\*significance at the 1% level.



Table 4: Part 1: Coefficients from testing model and estimated misweighting in PE risk assessment

	$\beta$ from testing equation (1)	Misweighting amount (2)	Std error of misweighting (3)	T statistic of misweighting (4)
<i>Underweighted risk factors</i>				
Prior hospital visit w/in 30 days	-0.0094	0.1070	0.0121	8.8430
Prior hospital visit w/in 7 days	-0.0041	0.1128	0.0130	8.6769
Prostate cancer (CCW)	0.0014	0.0298	0.0048	6.2083
Cancer metastasis (Elixhauser)	-0.0155	0.0726	0.0128	5.6719
History of deep vein thrombosis	0.0092	0.0571	0.0114	5.0088
History of pulmonary embolism	0.0315	0.0666	0.0145	4.5931
Rheumatoid arthritis, osteoarthritis (CCW)	0.0053	0.0091	0.0024	3.7917
Endometrial cancer (CCW)	-0.0011	0.0547	0.0153	3.5752
Obesity (Elixhauser)	0.0095	0.0218	0.0076	2.8684
Paralysis (Elixhauser)	-0.0026	0.0331	0.0117	2.8291
Other neurological conditions (Elixhauser)	-0.0043	0.0194	0.0075	2.5867
Any prior admission history	0.0028	0.0102	0.0041	2.4878
Alzheimer's disease (CCW)	-0.0023	0.0152	0.0064	2.3750
Colorectal cancer (CCW)	-0.0012	0.0136	0.0067	2.0299
<i>Overweighted risk factors</i>				
Ischemic heart disease (CCW)	0.0007	-0.0226	0.0023	-9.8261
Chronic obstructive pulmonary disease (CCW)	0.0132	-0.0182	0.0036	-5.0556
Atrial fibrillation (CCW)	-0.0066	-0.0156	0.0036	-4.3333
Depression (Elixhauser)	0.0033	-0.0208	0.0069	-3.0145
Peripheral vascular disease (Elixhauser)	-0.0013	-0.0214	0.0071	-3.0141
Diabetes (CCW)	-0.0055	-0.0087	0.0029	-3.0000
Osteoporosis (CCW)	0.0024	-0.0087	0.0033	-2.6364
Deficiency anemias (Elixhauser)	-0.0004	-0.0142	0.0056	-2.5357
Asthma (CCW)	0.0043	-0.0088	0.0040	-2.2000
Chronic pulmonary disease (Elixhauser)	-0.0042	-0.0094	0.0048	-1.9583
<i>Demographic factors</i>				
Black	-0.0074	0.0257	0.0044	5.8409
Asian	0.0005	-0.0386	0.0118	-3.2712
Hispanic	-0.0056	-0.0168	0.0097	-1.7320
Female	0.0014	0.0000	0.0024	0.0000
Age 65-69	-0.0012	0.0119	0.0037	3.2162
Age 70-74	-0.0089	0.0129	0.0052	2.4808
Age 75-79	-0.0024	0.0140	0.0038	3.6842
Age 80-84	-0.0033	0.0166	0.0039	4.2564
Age 85-89	-0.0043	0.0208	0.0042	4.9524
Age 90-94	-0.0127	0.0132	0.0078	1.6923

Notes: This table is continued in Appendix Table A.1, which reports results for the remaining comorbidities which show no significant evidence of under- or over-weighting. Column 1 reports marginal effects from coefficient estimates of the testing equation (i.e. equation 2); for example, patients who were admitted to the hospital within 30 days are 0.94 percentage points less likely to be tested, after controlling for included PE risk factors and physicians' testing thresholds. Column 2 reports estimates of physicians' misweighting of these PE risk factors estimated from equation 14; for example, physicians' observed testing patterns suggest they are underestimating the PE risk associated with a prior hospital visit in the past 30 days by 10.7 percentage points. Column 3 reports standard errors on these misweighting terms. Column 4 reports t-statistics. Variables are sorted by statistical significance, with the exception of demographic risk factors.

Table 5: Distribution of testing thresholds and misweighting under alternative estimation strategies

	Baseline parametric model, all comorbidities	Parametric model, Elixhauser comorbidities excluded	Parametric model, Elixhauser comorbidities and demographics excluded
	(1)	(2)	(3)
Mean of $\tau_d$	0.0563	0.0623	0.0662
Standard Deviation of $\tau_d$	0.0540	0.0396	0.0394
Average absolute value of misweighting	0.0226	0.0214	0.0200
Standard deviation of misweight	0.0347	0.0336	0.0329
Number of observations	1,890,660	1,890,660	1,890,660
	Heteroskedastic parametric model	Semiparametric model, linear polynomial	Semiparametric model, cubic polynomial
	(4)	(5)	(6)
Mean of $\tau_d$	0.0703	0.0672	0.0661
Standard Deviation of $\tau_d$	0.0514	0.0539	0.0541
Average absolute value of misweighting	0.0212	0.0207	0.0208
Standard deviation of misweight	0.0361	0.0357	0.0364
Number of observations	861,707	861,707	861,707

Notes: Panel 1 reports the estimated posterior mean and standard deviation of physician testing thresholds  $\tau_d$  from our baseline parametric model, after applying the Bayesian shrinkage described in Appendix F. Recall that  $\tau_d$  is the threshold probability of a positive test at which a physician determines it is worthwhile to test a patient. The average absolute value of misweighting calculates the absolute value of the difference between physicians' assessment of the patient's PE probability and the estimated risk associated with the patient's comorbidities, and then averages this value across all patients. The standard deviation of misweighting describes how the amount of misweighting varies across patients. Panel 2 reports results from the parametric model that excludes all Elixhauser comorbidities. Panel 3 reports results from the parametric model that excludes both Elixhauser comorbidities and demographic variables. Panel 4 reports results from the heteroskedastic model described in Section 6.2, which allows the variance of  $\eta_{id}$  to differ across physicians. Panels 5 and 6 report results from the semiparametric model described in Section 6.3, where Panel 5 fits the function  $\lambda(\cdot)$  with a linear function and Panel 6 applies a cubic polynomial. Models estimated in Panels 4, 5, and 6 exclude Elixhauser comorbidities and demographic variables and are estimated on a random subsample of half of the physicians for computational tractability.

Table 6: Calibration Parameters

<i>Definition</i>	<i>Value</i>	<i>Parameter</i>	<i>Source</i>
test sensitivity	0.83	$s$	Stein et al., 2006
baseline false positive rate	0.04	$f_p$	Stein et al., 2006
value of a statistical life	\$1,000,000	$VSL$	Murphy and Topel, 2006
medical benefit of treating PE	$0.025VSL$	$MB$	Lessler et al., 2009
medical cost of treating PE	$0.0017VSL$	$MC$	Lessler et al., 2009
financial cost of testing	\$300	$c$	estimated from Medicare claims
financial cost of PE treatment	\$2,800	$CT$	estimated from Medicare claims

Notes: Calibrated parameters of the model applied in welfare simulations reported in Section 7.

Table 7: Patient welfare with observed testing thresholds vs. in simulations with no overtesting

	<i>False positive rate of 4 percent</i>		<i>False positive rate of 3 percent</i>		<i>False positive rate of 0 percent</i>	
	$\tau^*=0.062$		$\tau^*=0.050$		$\tau^*=0.015$	
	Actual (1)	Simulation (2)	Actual (3)	Simulation (4)	Actual (5)	Simulation (6)
<i>Description of simulation results:</i>						
Fraction of doctors over-testing	83.7%	0%	67.2%	0%	10.4%	0%
Percent of patients tested	3.8%	1.9%	3.8%	2.6%	3.8%	3.7%
Number of patients tested	71,314	35,140	71,314	49,390	71,314	70,497
Test yield among tested patients	7.0%	9.0%	7.0%	8.3%	7.0%	7.1%
<i>Welfare analysis:</i>						
Total financial costs of testing (\$ millions)	35.6	19.5	35.6	26.4	35.6	35.3
Total medical cost of testing (\$ millions)	8.5	5.4	8.5	6.9	8.5	8.5
Total medical benefits of testing (\$ millions)	57.5	46.3	74.6	67.6	125.0	124.8
Net benefits of testing (\$ millions)	13.5	21.4	30.4	34.2	80.9	81.0
Total (financial + medical) costs per test (\$)	618.9	709.1	618.9	675.3	618.9	621.2
Total benefits per test (\$)	806.9	1318.7	1045.5	1368.3	1752.8	1770.5
Net benefits per test (\$)	188.1	609.6	426.7	693.0	1134.0	1149.3

Notes: We compare testing behavior and social welfare under the estimated posterior distribution of physician testing thresholds  $\tau_d$  (in odd numbered columns) to simulated behavior assuming all physicians with thresholds below the calibrated optimum are reassigned to the optimal testing threshold of  $\tau_d = \tau^*$  (in even numbered columns). The simulated results do not correct for misweighting. We report results under three different assumptions about the rate of false positive test results, described in the column headers.

Table 8: Patient welfare with observed misweighting vs. in simulations with no misweighting

	<i>False positive rate of 4%</i>	
	Actual (1)	No misweighting (2)
<i>Description of simulation results:</i>		
Percent of patients tested	3.8%	4.3%
Number of patients tested	71314	81410
Test yield among tested patients	7.0%	9.2%
Number of positive tests detected	5019	7526
<i>Welfare analysis:</i>		
Total financial costs of testing (\$ millions)	35.6	45.2
Total medical cost of testing (\$ millions)	8.5	12.4
Total medical benefits of testing (\$ millions)	57.5	106.8
Net benefits of testing (\$ millions)	13.5	49.1
Total (financial + medical) costs per test (\$)	618.9	707.8
Total benefits per test (\$)	806.9	1311.3
Net benefits per test (\$)	188.1	603.5

Notes: We compare testing behavior and social welfare under the observed physician weighting of patient risk factors (in column 1) to simulated behavior assuming that physicians target testing to patients with the highest expected probability of a positive test based on observable demographics and comorbidities (in column 2). The simulated results in Panel B allow  $\tau_d$  to follow the estimated posterior distribution (i.e. without correcting for overtesting).

## A Appendix: For Online Publication

Table A.1: Coefficients from testing model and estimated misweighting in PE risk assessment (continued)

	$\beta$ from testing equation (1)	Misweighting amount (2)	Std error of misweighting (3)	T statistic of misweighting (4)
<i>Other comorbidities</i>				
History of hip fracture (CCW)	-0.0035	0.0192	0.0116	1.6552
Alzheimer's related dementias (CCW)	-0.0060	0.0077	0.0047	1.6383
Anemia (CCW)	-0.0023	0.0038	0.0024	1.5833
Depression (CCW)	-0.0008	0.0042	0.0031	1.3548
Hypertension (CCW)	0.0008	0.0033	0.0025	1.3200
Solid tumor w/o metastasis (Elixhauser)	-0.0066	0.0145	0.0112	1.2946
Benign prostatic hyperplasia (CCW)	-0.0014	0.0046	0.0038	1.2105
Hypothyroidism (Elixhauser)	-0.0009	0.0068	0.0060	1.1333
Liver disease (Elixhauser)	-0.0066	0.0219	0.0195	1.1231
Prior surgery within 1 year	0.0136	0.0239	0.0215	1.1116
Blood loss anemia (Elixhauser)	-0.0044	0.0126	0.0118	1.0678
Breast cancer (CCW)	0.0066	0.0046	0.0049	0.9388
Stroke / Transient ischemic attack (CCW)	-0.0099	0.0035	0.0046	0.7609
Chronic kidney disease (CCW)	-0.0091	0.0024	0.0042	0.5714
Psychoses (Elixhauser)	-0.0057	0.0046	0.0126	0.3651
Congestive heart failure (Elixhauser)	-0.0022	0.0018	0.0056	0.3214
Congestive heart failure (CCW)	-0.0006	0.0008	0.0028	0.2857
Drug abuse (Elixhauser)	0.0059	0.0060	0.0304	0.1974
Alcohol abuse (Elixhauser)	0.0008	0.0020	0.0149	0.1342
Pulmonary circulation disease (Elixhauser)	-0.0035	0.0009	0.0107	0.0841
Acute myocardial infarction (CCW)	-0.0058	0.0002	0.0090	0.0222
Lymphoma (Elixhauser)	-0.0174	-0.0005	0.0220	-0.0227
Coagulation deficiency (Elixhauser)	-0.0001	-0.0006	0.0109	-0.0550
Weight loss (Elixhauser)	-0.0054	-0.0021	0.0119	-0.1765
Prior surgery within 30 days	0.0151	-0.0047	0.0191	-0.2461
Arthritis (Elixhauser)	0.0044	-0.0032	0.0096	-0.3333
Fluid & electrolyte disorders (Elixhauser)	-0.0013	-0.0022	0.0047	-0.4681
Acquired hypothyroidism (CCW)	0.0022	-0.0020	0.0035	-0.5714
Hyperlipidemia (CCW)	0.0054	-0.0017	0.0024	-0.7083
Hypertension (Elixhauser)	0.0012	-0.0051	0.0040	-1.2750
Diabetes w/chronic complications (Elixhauser)	-0.0080	-0.0176	0.0115	-1.5304
Glaucoma (CCW)	-0.0003	-0.0047	0.0029	-1.6207
Diabetes w/o chronic complications (Elixhauser)	-0.0023	-0.0085	0.0051	-1.6667
Lung cancer (CCW)	-0.0142	-0.0198	0.0113	-1.7522
Cataracts (CCW)	-0.0010	-0.0037	0.0021	-1.7619
Valvular disease (Elixhauser)	-0.0031	-0.0116	0.0060	-1.9333

Notes: Table continued from Table 4, which reported coefficients on all comorbidities with significant evidence of misweighting as well as key demographic variables. Column 1 reports marginal effects from coefficient estimates of the testing equation (i.e. equation 2); for example, patients who were admitted to the hospital within 30 days are 0.94 percentage points less likely to be tested, after controlling for included PE risk factors and physicians' testing thresholds. Column 2 reports estimates of physicians' misweighting of these PE risk factors estimated from equation 14; for example, physicians' observed testing patterns suggest they are underestimating the PE risk associated with a prior hospital visit in the past 30 days by 10.7 percentage points. Column 3 reports standard errors on these misweighting terms. Column 4 reports t-statistics. Variables are sorted by statistical significance.

Table A.2: Part 1: Assessment of misweighting with varying included covariates

	<i>All comorbidities</i>		<i>Excluding Elixhauser comorbidities</i>		<i>Excluding Elixhauser comorbidities and demographics</i>	
	Misweighting amount	Standard error	Misweighting amount	Standard error	Misweighting amount	Standard error
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Underweighted risk factors</i>						
Prior hospital visit w/in 30 days	0.1070	0.0121	0.1025	0.0125	0.1045	0.0125
Prior hospital visit w/in 7 days	0.1128	0.0130	0.1091	0.0133	0.1105	0.0133
Prostate cancer (CCW)	0.0298	0.0048	0.0311	0.0048	0.0318	0.0046
Cancer metastasis (Elixhauser)	0.0726	0.0128	0.0843	0.0134	0.0892	0.0134
History of deep vein thrombosis	0.0571	0.0114	0.0560	0.0113	0.0570	0.0113
History of pulmonary embolism	0.0666	0.0145	0.0800	0.0142	0.0827	0.0141
Rheumatoid arthritis, osteoarthritis (CCW)	0.0091	0.0024	0.0097	0.0025	0.0108	0.0024
Endometrial cancer (CCW)	0.0547	0.0153	0.0438	0.0154	0.0405	0.0153
Obesity (Elixhauser)	0.0218	0.0076				
Paralysis (Elixhauser)	0.0331	0.0117				
Other neurological conditions (Elixhauser)	0.0194	0.0075				
Any prior admission history	0.0102	0.0041	0.0033	0.0029	0.0028	0.0029
Alzheimer's disease (CCW)	0.0152	0.0064	0.0158	0.0065	-0.0036	0.0092
Colorectal cancer (CCW)	0.0136	0.0067	0.0166	0.0067	0.0163	0.0067
<i>Overweighted risk factors</i>						
Ischemic heart disease (CCW)	-0.0226	0.0023	-0.0233	0.0023	-0.0226	0.0023
Chronic obstructive pulmonary disease (CCW)	-0.0182	0.0036	-0.0158	0.0037	-0.0159	0.0037
Atrial fibrillation (CCW)	-0.0156	0.0036	-0.0172	0.0036	-0.0175	0.0036
Depression (Elixhauser)	-0.0208	0.0069				
Peripheral vascular disease (Elixhauser)	-0.0214	0.0071				
Diabetes (CCW)	-0.0087	0.0029	-0.0115	0.0028	-0.0105	0.0028
Osteoporosis (CCW)	-0.0087	0.0033	-0.0079	0.0033	-0.0075	0.0032
Deficiency anemias (Elixhauser)	-0.0142	0.0056				
Asthma (CCW)	-0.0088	0.0040	-0.0086	0.0040	-0.0072	0.0040
Chronic pulmonary disease (Elixhauser)	-0.0094	0.0048				
<i>Demographic factors</i>						
Black	0.0257	0.0044	0.0189	0.0045		
Asian	-0.0386	0.0118	-0.0392	0.0118		
Hispanic	-0.0168	0.0097	-0.0142	0.0100		
Female	0.0000	0.0024	0.0000	0.0024		
Age 65-69	0.0119	0.0037	0.0103	0.0037		
Age 70-74	0.0129	0.0052	0.0092	0.0053		
Age 75-79	0.0140	0.0038	0.0122	0.0038		
Age 80-84	0.0166	0.0039	0.0133	0.0039		
Age 85-89	0.0208	0.0042	0.0181	0.0042		
Age 90-94	0.0132	0.0078	0.0075	0.0081		

Notes: Table continued on next page. Column 1 reports estimates of physicians' misweighting of these PE risk factors estimated from equation 14 under the baseline specification with full set of included covariates. Column 2 reports standard errors on these misweighting terms. (Columns 1 and 2 replicate results reported in Table 4 for purposes of comparison.) Columns 3 and 4 also report misweighting terms and standard errors, now from the model that excludes the Elixhauser comorbidity set. Columns 5 and 6 report results from the model that excludes both Elixhauser comorbidities and demographic factors.

Table A2 Part 2: Assessment of misweighting with varying included covariates

	<i>All comorbidities</i>		<i>Excluding Elixhauser comorbidities</i>		<i>Excluding Elixhauser comorbidities and demographics</i>	
	Misweighting amount (1)	Standard error (2)	Misweighting amount (3)	Standard error (4)	Misweighting amount (5)	Standard error (6)
<i>Other comorbidities</i>						
History of hip fracture (CCW)	0.0192	0.0116	0.0025	0.0118	0.0042	0.0117
Alzheimer's related dementias (CCW)	0.0077	0.0047	0.0070	0.0048	0.0070	0.0049
Anemia (CCW)	0.0038	0.0024	0.0014	0.0024	0.0024	0.0024
Depression (CCW)	0.0042	0.0031	-0.0006	0.0029	-0.0010	0.0029
Hypertension (CCW)	0.0033	0.0025	0.0042	0.0024	0.0052	0.0024
Solid tumor w/o metastasis (Elixhauser)	0.0145	0.0112				
Benign prostatic hyperplasia (CCW)	0.0046	0.0038	0.0062	0.0038	0.0070	0.0035
Hypothyroidism (Elixhauser)	0.0068	0.0060				
Liver disease (Elixhauser)	0.0219	0.0195				
Prior surgery within 1 year	0.0239	0.0215	0.0352	0.0217	0.0293	0.0218
Blood loss anemia (Elixhauser)	0.0126	0.0118				
Breast cancer (CCW)	0.0046	0.0049	0.0089	0.0049	0.0095	0.0049
Stroke / Transient ischemic attack (CCW)	0.0035	0.0046	0.0027	0.0047	0.0050	0.0047
Chronic kidney disease (CCW)	0.0024	0.0042	0.0031	0.0044	0.0014	0.0044
Psychoses (Elixhauser)	0.0046	0.0126				
Congestive heart failure (Elixhauser)	0.0018	0.0056	-0.0053	0.0056	-0.0055	0.0056
Congestive heart failure (CCW)	0.0008	0.0028	0.0007	0.0028	0.0020	0.0028
Drug abuse (Elixhauser)	0.0060	0.0304				
Alcohol abuse (Elixhauser)	0.0020	0.0149				
Pulmonary circulation disease (Elixhauser)	0.0009	0.0107				
Acute myocardial infarction (CCW)	0.0002	0.0090	-0.0026	0.0092	0.0153	0.0066
Lymphoma (Elixhauser)	-0.0005	0.0220				
Coagulation deficiency (Elixhauser)	-0.0006	0.0109				
Weight loss (Elixhauser)	-0.0021	0.0119				
Prior surgery within 30 days	-0.0047	0.0191	-0.0066	0.0192	-0.0031	0.0192
Arthritis (Elixhauser)	-0.0032	0.0096				
Fluid & electrolyte disorders (Elixhauser)	-0.0022	0.0047				
Acquired hypothyroidism (CCW)	-0.0020	0.0035	0.0007	0.0030	0.0013	0.0030
Hyperlipidemia (CCW)	-0.0017	0.0024	-0.0005	0.0025	-0.0013	0.0025
Hypertension (CCW)	-0.0051	0.0040				
Diabetes w/complications (Elixhauser)	-0.0176	0.0115				
Glaucoma (CCW)	-0.0047	0.0029	-0.0043	0.0029	-0.0023	0.0029
Diabetes w/o complications (Elixhauser)	-0.0085	0.0051				
Lung cancer (CCW)	-0.0198	0.0113	-0.0219	0.0117	-0.0266	0.0116
Cataracts (CCW)	-0.0037	0.0021	-0.0029	0.0021	-0.0017	0.0020
Valvular disease (Elixhauser)	-0.0116	0.0060				

Notes: Table continued from previous page. Column 1 reports estimates of physicians' misweighting of these PE risk factors estimated from equation 14 under the baseline specification with full set of included covariates. Column 2 reports standard errors on these misweighting terms. (Columns 1 and 2 replicate results reported in Table 4 for purposes of comparison.) Columns 3 and 4 also report misweighting terms and standard errors, now from the model that excludes the Elixhauser comorbidity set. Columns 5 and 6 report results from the model that excludes both Elixhauser comorbidities and demographic factors.



Table A.3: Part 1: Assessing the costs of misweighting by variable

	<i>Net Benefits</i>	<i>Change in net benefits</i>
<b>Original</b>	<b>13.279</b>	
Age 65-69	12.323	-0.956
Age 70-74	12.078	-0.245
Age 75-79	11.580	-0.498
Age 80-84	11.988	0.408
Age 85-89	13.560	1.572
Age 90-94	13.695	0.135
Black	15.486	1.791
Asian	15.707	0.221
Hispanic	15.802	0.095
Acute myocardial infarction (CCW)	15.802	0.000
Alzheimer's disease (CCW)	16.712	0.910
Chronic obstructive pulmonary disease (CCW)	18.879	2.167
Congestive heart failure (CCW)	18.815	-0.064
History of hip fracture (CCW)	18.980	0.165
Anemia (CCW)	19.164	0.184
Asthma (CCW)	19.343	0.179
Hyperlipidemia (CCW)	19.516	0.173
Benign prostatic hyperplasia (CCW)	19.591	0.075
Hypertension (CCW)	19.432	-0.159
Acquired hypothyroidism (CCW)	19.426	-0.006
Alzheimer's related dementias (CCW)	19.644	0.218
Atrial fibrillation (CCW)	20.498	0.854
Cataracts (CCW)	20.625	0.127
Chronic kidney disease (CCW)	20.611	-0.014
Diabetes (CCW)	21.392	0.781
Glaucoma (CCW)	21.484	0.092
Ischemic heart disease (CCW)	23.516	2.032
Depression (CCW)	23.616	0.100
Osteoporosis (CCW)	23.677	0.061
Rheumatoid arthritis, osteoarthritis (CCW)	24.503	0.826
Stroke / Transient ischemic attack (CCW)	24.603	0.100
Breast cancer (CCW)	24.664	0.061
Colorectal cancer (CCW)	25.079	0.415
Prostate cancer (CCW)	26.588	1.509
Lung cancer (CCW)	26.541	-0.047
Endometrial cancer (CCW)	27.117	0.576

Notes: This table is continued on the next page. This table reports results of a series of simulation exercises where we test the welfare impact of correcting for physician misweighting of observed risk factors, one variable at a time. This exercise allows us to assess which specific risk factors are the biggest contributors to the welfare costs associated with misweighting. We proceed in the order listed in the table and show how the total net benefits of testing (in \$ millions) change from their observed value of 13.279 to the final value 49.132 in the absence of any misweighting, by correcting one additional variable in each row. Note that because we continue to allow physician thresholds to vary and do not correct for all risk factors at once, correcting a single additional risk factor occasionally leads to a small decline in net benefits. The results of this exercise may be sensitive to the order in which risk factors are corrected.

Table A3 Part 2: Assessing the costs of misweighting by variable

	<i>Net Benefits</i>	<i>Change in net benefits</i>
Prior surgery within 30 days	26.311	-0.806
Prior surgery within 1 year	30.794	4.483
Any prior admission history	32.632	1.838
Valvular disease (Elixhauser)	32.534	-0.098
Pulmonary circulation disease (Elixhauser)	32.546	0.012
Peripheral vascular disease (Elixhauser)	32.496	-0.050
Paralysis (Elixhauser)	32.927	0.431
Other neurological conditions (Elixhauser)	33.271	0.344
Diabetes w/o chronic complications (Elixhauser)	33.100	-0.171
Diabetes w/chronic complications (Elixhauser)	33.058	-0.042
Hypothyroidism (Elixhauser)	33.195	0.137
Liver disease (Elixhauser)	33.287	0.092
Lymphoma (Elixhauser)	33.286	-0.001
Solid tumor w/o metastasis (Elixhauser)	33.518	0.232
Arthritis (Elixhauser)	33.509	-0.009
Coagulation deficiency (Elixhauser)	33.504	-0.005
Obesity (Elixhauser)	33.840	0.336
Weight loss (Elixhauser)	33.825	-0.015
Fluid & electrolyte disorders (Elixhauser)	33.770	-0.055
Blood loss anemia (Elixhauser)	33.866	0.096
Deficiency anemias (Elixhauser)	33.668	-0.198
Alcohol abuse (Elixhauser)	33.673	0.005
Drug abuse (Elixhauser)	33.675	0.002
Psychoses (Elixhauser)	33.687	0.012
Depression (Elixhauser)	33.706	0.019
Hypertension (Elixhauser)	33.176	-0.530
History of deep vein thrombosis	34.174	0.998
History of pulmonary embolism	35.186	1.012
Prior hospital visit w/in 30 days	43.135	7.949
Prior hospital visit w/in 7 days	47.871	4.736
Female	47.871	0.000
Chronic pulmonary disease (Elixhauser)	47.903	0.032
Congestive heart failure (Elixhauser)	47.914	0.011
Cancer metastasis (Elixhauser)	49.132	1.218

Notes: This table is continued from the previous page. This table reports results of a series of simulation exercises where we test the welfare impact of correcting for physician misweighting of observed risk factors, one variable at a time. This exercise allows us to assess which specific risk factors are the biggest contributors to the welfare costs associated with misweighting. We proceed in the order listed in the table and show how the total net benefits of testing (in \$ millions) change from their observed value of 13.279 to the final value 49.132 in the absence of any misweighting. Note that because we continue to allow physician thresholds to vary and do not correct for all risk factors at once, correcting a single additional risk factor occasionally leads to a small decline in net benefits. The results of this exercise may also be sensitive to the order in which risk factors are corrected.

Table A.4: Sensitivity of welfare simulations to calibration parameters

<b>A. Counterfactual with no overtesting</b>			
	<i>Percent tested</i>	<i>Test yield</i>	<i>Change in net benefits</i>
<i>False positive rate</i>			
0.00	0.037	0.071	0.093
0.03	0.026	0.083	3.802
<b>0.04</b>	0.019	0.090	8.144
<i>Value of a statistical life</i>			
\$500,000	0.005	0.137	15.748
<b>\$1,000,000</b>	0.019	0.090	8.144
\$1,500,000	0.025	0.081	5.249
<i>Test sensitivity</i>			
0.75	0.019	0.090	8.080
<b>0.83</b>	0.019	0.090	8.144
0.90	0.018	0.090	8.191
<i>Financial cost of testing</i>			
\$0	0.033	0.075	0.725
<b>\$300</b>	0.019	0.090	8.144
\$500	0.012	0.104	16.872
<b>B. Counterfactual with no misweighting</b>			
	<i>Percent tested</i>	<i>Test yield</i>	<i>Change in net benefits</i>
<i>False positive rate</i>			
0.00	0.043	0.090	44.134
0.03	0.043	0.090	38.094
<b>0.04</b>	0.043	0.090	35.853
<i>Value of a statistical life</i>			
\$500,000	0.043	0.090	13.184
<b>\$1,000,000</b>	0.043	0.090	35.853
\$1,500,000	0.043	0.090	58.522
<i>Test sensitivity</i>			
0.75	0.043	0.090	36.120
<b>0.83</b>	0.043	0.090	35.853
0.90	0.043	0.090	35.660
<i>Financial cost of testing</i>			
\$0	0.043	0.090	38.882
<b>\$300</b>	0.043	0.090	35.853
\$500	0.043	0.090	33.834

Notes: This table supplements Tables 7 and 8 and displays the simulated welfare benefits of changing physician practice patterns under a range of calibration parameters. Each row represents a separate simulation exercise; bold rows indicate the baseline parameter values used for our main welfare analysis. The changes in net benefits (column 3) are reported in millions of dollars, compared to welfare under observed testing thresholds and misweighting. In any given row, all parameters aside from the one in question are kept constant at the values listed in Table 6. Panel A displays testing behavior and the improvement in social welfare under simulations assuming all physicians with thresholds below the calibrated optimum are reassigned to the optimal testing threshold of  $\tau_d = \tau^*$  (but maintaining the observed degree of misweighting). Panel B displays testing behavior and the improvement in social welfare under simulations assuming that physicians target testing to patients with the highest expected probability of a positive test based on observable demographics and comorbidities (but maintaining the observed degree of overtesting).

## B Physician decision tree & value of a negative CT scan

The flowchart depicted in Appendix Figure B.1 below shows a typical clinical pathway for a patient who may receive a chest CT to test for PE. The most common symptom that leads to the consideration of PE as a diagnosis is chest pain; this is a nonspecific symptom that could also indicate a cardiac problem, pneumonia, or a number of other conditions. Blood oxygen tests and an EKG are likely to be performed immediately at the bedside, and if they suggest a cardiac problem, the patient will receive a more complete cardiac workup.

If cardiac conditions are ruled out, the doctor may then be considering pneumonia, pleural effusion, and pulmonary embolism as possible diagnoses. A chest x-ray and D-dimer blood test would be the typical next steps. A chest x-ray is a low cost test with low levels of radiation exposure and little medical risk; it is highly effective at diagnosing pneumonia and pleural effusion, which are more common than PE. If the x-ray is negative, then the physician may become more concerned about the risk of PE, since other more common conditions causing chest pain have been ruled out. A chest x-ray is a commonplace and recommended antecedent to a CT scan; the popular Geneva risk scoring system for evaluating whether patient’s PE risk necessitates a CT scan includes chest X-ray findings among the seven risk factors used to calculate the score.

At this point, the physician may consider ordering a D-dimer, an inexpensive blood test that provides further information about a patient’s risk of PE. A low-risk result on the D-dimer suggests the patient does not have a PE and the physician may forego a CT scan. A positive D-dimer result is not diagnostic of PE, but suggests an elevated probability of this condition. At this point, the physician would consider ordering a CT scan. Over our study period, the popularity of the D-dimer as an additional screening tool for PE was on the rise. Although we cannot observe the use of the D-dimer in our data, variation in D-dimer utilization is one mechanism by which physician CT ordering behavior may vary.

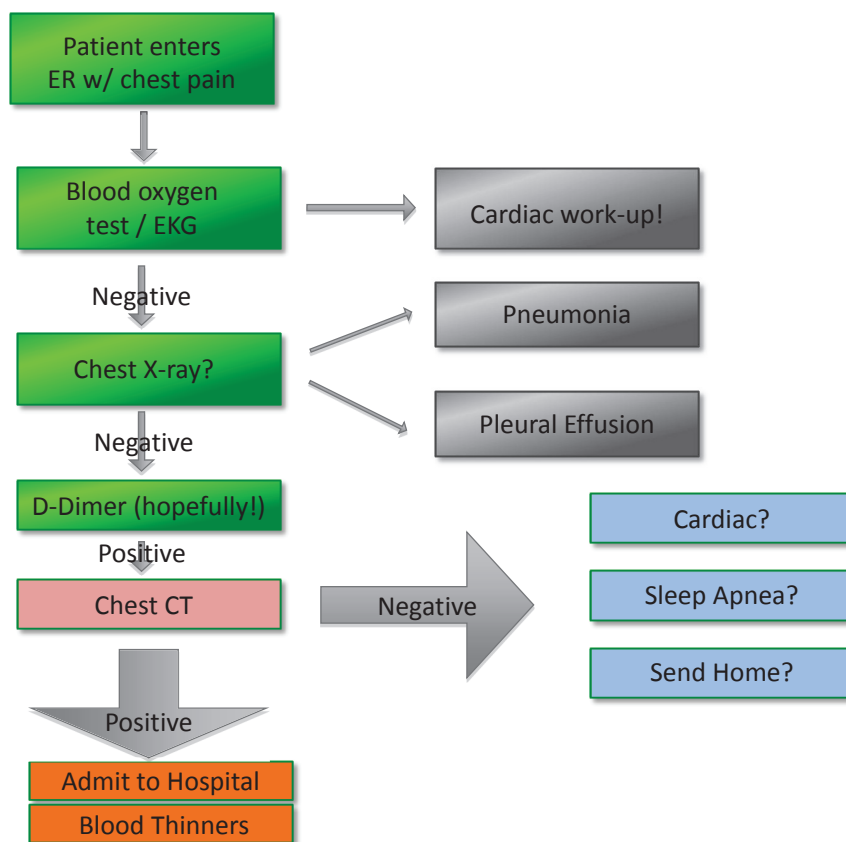
The physician will typically order a chest CT after ruling out these common causes of chest pain. A chest CT with contrast is useful for diagnosing pulmonary embolism, but otherwise adds little new information that may aid diagnosis of other possible acute conditions.<sup>21</sup> A positive test will typically lead to a hospital admission and treatment with blood thinners. Imaging is required for diagnosing PE; even high risk patients have a relatively low probability of PE and PE treatment is medically risky, so it is not a condition that would be treated presumptively without imaging.

A negative CT scan will leave the physician with a broad field of possible alternative diagnoses, including a more subtle cardiac condition, sleep apnea, infection, or a false alarm, and the CT scan result will not be helpful in distinguishing between these possibilities. Ruling out a chest CT has only a modest impact on the posterior probabilities of the other conditions that may be causing a patient’s symptoms, since the *ex ante* probability of PE is relatively low—even for higher risk patients. For these reasons, the informational value of a negative test is low.

---

<sup>21</sup>In Appendix C, we provide a detailed discussion of other conditions that can be diagnosed by chest CT and how we empirically address these possibilities.

Figure B.1: Clinical Assessment of Patient with Potential Pulmonary Embolism



## C Testing for Multiple Conditions

An important caveat to our above analysis is that claims data is only sufficient to identify CPT codes for “chest CT with contrast”; we cannot isolate CT scans that follow the PE testing protocol specifically. Although tests for PE are the primary indication for chest CTs in the emergency room setting, there are other possibilities. Because of this limitation, some of the tests we have labeled as “negative” since the patient is not diagnosed with pulmonary embolism may be tests performed for a different indication. There are five main alternative indications for CT scans in an emergency department setting: trauma, lung or chest cancers, aortic dissection, pleural effusion, and pneumonia. We discuss our approach to each of these alternative diagnoses in turn.

We exclude from the estimation sample patients with diagnosis codes related to trauma (such as fractures, injury, motor vehicle accidents), when these codes are associated with bills on the same day as the patient’s emergency department evaluation. Chest CTs for these patients are likely aiming to assess damage from a trauma rather than a pulmonary embolism. In a detailed sample of patient records from chest CT scans performed in the emergency room of a large hospital, diagnosis codes associated with the radiology bills readily distinguished traumas from other scanning indications.

Similarly, we exclude patients with a history of aortic aneurysm, aortic dissection, or other arterial dissection, in order to eliminate patients for whom chest CTs may be intended to evaluate for aortic dissection. Aortic dissections are extremely rare, with only approximately 9000 cases per year in the United States, making it over 30 times less common than pulmonary embolism (Meszaros et al. 2000).

It is unusual for a cancer diagnosis to be made for the first time in the ED, but patients with worsening symptoms as a result of tumor growth or metastasis and occasional new diagnoses may be seen. CT scanning is routinely used to diagnose and stage cancers. In our sample of detailed ED chest CT records from the academic medical center, fewer than 1% of the scans were used to diagnose or stage cancers. In the Medicare data, we exclude those patients with chest cancer indicated on their visit to the emergency room or associated inpatient visit from our preferred estimation sample.

Chest CTs can be used to guide a procedure to treat patients with pleural effusion, which is typically first diagnosed with a chest X-ray. Because a chest CT is not commonly a diagnostic test for pleural effusion but rather an input into the treatment of the disease, we can exclude patients from the sample with diagnoses of pleural effusion. Since some patients are diagnosed with both pleural effusion and pulmonary embolism, and in these patients the chest CT was likely serving a diagnostic role, we do not exclude pleural effusion patients with a diagnosis of pulmonary embolism. These sample restrictions will tend to overstate the rate of positive testing and bias us away from finding evidence of overtesting, since we may be excluding some pleural effusion patients who are being tested for pulmonary embolism but have a negative test result.

Together, these exclusions for patients with trauma, cancer, or pleural effusion remove 32% of patients receiving chest CTs from our sample. Results presented in the paper are qualitatively similar when these patients are included.

Finally, chest CTs can be used to diagnose pneumonia. Pneumonia can also be reliably diagnosed with cheaper and lower radiation technologies (David et al. 2012); the added value of a chest CT

with contrast in an ED setting for diagnosing these alternative conditions is very modest (Venkatesh et al. 2013). Technically, the value of a chest CT scan for diagnosing a condition that could otherwise be detected with an X-ray is bounded by the costs of the X-ray, which is about \$30 in our sample. Accounting for a \$30 additional net benefit from diagnosing pneumonia when indicated does not substantively change our results about the welfare costs of overtesting.

## D Validating our approach to coding test results in claims data

We identify positive tests on the basis of Medicare Part A hospital claims that include a diagnosis code for PE among any of the diagnoses associated with the hospital stay; we assume all other CT scans failed to detect PE. We have validated our approach to identifying positive tests by using cross-referenced patient chart and hospital billing data from two large academic medical centers. The evidence from these centers suggest that we are unlikely to understate physicians' testing thresholds due to undercounting of positive test results. In particular, we may undercount positive tests in the Medicare claims data for two reasons: if patients with PE are not admitted to the hospital; or if patients with PE are admitted but their inpatient bill does not include a diagnosis of pulmonary embolism.

At the two academic medical centers, we found that 90% of patients who test positive for PE in the emergency department were admitted within 1 day. Patients with very small PEs may occasionally be discharged after brief observation and treated with blood thinning agents as outpatients if the PE appeared small on the scan and the patient has no other complicating health conditions; this likely accounts for most of the cases where a test is coded as positive on the basis of patient chart data but no inpatient admission is recorded. Note that this suggests that we are undercounting positive tests precisely for the patient group for whom the benefits of treatment are the lowest.

Among patients with positive PE CT scans recorded in chart data who are subsequently admitted to the hospital, 87% have a diagnosis of pulmonary embolism recorded on the bill for their inpatient hospital stay. PE may not be recorded on the bill for two main reasons: the patient may have other medical conditions that are treated during the hospital stay and are reimbursed at a higher rate, such that there is no billing incentive to include PE among the inpatient diagnoses; or, the bill may simply be incorrectly coded. In total, 21% of patients diagnosed with PE in the emergency department (ED) do not have an inpatient claim with a PE diagnosis.

Of patients with a negative PE CT scan recorded in their emergency department chart, 1.5% have a diagnosis of pulmonary embolism recorded on the bill for an ensuing hospital stay. In the claims data, we would mistakenly attribute this diagnosis to the ED workup. This error could occur if the patient develops a PE later in his hospital course and receives a subsequent positive CT test, a plausible mechanism given that the immobilization frequently associated with hospital stays is a risk factor for PEs; alternatively, these PE diagnosis codes could indicate billing errors.

Taken together, these data suggest that of the 6% of CT tests that we code as positive in the Medicare data, 20% of the patients had negative findings on their initial ED PE CT. Of the 94% of tests we code as negative, 1.1% of the patients had positive ED PE CTs. The overall rate of positive tests is almost exactly equal to what it would be if no such coding mistakes were made, since



these two types of coding errors offset each other. This suggests that the limitations of this coding algorithm should not contribute to overstatements of the degree of overtesting in our Medicare sample.

## E Derivation and estimation of structural model

In this section, we describe the derivation and estimation of our structural model in more detail. This section is meant to complement the discussion in Section 4, by filling in additional algebraic steps needed to complete the estimation. We begin by outlining our parametric assumptions and describe the testing equation. Second, we derive the test outcome equation which is used to estimate the distribution of  $\tau_d$ , the degree of misweighting, and a scaling factor which relates the testing and test outcome equations.

Recall our assumption that doctor  $d$ 's ex ante belief about the probability of a positive test for patient  $i$  is given by  $q'_{id} = x_{id}\beta' + \alpha'_d + \eta_{id}$ . Although our baseline model assumes that  $\eta_{id}$  is independently and identically distributed across doctors and patients, in Section 6.2 we extend the model to allow for physician-specific heteroskedasticity. The motivation and results of this extension are discussed in more detail in that section. Because the heteroskedastic estimation procedure is a straightforward generalization of our baseline model, we use notation below that allows for heteroskedasticity and thus covers both the baseline model and its heteroskedastic extension.

We assume that the distribution of  $\eta_{id}$  follows a particular functional form, which is a mixture of a Uniform and a Bernoulli distribution; in particular,  $\eta_{id} \sim U(-\eta_d, \eta_d)$  with probability  $1 - p_d$  and  $\eta_{id} \sim U[v - \eta_d, v + \eta_d]$  with probability  $p_d$ . The baseline model in the text assumes homoskedasticity, so that  $p_d = p$  and  $\eta_d = \eta$  and we note below how this affects the estimation procedure.

Assume that doctors test a patient if and only if the patient's perceived probability of a positive test exceeds a physician-specific threshold, i.e.  $q'_{id} > \tau_d$ . Let  $I'_{id} \equiv x_{id}\beta' + \theta'_d$  where  $\theta'_d = \alpha'_d - \tau_d$ . Also as in the text,  $q_{id} = x_{id}\beta + \alpha_d + \eta_{id}$  gives the actual ex ante likelihood of a positive test. Let  $I_{id} \equiv x_{id}\beta + \theta_d$  denote the unprimed version of the propensity to test (i.e. the testing propensity we would observe if physicians correctly weighted observable comorbidities to maximize test yields).

$$\begin{aligned}
 Pr(Test_{id} = 1) &= Pr(q'_{id} > \tau_d) \\
 &= Pr(I'_{id} + \eta_{id} > 0) \\
 &= 1 - Pr(\eta_{id} < -I'_{id})
 \end{aligned} \tag{20}$$

Assume the distribution of  $\eta_{id}$  is such that  $I'_{id} + v < \eta_d$  for all  $I'_{id}$  and  $\eta_d$  so there is no testing propensity  $I'_{id}$  at which patients are always tested regardless of the value of  $\eta_{id}$ . Assume further that patients are never tested if the  $v$  shock is not realized. For example, the  $v$  shock could represent symptoms that would lead the physician to suspect PE, such as chest pain and shortness of breath. Then, given our distributional assumptions:  $Pr(\eta_{id} < -I'_{id}) = 1 - p_d + p_d \cdot \min\left\{1, \frac{\eta_d - (I'_{id} + v)}{2\eta_d}\right\}$ . Thus:

$$\begin{aligned}
Pr(Test_{id} = 1) &= p \left[ 1 - \min \left\{ 1, \frac{1}{2} - \frac{I'_{id} + v}{2\eta_d} \right\} \right] \\
&= \max \left\{ 0, \frac{p_d}{2} + \frac{p_d(I'_{id} + v)}{2\eta_d} \right\}
\end{aligned} \tag{21}$$

We estimate this equation by non-linear least squares. In the heteroskedastic model, we recover:  $\beta'$  (up to a scaling normalization),  $\hat{\eta}_d = C \frac{p_d}{2\eta_d}$  (where the value of the constant  $C$  depends on the normalization of  $\beta$ ), and  $\hat{\theta}'_d = \frac{p_d}{2} + \frac{p_d\theta'_d + v}{2\eta_d}$ . Intuitively, heteroskedasticity in  $\eta_d$  is identified by the fact that observables are less predictive of testing behavior for doctors with more private information. In the homoskedastic model where  $p_d = p$  and  $\eta_d = \eta$ , this simplifies so that we are estimating  $\hat{\beta}' = \frac{p\beta'}{2\eta}$  and  $\hat{\theta}'_d = \frac{p}{2} + \frac{p(\theta'_d + v)}{2\eta}$ .

In either the homoskedastic or heteroskedastic case, we can use the predicted values from estimation of equation 21 to construct an estimate of  $\tilde{I}'_{id} = \frac{p_d}{2} + \frac{p_d(I'_{id} + v)}{2\eta_d}$ . Estimating the heteroskedastic model requires an additional sample restriction at this stage. In theory,  $\eta_d$  is identified for all doctors. In practice, for a very small number of doctors, the estimated  $\eta_d$  would diverge to  $\infty$  because patients with larger  $x_{id}\beta'$  are less likely to be tested, due to random variation in a limited per-doctor sample. These doctors are excluded from the final sample for estimation when we turn to the heteroskedastic model.

Returning to the testing outcomes equation, our distributional assumptions imply that:  $E(\eta_{id} | \eta_{id} > -I'_{id}) = \frac{\eta_d - (I'_{id} + v)}{2}$ . Thus:

$$\begin{aligned}
E(q_{id} | Test_{id} = 1) &= \tau_d + I_{id} + E(\eta_{id} | \eta_{id} > -I'_{id}) \\
&= \tau_d + I_{id} + \frac{\eta_d - (I'_{id} + v)}{2} \\
&= \tau_d + (I_{id} - I'_{id}) + \frac{\eta_d + I'_{id} + v}{2} \\
&= \tau_d + \frac{\eta_d + I'_{id} + v}{2} + x_{id}(\beta - \beta') + (\alpha - \alpha') \\
&= \tau_d + \frac{\eta_d + I'_{id} + v}{2} + (x_{id} - E_d(x_{id}))(\beta - \beta')
\end{aligned} \tag{22}$$

where the last line follows from the assumption that  $E_d(q_{id} | Test_{id} = 1) = E_d(q'_{id} | Test_{id} = 1)$  so that doctors have overall unbiased beliefs about the average likelihood of a positive test across all their tested patients. From our definition of  $\tilde{I}'_{id}$  above, it follows that  $\frac{\eta_d + I'_{id} + v}{2} = \frac{\eta_d \tilde{I}'_{id}}{p_d}$  and so:

$$\begin{aligned}
E(Z_{id} | Test_{id} = 1) &= E(q_{id} | T_{id} = 1) \\
&= \tau_d + \frac{\eta_d \tilde{I}'_{id}}{p_d} + (x_{id} - E_d(x_{id}))(\beta - \beta')
\end{aligned} \tag{23}$$

where  $\tilde{I}'_{id}$  is the propensity estimated from the testing equation, and  $Z_{id}$  is the realized testing outcome (1 for a positive test, 0 for a negative test).

We can estimate this model by non-linear least squares but we need an additional exclusion

restriction so that the coefficient on  $\tilde{I}'_{id}$  is identified by more than just functional form. As discussed in Section 4.3, this restriction is that we effectively know  $\tau_d$  for high volume doctors who test marginal patients—i.e. patients who are very unlikely to be tested based on observables but are nonetheless tested—because we observe test outcomes among those patients. In practice, we also need to be careful about the misweighting term. If we average observed test outcomes  $Z_{id}$  among tested marginal patients (i.e. patients with  $\tilde{I}'_{id} = 0$ ) for doctors who have such patients, then for each of those doctors we obtain an estimate of:

$$QQ_d = \tau_d + (E_{m,d}(x_{id}) - E_d(x_{id}))(\beta - \beta') \quad (24)$$

where  $E_{m,d}(x_{id})$  gives the mean of  $x_{id}$  among only tested marginal patients for a given doctor. For doctors with marginal patients, we have:

$$E(Z_{id}|Test_{id} = 1) - QQ_d = \frac{\eta_d \tilde{I}'_{id}}{p_d} + (x_{id} - E_{m,d}(x_{id}))(\beta - \beta') \quad (25)$$

Because we observe only a small number of marginal patients for each doctor, we can construct:  $\widehat{QQ}_d = QQ_d + e_d$ , a noisy estimate of  $QQ_d$ . Thus, let  $Y_{id} = Z_{id}$  for doctors with no marginal tested patients and  $Y_{id} = Z_{id} - \widehat{QQ}_d$  for doctors with marginal tested patients. Further, let  $X_{id} = (x_{id} - E_{m,d}(x_{id}))$  for doctors with marginal tested patients and  $X_{id} = (x_{id} - E_d(x_{id}))$  for doctors with no marginal tested patients. Finally, let  $M_d$  denote an indicator for whether a doctor has marginal tested patients. This gives the estimating equation:

$$Y_{id} = (1 - M_d)\tau_d + \frac{\eta_d \tilde{I}'_{id}}{p_d} + X_{id}(\beta - \beta') + \epsilon_{id} \quad (26)$$

where  $\epsilon_{id} = M_d e_d + u_{id}$  includes both the noise in the estimation of  $QQ_d$  and the prediction error in  $Z_{id} = E(q_{id}|Test_{id} = 1) + u_{id}$ . This model can be estimated by least squares.

In the homoskedastic case,  $\frac{\eta_d}{p_d}$  is a constant which we recover from least squares estimation of equation 26. In the heteroskedastic model, we estimated  $\hat{\eta}_d = C \frac{p_d}{2\eta_d}$  in the testing equation, so the 2nd term in equation 26 is replaced by  $\frac{\tilde{I}'_{id}}{\hat{\eta}_d}$  and the recovered coefficient tells us  $\frac{C}{2}$ , which is sufficient given  $\hat{\eta}_d$  to recover  $\frac{p_d}{\eta_d}$ .

Following this procedure, we estimate the model and analyze the results described in Section 5. This model is also the basis of the welfare exercises reported in Section 7.

## F “Empirical Bayes” Estimates of $\tau_d$

In this section, we describe how we compute the distribution of the underlying  $\tau_d$  from the observed distribution of  $\hat{\tau}_d$  which includes both the underlying true variation and sampling error. We call this an “empirical Bayes” estimate because of the intuition that we are recovering the true underlying distribution of  $\tau_d$  from noisy estimates, but our specific model does not recover a posterior mean estimate of the parameter for each doctor. Results of this procedure are reported in Table 5. (Note that the welfare results reported in Section 7 require more restrictive assumptions of the empirical Bayes procedure and do recover a posterior estimate of  $\tau_d$  for each doctor. These additional

restrictions are described below and in Section 7.2.)

In order to form our estimate of the true distribution of  $\tau_d$ , we will proceed as follows:

1. Estimate the mean and variance of this distribution for doctors with no marginal tested patients.
2. Estimate the mean and variance of this distribution for doctors who do have marginal tested patients.
3. Apply the law of total variance to compute the mean and variance of the mixture distribution which combines the distributions for doctors with and without marginal tested patients.
4. Make a parametric assumption so that the mean and variance uniquely pin down the posterior distribution. (Required only for welfare simulations reported in Section 7.2.)

We start with our estimating equation from Appendix E, equation 26, reproduced below.

$$Y_{id} = (1 - M_d)\tau_d + \frac{\eta_d \tilde{I}'_{id}}{p_d} + X_{id}(\beta - \beta') + \epsilon_{id} \quad (27)$$

We can rewrite this equation in matrix form as:

$$Y = D\tau_{nm} + X\beta + \epsilon \quad (28)$$

where  $D$  includes the doctor fixed effects for all doctors who lack marginal tested patients (as indicated by the  $nm$  subscript) and  $X\beta$  includes the constant terms, the  $\tilde{I}'_{id}$  terms and the misweighting terms.

Our goal econometrically will be to relate the observed across doctor variance of  $\tau_{nm}$  (which includes estimation error) with the underlying true variance of  $\tau_{nm}$ .

Let  $M_x = I_n - X(X'X)^{-1}X'$  where  $I_n$  is the identity matrix. Partialing out gives:

$$M_x Y = M_x D \tau_{nm} + M_x \epsilon \quad (29)$$

Let  $S = M_x D$ . Then our estimator of  $\tau$  is given by:

$$\hat{\tau}_{nm} = \tau_{nm} + (S'S)^{-1}S'M_x\epsilon \quad (30)$$

For a vector  $x$ , define  $var(x) = E(xx') - E(x)E(x')$ . Define  $var_d(x) = E(x'x) - E_d(x)^2$ , i.e. the scalar generated by taking the variance across the observations in the vector. Taking the “outer product” variance of both sides of equation 30 gives:

$$\begin{aligned} var(\hat{\tau}_{nm}) &= var(\tau_{nm}) + (S'S)^{-1}S'M_x var(\epsilon)M_x S(S'S)^{-1} \\ &= var(\tau_{nm}) + (S'S)^{-1}S' var(\epsilon)S(S'S)^{-1} \end{aligned} \quad (31)$$

where the second line uses the fact that  $M_x M_x = M_x$ . Let  $S^{(i)'}$  denote the  $i$ th row of  $S$ . Assuming  $var(\epsilon)$  is a diagonal matrix,  $S_0 = \frac{1}{N} \sum_{i=1}^N e_i^2 S^{(i)} S^{(i)'} \rightarrow_p \frac{1}{N} \sum_{i=1}^N \epsilon_i^2 S^{(i)} S^{(i)'} = \frac{1}{N} S' var(\epsilon) S$ . This is

asymptotically equivalent to:

$$var(\tau_{nm}) = var(\hat{\tau}_{nm}) - (S'S)^{-1} \left( \sum_{i=1}^N e_i^2 S^{(i)} S^{(i)'} \right) (S'S)^{-1} \quad (32)$$

where  $e_i$  are the residuals from equation 28. Finally, using the fact that  $var_d(\tau_{nm}) = \frac{1}{N_{doc}} tr(var(\tau_{nm}))$  where  $N_{doc}$  is the number of doctors with no marginal tested patients (i.e. the docs for whom we are currently estimating  $\tau_d$ ), we have:

$$var_d(\tau_{nm}) = var_d(\hat{\tau}_{nm}) - \frac{1}{N_{doc}} tr \left( (S'S)^{-1} \left( \sum_{i=1}^N e_i^2 S^{(i)} S^{(i)'} \right) (S'S)^{-1} \right) \quad (33)$$

This equation allows us to recover  $var_d(\tau)$ , the variance of  $\tau_d$  for doctors who lack marginal tested patients. In order to recover  $\tau_d$  for doctors who do have marginal tested patients, we use the fact from equation 23 that:

$$E(Z_{id}|Test_{id} = 1) - (x_{id} - E_d(x_{id}))(\beta - \beta') = \tau_d \quad (34)$$

if we restrict to marginal tested patients of those doctors (meaning that  $\tilde{I}'_{id} = 0$ ). This equation can be written as a special case of equation 28, with  $Y_{id} = Z_{id} - (x_{id} - E_d(x_{id}))(\beta - \beta')$ . Note that  $D$  now denotes the matrix of doctor fixed effects for doctors *with* marginal tested patients,  $N_{marg}$  denotes the number of doctors with marginal tested patients, and  $X = 0$ . This simplification means that  $S = D$  and we have:

$$var_d(\tau_{marg}) = var_d(\hat{\tau}_{marg}) - \frac{1}{N_{marg}} tr \left( (D'D)^{-1} \left( \sum_{i=1}^N e_i^2 D^{(i)} D^{(i)'} \right) (D'D)^{-1} \right) \quad (35)$$

where in this case the residuals are computed from estimation of equation 34 by OLS on the sample of physicians with marginal tested patients and only those marginal tested patients included in the estimation.

To combine these distributions into a single distribution of  $\tau_d$ , we note that  $\tau_d$  is a random variable whose mean and variance are  $\mu_m = E(\tau_{marg})$  and  $\sigma_m^2 = Var_d(\tau_{marg})$  with probability  $P_m$  (the fraction of doctors who have some marginal tested patients) and  $\mu_{nm} = E(\tau_{nm})$  and  $\sigma_{nm}^2 = Var_d(\tau_{nm})$  respectively with probability  $1 - P_m$ . This implies:

$$\begin{aligned} E(\tau) &= P_m \mu_m + (1 - P_m) \mu_{nm} \\ var_d(\tau) &= P_m \sigma_m^2 + (1 - P_m) \sigma_{nm}^2 + P_m \mu_m^2 + (1 - P_m) \mu_{nm}^2 - (P_m \mu_m + (1 - P_m) \mu_{nm})^2 \end{aligned} \quad (36)$$

where the second equation follows from the law of total variance.

For simulations and welfare analyses, we further assume that  $\tau_d + M$  is log-normally distributed with mean  $E(\tau)$ , variance  $var_d(\tau)$  and minimum possible value  $M = fp$ .  $fp$  is the value we would estimate for patients in equation 26 if there were no PE incidence so that the only positive tests were false positives (implying  $E(Z_{id}|Test_{id} = 1) = fp$ , the rate of false positives). In order to recover an estimate of  $\tau_d$  for each doctor, we redraw values of  $\tau$  from the simulated distribution, order them

from least to greatest, and assign each doctor a  $\tau$  from the simulated distribution which matches that doctor’s rank among estimated  $\tau_d$ .

## G Simulations of testing behavior and test yields

This section describes how we apply our structural model to simulate the relationships plotted in Figure 3 and discussed in sections 5.1 and 5.4. The first exercise illustrates the hypothetical relationship between average physician testing propensities and positive test rates, if all doctors were to have the same testing threshold. We simulate testing decisions and test outcomes under a counterfactual where  $\tau_d$  is held constant across doctors, at the estimated average value  $E(\tau_d) = 0.056$ .

To calculate the new values of the testing propensities under this counterfactual where  $\tau_d = E(\tau_d)$  for all doctors, we start by considering the estimated testing propensity:  $\tilde{I}'_{id} = \frac{p}{2} + \frac{p(x_{id}\beta' + \theta'_d + v)}{2\eta}$ . To simulate the testing propensity under the counterfactual where testing thresholds are held constant at their mean,  $\tilde{I}'_{id}{}^{\tau_d=E(\tau_d)}$ , we need to add our estimate of  $(\hat{\tau}_d - E(\tau_d))\frac{p}{2\eta}$  back to original estimate of  $\tilde{I}'_{id}$ .

Because the estimated  $\hat{\tau}_d$  are noisy and overstate the true variance in the distribution, we calculate a posterior, shrunk estimate of each  $\tau_d$  before proceeding with this counterfactual exercise. At this stage, we need to make a distributional assumption about physician testing thresholds  $\tau_d$ . We assume they follow a log-normal distribution with mean and variance determined by the empirical Bayes estimates described above, and the same relative rank as in the raw estimated distribution (i.e. the doctor with the 20th largest estimated  $\hat{\tau}_d$  will also have the 20th largest posterior  $\tau_d$ ).

Plugging in our new, simulated estimates of  $\tilde{I}'_{id}{}^{\tau_d=E(\tau_d)}$  and setting  $\tau_d = E(\tau_d)$ , we calculate  $E(Z_{id}|Test_{id} = 1)$  for each patient following equation 13 and use these estimates to simulate average test yields. Results of this simulation exercise are reported in Section 5.1 and pictured in Figure 3.

The second simulation exercise considers the role of misweighting in determining the relationship between testing propensities and test yield. We simulate the counterfactual relationship between physicians’ average testing propensities and test yields that would be observed if there were no heterogeneity in testing thresholds *and* no misweighting of observable risk factors. Eliminating misweighting should increase the test yield for all values of the testing propensity by improving the targeting of PE CT tests to the highest risk patients.

First we simulate how testing propensities  $\tilde{I}'_{id}{}^{\tau_d=E(\tau_d)}$  would change if there were also no misweighting of patient risk factors. In particular, we add a correction factor  $(x - E(x))\frac{\beta - \beta'}{2(\eta/p)}$  to  $\tilde{I}'_{id}{}^{\tau_d=E(\tau_d)}$  to calculate new simulated testing propensities  $\tilde{I}_{id}^{sim}$  under the counterfactual with no misweighting. Based on these new values of  $\tilde{I}_{id}^{sim}$ , we calculate the expected test yield according to the formula  $E(Z_{id}^{sim}|Test_{id} = 1) = E(\tau_d) + \frac{\eta}{p}\tilde{I}_{id}^{sim}$  (from equation 13). Results of this simulation exercise are reported in Section 5.4 and pictured in Figure 3.

## H Computing the welfare costs of overtesting and misweighting

In order to calculate the welfare costs of overtesting and misweighting, we must first understand how false positive and false negative test results will affect the costs and benefits of testing, and the

calibrated optimal physician testing threshold. We begin by calculating the net utility of treatment, given that there are both false positive and false negative test results. Let  $PE_{id}$  denote the event that patient  $i$  truly has a PE. As before,  $Z_{id}$  is an indicator which is 1 if a test is positive.  $MB$  denotes the medical benefits of treatment if the patient has a PE,  $MC$  denotes the medical costs of treatment and  $CT$  denotes the financial cost of treatment. Then the net utility of a positive test is given by:

$$NU_{id} = Pr(PE_{id}|Z_{id} = 1)MB - MC - CT \quad (37)$$

The medical benefits of treatment accrue only if the positive test result is a “true positive,” i.e. the patient actually has a PE. If there are more false positives, the medical benefits of any observed positive test will be smaller. In contrast, the medical risks and financial costs of treatment are incurred for any treated patient regardless of whether he actually has a PE.

Let  $s$  denote the sensitivity of the test (one minus the probability of a false negative) and  $fp$  denote the probability of a false positive. Applying Bayes’ Rule and the law of total probability, we can rewrite net utility as:

$$NU_{id} = \frac{s(q_{id} - fp)}{q_{id}(s - fp)}MB - MC - CT \quad (38)$$

Given the net utility associated with treating a patient with a positive test, the net benefits of testing also depend on the probability of a positive test,  $q_{id}$  and the costs of testing  $c$ . We can therefore write the net benefits of testing as:

$$\begin{aligned} B_{id} &= q_{id}NU_{id} - c \\ &= \frac{s(q_{id} - fp)}{(s - fp)}MB - q_{id}MC - q_{id}CT - c \end{aligned} \quad (39)$$

Let  $\hat{N}U = \frac{s}{s-fp}MB - MC - CT$  and  $\hat{c} = c + \frac{s \cdot fp}{s-fp}MB$ . Then we can rewrite the net benefits of testing as:

$$B_{id} = q_{id}\hat{N}U - \hat{c} \quad (40)$$

The optimal testing threshold  $\tau^*$  will be the threshold at which the expected net benefits of testing are zero, or  $\tau^*\hat{N}U = \hat{c}$ .

Once we have recovered the optimal testing threshold, we can apply the structural model described in Section 4 and Appendix E, to compute the welfare cost of overtesting as follows. Let  $\hat{t}_{id}(\tau_d, \Delta\beta)$  denote the probability that consumer  $i$  is tested by doctor  $d$  as a function of  $\tau_d$  and the vector of weighting errors physicians make in assessing PE risk. The vector of misweighting errors is labeled as  $\Delta\beta = \beta - \beta'$ . Let  $\hat{Z}_{id}(\tau_d, \Delta\beta)$  denote the probability of a positive test conditional on testing.

To compute testing behavior under the counterfactual where all doctors utilize the optimal testing threshold  $\tau^*$ , we estimate  $\hat{t}_{id}(\tau^*, \Delta\beta)$  using the fact that  $I(\tau^*, \Delta\beta) = I(\tau_d, \Delta\beta) + (\tau_d - \tau^*)$  which implies  $\tilde{I}'(\tau^*, \Delta\beta) = \tilde{I}'(\tau_d, \Delta\beta) + \frac{p(\tau_d - \tau^*)}{2\eta}$ . Having adjusted the testing propensities, we can now calculate the expected probability of a positive test  $\hat{Z}_{id}(\tau^*, \Delta\beta) = \frac{\eta \tilde{I}'_{id}(\tau^*, \Delta\beta)}{p} + (x_{id} - E_d(x_{id}))(\beta - \beta')$ .



Welfare simulations to evaluate the costs of misweighting parallel the derivation above. In particular, to compute the propensity to test with no misweighting,  $\hat{t}_{id}(\tau_d, 0)$ , we use the fact that  $I(\tau_d, 0) = I(\tau_d, \Delta\beta) + (x_{id} - E_d(x_{id}))\Delta\beta$  which implies  $\tilde{I}'(\tau_d, 0) = \tilde{I}'(\tau_d, \Delta\beta) + \frac{p(x_{id} - E_d(x_{id}))\Delta\beta}{2\eta}$ . Given this adjustment to the testing propensities, we can calculate expected test outcomes according to the following formula:  $\hat{Z}_{id}(\tau_d, 0) = \tau_d + \frac{\eta\tilde{I}'_{id}(\tau_d, 0)}{p}$ .

To complete the welfare calculations, we must apply assumptions about the expected medical benefits, medical costs and financial costs associated with treatment of positive tests. Following the notation above, we have:

$$MB(\tau_d, \Delta\beta) = \sum_i Pr(Test_i = 1) \cdot Pr(PE_{id}|Test_i = 1)MB_{id} \quad (41)$$

$$= \sum_i \hat{t}_{id}(\tau_d, \Delta\beta) \frac{s(\hat{Z}_{id}(\tau_d, \Delta\beta) - fp)}{(s - fp)} MB_{id}$$

$$MC(\tau_d, \Delta\beta) = \sum_i Pr(Test_i = 1)Pr(Z_{id} = 1|Test_i = 1)MC_{id} \quad (42)$$

$$= \sum_i \hat{t}_{id}(\tau_d, \Delta\beta)\hat{Z}_{id}(\tau_d, \Delta\beta)MC_{id}$$

$$FC(\tau_d, \Delta\beta) = \sum_i Pr(Test_i = 1)(c + P(Z_{id} = 1|Test_i = 1)CT_{id}) \quad (43)$$

$$= \sum_i \hat{t}_{id}(\tau_d, \Delta\beta)(c + \hat{Z}_{id}(\tau_d, \Delta\beta)CT_{id})$$

$$NB(\tau_d, \Delta\beta) = MB(\tau_d, \Delta\beta) - MC(\tau_d, \Delta\beta) - FC(\tau_d, \Delta\beta) \quad (44)$$

where  $MB$  denote the medical benefits of testing (derived in Section 7.1),  $MC$  denotes the medical costs of testing,  $FC$  denotes the financial costs of testing and  $NB$  denotes the net benefits of testing as a function of these objects. The test sensitivity is given by  $s$ , and  $fp$  is the false positive rate. We define the welfare cost of overtesting as  $NB(\tau^*, \Delta\beta) - NB(\hat{\tau}_d, \Delta\beta)$  and the welfare cost from misweighting as  $NB(\hat{\tau}_d, 0) - NB(\hat{\tau}_d, \Delta\beta)$  where  $\hat{\tau}_d$  is drawn from the estimated underlying distribution of  $\tau_d$  which we recover using the methods outlined in Appendix F above.