

Online Appendix for  
“Selection with Variation in Diagnostic Skill:  
Evidence from Radiologists”

David C. Chan  
Matthew Gentzkow  
Chuan Yu  
December 2021

<b>A</b>	<b>Monotonicity Conditions</b>	<b>A.2</b>
<b>B</b>	<b>Identification of Preferences</b>	<b>A.3</b>
<b>C</b>	<b>Mapping Data to ROC Space</b>	<b>A.4</b>
<b>D</b>	<b>Tests of Monotonicity</b>	<b>A.6</b>
<b>E</b>	<b>Details of Structural Analysis</b>	<b>A.9</b>
	E.1 Optimal Diagnostic Thresholds	A.9
	E.2 Simulated Maximum Likelihood Estimation	A.13
	E.3 Empirical Bayes Posterior Means	A.14
<b>F</b>	<b>Robustness</b>	<b>A.14</b>
<b>G</b>	<b>Extensions</b>	<b>A.17</b>
	G.1 General Loss for False Negatives	A.17
	G.2 Incorrect Beliefs	A.23
	G.3 Simulation of Linear Risk Adjustment	A.24
	G.4 Controlling for Radiologist Skill	A.25

## A Monotonicity Conditions

We begin with the covariance object of interest under average monotonicity of Frandsen et al. (2019) (Condition 2). For a given case  $i$  and set of agents  $\mathcal{J}$ , define

$$\Psi_{i,\mathcal{J}} = \sum_{j \in \mathcal{J}} \rho_j (P_j - \bar{P}) (d_{ij} - \bar{d}_i),$$

where  $\rho_j$  is the share of cases assigned to agent  $j$ ,  $\bar{P} = \sum_j \rho_j P_j$  is the  $\rho$ -weighted average treatment propensity, and  $\bar{d}_i = \sum_j \rho_j d_{ij}$  is the  $\rho$ -weighted average potential treatment of case  $i$ .

To consider probabilistic monotonicity (Condition 3), which allows  $d_{ij}$  to be random, we consider the probability limit of  $\Psi_{i,\mathcal{J}}$  over random draws of  $d_{ij}$ , as the number of draws grows large:

$$\bar{\Psi}_{i,\mathcal{J}} = \sum_{j \in \mathcal{J}} \rho_j (P_j - \bar{P}) \left( \Pr(d_{ij} = 1) - E[\bar{d}_i] \right),$$

where  $E[\bar{d}_i] = \sum_j \rho_j \Pr(d_{ij} = 1)$ .

**Proposition A.1.** *Probabilistic monotonicity (Condition 3) in some set of agents  $\mathcal{J}$  implies  $\bar{\Psi}_{i,\mathcal{J}} \geq 0$  for all  $i$ .*

*Proof.* Under probabilistic monotonicity, for any  $j$  and  $j'$ ,  $P_j > P_{j'}$  implies that  $\Pr(d_{ij} = 1) \geq \Pr(d_{ij'} = 1)$  for all  $i$ . Thus, any ( $\rho$ -weighted) covariance between  $P_j$  and  $\Pr(d_{ij} = 1)$  must be weakly positive for all  $i$ , in any set of agents  $\mathcal{J}$  where probabilistic monotonicity holds.  $\bar{\Psi}_{i,\mathcal{J}}$  is in fact the  $\rho$ -weighted covariance between  $P_j$  and  $\Pr(d_{ij} = 1)$  for a given  $i$ , so  $\bar{\Psi}_{i,\mathcal{J}} \geq 0$  for all  $i$ .  $\square$

To analyze the implications of skill-propensity independence (Condition 4), we define the limit as the number of agents grows large. We assume that when the set of agents is  $\mathcal{J}$ , the skill  $\alpha_j$ , diagnosis rate  $P_j$ , an assignment weight  $\varsigma_j$  such that  $\rho_j = \varsigma_j / \sum_{j' \in \mathcal{J}} \varsigma_{j'}$ , and any other decision-relevant characteristics of each agent  $j \in \mathcal{J}$  are drawn independently from a distribution  $\mathcal{H}$ .

For a case  $i$ , let  $\mathcal{G}$  denote the distribution of  $(\alpha_{j(i)}, P_{j(i)})$  incorporating the uncertainty from both the draws from  $\mathcal{H}$  and the assignment process. Skill-propensity independence (Condition 4) implies that  $\alpha_{j(i)}$  and  $P_{j(i)}$  are independent under  $\mathcal{G}$ . We let  $\pi_i(\alpha, p)$  denote the probability that the case is diagnosed conditional on the assigned agent's skill  $\alpha$  and diagnosis rate  $p$ , and  $\pi_i(p)$  denote the probability conditional only on  $p$ . Probabilistic monotonicity (Condition 3) implies that  $\pi_i(\alpha, p)$  is increasing in  $p$ .

Let  $\bar{\Psi}_i$  denote the probability limit of  $\Psi_{i,\mathcal{J}}$  as the number of agents in  $\mathcal{J}$  grows large.

**Proposition A.2.** *Skill-propensity independence (Condition 4) implies  $\bar{\Psi}_i \geq 0$  for all  $i$ .*

*Proof.* Note that under skill-propensity independence we can write  $\mathcal{G}(\alpha, p) = \mathcal{G}_\alpha(\alpha) \mathcal{G}_p(p)$ , where  $\mathcal{G}_\alpha$  and  $\mathcal{G}_p$  are the marginal distributions of  $p$  and  $\alpha$ . By the law of large numbers, the probability

limit  $\bar{\Psi}_i$  is the expectation under the joint distribution  $\mathcal{G}$ :  $\bar{\Psi}_i = E_{\mathcal{G}} \left[ \left( p - \bar{P} \right) \left( \pi_i(\alpha, p) - \bar{d}_i \right) \right]$ . Moreover,

$$\begin{aligned} E_{\mathcal{G}} \left[ \left( p - \bar{P} \right) \left( \pi_i(\alpha, p) - \bar{d}_i \right) \right] &= \int_p \int_{\alpha} \left( p - \bar{P} \right) \pi_i(\alpha, p) d\mathcal{G}(\alpha, p) \\ &= \int_p \left( p - \bar{P} \right) \pi_i(p) d\mathcal{G}_p(p) \\ &\geq 0 \end{aligned}$$

The first equality uses the fact that  $E_{\mathcal{G}} \left[ \left( P_j - \bar{P} \right) \bar{d}_i \right] = 0$ , the second equality uses skill-propensity independence, and the final inequality uses  $\bar{P} = E_{\mathcal{G}} [P_j]$  and the fact that  $\pi_i(\alpha, p)$  increasing in  $p$  implies  $\pi_i(p)$  increasing in  $p$ .  $\square$

## B Identification of Preferences

**Proposition B.3.** *If the posterior probability of  $s_i = 1$  is continuously increasing in  $w_{ij}$  for any signal, ROC curves must be smooth and concave.*

*Proof.* Without loss of generality, consider a uniform signal  $w \sim U(0, 1)$ . Then under the threshold rule noted in Section 2.1,  $P_j = 1 - \tau_j$ . Furthermore,

$$\begin{aligned} TPR_j &= \frac{1}{S} \int_{1-P_j}^1 \Pr(s = 1 | w, \alpha_j) dw; \\ FPR_j &= \frac{1}{1-S} \int_{1-P_j}^1 1 - \Pr(s = 1 | w, \alpha_j) dw. \end{aligned}$$

This implies a slope in ROC space of  $\frac{1-S}{S} \frac{\Pr(s=1|1-P_j, \alpha_j)}{1-\Pr(s=1|1-P_j, \alpha_j)}$  at  $P_j$ , which is decreasing in  $P_j$  if  $\Pr(s = 1 | w, \alpha_j)$  is increasing in  $w$ .  $\square$

**Proposition B.4.** *Knowing the cost of a false negative relative to a false positive,  $\beta_j \equiv \frac{u_j(1,1) - u_j(0,1)}{u_j(0,0) - u_j(1,0)} \in (0, \infty)$ , is sufficient to identify the function  $u_j(\cdot, \cdot)$  up to normalizations.*

*Proof.* The agent's expected loss from choosing  $d = 1$  rather than  $d = 0$  is

$$E[u(1, s) - u(0, s) | w, \alpha] = [u(1, 1) - u(0, 1)] \Pr(s = 1 | w, \alpha) + [u(1, 0) - u(0, 0)] \Pr(s = 0 | w, \alpha).$$

The optimal decision is thus  $d = 1$  if and only if

$$\frac{u(1, 1) - u(0, 1)}{u(0, 0) - u(1, 0)} \geq \frac{\Pr(s = 0 | w, \alpha)}{\Pr(s = 1 | w, \alpha)}.$$

$\square$

## C Mapping Data to ROC Space

In this appendix, we detail parameters that map the observed data on diagnoses ( $d_i$ ) and false negatives ( $m_i$ ) for each patient to the key objects of the true positive rate ( $TPR_j$ ) and the false positive rate ( $FPR_j$ ) for each radiologist in ROC space. As discussed in Section 4.1, this mapping requires a parameter for the prevalence of pneumonia, or  $S = 1 - \Phi(\bar{v})$ . Under quasi-random assignment, this prevalence of pneumonia is (conditionally) the same across radiologists.

In addition, we allow for two additional parameters to address practical concerns. First, some chest X-rays are ordered for reasons completely unrelated to pneumonia (e.g., rib fractures). We thus consider a proportion of cases  $\kappa$  that are not at risk for pneumonia and are recognized as such by all radiologists. Second, we do not observe false negatives immediately at the same time that the chest X-ray is read. So we allow for a share  $\lambda$  of undiagnosed cases that do not have pneumonia to develop it and be diagnosed subsequently, thus being incorrectly observed as false negatives.

We begin with the observed radiologist-specific diagnosis and miss rates  $P_j^{\text{obs}}$  and  $FN_j^{\text{obs}}$ , which are population values of the estimates  $\widehat{P}_j^{\text{obs}}$  and  $\widehat{FN}_j^{\text{obs}}$  defined in the main text. They relate to true shares  $FN_j$ ,  $TN_j$ ,  $FP_j$ , and  $TP_j$  as follows:

$$P_j^{\text{obs}} = (1 - \kappa)(TP_j + FP_j) = (1 - \kappa)P_j; \quad (\text{C.1})$$

$$FN_j^{\text{obs}} = (1 - \kappa)(FN_j + \lambda TN_j). \quad (\text{C.2})$$

Using Equations (C.1) and (C.2) above and the fact that  $TN_j = 1 - P_j - FN_j$ , we derive

$$FN_j = \frac{\lambda P_j^{\text{obs}} + FN_j^{\text{obs}}}{(1 - \kappa)(1 - \lambda)} - \frac{\lambda}{1 - \lambda}. \quad (\text{C.3})$$

We can derive the remaining shares by using  $TN_j = 1 - P_j - FN_j$ ,  $TP_j = S - FN_j$ , and  $FP_j = P_j - TP_j$ :

$$\begin{aligned} TN_j &= \frac{1}{1 - \lambda} - \frac{P_j^{\text{obs}} + FN_j^{\text{obs}}}{(1 - \kappa)(1 - \lambda)}; \\ TP_j &= S - \left( \frac{\lambda P_j^{\text{obs}} + FN_j^{\text{obs}}}{(1 - \kappa)(1 - \lambda)} - \frac{\lambda}{1 - \lambda} \right); \\ FP_j &= \frac{P_j^{\text{obs}} + FN_j^{\text{obs}}}{(1 - \kappa)(1 - \lambda)} - \frac{\lambda}{1 - \lambda} - S. \end{aligned}$$

The underlying true positive rates and false positive rates are thus

$$\begin{aligned} TPR_j &\equiv \frac{TP_j}{TP_j + FN_j} = 1 - \frac{1}{S} \left( \frac{\lambda P_j^{\text{obs}} + FN_j^{\text{obs}}}{(1 - \kappa)(1 - \lambda)} - \frac{\lambda}{1 - \lambda} \right); \\ FPR_j &\equiv \frac{FP_j}{FP_j + TN_j} = \frac{1}{1 - S} \left( \frac{P_j^{\text{obs}} + FN_j^{\text{obs}}}{(1 - \kappa)(1 - \lambda)} - \frac{\lambda}{1 - \lambda} - S \right). \end{aligned}$$

Conditional on  $S$ ,  $\kappa$ , and  $\lambda$ , we can thus transform data for a given radiologist in reduced-form space to the relevant radiologist-specific rates in ROC space:

$$\left(P_j^{\text{obs}}, FN_j^{\text{obs}}\right) \xrightarrow{S, \kappa, \lambda} (FPR_j, TPR_j).$$

In Figure V, we show the implied  $(FPR_j, TPR_j)$  based on  $(\widehat{P}_j^{\text{obs}}, \widehat{FN}_j^{\text{obs}})$  and model estimates of  $S$ ,  $\kappa$ , and  $\lambda$ . This figure does not account for the fact that  $(\widehat{P}_j^{\text{obs}}, \widehat{FN}_j^{\text{obs}})$  are measured in finite sample, and we simply impose that  $TPR_j \leq 1$ ,  $FPR_j \geq 0$ , and  $TPR_j \geq FPR_j$ , sequentially. The first step of  $TPR_j \leq 1$  truncates 597 out of 3,199 radiologists (or 18.7% of radiologists), which mainly comes from the radiologists whose observed miss rate,  $\widehat{FN}_j^{\text{obs}}$ , is smaller than  $\lambda$ . The second step of  $FPR_j \geq 0$  truncates 44 radiologists. The third step of  $TPR_j \geq FPR_j$  truncates 68 radiologists. In Appendix Figure A.14, we plot empirical Bayes posterior means of  $(FPR_j, TPR_j)$  based on  $(\widehat{P}_j^{\text{obs}}, \widehat{FN}_j^{\text{obs}})$  and all estimated model parameters.

While ROC-space radiologist rates depend on  $S$ ,  $\kappa$ , and  $\lambda$ , it is important to note that two key findings are invariant to these parameters. First, Figure VI and Appendix Figure A.9 imply an upward-sloping relationship between  $P_j^{\text{obs}}$  and  $FN_j^{\text{obs}}$ . By Equations (C.1) and (C.3), we can see that this violates the prediction that  $\Delta \in [-1, 0]$ , based on  $P_j$  and  $FN_j$ . Specifically, comparing two radiologists  $j$  and  $j'$ , Equations (C.1) and (C.3) imply that

$$\frac{FN_j^{\text{obs}} - FN_{j'}^{\text{obs}}}{P_j^{\text{obs}} - P_{j'}^{\text{obs}}} = (1 - \lambda) \frac{FN_j - FN_{j'}}{P_j - P_{j'}} - \lambda \in [-1, -\lambda].$$

So the coefficient estimand  $\Delta^{\text{obs}} > 0$  from a regression of  $FN_j^{\text{obs}}$  on  $P_j^{\text{obs}}$  implies that  $\Delta > 0$  for any  $\lambda \in [0, 1)$ .

Second, by Remark 2, an upward sloping relationship between  $P_j$  and  $FN_j$  contradicts uniform skill regardless of  $S$ . Therefore, regardless of  $S$ , the pattern of  $(FPR_j, TPR_j)$  across radiologists in ROC space, as in Figure V, should remain downward-sloping and inconsistent with the assumption of uniform skill.<sup>1</sup>

To illustrate the second point, we show in Appendix Figure A.6 that the pattern of  $(FPR_j, TPR_j)$  across radiologists remains inconsistent with uniform skill, at lower and upper bounds for  $S$ . To construct these bounds, we first divide all radiologists into ten bins based on their diagnosed shares  $\widehat{P}_j$ . For each bin  $q$ , we set a lower bound for  $S$  at the weighted-average (underlying) miss rate, or  $\underline{S}_q = \overline{FN}_q = \frac{\sum_{j \in \mathcal{J}_q} n_j \widehat{FN}_j}{\sum_{j \in \mathcal{J}_q} n_j}$ , where  $\mathcal{J}_q$  is the set of agents in bin  $q$ . In other words, we assume that all diagnoses are false positives. We set an upper bound for  $S$  at the weighted-average sum of the (underlying) miss rate and diagnosis rate, or  $\overline{S}_q = \overline{FN}_q + \overline{P}_q = \frac{\sum_{j \in \mathcal{J}_q} n_j (\widehat{FN}_j + \widehat{P}_j)}{\sum_{j \in \mathcal{J}_q} n_j}$ . Finally, we take the intersection of these bounds from all bins as the bounds in the full sample, which gives us

<sup>1</sup>Consider two agents  $j$  and  $j'$ . Let  $\Delta TPR \equiv TPR_j - TPR_{j'}$ ;  $\Delta FPR \equiv FPR_j - FPR_{j'}$ ;  $\Delta P \equiv P_j - P_{j'}$ ; and  $\Delta FN \equiv FN_j - FN_{j'}$ . It is easy to show that  $\Delta TPR = -\frac{1}{S} \Delta FN$  and  $\Delta FPR = \frac{1}{1-S} (\Delta P + \Delta FN)$ . So  $\frac{\Delta TPR}{\Delta FPR} = -\frac{1-S}{S} \frac{\Delta FN}{\Delta P + \Delta FN}$ . The condition that  $\frac{\Delta FN}{\Delta P} \in (-1, 0)$  is equivalent to the condition that  $\frac{\Delta TPR}{\Delta FPR} > 0$ , as long as  $S \in (0, 1)$ .

$\underline{S} = \max_{1 \leq q \leq 10} \underline{S}_q = 0.015$  and  $\bar{S} = \min_{1 \leq q \leq 10} \bar{S}_q = 0.073$ .

Further, as we discuss in Section 4.4, our overall results remain robust to alternative values for  $\kappa$ . As shown in Appendix Table A.10, model parameters are stable and suggest wide variation in diagnostic skill. Model implications for reducing variation by uniform preferences or uniform skill similarly remain robust.

## D Tests of Monotonicity

Under the standard monotonicity assumption (Condition 1(iii)), when comparing a radiologist  $j'$  who diagnoses more cases than radiologist  $j$ , there cannot be a case  $i$  such that  $d_{ij} = 1$  and  $d_{ij'} = 0$ . In this appendix, we conduct informal tests of this assumption that are standard in the judges-design literature, along the lines of tests in Bhuller et al. (2020) and Dobbie et al. (2018). These monotonicity tests confirm whether the first-stage estimates are non-negative in subsamples of cases. We first present results of implementing these standard tests. We then draw relationships between these tests, which do not reject monotonicity, and our analysis in Section 4, which strongly rejects monotonicity.

### Results

We define subsamples of cases based on patient characteristics. We consider four characteristics: probability of diagnosis (based on patient characteristics), age, arrival time, and race. We define two subsamples for each of the characteristics, for a total of eight subsamples: (i) above-median age, (ii) below-median age, (iii) above-median probability of diagnosis, (iv) below-median probability of diagnosis, (v) arrival time during the day (between 7 a.m. and 7 p.m.), (vi) arrival time at night (between 7 p.m. and 7 a.m.), (vii) white race, and (viii) non-white race.

The first testable implication follows from the following intuition: Under monotonicity, a radiologist who generally increases the probability of diagnosis should increase the probability of diagnosis in any subsample of cases. Following the judges-design literature, we construct leave-out propensities for pneumonia diagnosis and use these propensities as instruments for whether an index case is diagnosed with pneumonia, as in Equation (4).

In each of the eight subsamples indexed by  $r$ , we estimate the following first-stage regression, using observations in subsample  $\mathcal{I}_r$ :

$$d_i = \alpha_r Z_{j(i)} + \mathbf{X}_i \pi_r + \mathbf{T}_i \eta_r + \varepsilon_i. \quad (\text{D.4})$$

Consistent with our quasi-experiment in Assumption 1, we control for time categories interacted with station identities, or  $\mathbf{T}_i$ . We also control for patient characteristics  $\mathbf{X}_i$ , as in our baseline first-stage regression. Under monotonicity, we should have  $\alpha_r \geq 0$  for all  $r$ .

The second testable implication is slightly stronger: Under monotonicity, an increase in the probability of diagnosis by changing radiologists in any subsample of patients should correspond to increases in the probability of diagnosis in all other subsamples of patients. To capture this intuition,

we construct “reverse-sample” instruments that exclude any case in subsample  $r$ :

$$Z_j^{-r} = \frac{1}{|I_j \setminus \mathcal{I}_r|} \sum_{i \in I_j \setminus \mathcal{I}_r} d_i,$$

We estimate the first-stage regression, using observations in subsample  $\mathcal{I}_r$ :

$$d_i = \alpha_r Z_{j(i)}^{-r} + \mathbf{X}_i \pi_r + \mathbf{T}_i \eta_r + \varepsilon_i. \quad (\text{D.5})$$

As before, we control for patient characteristics  $\mathbf{X}_i$  and time categories interacted with station dummies  $\mathbf{T}_i$ , and we check whether  $\alpha_r \geq 0$  for all  $r$ .

In Appendix Table A.6, we show results for these informal monotonicity tests, based on Equations (D.4) and (D.5). Panel A shows results corresponding to the standard leave-out instrument, or  $\alpha_r$  from the Equation (D.4). Panel B shows results corresponding to the reverse-sample instrument, or  $\alpha_r$  from Equation (D.5). Each column corresponds to a different subsample. All 16 regressions yield strongly positive first-stage coefficients.

### Relationship with Reduced-Form Analysis

At a high level, the informal tests of monotonicity in the judges-design literature use information about observable case characteristics and treatment decisions, while our analysis in Section 4 exploits additional information about outcomes tied to an underlying state that is relevant for the classification decision. In this subsection, we will clarify the relationship between these analyses.

We begin with the standard condition for IV validity, Condition 1. Following Imbens and Angrist (1994), we abstract from covariates, assuming unconditional random assignment in Condition 1(ii), and consider a discrete multivalued instrument  $Z_i$ . In the judges design, the instrument can be thought of as the agent’s treatment propensity, or  $Z_i = P_{j(i)} \in \{p_1, p_2, \dots, p_K\}$ , which the leave-out instrument approaches with infinite data. We assume that  $p_1 < p_2 < \dots < p_K$ . We also introduce the notation  $d_i(Z_i) \in \{0, 1\}$  to denote potential treatment decisions as a function of the instrument; in our main framework, this amounts to  $d_{ij} = d_i(p)$  for all  $j$  such that  $P_j = p$ .

Now consider some binary characteristic  $x_i \in \{0, 1\}$ . We first note that the following Wald estimand between two consecutive values  $p_k$  and  $p_{k+1}$  of the instrument characterizes the probability that  $x_i = 1$  among compliers  $i$  such that  $d_i(p_{k+1}) > d_i(p_k)$ :

$$\frac{E[x_i d_i | Z_i = p_{k+1}] - E[x_i d_i | Z_i = p_k]}{E[d_i | Z_i = p_{k+1}] - E[d_i | Z_i = p_k]} = E[x_i | d_i(p_{k+1}) > d_i(p_k)].$$

Since  $x_i$  is binary, this Wald estimand gives us  $\Pr(x_i = 1 | d_i(p_{k+1}) > d_i(p_k)) \in [0, 1]$ .

Under Imbens and Angrist (1994), 2SLS of  $x_i d_i$  as an “outcome variable,” instrumenting  $d_i$  with all values of  $Z_i$ , will give us a weighted average of the Wald estimands over  $k \in \{1, \dots, K-1\}$ .



Specifically, consider the following equations:

$$x_i d_i = \Delta^x d_i + u_i^x; \quad (\text{D.6})$$

$$d_i = \alpha^x Z_i + v_i^x. \quad (\text{D.7})$$

The 2SLS estimator of  $\Delta^x$  in this set of equations should converge to a weighted average:

$$\Delta^x = \sum_{k=1}^{K-1} \Omega_k \Pr(x_i = 1 | d_i(p_{k+1}) > d_i(p_k)),$$

where weights  $\Omega_k$  are positive and sum to 1. Therefore, we would expect that  $\hat{\Delta}^x \in [0, 1]$ .

The informal monotonicity tests we conducted above ask whether some weighted average of  $\Pr(d_i(p_{k+1}) > d_i(p_k) | x_i = 1)$  is greater than 0. Since  $\Pr(x_i = 1) > 0$  and  $\Pr(d_i(p_{k+1}) > d_i(p_k)) > 0$ , the two conditions— $\Pr(d_i(p_{k+1}) > d_i(p_k) | x_i = 1) > 0$  and  $\Pr(x_i = 1 | d_i(p_{k+1}) > d_i(p_k)) > 0$ —are equivalent. Therefore, if we were to estimate Equations (D.6) and (D.7) by 2SLS, we would in essence be evaluating the same implication as the informal monotonicity tests standard in the literature.

In contrast, in a stylized representation of Section 4, we are performing 2SLS on the following equations:

$$m_i = \Delta d_i + u_i; \quad (\text{D.8})$$

$$d_i = \alpha Z_i + v_i. \quad (\text{D.9})$$

Recall that  $m_i = \mathbf{1}(d_i = 0, s_i = 1) = s_i(1 - d_i)$ . Following the same reasoning above, we can state the estimand  $\Delta$  as follows:

$$\Delta = - \sum_{k=1}^{K-1} \Omega_k \Pr(s_i = 1 | d_i(p_{k+1}) > d_i(p_k)),$$

which is a negative weighted average of conditional probabilities. This yields the same prediction that we stated in Remark 3 (i.e.,  $\Delta \in [-1, 0]$ ). As we discuss in Section 2.3, weaker conditions of monotonicity would leave this prediction unchanged.

More generally, we could apply the same reasoning to any binary potential outcome  $y_i(d) \in \{0, 1\}$  under treatment choice  $d \in \{0, 1\}$ . It is straightforward to show that, if we replace  $m_i$  with  $y_i d_i$  in Equation (D.8), the 2SLS system of Equations (D.8) and (D.9) would yield

$$\Delta = \sum_{k=1}^{K-1} \Omega_k \Pr(y_i(1) = 1 | d_i(p_{k+1}) > d_i(p_k)) \in [0, 1].$$

Alternatively, replacing  $m_i$  with  $-y_i(1 - d_i)$  in Equation (D.8) would imply

$$\Delta = \sum_{k=1}^{K-1} \Omega_k \Pr(y_i(0) = 1 | d_i(p_{k+1}) > d_i(p_k)) \in [0, 1].$$

How might we interpret our results together in Section 4 and in this appendix? We show above that the informal monotonicity tests are necessary for demonstrating that binary observable characteristics have admissible probabilities (i.e.,  $\Pr(x_i = 1) \in [0, 1]$ ) among compliers. On the other hand, our analysis in Section 4 strongly rejects that the key underlying state  $s_i$  has admissible probabilities among compliers. Specifically, our finding that  $\Delta \notin [-1, 0]$  is equivalent to showing that  $\Pr(s_i = 1) \notin [0, 1]$  among compliers, weighted by the probability that they contribute to the LATE. Observable characteristics may be correlated with  $s_i$ , but  $s_i$  is undoubtedly related to characteristics that are unobservable to the econometrician but, importantly, observable to radiologists. The importance of these unobservable characteristics will drive the difference between our analysis and the standard informal tests for monotonicity.

If monotonicity violations are more likely to occur between cases based on an underlying state than they to occur between cases based on observable characteristics, as would be plausible in classification decisions with variation in skill, then an analysis based on the underlying state should be stronger than an analysis based only on observable characteristics.

Finally, we note in Section 2.3 that our analysis in Section 4 is strongly connected to the conceptual intuition for testing IV validity described in Kitagawa (2015). Kitagawa (2015) shows that with data on treatment  $d_i$ , outcome  $y_i$ , and instrument  $Z_i$ , the *strongest* testable implication of IV validity is that potential outcomes should have positive density among compliers. Kitagawa (2015) and Mourifié and Wan (2017) extend this intuition when we also have access to some observable characteristic  $x_i$ . In this case, the implication of IV validity can be strengthened to requiring potential outcomes to have positive density among compliers *within each bin of  $x_i$* . Thus, to implement a stronger test of IV validity (including monotonicity), we could undertake a similar test of  $\Delta \in [-1, 0]$  using observations within each bin of  $x_i$ .

## E Details of Structural Analysis

### E.1 Optimal Diagnostic Thresholds

We provide a derivation of the optimal diagnostic threshold, given by Equation (7) in Section 5.1. We start with a general expression for the joint distribution of the latent index for each patient, or  $v_i$ , and radiologist signals, or  $w_{ij}$ . These signals determine each patient’s true disease status and diagnosis status:

$$\begin{aligned} s_i &= \mathbf{1}(v_i > \bar{v}); \\ d_{ij} &= \mathbf{1}(w_{ij} > \tau_j). \end{aligned}$$

We then form expectations of unconditional rates of false positives and false negatives, or  $FP_j \equiv \Pr(d_{ij} = 1, s_i = 0)$  and  $FN_j \equiv \Pr(d_{ij} = 0, s_i = 1)$ , respectively. Consider the radiologist-specific joint

distribution of  $(w_{ij}, v_i)$  as  $f_j(x, y)$ . Then

$$FN_j = \Pr(w_{ij} < \tau_j, v_i > \bar{v}) = \int_{-\infty}^{\tau_j} \int_{\bar{v}}^{+\infty} f_j(x, y) dy dx;$$

$$FP_j = \Pr(w_{ij} > \tau_j, v_i < \bar{v}) = \int_{\tau_j}^{+\infty} \int_{-\infty}^{\bar{v}} f_j(x, y) dy dx.$$

The joint distribution  $f_j(x, y)$  and  $\bar{v}$  are known to the radiologist. Given her expected utility function in Equation (6),

$$E[u_{ij}] = -(FP_j + \beta_j FN_j),$$

where  $\beta_j$  is the disutility of a false negative relative to a false positive, the radiologist sets  $\tau_j$  to maximize her expected utility.

The first order condition from expected utility is

$$-\frac{\partial FP_j}{\partial \tau_j} - \beta_j \frac{\partial FN_j}{\partial \tau_j} = 0.$$

Denote the marginal density of  $w_{ij}$  as  $g_j$ . Denote the conditional density of  $v_i$  given  $w_{ij}$  as  $f_j(y|x) = \frac{f_j(x, y)}{g_j(x)}$  and the conditional cumulative distribution as  $F_j(y|x) = \int_{-\infty}^y f_j(t|x) dt$ . Then solving this first order condition for the optimal threshold yields

$$\begin{aligned} -\frac{\partial FP_j}{\partial \tau_j} - \beta_j \frac{\partial FN_j}{\partial \tau_j} &= \int_{-\infty}^{\bar{v}} f_j(\tau_j, y) dy - \beta_j \int_{\bar{v}}^{+\infty} f_j(\tau_j, y) dy \\ &= \int_{-\infty}^{\bar{v}} f_j(y|\tau_j) g_j(\tau_j) dy - \beta_j \int_{\bar{v}}^{+\infty} f_j(y|\tau_j) g_j(\tau_j) dy \\ &= F_j(\bar{v}|\tau_j) g_j(\tau_j) - \beta_j (1 - F_j(\bar{v}|\tau_j)) g_j(\tau_j) = 0. \end{aligned}$$

The solution to the first order condition  $\tau_j^*$  satisfies

$$F_j(\bar{v}|\tau_j^*) = \frac{\beta_j}{1 + \beta_j}. \quad (\text{E.10})$$

Equation (E.10) can alternatively be stated as

$$\beta_j = \frac{F_j(\bar{v}|\tau_j^*)}{1 - F_j(\bar{v}|\tau_j^*)}.$$

This condition intuitively states that at the optimal threshold, the likelihood ratio of a false positive over a false negative is equal to the relative disutility of a false negative.

As a special case, when  $(w_{ij}, v_i)$  follows a joint-normal distribution, as in Equation (5), we know that  $v_i|w_{ij} \sim N(\alpha_j w_{ij}, 1 - \alpha_j^2)$ , or  $(v_i - \alpha_j w_{ij}) / \sqrt{1 - \alpha_j^2} | w_{ij} \sim N(0, 1)$ . This implies that

$F_j(\bar{v}|\tau_j^*) = \Phi\left(\left(\bar{v} - \alpha_j \tau_j^*\right) / \sqrt{1 - \alpha_j^2}\right)$ . Plugging in Equation (E.10) and rearranging, we obtain Equation (7):

$$\tau_j^*(\alpha_j, \beta_j) = \frac{\bar{v} - \sqrt{1 - \alpha_j^2} \Phi^{-1}\left(\frac{\beta_j}{1 + \beta_j}\right)}{\alpha_j}.$$

Below we verify that  $\partial^2 E[u_{ij}] / \partial \tau_j^2 < 0$  at  $\tau_j^*$  in a more general case, so  $\tau_j^*$  is the optimal threshold that maximizes expected utility.

### Comparative Statics

Returning to the general case, we need to impose a monotone likelihood ratio property to ensure that Equation (E.10) implies a unique solution and to analyze comparative statics.

**Assumption E.1 (Monotone Likelihood Ratio Property).** *The joint distribution  $f_j(x, y)$  satisfies*

$$\frac{f_j(x_2, y_2)}{f_j(x_2, y_1)} > \frac{f_j(x_1, y_2)}{f_j(x_1, y_1)}, \forall x_2 > x_1, y_2 > y_1, j.$$

We can rewrite the property using the conditional density:

$$\frac{f_j(y_2|x_2)}{f_j(y_1|x_2)} > \frac{f_j(y_2|x_1)}{f_j(y_1|x_1)}, \forall x_2 > x_1, y_2 > y_1, j.$$

That is, the likelihood ratio  $f_j(y_2|x_2) / f_j(y_1|x_2)$ , for  $y_2 > y_1$  and any  $j$ , always increases with  $x$ . In the context of our model, when a higher signal  $w_{ij}$  is observed, the likelihood ratio of a higher  $v_i$  over a lower  $v_i$  is higher than when a lower  $w_{ij}$  is observed. Intuitively, this means that the signal a radiologist receives is informative of the patient's true condition. As a special case, if  $f(x, y)$  is a bivariate normal distribution, the monotone likelihood ratio property is equivalent to a positive correlation coefficient.

Assumption E.1 implies *first-order stochastic dominance*. Fixing  $x_2 > x_1$  and considering any  $y_2 > y_1$ , Assumption E.1 implies

$$f_j(y_2|x_2) f_j(y_1|x_1) > f_j(y_2|x_1) f_j(y_1|x_2). \quad (\text{E.11})$$

Integrating this expression with respect to  $y_1$  from  $-\infty$  to  $y_2$  yields

$$\int_{-\infty}^{y_2} f_j(y_2|x_2) f_j(y_1|x_1) dy_1 > \int_{-\infty}^{y_2} f_j(y_2|x_1) f_j(y_1|x_2) dy_1.$$

Rearranging, we have

$$\frac{f_j(y_2|x_2)}{f_j(y_2|x_1)} > \frac{F_j(y_2|x_2)}{F_j(y_2|x_1)}, \forall y_2.$$

Similarly, integrating Equation (E.11) with respect to  $y_2$  from  $y_1$  to  $\infty$  yields

$$\int_{y_1}^{+\infty} f_j(y_2|x_2) f_j(y_1|x_1) dy_2 > \int_{y_1}^{+\infty} f_j(y_2|x_1) f_j(y_1|x_2) dy_2.$$

Rearranging, we have

$$\frac{1 - F_j(y_1|x_2)}{1 - F_j(y_1|x_1)} > \frac{f_j(y_1|x_2)}{f_j(y_1|x_1)}, \forall y_1.$$

Combining the two inequalities, we have

$$F_j(y|x_1) > F_j(y|x_2), \forall y. \quad (\text{E.12})$$

Under Equation (E.12), for a fixed  $\bar{v}$ ,  $F_j(\bar{v}|\tau_j)$  decreases with  $\tau$ , i.e.,  $\partial F_j(\bar{v}|\tau_j)/\partial \tau_j < 0$ . We can now verify that

$$\left. \frac{\partial^2 E[u_{ij}]}{\partial \tau_j^2} \right|_{\tau_j = \tau_j^*} = (1 + \beta_j) g_j(\tau_j^*) \left. \frac{\partial F_j(\bar{v}|\tau_j)}{\partial \tau_j} \right|_{\tau_j = \tau_j^*} < 0.$$

Therefore,  $\tau_j^*$  represents an optimal threshold that maximizes expected utility.

Using Equation (E.12) and the Implicit Function Theorem, we can also derive two reasonable comparative static properties of the optimal threshold. First,  $\tau_j^*$  decreases with  $\beta_j$ :

$$\left. \frac{\partial \tau_j^*}{\partial \beta_j} \right|_{\tau_j = \tau_j^*} = \frac{1}{(1 + \beta_j)^2} \left( \left. \frac{\partial F_j(\bar{v}|\tau_j)}{\partial \tau_j} \right|_{\tau_j = \tau_j^*} \right)^{-1} < 0.$$

Second,  $\tau_j^*$  increases with  $\bar{v}$ :

$$\left. \frac{\partial \tau_j^*}{\partial \bar{v}} \right|_{\tau_j = \tau_j^*} = -f_j(\bar{v}|\tau_j^*) \left( \left. \frac{\partial F_j(\bar{v}|\tau_j)}{\partial \tau_j} \right|_{\tau_j = \tau_j^*} \right)^{-1} > 0.$$

In other words, holding fixed the signal structure, a radiologist will increase her diagnosis rate when the relative disutility of false negatives increases and will decrease her diagnosis rate when pneumonia is less prevalent.

We next turn to analyzing the comparative statics of the optimal threshold with respect to skill. For a convenient specification with single-dimensional skill, we return to the specific case of joint-normal signals:

$$\begin{pmatrix} v_i \\ w_{ij} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \alpha_j \\ \alpha_j & 1 \end{pmatrix} \right).$$

Taking the derivative of the optimal threshold with respect to  $\alpha_j$  in Equation (7), we have

$$\frac{\partial \tau_j^*}{\partial \alpha_j} = \frac{\Phi^{-1}\left(\frac{\beta_j}{1+\beta_j}\right) - \bar{\nu} \sqrt{1 - \alpha_j^2}}{\alpha_j^2 \sqrt{1 - \alpha_j^2}}.$$

These relationships yield the following observations. When  $\alpha_j = 1$ ,  $\tau_j^* = \bar{\nu}$ . When  $\alpha_j = 0$ , the radiologist diagnoses no one if  $\beta_j < \frac{\Phi(\bar{\nu})}{1-\Phi(\bar{\nu})}$  (i.e.,  $\tau_j^* = \infty$ ), and the radiologist diagnoses everyone if  $\beta_j > \frac{\Phi(\bar{\nu})}{1-\Phi(\bar{\nu})}$  (i.e.,  $\tau_j^* = -\infty$ ). When  $\alpha_j \in (0, 1)$ , the relationship between  $\tau_j^*$  and  $\alpha_j$  depends on the prevalence parameter  $\bar{\nu}$ . Generally, if  $\beta_j$  is greater than some upper threshold  $\bar{\beta}$ ,  $\tau_j^*$  will always increase with  $\alpha_j$ ; if  $\beta_j$  is less than some lower threshold  $\underline{\beta}$ ,  $\tau_j^*$  will always decrease with  $\alpha_j$ ; if  $\beta_j \in (\underline{\beta}, \bar{\beta})$  is in between the lower and upper thresholds,  $\tau_j^*$  will first decrease then increase with  $\alpha_j$ . The thresholds for  $\beta_j$  depend on  $\bar{\nu}$ :

$$\begin{aligned} \underline{\beta} &= \min\left(\frac{\Phi(\bar{\nu})}{1-\Phi(\bar{\nu})}, 1\right); \\ \bar{\beta} &= \max\left(\frac{\Phi(\bar{\nu})}{1-\Phi(\bar{\nu})}, 1\right). \end{aligned}$$

The closer  $\bar{\nu}$  is to 0, the less space there will be between the thresholds. The range of  $\beta_j$  between the thresholds generally decreases as  $\bar{\nu}$  decreases.

Intuitively, there are two forces that drive the relationship between  $\tau_j^*$  and  $\alpha_j$ . First, the threshold of radiologists with low skill will depend on the overall prevalence of pneumonia. If pneumonia is uncommon, then radiologists with low skill will tend to diagnose fewer patients; if pneumonia is common, then radiologists with low skill will tend to diagnose more patients. Second, the threshold will depend on the relative disutility of false negatives,  $\beta_j$ . If  $\beta_j$  is high enough, then radiologists with lower skill will tend to diagnose more patients with pneumonia. Depending on the size of  $\beta_j$ , this mechanism may not be enough to have  $\tau_j^*$  always increasing in  $\alpha_j$ .

## E.2 Simulated Maximum Likelihood Estimation

In Section 5.2, we estimate the hyperparameter vector  $\theta \equiv (\mu_\alpha, \mu_\beta, \sigma_\alpha, \sigma_\beta, \lambda, \bar{\nu})$  by maximum likelihood:

$$\hat{\theta} = \arg \max_{\theta} \sum_j \log \int \mathcal{L}_j(\tilde{n}_j^d, \tilde{n}_j^m, n_j | \gamma_j) f(\gamma_j | \theta) d\gamma_j.$$

To calculate the radiologist-specific likelihood,

$$\mathcal{L}_j(\tilde{n}_j^d, \tilde{n}_j^m, n_j | \theta) = \int \mathcal{L}_j(\tilde{n}_j^d, \tilde{n}_j^m, n_j | \gamma_j) f(\gamma_j | \theta) d\gamma_j,$$

we need to evaluate the integral numerically. We approximate the integral using multiple-dimensional sparse grids as introduced in Heiss and Winschel (2008), which generates  $R$  nodes  $\gamma_j^r$  following the density  $f(\gamma_j | \theta)$ , given any hyperparameter vector  $\theta$ . These nodes are chosen based on Gaussian

quadratures and are assigned weights  $w^r$  such that  $\sum_r w^r = 1$ . We use a high accuracy level, which leads to  $R = 921$  nodes in a two-dimensional integral. Then we take the weighted average across all nodes of the likelihood as an approximation of the integral:

$$\mathcal{L}_j\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\theta}\right) \approx \sum_{r=1}^R w^r \mathcal{L}_j\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j^r\right).$$

The overall log-likelihood becomes

$$\log \mathcal{L}\left(\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j\right)_{j=1}^J \mid \boldsymbol{\theta}\right) \approx \sum_{j=1}^J \log \left( \sum_{r=1}^R w^r \mathcal{L}_j\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j^r\right) \right).$$

### E.3 Empirical Bayes Posterior Means

After estimating  $\hat{\boldsymbol{\theta}}$ , we want to find the empirical Bayes posterior mean  $\hat{\boldsymbol{\gamma}}_j = (\hat{\alpha}_j, \hat{\beta}_j)$  for each radiologist  $j$ . Using Bayes' theorem, the empirical conditional posterior distribution of  $\boldsymbol{\gamma}_j$  is

$$f\left(\boldsymbol{\gamma}_j \mid \tilde{n}_j^d, \tilde{n}_j^m, n_j; \hat{\boldsymbol{\theta}}\right) = \frac{f\left(\boldsymbol{\gamma}_j, \tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \hat{\boldsymbol{\theta}}\right)}{f\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \hat{\boldsymbol{\theta}}\right)} = \frac{f\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j\right) f\left(\boldsymbol{\gamma}_j \mid \hat{\boldsymbol{\theta}}\right)}{\int f\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j\right) f\left(\boldsymbol{\gamma}_j \mid \hat{\boldsymbol{\theta}}\right) d\boldsymbol{\gamma}_j},$$

where  $f\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j\right)$  is equivalent to  $\mathcal{L}_j\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j\right)$ . The denominator is then equivalent to the likelihood  $\mathcal{L}_j\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\theta}\right)$ . The empirical Bayes predictions are the following posterior means:

$$\hat{\boldsymbol{\gamma}}_j = \int \boldsymbol{\gamma}_j f\left(\boldsymbol{\gamma}_j \mid \tilde{n}_j^d, \tilde{n}_j^m, n_j; \hat{\boldsymbol{\theta}}\right) d\boldsymbol{\gamma}_j = \frac{\int \boldsymbol{\gamma}_j f\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j\right) f\left(\boldsymbol{\gamma}_j \mid \hat{\boldsymbol{\theta}}\right) d\boldsymbol{\gamma}_j}{\int f\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j\right) f\left(\boldsymbol{\gamma}_j \mid \hat{\boldsymbol{\theta}}\right) d\boldsymbol{\gamma}_j}.$$

As above, the integrals are evaluated numerically using sparse grids. We generate  $R$  nodes  $\boldsymbol{\gamma}_j^r$  following the density  $f\left(\boldsymbol{\gamma}_j \mid \hat{\boldsymbol{\theta}}\right)$  and calculate the empirical Bayes posterior means as

$$\hat{\boldsymbol{\gamma}}_j = \frac{\sum_{r=1}^R w^r \boldsymbol{\gamma}_j^r f\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j^r\right)}{\sum_{r=1}^R w^r f\left(\tilde{n}_j^d, \tilde{n}_j^m, n_j \mid \boldsymbol{\gamma}_j^r\right)}.$$

## F Robustness

In this appendix, we discuss alternative empirical implementations from the baseline approach. Appendix Table A.8 presents results for the following empirical approaches:

1. **Baseline.** This column presents results for the baseline empirical approach. This approach uses observations from all stations; the sample selection procedure is given in Appendix Table A.1. We risk-adjust diagnosis and false negative status by 77 patient characteristic variables,

described in Section 4.2, in addition to the controls for time dummies interacted with stations dummies required for plausible quasi-random assignment in Assumption 1. We define a false negative as a case that was not diagnosed initially with pneumonia but returned within 10 days and was diagnosed at that time with pneumonia.

2. **Balanced.** This approach modifies the baseline approach by restricting to 44 stations we select in Section 4.2 with stronger evidence for quasi-random assignment. Risk-adjustment and the definition of a false negative are unchanged from baseline.
3. **VA users.** This approach restricts attention to a sample of veterans who use VA care more than non-VA care. We identify this sample among dual enrollees in Medicare and the VA. We access both VA and Medicare records of care inside and outside the VA, respectively. We count the number of outpatient, ED, and inpatient visits in the VA and in Medicare, and keep veterans who have more total visits in the VA than in Medicare. The risk-adjustment and outcome definition are unchanged from baseline.
4. **Admission.** This approach redefines a false negative to only occur among patients with a greater than 50% predicted chance of admission. Patients with a lower predicted probability of admission are all coded to have  $m_i = 0$ . The sample selection and risk adjustment are the same as in baseline.
5. **Minimum controls.** This approach only controls for time dummies interacted with station dummies,  $\mathbf{T}_i$ , as specified by Assumption 1, without the 77 patient characteristic variables. The sample and outcome definition are unchanged from baseline.
6. **No controls.** This approach includes no controls. That is, we bypass the risk-adjustment procedure and use raw counts  $(n_j^d, n_j^m, n_j)$  in the likelihood, rather than the risk-adjusted counts  $(\tilde{n}_j^d, \tilde{n}_j^m, n_j)$ .
7. **Fix  $\lambda$ , flexible  $\rho$ .** This approach allows for flexible estimation of  $\rho$  in the structural model (whereas we assume that  $\rho = 0$  in the baseline structural model). Using results from our baseline estimation, we fix  $\lambda = 0.026$  instead.

## Rationale

Relative to the baseline approach, the “balanced” and “minimum controls” approaches respectively evaluate the importance of selecting stations with stronger evidence of quasi-random assignment and of controlling for rich patient observable characteristics. If results are robust under these approaches, then it is less likely that potential non-random assignment could be driving our results.

We evaluate results under the “VA users” approach in order to assess the potential threat that false negatives may be unobserved if patients fail to return to the VA. Although the process of returning to the VA is endogenous, it is only a concern under non-random assignment of patients to radiologists or under exclusion violations in which radiologists may influence the likelihood that a patient returns



to the VA, separate of incurring a false negative. Veterans who predominantly use the VA relatively to non-VA options are more likely to return to the VA for unresolved symptoms. Therefore, if results are robust under this approach, then exclusion violations and endogenous return visits are unlikely to explain our key findings.

Similarly, we assess an alternative definition of a false negative in the “admission” approach, requiring that patients are highly likely to be admitted as an inpatient based on their observed characteristics. Admitted patients have a built-in pathway for re-evaluation if signs and symptoms persist, worsen, or emerge; they need not decide to return to the VA. This approach also addresses a related threat that fellow ED radiologists may be more reluctant to contradict some radiologists than others, since admitted patients typically receive radiological evaluation from other divisions of radiology.

We take the “no controls” approach in order to assess the importance of linear risk-adjustment for our structural results. Although linear risk adjustment may be inconsistent with our nonlinear structural model, we expect that structural results should be qualitatively unchanged if risk-adjustment is relatively unimportant. In “fix  $\lambda$ , flexible  $\rho$ ,” we examine whether our structural model can rationalize the slight negative correlation between  $\alpha_j$  and  $\beta_j$  implied by the data in Appendix Figure A.13.

## Results

Appendix Table A.8 shows the robustness of key results under alternative implementations. Panel A reports sample statistics and reduced-form moments. All empirical implementations result in large variation in diagnosis and miss rates across radiologists. Standard deviations for both rates are weighted by the number of cases. The standard deviation of residual miss rates, after controlling for radiologist diagnosis rates, reveals that substantial heterogeneity in outcomes remains even after controlling for heterogeneity in decisions. This suggests violations, under all approaches, in the strict version of monotonicity in Condition 1(iii). Most importantly, the IV slope remains similarly positive across approaches. This suggests consistently strong violations in the weaker monotonicity conditions in Conditions 2-4.

Panel B of Appendix Table A.8 summarizes policy implications from decomposing variation into skill and preference components, as described in Section 6. In most implementations, more variation in diagnosis can be explained by heterogeneity in skill than by heterogeneity in preferences. An even larger proportion of variation in false negatives can be explained by heterogeneity in skill; essentially none of the variation in false negatives can be explained by heterogeneity in preferences.

Appendix Table A.9 shows corresponding structural model results under each of these alternative implementations. Panel A reports parameter estimates, and Panel B reports moments in the distribution of  $(\alpha_j, \beta_j)$  implied by the model parameters. The implementations again suggest qualitatively similar distributions of  $\alpha$ ,  $\beta$ , and  $\tau$ .

## G Extensions

### G.1 General Loss for False Negatives

Our baseline specification of utility in Equation (6) considers a fixed loss for any false negative relative to the loss for a false positive. In reality, some cases of pneumonia (e.g., those involving particularly virulent strains or vulnerable patients) may be much more costly to miss. In this appendix, we show that implications are qualitatively unchanged under a more general model with losses for false negatives that may be higher for these more severe cases.

We consider the following utility function:

$$u_{ij} = \begin{cases} -1, & \text{if } d_{ij} = 1, s_i = 0, \\ -\beta_j h(v_i), & \text{if } d_{ij} = 0, s_i = 1, \\ 0, & \text{otherwise,} \end{cases}$$

where  $h(v_i)$  is bounded, differentiable, and weakly increasing in  $v_i$ .<sup>2</sup> As before,  $s_i \equiv \mathbf{1}(v_i > \bar{v})$ , and  $\beta_j > 0$ . Without loss of generality, we assume  $h(\bar{v}) = 1$ , so  $h(v_i) \geq 1, \forall v_i$ .

Denote the conditional density of  $v_i$  given  $w_{ij}$  as  $f_j(v_i|w_{ij})$  and the corresponding conditional cumulative density as  $F_j(v_i|w_{ij})$ . Expected utility, conditional on  $w_{ij}$  and  $d_{ij} = 0$ , is

$$\begin{aligned} E_{v_i} [u_{ij}(v_i, d_{ij} = 0)|w_{ij}] &= -\beta_j E_{v_i} [h(v_i) \mathbf{1}(d_{ij} = 0, s_i = 1)|w_{ij}] \\ &= -\beta_j \int_{\bar{v}}^{+\infty} h(v_i) f_j(v_i|w_{ij}) dv_i. \end{aligned}$$

The corresponding expectation when  $d_{ij} = 1$  is

$$\begin{aligned} E_{v_i} [u_{ij}(v_i, d_{ij} = 1)|w_{ij}] &= -\Pr(s_i = 0, d_{ij} = 1|w_{ij}) \\ &= -\int_{-\infty}^{\bar{v}} f_j(v_i|w_{ij}) dv_i = \int_{\bar{v}}^{+\infty} f_j(v_i|w_{ij}) dv_i - 1. \end{aligned}$$

The radiologist chooses  $d_{ij} = 1$  if and only if  $E_{v_i} [u_{ij}(v_i, d_{ij} = 1)|w_{ij}] > E_{v_i} [u_{ij}(v_i, d_{ij} = 0)|w_{ij}]$ , or

$$\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i|w_{ij}) dv_i > 1.$$

If  $h(v_i) = 1$  for all  $v_i$ , then this condition reduces to  $\Pr(v_i > \bar{v}|w_{ij}) = 1 - F_j(\bar{v}|w_{ij}) > \frac{1}{1 + \beta_j}$ . In the general form, if the radiologist is indifferent in diagnosing or not diagnosing, we have

<sup>2</sup>The boundedness assumption ensures that the integrals below are well-defined. This is a sufficient condition but not necessary. The differentiability assumption simplifies calculation.

$$\begin{aligned}
1 &= \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i | w_{ij}) dv_i \\
&= \int_{\bar{v}}^{+\infty} (1 + \beta_j) f_j(v_i | w_{ij}) dv_i + \int_{\bar{v}}^{+\infty} \beta_j (h(v_i) - 1) f_j(v_i | w_{ij}) dv_i \\
&\geq (1 + \beta_j)(1 - F_j(\bar{v} | w_{ij})),
\end{aligned}$$

as we assume  $h(v_i) \geq 1$ . Now the marginal patient may have a lower conditional probability of having pneumonia than the case where  $h(v_i) = 1, \forall v_i$ , as false negatives may be more costly.

Define the optimal diagnosis rule as

$$d_j(w_{ij}) = \mathbf{1} \left( \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i | w_{ij}) dv_i > 1 \right).$$

Proposition G.5 shows conditions under which the optimal diagnosis rule satisfies the threshold crossing property.

**Proposition G.5.** *Suppose the following two conditions hold:*

1. *For any  $w'_{ij} > w_{ij}$ , the conditional distribution of  $v_i$  given  $\epsilon'_{ij}$  first-order dominates (FOSD) the conditional distribution of  $v_i$  given  $\epsilon_{ij}$ , i.e.,  $F_j(v_i | w'_{ij}) < F_j(v_i | w_{ij})$ ,  $\forall v_i$ ,*

2.  $0 < F_j(\bar{v} | w_{ij}) < 1, \forall w_{ij}$ .  $\lim_{w_{ij} \rightarrow -\infty} F_j(\bar{v} | w_{ij}) = 1$  and  $\lim_{w_{ij} \rightarrow +\infty} F_j(\bar{v} | w_{ij}) = 0$ .

*Then the optimal diagnosis rule satisfies the threshold-crossing property, i.e., for any radiologist  $j$ , there exists  $\tau_j^*$  such that*

$$d_j(w_{ij}) = \begin{cases} 0, & w_{ij} < \tau_j^*, \\ 1, & w_{ij} \geq \tau_j^*. \end{cases}$$

We first prove the following lemma.

**Lemma G.6.** *Suppose  $w'_{ij} > w_{ij}$ . If  $F_j(v_i | w'_{ij}) < F_j(v_i | w_{ij})$ , for each  $v_i$ , then  $d_j(w_{ij}) = 1$  implies  $d_j(w'_{ij}) = 1$ .*

*Proof.* Using integration by parts, we have

$$\begin{aligned}
&\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \left( f_j(v_i | w'_{ij}) - f_j(v_i | w_{ij}) \right) dv_i \\
&= (1 + \beta_j h(v_i)) \left( F_j(v_i | w'_{ij}) - F_j(v_i | w_{ij}) \right) \Big|_{\bar{v}}^{+\infty} - \int_{\bar{v}}^{+\infty} \beta_j h'(v_i) \left( F_j(v_i | w'_{ij}) - F_j(v_i | w_{ij}) \right) dv_i \\
&= -(1 + \beta_j) \left( F_j(\bar{v} | w'_{ij}) - F_j(\bar{v} | w_{ij}) \right) - \int_{\bar{v}}^{+\infty} \beta_j h'(v_i) \left( F_j(v_i | w'_{ij}) - F_j(v_i | w_{ij}) \right) dv_i > 0,
\end{aligned}$$

since  $F_j(v_i | w'_{ij}) < F_j(v_i | w_{ij}), \forall v_i$ ,  $h(v_i)$  is bounded,  $h(\bar{v}) = 1$ , and  $h'(v_i) \geq 0$ .

We now proceed to the proof of Proposition G.5. □

*Proof.* The second condition of Proposition G.5 ensures that

$$\lim_{w_{ij} \rightarrow -\infty} \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i | w_{ij}) dv_i \leq (1 + M\beta_j)(1 - \lim_{w_{ij} \rightarrow -\infty} F_j(\bar{v} | w_{ij})) = 0 < 1;$$

$$\lim_{w_{ij} \rightarrow +\infty} \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i | w_{ij}) dv_i \geq (1 + \beta_j)(1 - \lim_{w_{ij} \rightarrow +\infty} F_j(\bar{v} | w_{ij})) = 1 + \beta_j > 1,$$

where  $M = \sup h(v_i)$ . So  $\lim_{w_{ij} \rightarrow -\infty} d_j(w_{ij}) = 0$  and  $\lim_{w_{ij} \rightarrow +\infty} d_j(w_{ij}) = 1$ . Using Lemma G.6, the optimal diagnosis rule satisfies the threshold-crossing property. In particular, the optimal threshold  $\tau_j^*$  satisfies

$$\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i | \tau_j^*) dv_i = 1.$$

□

**Proposition G.7.** *Suppose the conditions in Proposition G.5 hold and  $f_j$  is fixed. Then the optimal threshold  $\tau_j^*$  decreases with  $\beta_j$ . In particular,  $\tau_j^* \rightarrow +\infty$  as  $\beta_j \rightarrow 0^+$  and  $\tau_j^* \rightarrow -\infty$  as  $\beta_j \rightarrow +\infty$ .*

*Proof.* Consider radiologists  $j$  and  $j'$  with  $\beta_j > \beta_{j'}$ . Denote their optimal thresholds as  $\tau_j^*$  and  $\tau_{j'}^*$ , respectively. We have  $\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i | \tau_j^*) dv_i = 1$  and

$$\begin{aligned} & \int_{\bar{v}}^{+\infty} (1 + \beta_{j'} h(v_i)) f_j(v_i | \tau_j^*) dv_i - \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i | \tau_j^*) dv_i \\ &= (\beta_{j'} - \beta_j) \int_{\bar{v}}^{+\infty} h(v_i) f_j(v_i | \tau_j^*) dv_i < 0. \end{aligned}$$

So  $\int_{\bar{v}}^{+\infty} (1 + \beta_{j'} h(v_i)) f_j(v_i | \tau_j^*) dv_i < 1$ , or  $d_{j'}(\tau_j^*) = 0$ . By Proposition G.5, we know that  $\tau_j^* < \tau_{j'}^*$ .

Since  $\tau_j^*$  decreases with  $\beta_j$ , if bounded below or above, it must have limits as  $\beta_j$  approaches  $+\infty$  or  $0^+$ . We can confirm that this is not the case. For example, suppose  $\tau_j^*$  is bounded below. The limit

exists and is denoted by  $\underline{\tau}$ . Take  $\beta_j \geq \frac{1}{1 - F(\bar{v} | \underline{\tau})}$ . Then

$$\begin{aligned} \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i | \tau_j^*) dv_i &\geq \left(1 + \frac{1}{1 - F(\bar{v} | \underline{\tau})}\right) (1 - F_j(\bar{v} | \tau_j^*)) \\ &> \left(1 + \frac{1}{1 - F(\bar{v} | \underline{\tau})}\right) (1 - F_j(\bar{v} | \underline{\tau})) = 2 - F_j(\bar{v} | \underline{\tau}). \end{aligned}$$

The second inequality holds since  $\tau_j^* > \underline{\tau}$ . Take the limit and we have

$$\lim_{\beta_j \rightarrow +\infty} \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) f_j(v_i | \tau_j^*) dv_i \geq 2 - F_j(\bar{v} | \underline{\tau}) > 1.$$

This is a contraction, so  $\tau_j^*$  is not bounded below. Similarly, we can show  $\tau_j^*$  is not bounded above. □

From now on, we assume  $w_{ij}$  and  $v_i$  follow a bivariate normal distribution:

$$\begin{pmatrix} w_{ij} \\ v_i \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \alpha_j \\ \alpha_j & 1 \end{pmatrix}\right).$$

Conditional on observing  $w_{ij}$ , the true signal  $v_i$  follows a normal distribution  $\mathcal{N}(\alpha_j w_{ij}, 1 - \alpha_j^2)$ . So

$$F_j(v_i|w_{ij}) = \Phi\left(\frac{v_i - \alpha_j w_{ij}}{\sqrt{1 - \alpha_j^2}}\right),$$

where  $\Phi(\cdot)$  is the CDF of the standard normal distribution.

**Corollary G.8.** *Suppose  $w_{ij}$  and  $v_i$  follow the bivariate normal distribution specified above. Then if  $\alpha_j > 0$ , the optimal diagnosis rule satisfies the threshold-crossing property.*

*Proof.* When  $w_{ij}$  and  $v_i$  follow the bivariate normal distribution with the correlation coefficient being  $\alpha_j$ , we have  $F_j(v_i|w_{ij}) = \Phi\left(\frac{v_i - \alpha_j w_{ij}}{\sqrt{1 - \alpha_j^2}}\right)$ . It is easy to verify that the two conditions in Proposition G.5 hold if  $\alpha_j > 0$ .

Define the optimal threshold  $\tau_j^* = \tau_j(\alpha_j, \beta_j; \bar{h}(\cdot))$  by

$$\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \frac{1}{\sqrt{1 - \alpha_j^2}} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i = 1,$$

where  $\phi(\cdot)$  is the density of the standard normal distribution. □

**Corollary G.9.** *The optimal threshold satisfies*

$$\frac{\bar{v} - \sqrt{1 - \alpha_j^2} \Phi^{-1}\left(\frac{\beta_j M}{1 + \beta_j M}\right)}{\alpha_j} \leq \tau_j^* \leq \frac{\bar{v} - \sqrt{1 - \alpha_j^2} \Phi^{-1}\left(\frac{\beta_j}{1 + \beta_j}\right)}{\alpha_j},$$

where  $M = \sup h(v_i)$ .

*Proof.* Since  $h(v_i) \geq 1$ , we have

$$\begin{aligned} 1 &= \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \frac{1}{\sqrt{1 - \alpha_j^2}} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i \\ &\geq (1 + \beta_j) \int_{\bar{v}}^{+\infty} \frac{1}{\sqrt{1 - \alpha_j^2}} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i \\ &= (1 + \beta_j) \left(1 - \Phi\left(\frac{\bar{v} - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right)\right). \end{aligned}$$

Rearrange and we can get the upper bound of  $\tau_j^*$ . Similarly, we can derive the lower bound of  $\tau_j^*$ .

The proposition below summarizes the relation between the general case and case where  $h(v_i) = 1, \forall v_i$ . □

**Proposition G.10.** Let  $\tau_j^* = \tau_j(\alpha_j, \beta_j; h(\cdot))$ . Define

$$\beta'_j = \beta'_j(\alpha_j, \beta_j; h(\cdot)) = \beta_j \frac{\int_{\bar{v}}^{+\infty} h(v_i) \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i}{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i}.$$

Then we can use the new  $\beta'_j$  to characterize the optimal threshold:

$$\tau_j(\alpha_j, \beta_j; h(\cdot)) = \tau_j(\alpha_j, \beta'_j; h(\cdot) = 1).$$

*Proof.* Let  $\tau_j^* = \tau_j(\alpha_j, \beta_j; h(\cdot))$  and  $\tau_j^{*'} = \tau_j(\alpha_j, \beta'_j; h(\cdot) = 1)$ . Then

$$\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \frac{1}{\sqrt{1 - \alpha_j^2}} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i = \int_{\bar{v}}^{+\infty} (1 + \beta'_j) \frac{1}{\sqrt{1 - \alpha_j^2}} \phi\left(\frac{v_i - \alpha_j \tau_j^{*'}}{\sqrt{1 - \alpha_j^2}}\right) dv_i = 1.$$

Substitute the expression of  $\beta'_j$  into the second equality and we have

$$\begin{aligned} & \int_{\bar{v}}^{+\infty} \left( 1 + \beta_j \frac{\int_{\bar{v}}^{+\infty} h(v_i) \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i}{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i} \right) \frac{1}{\sqrt{1 - \alpha_j^2}} \phi\left(\frac{v_i - \alpha_j \tau_j^{*'}}{\sqrt{1 - \alpha_j^2}}\right) dv_i = 1 \\ \Rightarrow & \int_{\bar{v}}^{+\infty} \frac{\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i}{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i} \frac{1}{\sqrt{1 - \alpha_j^2}} \phi\left(\frac{v_i - \alpha_j \tau_j^{*'}}{\sqrt{1 - \alpha_j^2}}\right) dv_i = 1 \\ \Rightarrow & \underbrace{\frac{1}{\sqrt{1 - \alpha_j^2}} \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i}_{=1} \frac{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^{*'}}{\sqrt{1 - \alpha_j^2}}\right) dv_i}{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i} = 1 \\ \Rightarrow & \int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^{*'}}{\sqrt{1 - \alpha_j^2}}\right) dv_i = \int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i. \end{aligned}$$

So we have  $\tau_j^{*'} = \tau_j^*$ . □

**Proposition G.11.** For fixed  $\beta_j$  and  $h(\cdot)$ ,  $\beta'_j = \beta'_j(\alpha_j, \beta_j; h(\cdot))$  decreases with  $\alpha_j$ .

*Proof.* The optimal threshold  $\tau_j^* = \tau_j(\alpha_j, \beta_j; h(\cdot))$  is given by

$$\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \frac{1}{\sqrt{1 - \alpha_j^2}} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i = 1.$$

By Proposition G.10, we can write

$$\begin{aligned} \beta'_j = \beta_j & \frac{\int_{\bar{v}}^{+\infty} h(v_i) \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i}{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i} = \frac{\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i) - 1) \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i}{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i} \\ & = \frac{\int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i - \int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i}{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i} = \frac{\sqrt{1 - \alpha_j^2}}{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i} - 1. \end{aligned}$$

Define  $x_i = \frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}$ . Then  $dv_i = \sqrt{1 - \alpha_j^2} dx_i$ . Using variable transformation, we have

$$\beta'_j = \frac{\sqrt{1 - \alpha_j^2}}{\int_{\bar{v}}^{+\infty} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i} - 1 = \frac{1}{1 - \Phi\left(\frac{\bar{v} - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right)} - 1.$$

Denote  $Q(v_i, \alpha_j, \beta_j) = \frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}$ . For fixed  $\beta_j$ , the relationship between  $\beta'_j$  and  $\alpha_j$  reduces the relationship between  $Q(\bar{v}, \alpha_j, \beta_j)$  and  $\alpha_j$ . Using integration by parts for the formula of the optimal threshold, we have

$$\begin{aligned} 1 & = \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \frac{1}{\sqrt{1 - \alpha_j^2}} \phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i = \int_{\bar{v}}^{+\infty} (1 + \beta_j h(v_i)) \frac{\partial \Phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right)}{\partial v_i} dv_i \\ & = (1 + \beta_j h(v_i)) \Phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) \Big|_{\bar{v}}^{+\infty} - \int_{\bar{v}}^{+\infty} \beta_j h'(v_i) \Phi\left(\frac{v_i - \alpha_j \tau_j^*}{\sqrt{1 - \alpha_j^2}}\right) dv_i \\ & = 1 + \beta_j M - (1 + \beta_j) \Phi(Q(\bar{v}, \alpha_j, \beta_j)) - \beta_j \int_{\bar{v}}^{+\infty} h'(v_i) \Phi(Q(v_i, \alpha_j, \beta_j)) dv_i, \end{aligned}$$

where  $M = \sup h(v_i)$ . Take the derivative with respect to  $\alpha_j$ ,

$$0 = -(1 + \beta_j)\phi(Q(\bar{v}, \alpha_j, \beta_j))\frac{\partial Q(\bar{v}, \alpha_j, \beta_j)}{\partial \alpha_i} - \beta_j \int_{\bar{v}}^{+\infty} h'(v_i)\phi(Q(v_i, \alpha_j, \beta_j))\frac{\partial Q(v_i, \alpha_j, \beta_j)}{\partial \alpha_j} dv_i. \quad (\text{G.13})$$

We want to show that  $\frac{\partial Q(\bar{v}, \alpha_j, \beta_j)}{\partial \alpha_i} \leq 0$  for all  $\alpha_j \in (0, 1)$ . We prove this by contradiction. Assume

that for some  $\alpha'_j \in (0, 1)$ , we have  $\left. \frac{\partial Q(\bar{v}, \alpha_j, \beta_j)}{\partial \alpha_i} \right|_{\alpha_j = \alpha'_j} > 0$ . Since  $\frac{\partial^2 Q(v_i, \alpha_j, \beta_j)}{\partial \alpha_j \partial v_i} = \frac{\alpha_j}{(1 - \alpha_j)^{3/2}} > 0$ ,

we know that  $\frac{\partial Q(\bar{v}, \alpha_j, \beta_j)}{\partial \alpha_i}$  increases with  $v_i$  for any fixed  $\alpha_j \in (0, 1)$ , in particular for  $\alpha_j = \alpha'_j$ . Then  $\left. \frac{\partial Q(v_i, \alpha_j, \beta_j)}{\partial \alpha_i} \right|_{\alpha_j = \alpha'_j} \geq \left. \frac{\partial Q(\bar{v}, \alpha_j, \beta_j)}{\partial \alpha_i} \right|_{\alpha_j = \alpha'_j} > 0$  for any  $v_i \geq \bar{v}$ . Since  $h'(v_i) \geq 0$ , we have

$$\left. \frac{\partial Q(\bar{v}, \alpha_j, \beta_j)}{\partial \alpha_i} \right|_{\alpha_j = \alpha'_j} > 0, \int_{\bar{v}}^{+\infty} h'(v_i)\phi(Q(v_i, \alpha_j, \beta_j))\frac{\partial Q(v_i, \alpha_j, \beta_j)}{\partial \alpha_j} dv_i \Big|_{\alpha_j = \alpha'_j} \geq 0.$$

Then Equation (G.13) cannot hold for  $\alpha_j = \alpha'_j$ , as the right hand is strictly negative, a contradiction.

So, we must have  $\frac{\partial Q(\bar{v}, \alpha_j, \beta_j)}{\partial \alpha_i} \leq 0, \forall \alpha_j \in (0, 1)$ . Therefore,

$$\frac{\partial \beta'_j}{\partial \alpha_j} = \frac{\phi(Q(\bar{v}, \alpha_j, \beta_j))\frac{\partial Q(\bar{v}, \alpha_j, \beta_j)}{\partial \alpha_j}}{(1 - \Phi(Q(\bar{v}, \alpha_j, \beta_j)))^2} \leq 0.$$

□

## G.2 Incorrect Beliefs

Under the model of radiologist signals implied by Equation (5), we can identify each radiologist's skill  $\alpha_j$  and her diagnostic threshold  $\tau_j$ . The utility in Equation (6) implies the optimal threshold in Equation (7), as a function of skill  $\alpha_j$  and preference  $\beta_j$ . If radiologists know their skill, then this allows us to infer  $\beta_j$  from  $\alpha_j$  and  $\tau_j$ .

In this appendix, we allow for the possibility that radiologists may be misinformed about their skill: A radiologist may believe she has skill  $\alpha'_j$  even though her true skill is  $\alpha_j$ . Since only (true)  $\alpha_j$  and  $\tau_j$  are identified, we cannot separately identify  $\alpha'_j$  and  $\beta_j$  from Equation (7). In this exercise, we therefore assume  $\beta_j$ , in order to infer  $\alpha'_j$  for each radiologist.

We start with our baseline model and form an empirical Bayes posterior mean of  $(\alpha_j, \beta_j)$  for each radiologist. We use Equation (7) to impute the empirical Bayes posterior mean of  $\tau_j$ . Thus, for each radiologist, we have an empirical Bayes posterior mean of  $(\alpha_j, \beta_j, \tau_j)$  from our baseline model; the distributions of the posterior means for  $\alpha_j$ ,  $\beta_j$ , and  $\tau_j$  are shown in separate panels of Appendix Figure A.13.

To extend this analysis to impute each radiologist's belief about her skill,  $\alpha'_j$ , we perform the



following two additional steps: First, we take the mean of the distribution of empirical Bayes posterior means  $\{\beta_j\}_{j \in \mathcal{J}}$ , which we calculate as 6.71. Second, we set all radiologists to have  $\beta_j = 6.71$ . We use each radiologist's empirical Bayes posterior mean of  $\tau_j$  and the formula for the optimal threshold in Equation (7) to infer her belief about her skill,  $\alpha'_j$ .

The relationship between  $\alpha'_j$ ,  $\beta_j$ , and  $\tau_j$  is shown in Figure IX. As shown in the figure, for  $\beta_j = 6.71$ , the comparative statics of  $\tau_j^*$  are first decreasing and then increasing with a radiologist's perceived  $\alpha'_j$ . Thus, holding fixed  $\beta_j = 6.71$ , an observed  $\tau_j$  does not generally imply a single value of  $\alpha'_j$ . If  $\tau_j$  is too low, then there will not be a value of  $\alpha'_j$  to generate  $\tau_j$  with  $\beta_j = 6.71$ ; this case occurs only for a minority of radiologists. Other  $\tau_j$  generally can be consistent with either a value of  $\alpha'_j$  on the downward-sloping part of the curve or with a value of  $\alpha'_j$  on the upward-sloping part of the curve. In this case, we take the higher value of  $\alpha'_j$ , since the vast majority of empirical Bayes posterior means of  $\alpha_j$  are on the upward-sloping part of Figure IX.

Appendix Figure A.19 plots each radiologist's perceived skill, or  $\alpha'_j$ , on the  $y$ -axis and her actual skill, or  $\alpha_j$ , on the  $x$ -axis. The plot shows that the radiologists' perceptions of their skill generally correlate well with their actual skill, particularly among higher-skilled radiologists. Lower-skilled radiologists, however, tend to over-estimate their skill relative to the truth.

### G.3 Simulation of Linear Risk Adjustment

As described in Section 5.2, we estimate our structural model using moments for each radiologist that are risk-adjusted by linear regressions. An alternative approach would be to explicitly incorporate heterogeneity in  $\Pr(s_i = 1)$ , by station, time, and patient characteristics, into the structural model. While this approach is more consistent with the structural model, it is often computationally prohibitive.

In this appendix section, we use Monte Carlo simulations to examine the effectiveness of linear risk adjustment in recovering the underlying structural parameters of our model. Specifically, we fix the set of radiologists at each station and the number of patients that each radiologist examines, or  $n_j$ , to match the actual data. Assuming that parameter estimates in Table I are the truth, we simulate primitives  $\{\alpha_j, \beta_j\}_{j \in \mathcal{J}}$  independent of  $n_j$ . We also simulate at-risk patients from a binomial distribution with the probability of being at risk of  $1 - \kappa$ .

For patients at risk, we simulate their latent index  $v_i$  and the radiologist-observed signal  $w_{ij}$  using  $\alpha_j$  of the assigned radiologist  $j$ . Importantly, in this simulation, we model *conditional* random assignment of patients to radiologists within station. For  $v_i$  and  $w_{ij}$  that are jointly normally distributed, as in Equation (5),

$$\begin{pmatrix} v_i \\ w_{ij} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \alpha_j \\ \alpha_j & 1 \end{pmatrix} \right),$$

we have

$$s_i = \mathbf{1}(v_i > \bar{v}_{\ell(j)}),$$

where  $\bar{v}_{\ell(j)}$  depends on the station  $\ell(j)$  in which radiologist  $j$  works. Radiologists know  $\bar{v}_{\ell(j)}$ . The

optimal threshold is then

$$\tau^*(\alpha_j, \beta_j; \ell(j)) = \frac{\bar{v}_{\ell(j)} - \sqrt{1 - \alpha_j^2} \Phi^{-1}\left(\frac{\beta_j}{1 + \beta_j}\right)}{\alpha_j},$$

which generates  $d_{ij} = \mathbf{1}(w_{ij} > \tau^*(\alpha_j, \beta_j; \ell(j)))$ . We finally simulate patients who did not initially have pneumonia but later developed it with  $\lambda$ .

Each simulated dataset has the same number of observations as in the original dataset, with four variables for each patient  $i$ : the radiologist identifier  $j$ , the station identifier  $\ell$ , the diagnosis indicator  $d_i = \sum_j \mathbf{1}(j = j(i)) d_{ij}$ , and the (observed) false negative indicator  $m_i = \mathbf{1}(d_i = 0, s_i = 1)$ . We obtain risk-adjusted radiologist moments from the simulated data by regressing diagnosis or false negative indicators on radiologist dummies and station dummies.

The key object of confounding risk across groups of observations is the distribution of  $\bar{v}_\ell$ . We assume that this distribution is normal and calibrate its standard deviation based on the following target: the ratio of the standard deviation of unadjusted radiologist diagnosis rates to the standard deviation of adjusted radiologist diagnosis rates. In the actual data, these standard deviations are shown in Appendix Table A.8, as 1.966 and 1.023, respectively. Conceptually, the ratio of these standard deviations captures the net effect of risk adjustment on reduced-form radiologist diagnosis rates. In each of five simulated datasets, we calculate a similar ratio. In our calibration, we aim to match the average of these ratios across the five simulations, holding the random-generating seed fixed in each simulation.

In each of the simulations, we redo three sets of results based on unadjusted or adjusted radiologist moments. First, we re-estimate the model parameters. Second, we re-compute counterfactual variation in diagnoses and false negatives when either variation in skill or variation in preferences is eliminated, as described in Section 6.1. Third, we re-compute welfare under policy counterfactuals, as described in Section 6.2. As shown in Appendix Figure A.20, the results of this exercise suggest that linear risk adjustment eliminates most of the bias due to confounding variation in risk across groups of observations. For many estimated parameters and counterfactual results, the bias is almost eliminated by linear risk adjustment.

#### G.4 Controlling for Radiologist Skill

Intuitively, monotonicity should hold within bins of skill. In this appendix section, we explore a Monte Carlo proof of concept for whether controlling for agent skill in a judges-design regression can recover complier-weighted treatment effects. Specifically, we simulate data that match our observed data, taking structural estimates as the truth. We then evaluate whether we can recover the complier-weighted “treatment effect,” or  $-\Pr(s = 1)$  in our case, that one should obtain under IV validity when regressing  $m_i$  on  $d_i$ , instrumenting  $d_i$  with  $Z_i$ .

As in Appendix G.3, we take parameter estimates in Table I as the truth and simulate true primitives  $\{\alpha_j, \beta_j\}_{j \in \mathcal{J}}$ . We similarly fix observations per radiologist and simulate patients at risk. Among

these patients, we simulate  $v_i$  and  $w_{ij}$ . We determine which patients are diagnosed with pneumonia and which patients are false negatives based on  $\tau_j^*(\alpha_j, \beta_j)$ , in Equation (7), and  $\bar{v}$ . This implies that, unlike the simulations in Appendix G.3, patients are unconditionally randomly assigned. Finally, we simulate patients who did not initially have pneumonia but later developed it with  $\lambda$ .

In the remainder of this appendix section, we will derive the target LATE and then compare whether we can estimate it using various strategies to control for skill.

**Derivation of the Properly Specified Estimand.** The ideal experiment would be to compare radiologists with the same  $\alpha_j$ . However, we have a continuous distribution of  $\alpha_j$  and a finite number of radiologists. We therefore derive an approximation of the true relationship between  $FN_j^{\text{obs}}$  and  $P_j^{\text{obs}}$ , conditional on skill  $\alpha_j$ , under a large number of radiologists with the same skill and a large number of patients per radiologist. We then integrate this approximation over the distribution of skill.

Specifically,

$$P_j^{\text{obs}}(\alpha_j, \beta_j) = (1 - \kappa) \Pr(w_{ij} > \tau_j^*) = (1 - \kappa) \left(1 - \Phi(\tau_j^*)\right); \quad (\text{G.14})$$

$$FN_j^{\text{obs}}(\alpha_j, \beta_j) = (1 - \kappa) \left(\Pr(w_{ij} < \tau_j^*, v_i > \bar{v} \mid \alpha_j) + \lambda \Pr(w_{ij} < \tau_j^*, v_i < \bar{v} \mid \alpha_j)\right), \quad (\text{G.15})$$

where  $\tau_j^* = \tau^*(\alpha_j, \beta_j)$  in Equation (7). Conditional on  $\alpha_j$ , there exists a one-to-one mapping in the reduced-form space between  $FN_j^{\text{obs}}$  and  $P_j^{\text{obs}}$ .

Conditional on the realization of skill  $\alpha$ , we draw  $J + 1$  radiologists with varying  $\beta_j$  from the true distribution and derive their optimal thresholds  $\tau_j^*$ . We calculate their population diagnosis and miss rates as  $p_j = E[d_i \mid j(i) = j] = P_j^{\text{obs}}(\alpha_j, \beta_j)$  and  $\bar{m}_j = E[m_i \mid j(i) = j] = FN_j^{\text{obs}}(\alpha_j, \beta_j)$ , respectively. We consider the LATE when we use  $p_j$  as the scalar instrument for diagnosis  $d_i$ . We rank radiologists based on  $p_j$  from smallest to largest, so that  $p_0 < p_1 < \dots < p_J$ . From Theorem 2 of Imbens and Angrist (1994), the LATE conditional on skill  $\alpha$  is

$$\Delta^*(\alpha) = \sum_{j=1}^J \psi_j \delta_{j,j-1},$$

where

$$\psi_j = \frac{(p_j - p_{j-1}) \sum_{l=j}^J \rho_l (p_l - \bar{p})}{\sum_{m=1}^J (p_m - p_{m-1}) \sum_{l=j}^J \rho_l (p_l - \bar{p})},$$

$$\delta_{j,j-1} = \frac{\bar{m}_j - \bar{m}_{j-1}}{p_j - p_{j-1}}.$$

$\psi_j$  is a non-negative weight, which depends on the first-stage difference in diagnosis rates between radiologists and the probability of assignment to  $j$ , or  $\rho_j$ .  $\delta_{j,j-1}$  is the Wald estimand based on random assignment between  $j$  and  $j - 1$ . Note that  $\rho_j = (J + 1)^{-1}$  for all  $j$ , by random assignment, and  $\bar{p} = \frac{1}{J+1} \sum_{j=0}^J p_j$ .

We then simulate  $K$  values of  $\alpha_k$  from the true distribution to derive the LATE (unconditional on

skill) as

$$\Delta^* = \frac{1}{K} \sum_{k=1}^K \Delta^*(\alpha_k).$$

We choose reasonably large  $J = 1,000$  and  $K = 1,000$ . This can be seen as the approximation of the expectation of the LATE across many realizations of skill. We compute  $\Delta^* = -0.154$ .

**Estimation Results.** We then estimate the effect of diagnosis  $d_i$  on the false negative indicator  $m_i$  and present results in Appendix Table A.11. As in the main text, we estimate this effect by judges-design IV, exploiting the relationship between radiologist diagnosis and miss rates.

The standard specification is shown in Column 1 of all panels. Specifically, we perform 2SLS of  $m_i$  on  $d_i$ , instrumenting  $d_i$  by the leave-out diagnosis propensity  $Z_i$ , given in Equation (4). Since cases are randomly assigned unconditionally in this simulation, we include no further controls. This result is significantly positive, at 0.096, despite the true negative LATE of  $\Delta^* = -0.154$ .

In Panel A, we show results of regressions that control for true skill,  $\alpha_j$ . For Column 2 of this panel, we control for  $\alpha_j$  linearly in the 2SLS regression. For Columns 3-6, we divide  $\alpha_j$  into 5, 10, 20, and 50 bins, respectively, and include indicators for bins of  $\alpha_j$  as controls in the regression. The results in these columns encompass the true LATE.

In Panel B, we show results of similar regressions that replace functions of true skill  $\alpha_j$  with corresponding functions of the empirical Bayes posterior mean of  $\alpha_j$ , or  $\hat{\alpha}_j$ . Specifically, for Column 2, we control for  $\hat{\alpha}_j$  linearly; for Columns 3-6, we divide  $\hat{\alpha}_j$  into 5, 10, 20, and 50 bins, respectively, and include indicators for bins of  $\alpha_j$  as controls in the regression. To account for the fact that  $\hat{\alpha}_j$  is a generated regressor, we construct standard errors by 50 bootstrapped samples, drawing observations by radiologist with replacement and keeping the total number of radiologists fixed. These results are also strongly negative, but they are more negative than the true LATE. The confidence intervals are also substantially wider.

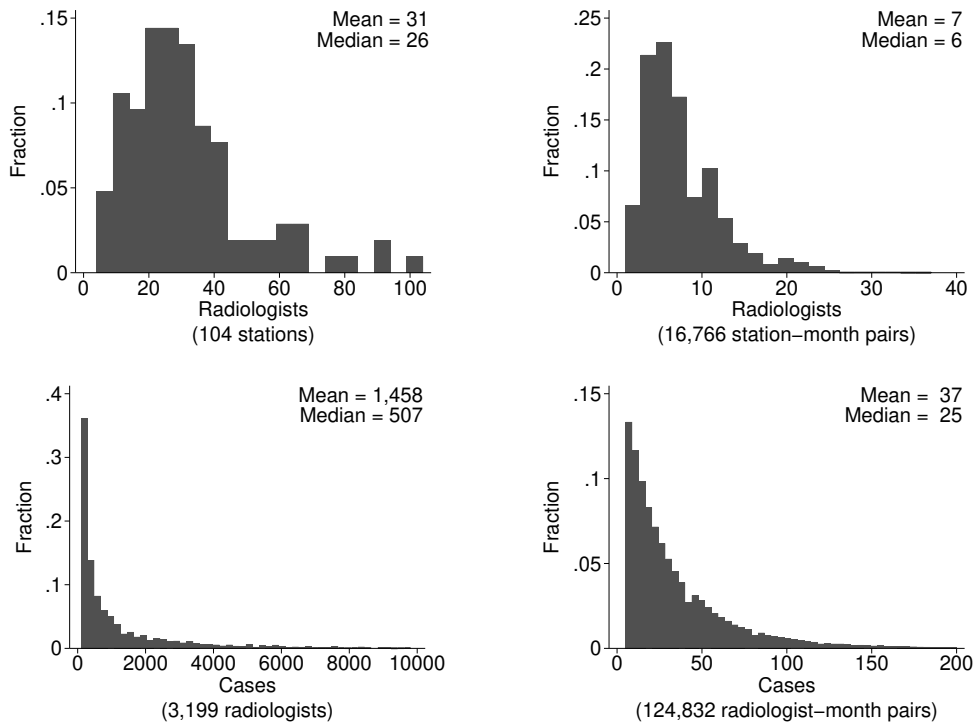
In Panel C, we show results from indirect least squares regressions of  $m_i$  on empirical Bayes posteriors of  $P_j$  and  $\alpha_j$ . For Column 2, we control for the posterior mean  $\hat{\alpha}_j$  linearly; for Columns 3-6, we control for posterior probabilities that  $\alpha_j$  resides in each of 5, 10, 20, and 50 bins, respectively. We construct standard errors by the same bootstrap procedure that we use for Panel B. The estimates of the LATE are negative and less biased than in Panel B. Nevertheless, they are still generally larger in magnitude than the true LATE.

These results suggest that we can recover the true LATE when we control for true skill. However, estimates are biased, albeit in the opposite direction in our simulation, when we use empirical Bayes posteriors of skill. In Appendix Figure A.21, we confirm that estimates from regressions that use empirical Bayes posteriors for radiologists with a very large number of cases approach the true LATE. Even so, the number of cases per radiologist is already high in our simulated sample. By construction, each radiologist has at least 100 cases, and we match the distribution of cases for each radiologist to the actual distribution, shown in Appendix Figure A.1. We leave further refinement of this approach in finite samples to future work.

## References

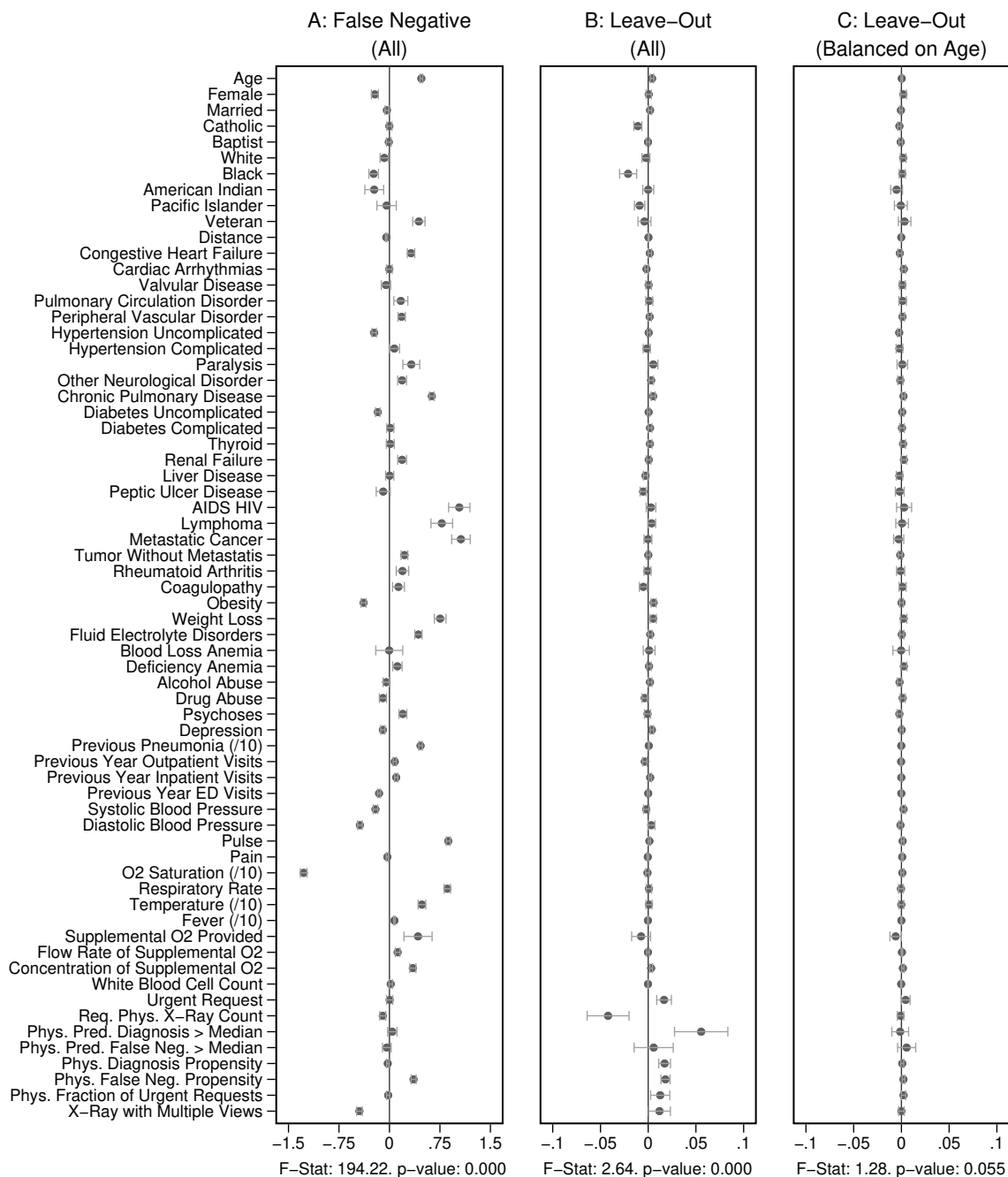
ANDREWS, M. J., L. GILL, T. SCHANK, AND R. UPWARD (2008): “High Wage Workers and Low Wage Firms: Negative Assortative Matching or Limited Mobility Bias?” *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 171, 673-697.

Figure A.1: Distribution of Radiologists and Cases



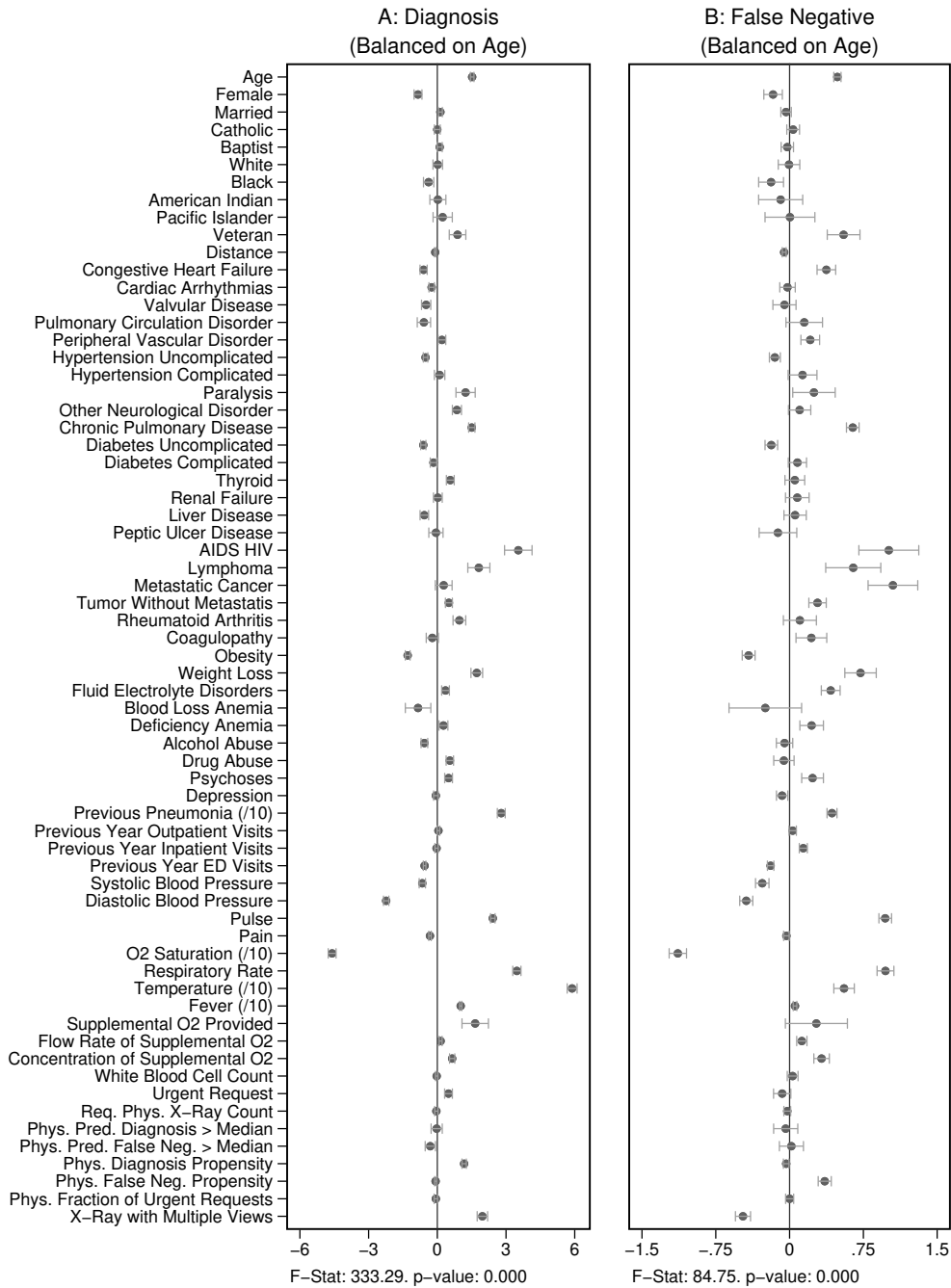
*Note:* This figure shows the distributions of radiologists across stations, of radiologists across station-months, of cases across radiologists, and of cases across radiologist-months. As shown in Appendix Table A.1, the minimum number of cases for a radiologist is 100, and the minimum number of cases for a radiologist-month pair is 5. In this figure, we truncate the number of cases per radiologist at 10,000; 57 radiologists, or 1.78% of the total, have more cases than this limit. We truncate the number of cases per radiologist-month at 200; 1,274 radiologist-months, or 1.02% of the total, have more cases than this limit.

Figure A.2: Covariate Balance (Miss Rate)



Note: This figure shows coefficients and 95% confidence intervals from regressions of the false-negative indicator  $m_i$  (left column) or the assigned radiologist's leave-out miss rate (middle and right columns) on covariates  $\mathbf{X}_i$ , controlling for time-station interactions  $\mathbf{T}_i$ . The 66 covariates are the variables listed in Appendix A.2, less the 11 variables that are indicators for missing values. The leave-out miss rate is calculated analogously to the leave-out diagnosis propensity  $Z_i$ . The left and middle panels use the full sample of stations. The right panel uses 44 stations with balance on age, defined in Section 4.2. The outcome variables are multiplied by 100. Continuous covariates are standardized so that they have standard deviations equal to 1. For readability, a few coefficients (and their standard errors) are divided by 10, as indicated by "/10" in the covariate labels. At the bottom of each panel, we report the  $F$ -statistic and  $p$ -value from the joint  $F$ -test of all covariates.

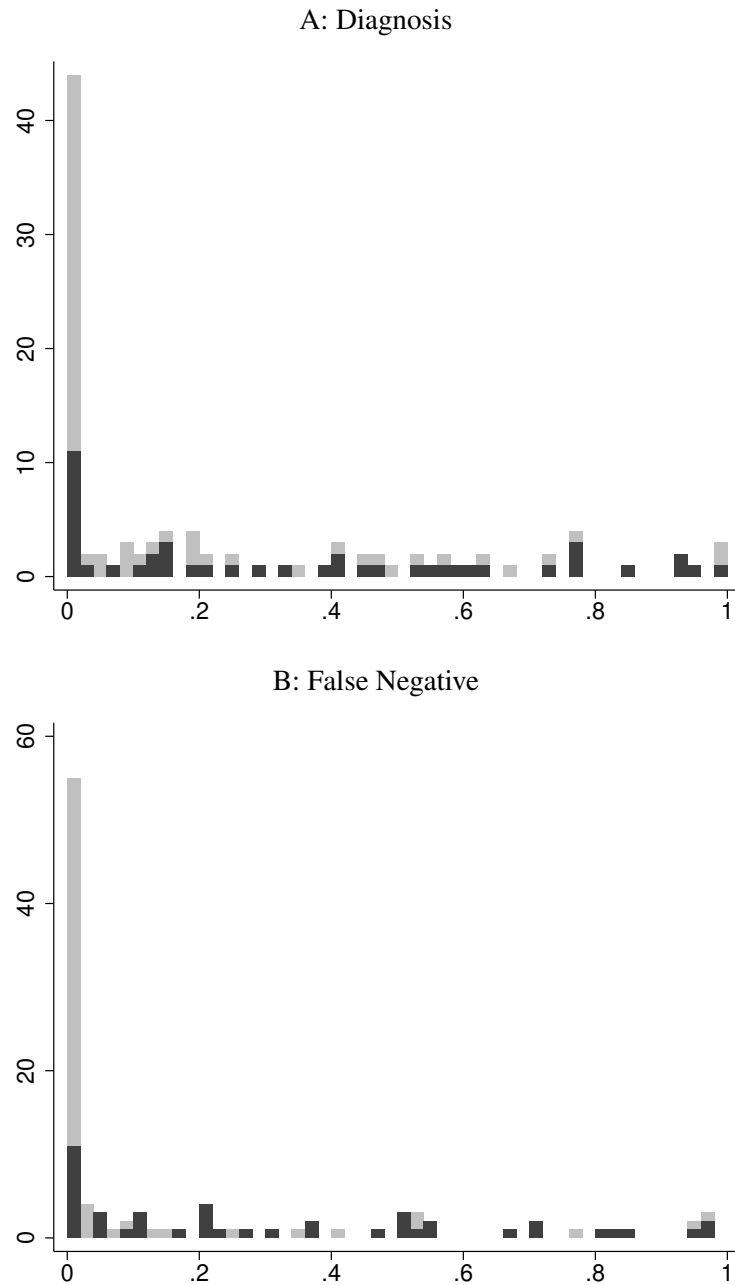
Figure A.3: Predicting Diagnosis and False Negatives (Stations with Balance on Age)



Note: This figure shows coefficients and 95% confidence intervals from regressions of diagnosis status  $d_i$  (left column) or the false negative indicator  $m_i$  (right column) on covariates  $\mathbf{X}_i$ , controlling for time-station interactions  $\mathbf{T}_i$  in the sample of 44 stations with balance on age (defined in Section 4.2). This is analogous to the left-hand columns of Figure VI and Appendix Figure A.2 respectively, with the restricted sample of stations. The outcome variables are multiplied by 100. The 66 covariates are the variables listed in Appendix A.2, less the 11 variables that are indicators for missing values. Continuous covariates are standardized so that they have standard deviations equal to 1. For readability, a few coefficients (and their standard errors) are divided by 10, as indicated by “/10” in the covariate labels. At the bottom of each panel, we report the  $F$ -statistic and  $p$ -value from the joint  $F$ -test of all covariates.

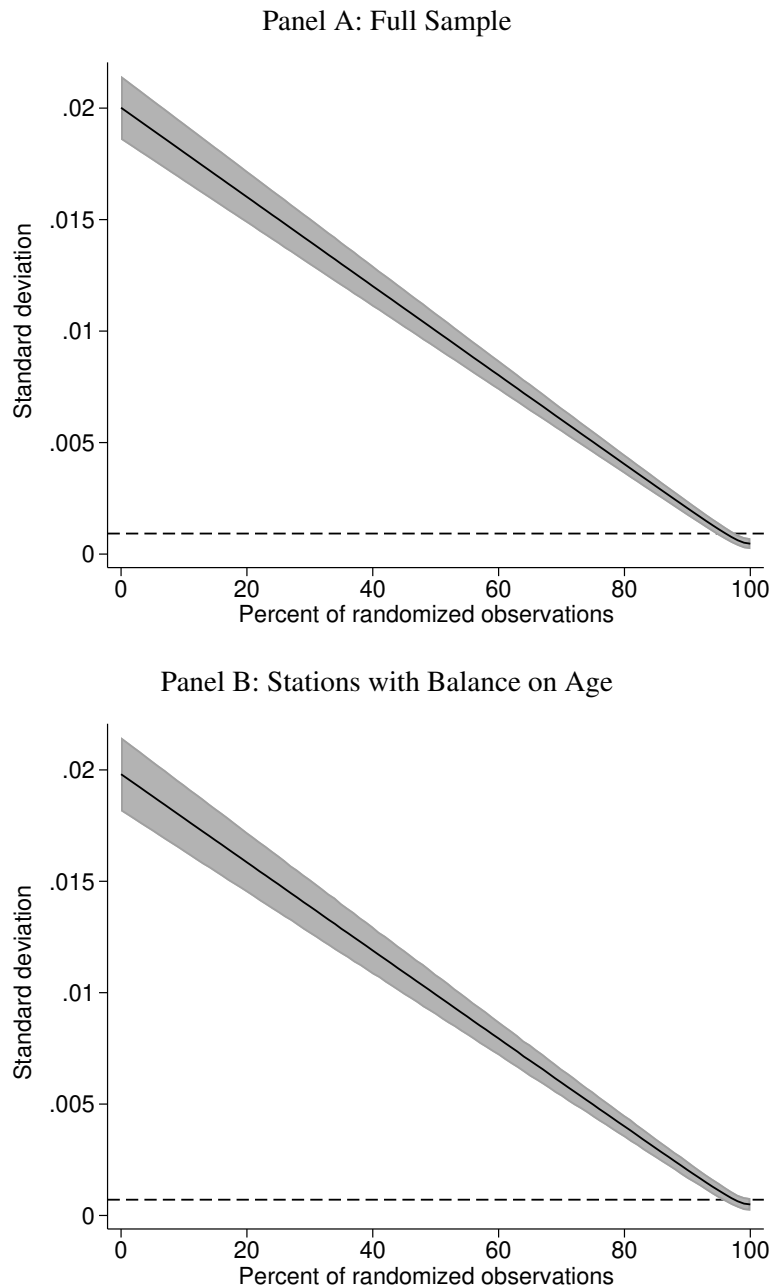


Figure A.4: Randomization Inference



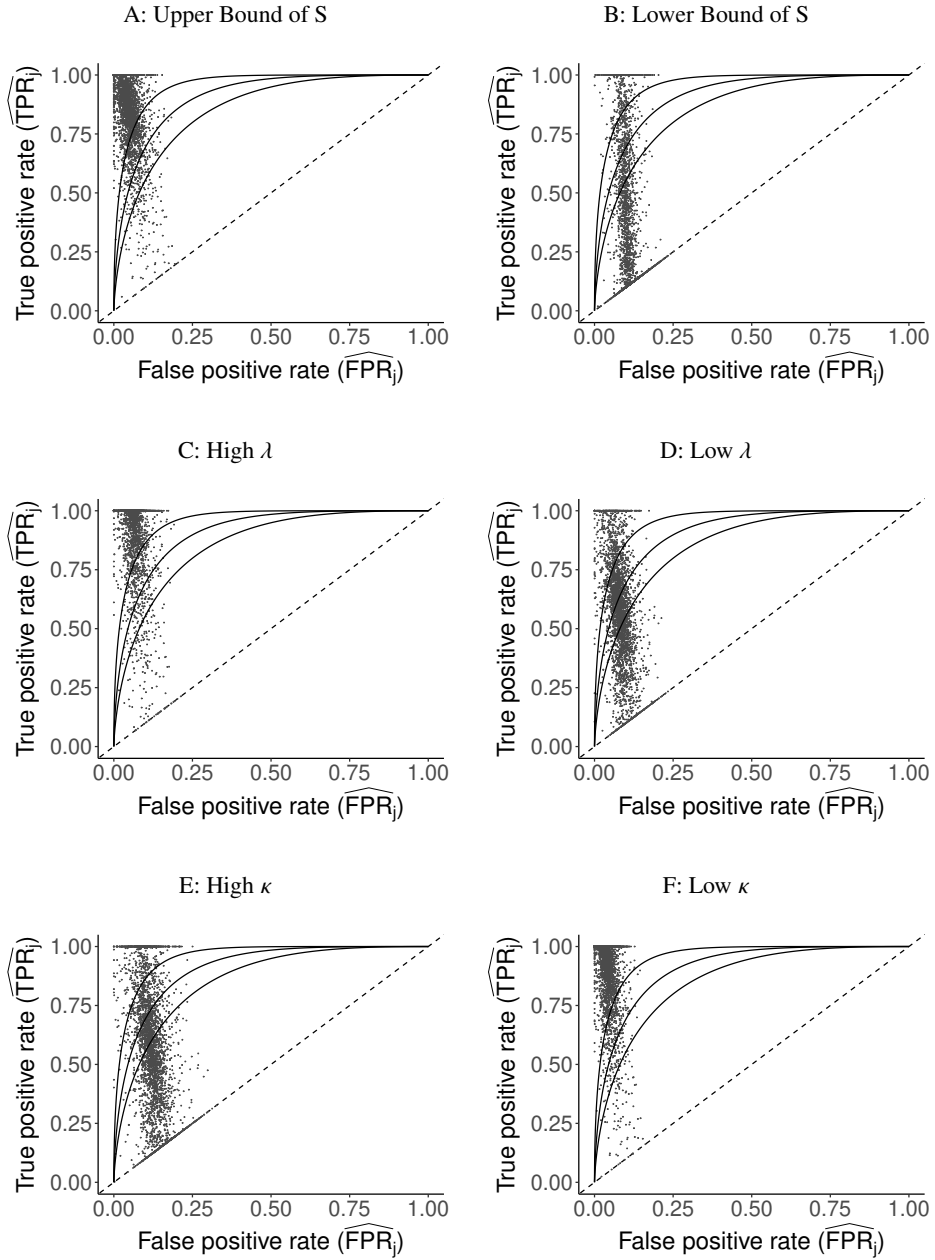
*Note:* This figure plots histograms of station-level  $p$ -values for quasi-random assignment computed using randomization inference. We first residualize predicted diagnosis and false negative indicators  $\hat{d}_i$  and  $\hat{m}_i$  by minimal controls  $\mathbf{T}_i$ . We then create 100 samples in each of which we randomly reassign the residualized values to patients within each station. For each of these samples as well as the baseline sample we regress the residualized values on radiologist dummies, and calculate the case-weighted standard deviation of estimated radiologist fixed effects. We then define the  $p$ -value for each station to be the share of the 100 samples that yield a larger standard deviation than the baseline sample. In each panel, light gray bars represent station counts among the 60 stations that fail the test according to age; dark gray bars represent station counts out of the 44 stations that pass the test according to age.

Figure A.5: Variation in Radiologist Miss Rates Under Counterfactual Sorting



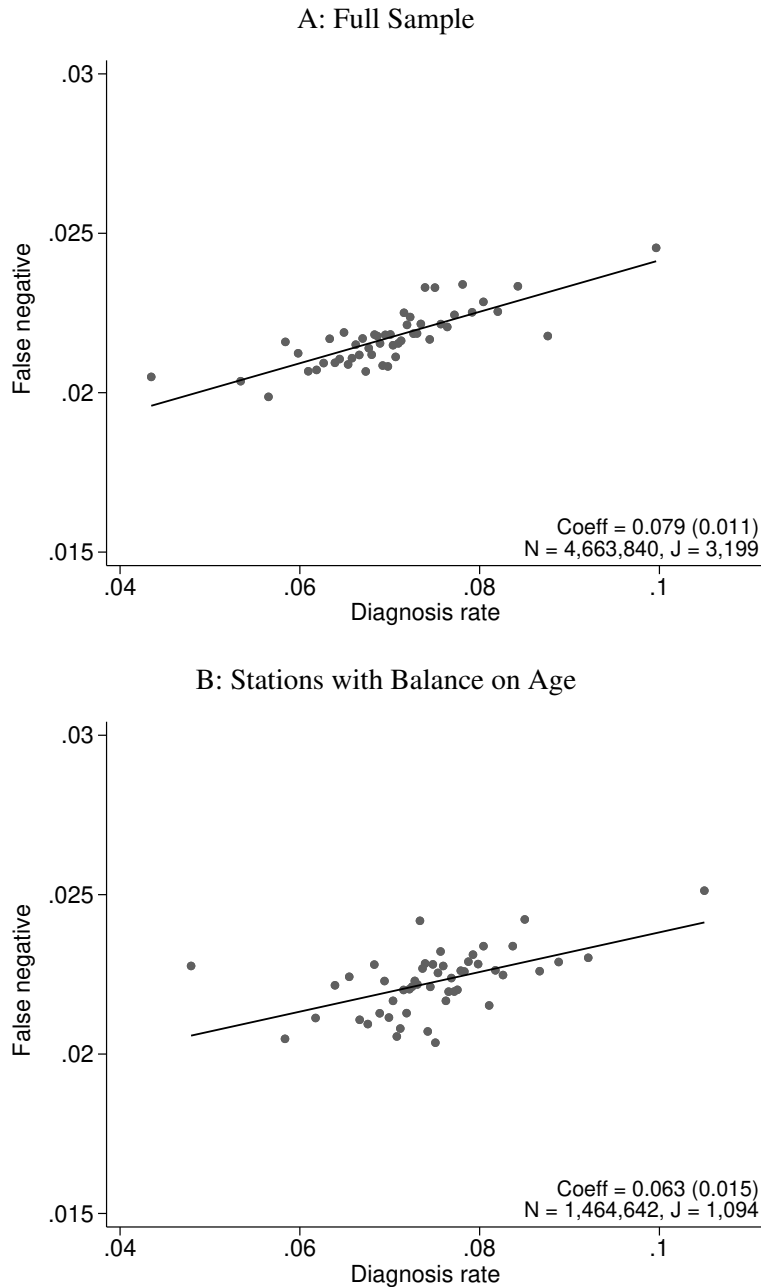
*Note:* This figure plots the standard deviation of radiologist fixed effects in simulations on the y-axis in resorted data where  $\iota \in [0, 100]$  percent of patients are randomly assigned to radiologists. The dashed line indicates the standard deviation in the observed data. Panel A shows results for the full sample. Panel B shows results for the sample of 44 stations selected for balance on age, as defined in Section 4.2. To construct the figure, we first residualize  $\hat{m}_i$  by minimal controls  $\mathbf{T}_i$ . We then create 101 samples. In each, we first reassign  $\iota \in \{0, 1, \dots, 100\}$  percent of cases randomly and the remaining cases perfectly sorted by  $\hat{m}_i$  to radiologists within the same station (holding the total number of cases for each radiologist constant). For each of these samples and the baseline sample, we regress the reassigned values on radiologist fixed effects and display the standard deviation of the estimated values. The shaded gray regions reflect 95% confidence intervals across 50 bootstrapped samples, drawn by radiologist blocks. The confidence interval corresponding to the dashed line in Panel A is  $\iota \in [96, 99]$ ; in Panel B, it is  $\iota \in [97, 100]$ .

Figure A.6: Projecting Data on ROC Space Using Alternative Parameter Values



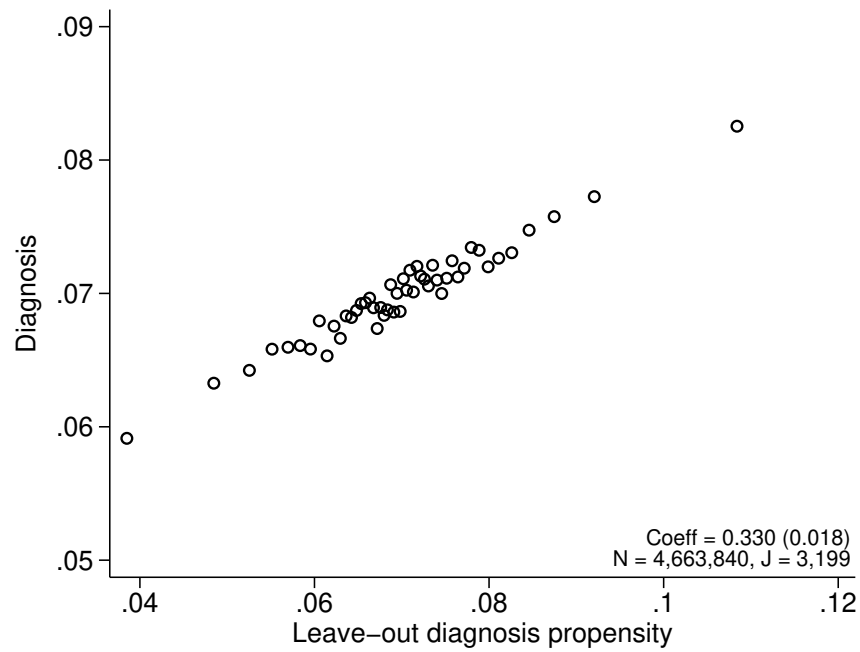
*Note:* This figure plots the true positive rate ( $\widehat{TPR}_j$ ) and false positive rate ( $\widehat{FPR}_j$ ) analogously to Figure V, under alternative values of prevalence ( $S$ ), the share of X-rays not at risk for pneumonia ( $\kappa$ ), and the share of cases in which pneumonia first manifests after the initial visit ( $\lambda$ ). In Panels A and B, we consider upper and lower bounds for  $S$ , as defined in Section 4.1. In Panels C and D, we increase and decrease  $\lambda$  by 50% relative to the baseline value  $\lambda = 0.026$ . In Panels E and F, we increase and decrease  $\kappa$  by 50% relative to its baseline value  $\kappa = 0.336$ . Appendix C provides details on this projection.

Figure A.7: Diagnosis and Miss Rates, Fixed Effects Specification



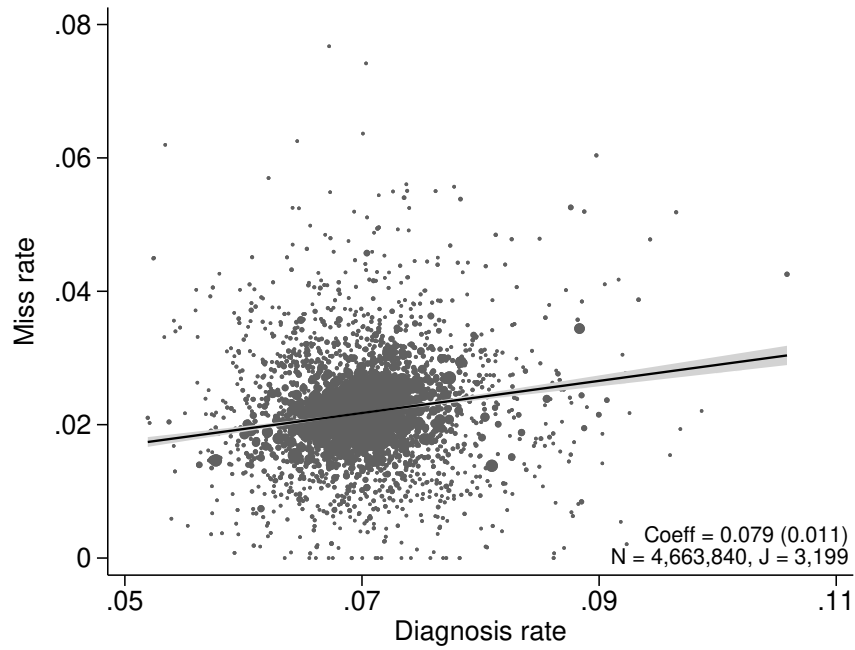
*Note:* This figure plots the relationship between miss rates and diagnosis rates across radiologists, using radiologist dummies as instruments. Plots are analogous to Figure VI. The  $x$ -axis plots  $\widehat{P}_j^{\text{obs}}$  and the  $y$ -axis plots  $\widehat{FN}_j^{\text{obs}}$ , defined in Section 4.3, both residualized by minimal controls of station-time interactions. Panel A shows results in the full sample of stations, and Panel B shows results in the subsample comprising 44 stations with balance on age, as defined in Section 4.2. The coefficient in each panel corresponds to the 2SLS estimate and standard error (in parentheses) for the corresponding IV regression, as well as the number of cases ( $N$ ) and the number of radiologists ( $J$ ). To account for clustering by radiologist, we test for first-stage joint significance by comparing an  $F$ -statistic of the radiologist dummies with  $F$ -statistics in 100 bootstrapped samples, drawn by a two-step procedure by radiologist and then by patient (both with replacement). The  $p$ -value for the joint significance is less than 0.01.

Figure A.8: First Stage



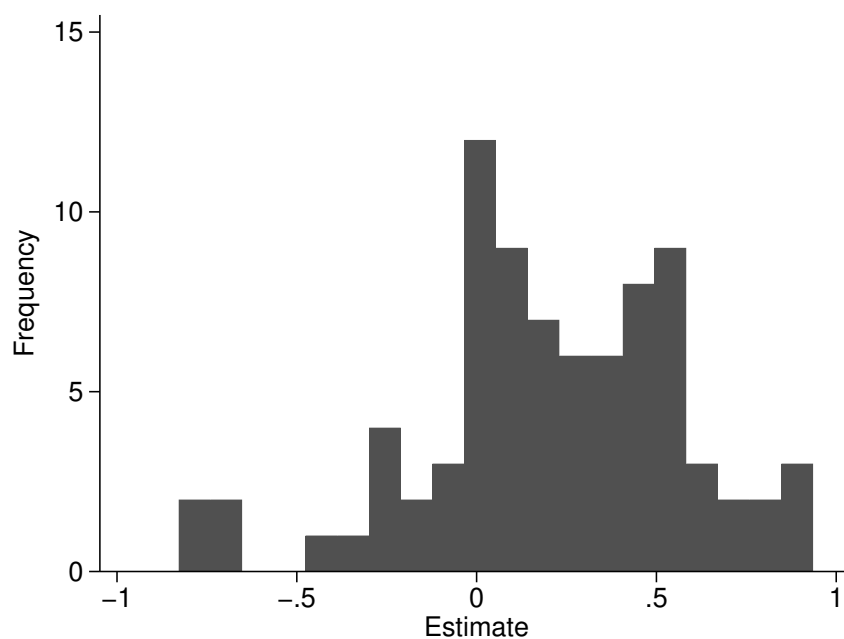
*Note:* This figure shows a binned scatter plot illustrating the first-stage relationship corresponding to Panel A of Figure VI. The  $y$ -axis shows residuals from a regression of diagnosis  $d_i$  on the covariates  $\mathbf{X}_i$  and minimal controls  $\mathbf{T}_i$ . The  $x$ -axis shows residuals from a regression of the leave-out propensity instrument  $Z_i$  on the same controls. The overall probability of diagnosis is added to residuals on the  $y$ -axis, and the average case-weighted  $Z_i$  is added to residuals on the  $x$ -axis. We report the first-stage coefficient as well as the number of cases ( $N$ ) and the number of radiologists ( $J$ ). The standard error is clustered at the radiologist level and shown in parentheses.

Figure A.9: Radiologist-Level Variation



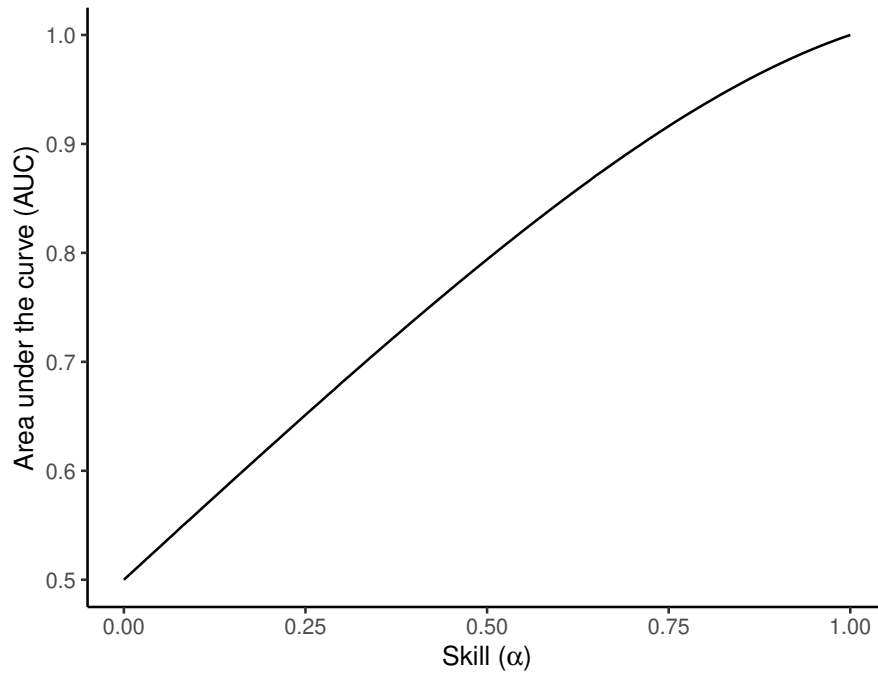
*Note:* This figure shows the relationship between radiologists' miss rates and diagnosis rates. We collapse the underlying data in Panel A of Figure VI to the radiologist level by taking the average. Each dot represents a radiologist, weighted by the number of cases. The coefficient and standard error are identical to those shown in Panel A of Figure VI. A radiologist in the case-weighted 90th percentile of miss rates has a miss rate 0.7 percentage points higher than that of a radiologist in the case-weighted 10th percentile. We calculate this by subtracting the case-weighted 10th percentile residual from the case-weighted 90th percentile residual from the underlying case-weighted regression.

Figure A.10: Distribution of Slope Estimates Across Stations



*Note:* This figure shows the distribution of station-level estimates of the slope  $\Delta$  relating radiologists' miss rates to their diagnosis rates. Each estimate is computed using the analogous IV procedure to that used to produce Figure VI with data from a single station. In the figure, 73 out of 104 stations have an estimate of the coefficient greater than zero.

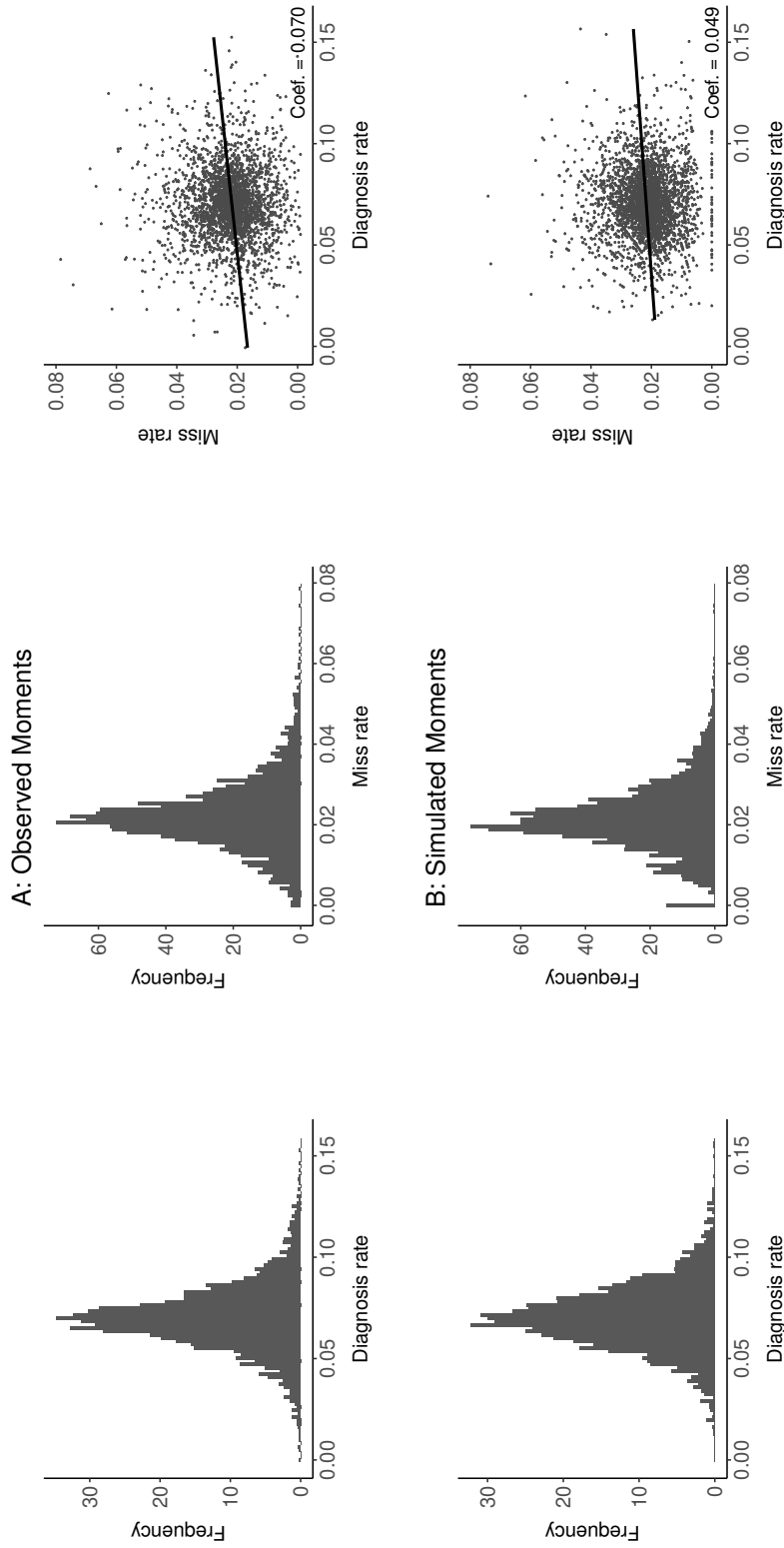
Figure A.11: Area Under the Curve (AUC) and Skill ( $\alpha$ )



*Note:* The Area Under the Curve (AUC) is the integral of an ROC curve. This figure shows the one-to-one mapping between AUC and the measure of skill  $\alpha$  under the assumptions of our structural model. When  $\alpha = 0$ , the ROC curve coincides with the 45-degree line and  $AUC = 0.5$ . When  $\alpha = 1$ , the ROC curve reduces to the left and top lines and  $AUC = 1$ .

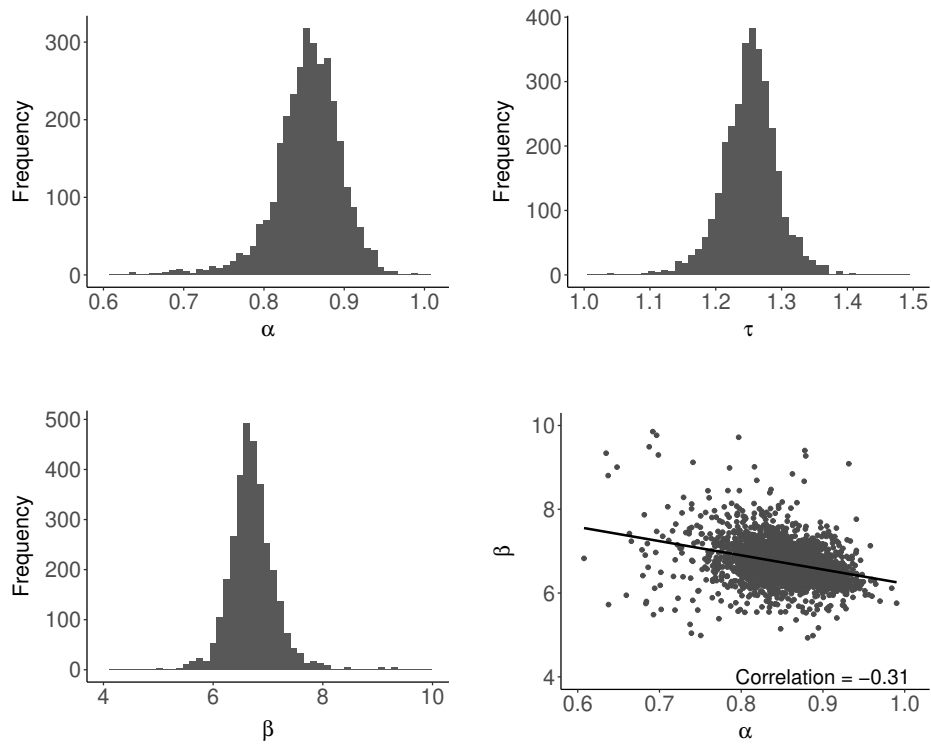


Figure A.12: Model Fit



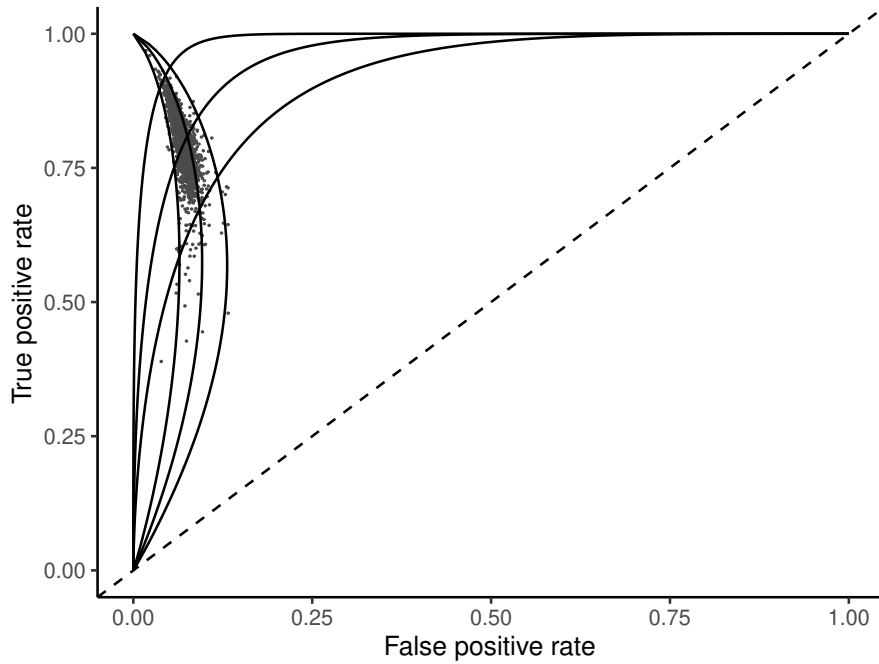
*Note:* This figure compares the actual moments observed in the data (the first row) with the moments simulated using the estimated parameters and simulated primitives from our main model estimates (the second row). To arrive at simulated moments in the second row, we first draw primitives for each radiologist,  $\alpha_j$  and  $\beta_j$ . We then simulate patients equal to the number assigned to the radiologist in the data, first drawing an indicator for whether the patient is at risk of pneumonia from a binomial distribution with the probability of being at risk  $1 - \kappa$ , then simulating their  $v_i$  and  $w_{ij}$  to determine their pneumonia status and the radiologist's diagnosis decision, given the threshold  $\bar{v}$  for pneumonia and the radiologist's diagnostic threshold  $\tau_j$ . For patients who are at risk, not diagnosed, and do not have pneumonia, we assign cases in which pneumonia first manifests after the initial visit with probability  $\lambda$ . Finally, we calculate the diagnosis and miss rate for each radiologist.

Figure A.13: Distributions of Radiologist Posterior Means



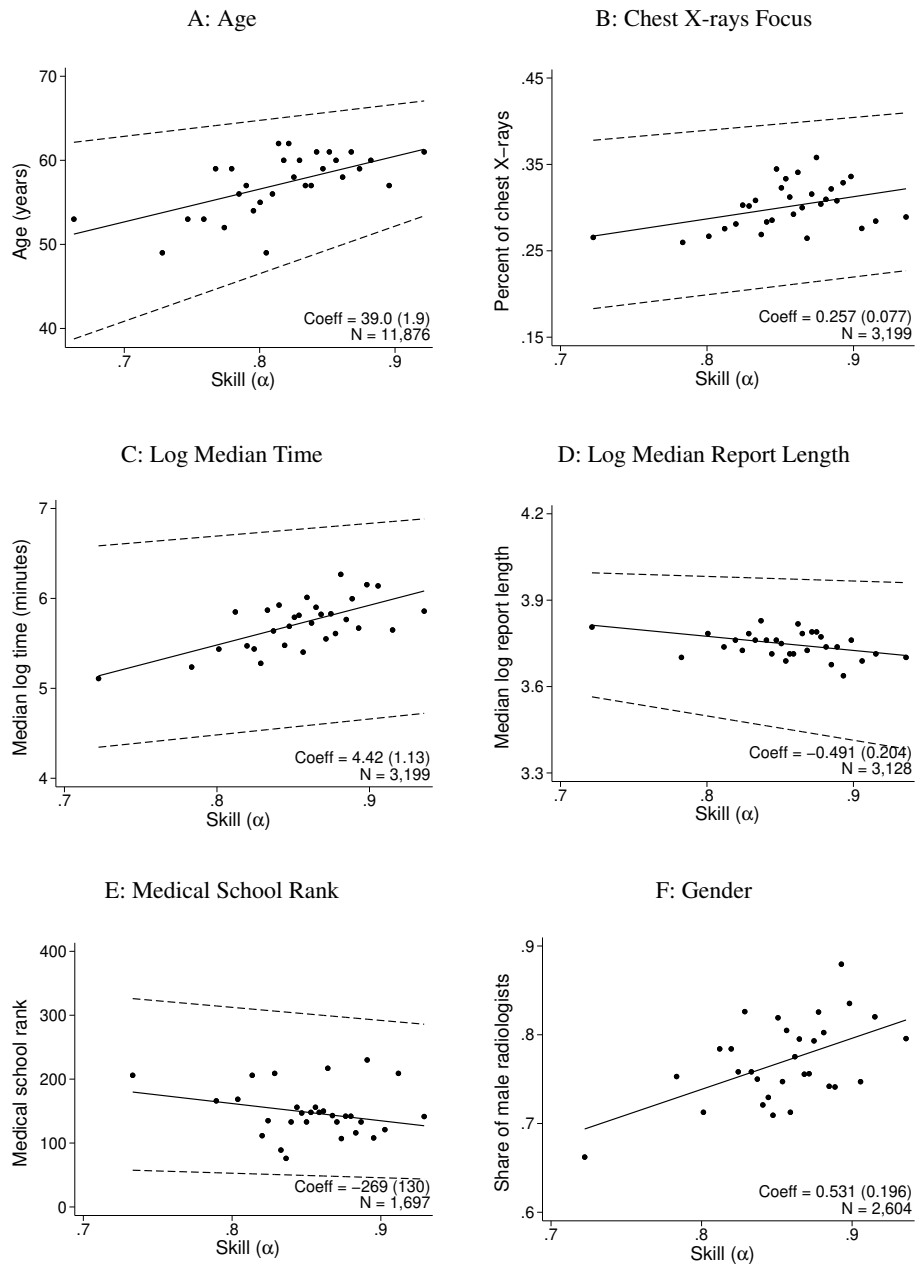
*Note:* This figure plots the distributions of radiologist empirical Bayes posterior means of our main specification. The first three subfigures plot the distributions of skill  $\hat{\alpha}_j$ , diagnostic thresholds  $\tau^*(\hat{\alpha}_j, \hat{\beta}_j)$ , and preferences  $\hat{\beta}_j$ . The last subfigure plots the joint distribution of skill and preferences. The method to calculate empirical Bayes posterior means is described in Appendix E.3.

Figure A.14: ROC Curve with Model-Generated Moments



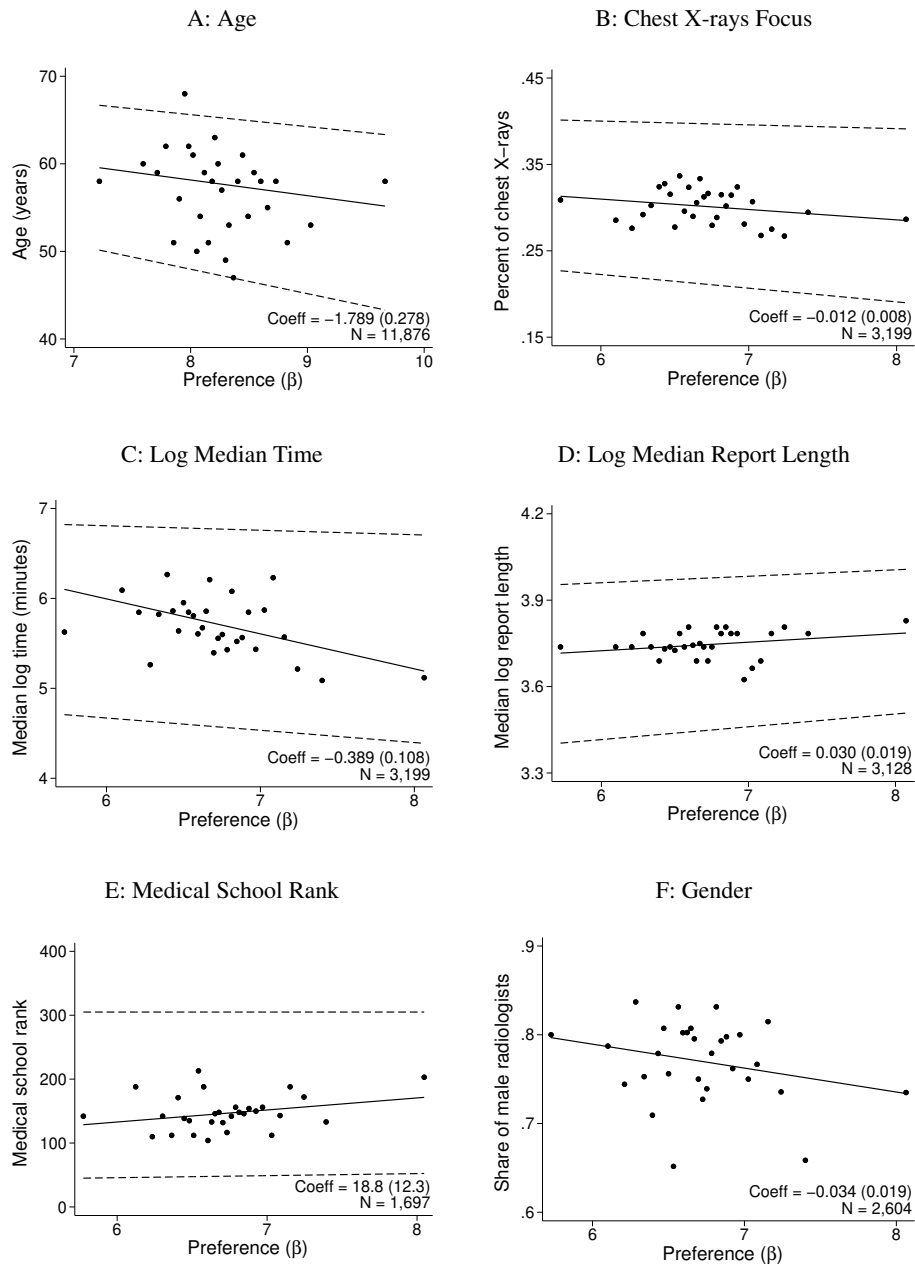
*Note:* This figure presents, for each radiologist, the true positive rate ( $TPR_j$ ) and false positive rate ( $FPR_j$ ) implied by radiologist posterior means of our main structural specification. Radiologist posterior means  $\hat{\gamma}_j = (\hat{\alpha}_j, \hat{\beta}_j)$  are calculated after estimating the model, described in Appendix E.3, and are the same as shown in Appendix Figure A.13. Large-sample  $P_j$  and  $FN_j$  are functions of radiologist primitives, given by  $p_{1j}(\gamma_j) \equiv \Pr(w_{ij} > \tau_j^* | \gamma_j)$  and  $p_{2j}(\gamma_j) \equiv \Pr(w_{ij} < \tau_j^*, v_i > \bar{v} | \gamma_j)$ , given in Section 5. As in Figure V,  $TPR_j = 1 - FN_j/S$  and  $FPR_j = (P_j + FN_j - S)/(1 - S)$ . This figure also plots the iso-preference curves for  $\beta \in \{5, 7, 9\}$  from (0,0) to (0,1) in ROC space. Each iso-preference curve illustrates how the optimal point in ROC space varies with skill for a fixed preference.

Figure A.15: Heterogeneity in Skill



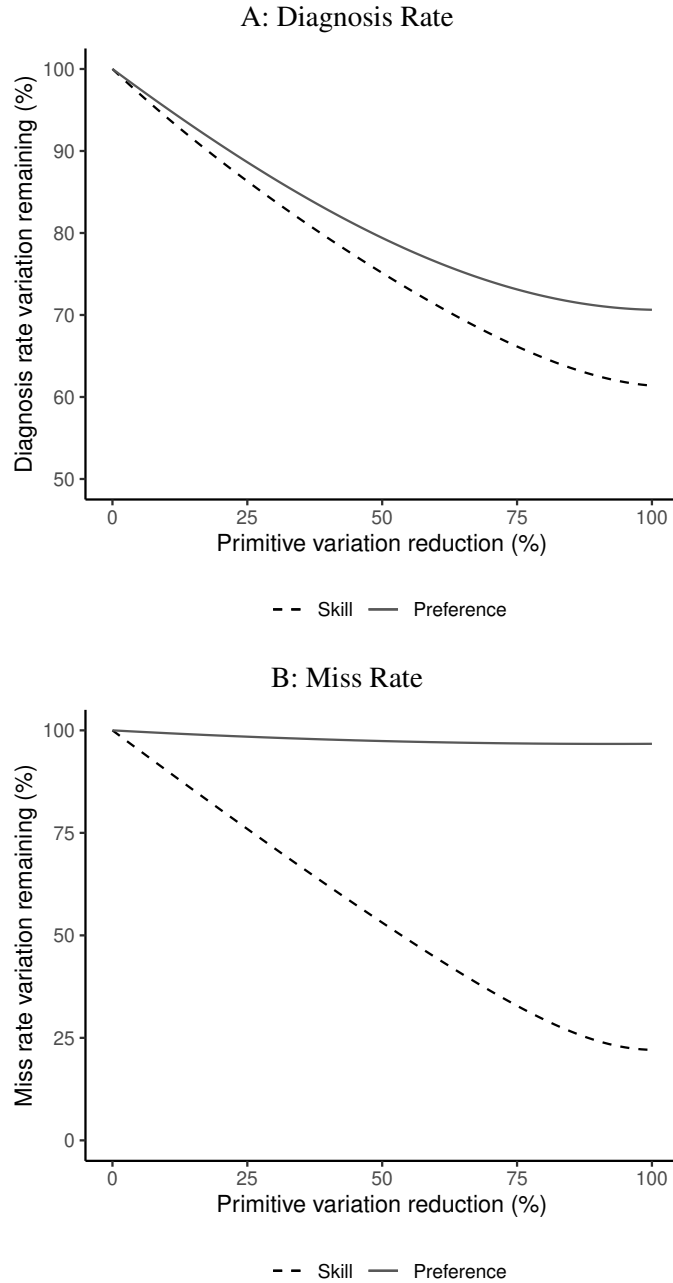
*Note:* This figure shows the relationship between the empirical Bayes posterior mean of a radiologist’s skill ( $\alpha$ ) on the x-axis and the following variables on the y-axis: (i) the radiologist’s age; (ii) the proportion of the radiologist’s exams that are chest X-rays; (iii) the log median time that the radiologist spends to generate a chest X-ray report; (iv) the log median length of the issue reports; (v) the rank of the medical school that the radiologist attended according to U.S. News & World Report; and (vi) gender. Except for gender, the three lines show the fitted values from the 25th, 50th, and 75th quantile regressions. For gender, the line shows the fitted values from an OLS regression. The dots are the median values of the variables on the y-axis within 30 bins of  $\alpha$ . Appendix Figure A.16 shows the corresponding plots with preferences ( $\beta$ ) on the x-axis. Some variables are missing for a subset of radiologists. For age, the result is based on a model that allows underlying primitives to vary by radiologist and age bin (we group five years as an age bin). See Section 5.5 for more details. Each panel reports the slope as well as the number of observations ( $N$ ). The standard error is shown in parentheses.

Figure A.16: Heterogeneity in Preferences



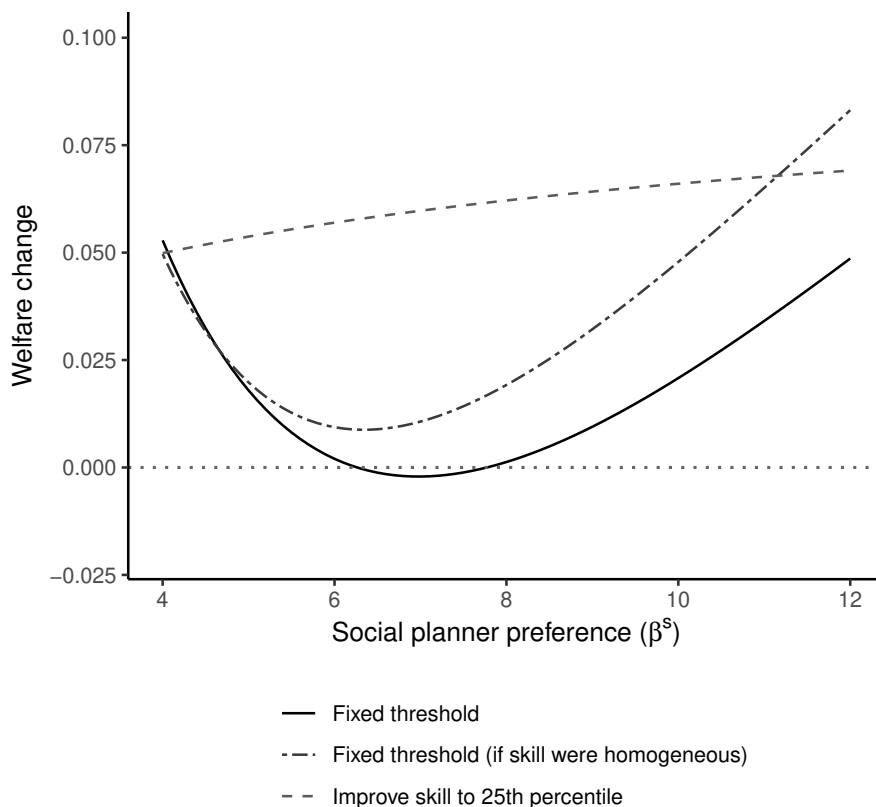
*Note:* This figure shows the relationship between a radiologist’s empirical Bayes posterior mean of her preference ( $\beta$ ) on the  $x$ -axis and the following variables on the  $y$ -axis: (i) the radiologist’s age; (ii) the proportion of the radiologist’s exams that are chest X-rays; (iii) the log median time that the radiologist spends to generate a chest X-ray report; (iv) the log median length of the issue reports; (v) the rank of the medical school that the radiologist attended according to U.S. News & World Report; and (vi) gender. Except for gender, the three lines show the fitted values from the 25th, 50th, and 75th quantile regressions. For gender, the line shows the fitted values from an OLS regression. The dots are the median values of the variables on the  $y$ -axis within each bin of  $\beta$ . 30 bins are used. Figure A.15 shows the corresponding plots with diagnostic skill ( $\alpha$ ) on the  $x$ -axis. Some variables are missing for a subset of radiologists. For age, the result is based on a model that allows underlying primitives to vary by radiologist and age bin (we group five years as an age bin). See Section 5.5 for more details. Each panel reports the slope as well as the number of observations ( $N$ ). The standard error is shown in parentheses.

Figure A.17: Variation Decomposition



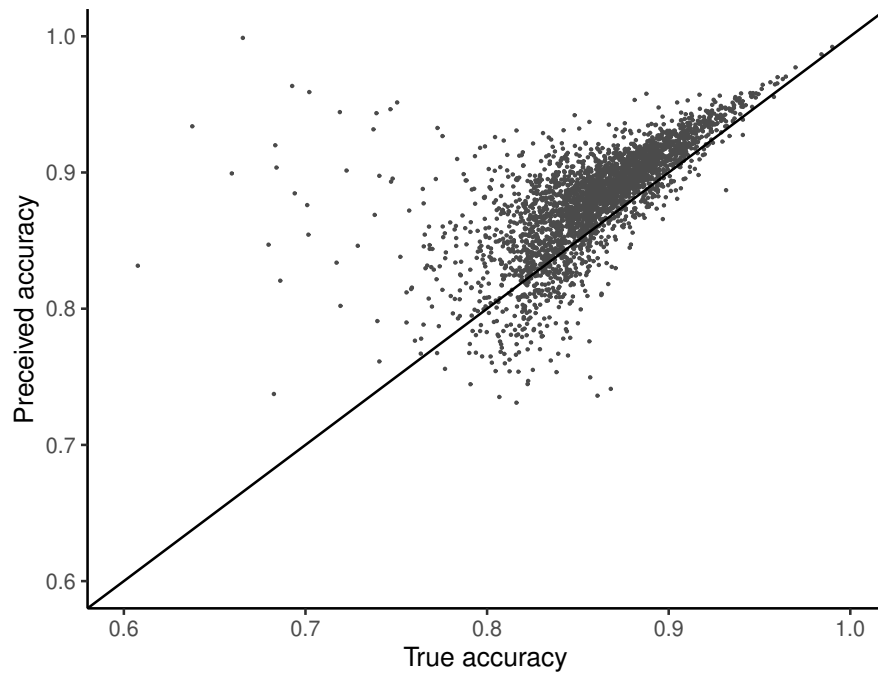
*Note:* This figure illustrates our method of calculating the variation in diagnosis and miss rates due to variation in skill and preferences. For  $x \in [0, 1]$ , we first keep  $\beta_j$  unchanged and replace  $\alpha_j$  by  $(1-x)\alpha_j + x \cdot \bar{\alpha}$ , where  $\bar{\alpha}$  is the median value of  $\alpha_j$ . When  $x = 0$ , this step simply gives  $\alpha_j$ . When  $x = 1$ , this step replaces all  $\alpha_j$  with  $\bar{\alpha}$  and thus eliminates all variation in  $\alpha_j$ . We derive the new diagnosis and miss rates under different  $x$ , calculate their standard deviations, and divide them by the original standard deviation with  $x = 0$ . We perform a similar calculation by shrinking  $\beta_j$  to the median value  $\bar{\beta}$  as  $x$  approaches 1 and keeping  $\alpha_j$  unchanged. Panel A shows the effect of reducing variation in skill or variation in preferences on the variation in diagnosis rates. Panel B shows the effect on the variation in miss rates. We report numbers that correspond to  $x = 1$  in Section 6.1.

Figure A.18: Counterfactual Policies



*Note:* This figure plots the counterfactual welfare gains of different policies. Welfare is defined in Equation (10) and is normalized to 0 for the status quo and 1 for the first best (no false positive or false negative outcomes). The  $x$ -axis represents different possible disutility weights that the social planner may place on false negatives relative to false positives, or  $\beta^s$ . The first policy imposes a common diagnostic threshold to maximize welfare. The second policy also imposes a common diagnostic threshold to maximize welfare but incorrectly computes welfare under the assumption that radiologists have the same diagnostic skill. The third policy trains radiologists to the 25th percentile of diagnostic skill (if their skill is below the 25th percentile) and allows them to choose their own diagnostic thresholds based on their preferences.

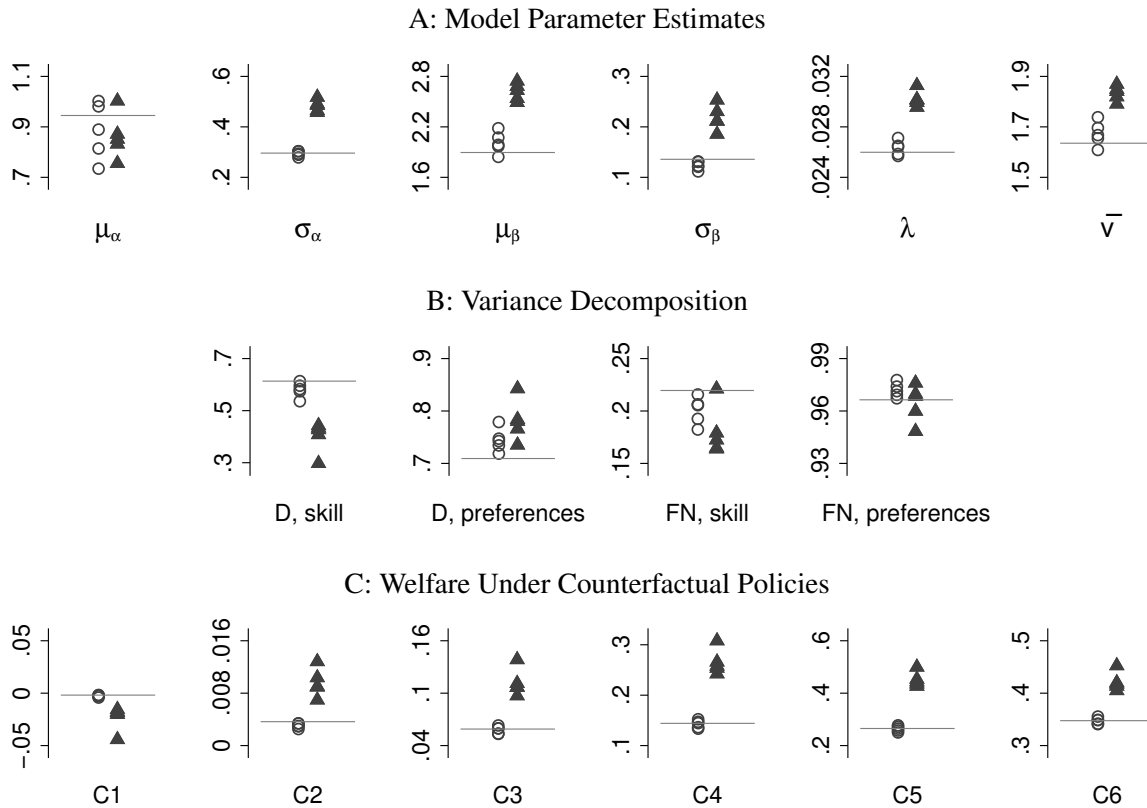
Figure A.19: Possibly Incorrect Beliefs about Accuracy



*Note:* This figure plots the relationship between radiologists' true accuracy and perceived accuracy, in an alternative model in which variation in diagnostic thresholds for a given skill is driven by variation in perceived skill, holding preferences fixed. This contrasts with the baseline model in which radiologists perceive their true skill but may vary in their preferences. We calculate the mean preference from our benchmark estimation results at  $\beta = 6.71$ , and we assign this preference parameter to all radiologists. We then use the formula for the optimal threshold as a function of  $\beta = 6.71$  and (perceived) accuracy to calculate perceived accuracy. Appendix G.2 describes this procedure to calculate perceived accuracy in further detail.

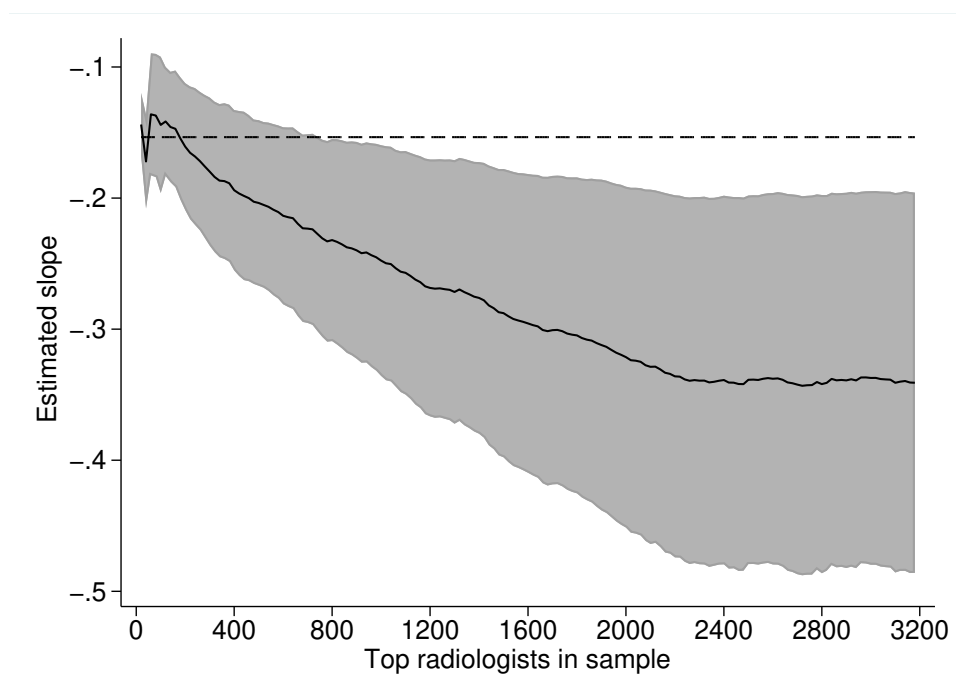


Figure A.20: Comparing Results with and without Risk Adjustment



*Note:* This figure shows structural results from simulated data with heterogeneity in pneumonia risk across stations. We simulate data to match the actual data in the number of radiologists in each station and the number of patients assigned to each radiologist. The simulated data come from the data generating process described in Appendix G.3, which matches the baseline model in Section 5.1 but allows for heterogeneity in pneumonia risk across stations. We take model parameter estimates in Table I as the truth and additionally include station-specific thresholds  $\bar{v}_\ell$  to model heterogeneity in pneumonia risk across stations. In each simulated dataset, we re-estimate structural parameters using radiologist diagnosis and miss rates that are either unadjusted (shown in triangles) or adjusted by linear regressions controlling for station dummies (shown in circles). Panel A shows model parameter estimates, as defined in Table I. Panel B shows variance decomposition results that follow from the model parameter estimates, as described in Section 6.1. Panel C similarly shows welfare under counterfactual policies, as described in Section 6.2. Horizontal lines denote true values of each object.

Figure A.21: Slope Estimates with Skill Controls, Radiologists Ordered by Volume



*Note:* This figure shows 2SLS estimates in simulated data of  $\Delta^*$  in subsamples of radiologists ordered by volume.  $\Delta^*$  is the LATE of diagnosis  $d_i$  on false negative  $m_i$  (i.e.,  $-\Pr(s_i)$ ), which we should obtain in valid judges-design (IV) regressions examining relationship between radiologist diagnosis and miss rates. We regress  $m_i$  on  $d_i$ , instrument  $d_i$  with the leave-out diagnosis propensity  $Z_i$  in Equation (4), and control for the empirical Bayes posterior mean of radiologist skill. Each estimate is based on a subsample of radiologists included in order of volume (from highest to lowest volume). The far-right end of the  $x$ -axis shows the estimate from the full sample; that estimate corresponds to Column 2 of Panel B in Appendix Table A.11. The 95% confidence interval is shaded in gray; standard errors are clustered by radiologist. The true estimand,  $\Delta^* = -0.154$ , is shown in the dashed line. Appendix G.4 provides further details.

Table A.1: Sample Selection

Sample step	Description	Cases	Radiologists
1. Select all chest X-ray observations from October 1999 to September 2015, inclusive	We define chest X-rays by the Current Procedural Terminology (CPT) codes of 71010 and 71020, and we require the status of the chest X-ray to be “complete”	5,523,995	6,330
2. Collapse multiple chest X-rays in a patient-day into one observation	If there are multiple radiologists among the chest X-rays, we assign the patient-day to the radiologist corresponding to the first chest X-ray in the patient-day	5,427,841	6,324
3. Retain patient-days that are at least 30 days from the last chest X-ray	Since we are interested in subsequent outcomes (e.g., return visits), we focus on initial chest X-rays with no prior chest X-rays within 30 days	4,828,550	6,283
4. Drop observations with missing radiologist identity or patient age or gender		4,823,985	6,283
5. Drop patients with age greater than 100 or less than 20		4,817,787	6,283
6. Drop radiologist-month pairs with fewer than 5 observations	This mitigates against limited mobility bias (Andrews et al. 2008), since we include month-year interactions as part of $\mathbf{T}_i$ in all our regression specifications of risk-adjustment	4,742,526	5,277
7. Drop radiologists with fewer than 100 remaining cases		4,663,840	3,199

*Note:* This table describes key sample selection steps, the number of cases, and the number of radiologists after each step.

Table A.2: Patient and Order Characteristic Variables

Category	Variables
Demographics (13 variables)	Age, indicator for male gender, indicator for married, 2 indicators for religion (Roman Catholic, Baptist, other religion as omitted), 4 indicators for race* (Black, White, American Indian, Pacific Islander, Asian/other race as omitted), indicator for veteran, distance between home and VA station performing X-ray*
Prior utilization (3 variables)	Previous year outpatient visits, previous year inpatient visits, previous year ED visits
Prior diagnoses (32 variables)	31 Elixhauser indicators (dividing hypertension indicator into 2 indicators for complicated and uncomplicated hypertension), indicator for prior pneumonia
Vital signs and WBC count (21 variables)	Systolic blood pressure*, diastolic blood pressure*, pulse*, pain*, O2 saturation*, respiratory rate*, temperature*, indicator for fever, indicator for supplemental O2 provided*, flow rate of supplemental O2, concentration of supplemental O2, white blood cell (WBC) count*
X-ray order (8 variables)	Indicator for urgent order, indicator for X-ray with multiple views (CPT 71020), number of X-rays by requesting physician, indicator for above-median average predicted diagnosis (based on the 13 demographic variables) of requesting physician, indicator for above-median average predicted false negative (based on the 13 demographic variables) of requesting physician, requesting physician leave-out share of pneumonia diagnosis, requesting physician leave-out share of false negatives, requesting physician leave-out share of urgent orders.

*Note:* This table describes 77 patient and X-ray order characteristic variables used as controls. \* behind a variable denotes that we include an additional variable to indicate missing values; there are 11 such variables. Predicted diagnosis and predicted false negative are predicted probabilities formed by running a linear probability regression of diagnosis indicator  $d_i$  and false negative indicator  $m_i$ , respectively, on demographic variables to calculate a linear fit for each patient. These predicted probabilities are averaged within each requesting physician.

Table A.3: Covariate Balance

	All Stations			Stations with Balance on Age		
Panel A: Diagnosis and Leave-Out Diagnosis Propensity						
	$d_1$	$d_2$	Diagnosis	Leave-Out Diagnosis Propensity	$d_2$	Leave-Out Diagnosis Propensity
Demographics	13	3,198	458.62 [0.000]	4.63 [0.000]	1,093	0.91 [0.538]
Prior diagnosis	32	3,198	550.12 [0.000]	3.60 [0.000]	1,093	1.44 [0.055]
Prior utilization	3	3,198	833.74 [0.000]	11.00 [0.000]	1,093	1.79 [0.147]
Vitals and WBC count	21	3,198	1341.36 [0.000]	4.01 [0.000]	1,093	1.00 [0.463]
Ordering characteristics	8	3,198	238.20 [0.000]	7.61 [0.000]	1,093	4.32 [0.000]
All variables	77	3,198	608.20 [0.000]	2.28 [0.000]	1,093	1.40 [0.015]
Panel B: False Negative and Leave-Out Miss Rate						
	$d_1$	$d_2$	False Negative	Leave-Out Miss Rate	$d_2$	Leave-Out Miss Rate
Demographics	13	3,198	456.37 [0.000]	4.43 [0.000]	1,093	1.98 [0.019]
Prior diagnosis	32	3,198	318.08 [0.000]	2.84 [0.000]	1,093	1.45 [0.053]
Prior utilization	3	3,198	1044.72 [0.000]	9.57 [0.000]	1,093	0.25 [0.863]
Vitals and WBC count	21	3,198	516.95 [0.000]	4.21 [0.000]	1,093	1.23 [0.213]
Ordering characteristics	8	3,198	304.37 [0.000]	11.26 [0.000]	1,093	2.32 [0.018]
All variables	77	3,198	194.22 [0.000]	2.64 [0.000]	1,093	1.28 [0.055]

*Note:* This table presents results of joint statistical significance from regressions of different outcomes on groups of patient characteristics. Each cell presents the  $F$ -statistic of the joint significance of a group of patient characteristics in a regression of an outcome, controlling for minimal controls  $\mathbf{T}_i$ . Panel A mirrors Figure IV, where Column 1 uses the diagnosis indicator as the outcome and Columns 2-3 use assigned radiologist's leave-out diagnosis propensity. Panel B mirrors Appendix Figure A.2, where Column 1 uses the false negative indicator as the outcome and Columns 2-3 use assigned radiologist's leave-out miss rate. In both panels, Columns 1 and 2 show regressions using the full sample of stations with 4,663,840 observations and Column 3 shows regressions using the sample of 44 stations with balance on age with 1,464,642 observations, described in Section 4.2.  $d_1$ , the first degree of freedom of the  $F$ -statistic, corresponds to the number of covariates;  $d_2$ , the second degrees of freedom, corresponds to the number of radiologists minus 1. The  $p$ -value corresponding to each  $F$ -statistic is displayed in brackets. Patient characteristics are described in further detail in Section 3 and Appendix Table A.2. Appendix Figure IV shows estimated coefficients and 95% confidence intervals for regressions with "all variables" in Panel A; Appendix Figure A.2 shows estimated coefficients and 95% confidence intervals for regressions with "all variables" in Panel B.

Table A.4: Balance

	Diagnosis Rate		Miss Rate	
	Below-Median	Above-Median	Difference	Difference
	Panel A: Full Sample			
Diagnosis	6.318 (0.029)	7.658 (0.030)	1.340 (0.045)	7.179 (0.032)
Predicted diagnosis	6.926 (0.017)	7.050 (0.015)	0.124 (0.022)	7.047 (0.014)
False negative	2.098 (0.013)	2.246 (0.012)	0.149 (0.017)	2.467 (0.011)
Predicted false negative	2.149 (0.005)	2.195 (0.004)	0.046 (0.006)	2.190 (0.004)
Number of cases	2,331,925	2,331,915	2,331,930	2,331,910
	Panel B: Stations with Balance on Age			
Diagnosis	6.901 (0.031)	8.085 (0.035)	1.185 (0.056)	7.613 (0.040)
Predicted diagnosis	7.402 (0.010)	7.414 (0.010)	0.012 (0.015)	7.408 (0.010)
False negative	2.179 (0.016)	2.273 (0.016)	0.094 (0.022)	2.480 (0.014)
Predicted false negative	2.214 (0.003)	2.222 (0.003)	0.008 (0.004)	2.217 (0.003)
Number of cases	732,322	732,320	732,321	732,321

*Note:* This table presents results assessing balance in patient characteristics. We divide patients into two groups with above- and below-median values of their assigned radiologist's diagnosis rates  $\hat{P}_j^{\text{obs}}$  (Columns 1-3) or miss rates  $\overline{FN}_j^{\text{obs}}$  (Columns 4-6) defined in Section 4.3, further risk-adjusted by minimal controls  $\mathbf{T}_i$ . In each panel, the patient groups are compared by actual diagnosis  $d_i$ , predicted diagnosis  $\hat{d}_i$ , actual false negative  $m_i$ , and predicted false negative  $\hat{m}_i$ . Predicted diagnosis and predicted false negative are formed by regressions using 77 patient characteristic variables, described in further detail in Section 3 and Appendix Table A.2. These outcomes are risk-adjusted by  $\mathbf{T}_i$ . Columns 1-2 and 4-5 show the mean of each residualized outcome across patients in each group; differences between groups are given in Columns 3 and 6. Standard errors shown in parentheses are computed by regressing the outcome on an above-median indicator and a below-median indicator, without a constant, and clustering by radiologist. Panel A shows results in all stations; Panel B shows results in stations with balance on age, described further in Section 4.2. In the last row of each panel, we display the number of cases in each group.

Table A.5: Statistics on Radiologist-Level Moments

	Mean	SD	Percentiles			
			10th	25th	75th	90th
Panel A: Observed, Risk-Adjusted						
Diagnosis rate $\widehat{P}_j^{\text{obs}}$	0.070	0.010	0.059	0.065	0.074	0.082
Miss rate $\widehat{FN}_j^{\text{obs}}$	0.022	0.005	0.017	0.019	0.024	0.027
Panel B: Also Adjusted for $\hat{\kappa} = 0.336$ and $\hat{\lambda} = 0.026$						
Diagnosis rate $\widehat{P}_j$	0.105	0.015	0.089	0.097	0.112	0.123
Miss rate $\widehat{FN}_j$	0.010	0.007	0.002	0.006	0.013	0.018
False positive rate $\widehat{FPR}_j$	0.068	0.019	0.048	0.057	0.078	0.090
True positive rate $\widehat{TPR}_j$	0.802	0.131	0.654	0.748	0.878	0.959

*Note:* This table presents statistics for various radiologist-level moments. Panel A shows raw risk-adjusted diagnosis and miss rates, which are fitted radiologist fixed effects from regressions of  $d_i$  and  $m_i$  on radiologist fixed effects, patient characteristics  $\mathbf{X}_i$ , and minimal controls  $\mathbf{T}_i$ , respectively. Panel B adjusts for the share of X-rays not at risk of pneumonia ( $\hat{\kappa} = 0.336$ ), calibrated in Section 3, and the share of cases whose pneumonia manifests after the first visit ( $\hat{\lambda} = 0.026$ ), estimated in Section 5.2. False positive rates and true positive rates are then computed using the estimated prevalence rate ( $\hat{S} = 0.051$ ). All statistics are weighted using the number of cases. See Appendix C for more details.

Table A.6: Informal Monotonicity Tests

Subsample	Outcome: Diagnosed, $d_i$							
	Older	Younger	High Pr ( $d_i$ )	Low Pr ( $d_i$ )	White	Non-White	Daytime	Nighttime
Panel A: Baseline								
Instrument, $Z_j$	0.230 (0.013)	0.413 (0.015)	0.149 (0.009)	0.482 (0.018)	0.346 (0.012)	0.280 (0.017)	0.353 (0.011)	0.233 (0.021)
Mean outcome	0.051	0.089	0.021	0.119	0.075	0.059	0.069	0.073
Observations	2,331,962	2,331,860	2,331,896	2,331,906	3,088,650	1,575,015	3,456,470	1,207,246
Panel B: Reverse-Sample								
Instrument, $Z_j^{-r}$	0.168 (0.009)	0.384 (0.016)	0.108 (0.006)	0.741 (0.032)	0.189 (0.010)	0.253 (0.014)	0.126 (0.008)	0.244 (0.019)
Mean outcome	0.051	0.089	0.021	0.119	0.075	0.059	0.069	0.073
Observations	2,331,962	2,331,860	2,331,896	2,331,906	3,046,649	1,570,742	3,321,569	1,200,498
Time $\times$ station fixed effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Patient controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes

*Note:* This table shows results from informal tests of monotonicity that are standard in the judges-design literature. Each column corresponds to a different subsample of observations. In each subsample, we run first stage regressions of the effect of a leave-out instrument on diagnosis, controlling for 77 variables for patient characteristics, described in Section 3 and Appendix Table A.2, and time dummies interacted with location dummies. Panel A shows results from Equation (D.4), using a standard leave-out instrument. Panel B shows results from Equation (D.5), using a reverse-sample instrument. See Appendix D for more details.



Table A.7: Judges-Design Estimates of the Effect of Diagnosis on Other Outcomes

Outcome	All Stations		Stations with Balance on Age	
Admissions within 30 days	1.114 (0.338)	0.633	-0.076 (0.219)	0.587
ED visits within 30 days	0.146 (0.121)	0.290	-0.385 (0.201)	0.290
ICU visits within 30 days	0.201 (0.051)	0.044	-0.088 (0.067)	0.042
Inpatient-days in initial admission	10.695 (2.317)	2.530	0.588 (2.193)	2.209
Inpatient-days within 30 days	11.383 (2.059)	3.330	-1.123 (1.879)	3.043
Mortality within 30 days	0.150 (0.032)	0.033	-0.126 (0.057)	0.033

*Note:* This table presents results using the assigned radiologist’s leave-out diagnosis propensity in Equation (4) as the instrument to calculate the effect of diagnosis on other outcomes, similar to the benchmark outcome of false negative status in Figure VI. All regressions control for 77 variables of patient characteristics, described in Section 3 and Appendix Table A.2, and time dummies interacted with location dummies. Columns 1 and 3 give results of the IV estimates. Standard errors are given in parentheses. Columns 2 and 4 report mean outcomes. Columns 1 and 2 show regressions using the full sample of stations; Columns 3 and 4 show regressions using the sample of 44 stations with balance on age, described in Section 4.2.

Table A.8: Alternative Specifications

	Baseline	Balanced	VA users	Admission	Minimum		Fix $\lambda$ , flexible $\rho$
					controls	No controls	
Panel A: Data and Reduced-Form Moments							
SD of diagnosis	1.023	1.031	1.095	1.027	1.231	1.966	1.023
SD of false negative status	0.499	0.461	0.580	0.427	0.532	0.752	0.499
SD of false negative residual	0.494	0.457	0.577	0.426	0.510	0.680	0.494
Slope, IV	0.291	0.344	0.357	0.201	0.270	0.189	0.291
Number of observations	4,663,840	1,464,642	3,099,211	4,663,601	4,663,840	4,663,840	4,663,840
Number of radiologists	3,199	1,094	3,199	3,199	3,199	3,199	3,199
Panel B: Variation Decomposition							
Diagnosis							
Uniform skill	0.613 (0.056)	0.634 (0.163)	0.619 (0.069)	0.715 (0.057)	0.515 (0.054)	0.350 (0.058)	0.615 (0.044)
Uniform preference	0.709 (0.079)	0.725 (0.120)	0.686 (0.103)	0.614 (0.086)	0.766 (0.058)	0.812 (0.051)	0.710 (0.071)
False negative							
Uniform skill	0.220 (0.046)	0.177 (0.074)	0.174 (0.048)	0.217 (0.050)	0.170 (0.029)	0.112 (0.016)	0.212 (0.040)
Uniform preference	0.966 (0.019)	0.981 (0.059)	0.981 (0.016)	0.971 (0.019)	0.977 (0.010)	0.992 (0.024)	0.969 (0.016)

*Note:* This table shows robustness of results under alternative implementations. “Baseline” presents our baseline results. “Balanced” presents results estimated only on the 44 stations we identify with quasi-random assignment. “VA users” restricts to a sample of veterans with more total visits in the VA than in Medicare. “Admission” defines false negatives only in patients with a high probability of admission. “Minimum controls” performs risk-adjustment only using time and stations. “No controls” presents results estimated using the raw diagnosis and miss rates without adjusting for stations, time, and patient characteristics. “Fix  $\lambda$ , flexible  $\rho$ ” presents results estimated by fixing  $\lambda$  at the estimated value in the baseline specification, but allowing  $\rho$ , the correlation between  $\alpha_j$  and  $\beta_j$ , to vary flexibly. Appendix F provides rationale for each of these implementations and further discussion. Standard errors for Panel B, shown in parentheses, are computed by block bootstrap, with replacement, at the radiologist level.

Table A.9: Alternative Specifications (Additional Detail)

	Baseline	Balanced	VA users	Admission	Minimum controls	No controls	Fix $\lambda$ , flexible $\rho$
Panel A: Model Parameter Estimates							
$\mu_\alpha$	0.945 (0.219)	0.516 (0.960)	0.809 (0.156)	0.820 (0.206)	0.890 (0.135)	1.091 (0.148)	0.911 (0.304)
$\sigma_\alpha$	0.296 (0.029)	0.227 (0.253)	0.421 (0.036)	0.246 (0.030)	0.383 (0.032)	0.784 (0.070)	0.294 (0.032)
$\mu_\beta$	1.895 (0.249)	2.564 (0.632)	1.900 (0.231)	2.066 (0.253)	2.059 (0.127)	1.938 (0.152)	1.928 (0.349)
$\sigma_\beta$	0.136 (0.044)	0.084 (0.193)	0.159 (0.047)	0.138 (0.034)	0.143 (0.031)	0.220 (0.064)	0.130 (0.055)
$\lambda$	0.026 (0.001)	0.029 (0.006)	0.022 (0.002)	0.016 (0.001)	0.027 (0.001)	0.025 (0.002)	- -
$\bar{\nu}$	1.635 (0.091)	1.873 (0.261)	1.678 (0.074)	1.704 (0.096)	1.681 (0.050)	1.597 (0.045)	1.649 (0.125)
$\rho$	-	-	-	-	-	-	-0.056 (0.168)
$\kappa$	0.336	0.336	0.336	0.336	0.336	0.336	0.336
Panel B: Radiologist Primitives							
Mean $\alpha$	0.855	0.728	0.806	0.769	0.832	0.826	0.847
10th percentile	0.756	0.610	0.631	0.647	0.689	0.542	0.744
90th percentile	0.934	0.833	0.937	0.874	0.940	0.985	0.929
Mean $\beta$	6.713	13.034	6.766	9.723	7.920	7.110	6.928
10th percentile	5.596	11.673	5.455	8.284	6.534	5.247	5.819
90th percentile	7.909	14.456	8.186	11.253	9.410	9.188	8.112
Mean $\tau$	1.252	1.213	1.307	1.253	1.252	1.307	1.253
10th percentile	1.165	1.138	1.193	1.165	1.139	1.075	1.167
90th percentile	1.336	1.290	1.412	1.339	1.364	1.461	1.336

*Note:* This table shows additional details of the robustness results under alternative specifications. The columns, each corresponding to an alternative specification, are the same as Appendix Table A.8. The parameters in Panel A are the same as discussed in Table I.

Table A.10: Model Results Under Alternative Values of  $\kappa$

Panel A: Value of $\kappa$			
$\kappa$	0.168	0.336	0.504
Panel B: Model Parameter Estimates			
$\mu_\alpha$	1.023	0.945	0.798
$\sigma_\alpha$	0.291	0.296	0.311
$\mu_\beta$	1.916	1.895	1.863
$\sigma_\beta$	0.143	0.136	0.129
$\lambda$	0.020	0.026	0.035
$\bar{v}$	1.740	1.635	1.499
Panel C: Variation Decomposition			
Diagnosis, Uniform skill	0.627	0.613	0.618
Diagnosis, Uniform preference	0.698	0.709	0.694
False negative, Uniform skill	0.224	0.220	0.216
False negative, Uniform preference	0.965	0.966	0.967

*Note:* This table presents the analogous results in Table I under different values of  $\kappa$ . In the baseline estimation,  $\kappa=0.336$  is calibrated as the fraction of patients whose probability of having pneumonia predicted by a machine learning algorithm is smaller than 0.01. We use two other values of  $\kappa$  that represent a 50% decrease (Column 1) and 50% increase (Column 3) around the calibrated value (Column 2). Panel A shows model parameter estimates corresponding to these alternative thresholds. Panel B shows the variation decomposition under these alternative thresholds. Parameters are described in further detail in Sections 5.1 and 5.2, and counterfactual variation exercise is described in further detail in Section 6.1.

Table A.11: Slope Estimates Controlling for Radiologist Skill

	(1)	(2)	(3)	(4)	(5)	(6)
Panel A: True Skill						
Diagnosis	0.096 (0.016)	-0.124 (0.014)	-0.132 (0.019)	-0.147 (0.019)	-0.155 (0.017)	-0.156 (0.017)
Panel B: Skill Posteriors						
Diagnosis	0.096 (0.016)	-0.342 (0.084)	-0.575 (0.084)	-0.668 (0.119)	-0.698 (0.143)	-0.752 (0.237)
Panel C: Indirect Least Squares						
Diagnosis	0.096 (0.016)	-0.251 (0.043)	-0.364 (0.034)	-0.369 (0.036)	-0.208 (0.058)	-0.051 (0.119)

*Note:* This table presents slope estimates in simulated data of  $\Delta^*$ , or the LATE of diagnosis  $d_i$  on false negative  $m_i$ , based on IV regressions identified by the judges-design relationship between radiologist diagnosis and miss rates. Column 1 in all panels presents the same specification, akin to the benchmark IV regression in the paper, instrumenting  $d_i$  with the leave-out diagnosis propensity  $Z_i$  in Equation (4), with no further controls. For Panel A, we additionally control for true (simulated) radiologist skill  $\alpha_j$ . For Column 2 of this panel, we control for linear  $\alpha_j$ ; for Columns 3-6, we control for indicators for each of 5, 10, 20, and 50 bins of  $\alpha_j$ , respectively. For Panel B, we use the empirical Bayes posteriors instead of true skill, defined in Appendix E.3. For Column 2 of this panel, we linearly control for the posterior mean of  $\alpha_j$ ; for Columns 3-6, we control for indicators for each of 5, 10, 20, and 50 bins of this posterior mean, respectively. Panel C shows results from indirect least squares, regressing  $m_i$  on posteriors of  $P_j$  and  $\alpha_j$  by OLS. For Column 2 of this panel, we control for the posterior mean of  $\alpha_j$ ; for Columns 3-6, we control for posterior probabilities that  $\alpha_j$  resides in each of 5, 10, 20, and 50 bins, respectively. Standard errors, shown in parentheses, are clustered by radiologist. In Panels B and C, standard errors are computed by 50 samples drawn by block bootstrap with replacement, at the radiologist level. We compute the true estimand  $\Delta^* = -0.154$ . Appendix G.4 provides further details.