## NBER WORKING PAPER SERIES

## MEASURING RACIAL DISCRIMINATION IN BAIL DECISIONS

David Arnold Will S. Dobbie Peter Hull

Working Paper 26999 http://www.nber.org/papers/w26999

NATIONAL BUREAU OF ECONOMIC RESEARCH 1050 Massachusetts Avenue Cambridge, MA 02138 April 2020

We thank Tim Armstrong, Leah Boustan, Sydnee Caldwell, Raj Chetty, Joseph Doyle, Matt Gentzkow, Ed Glaeser, Paul Goldsmith-Pinkham, Damon Jones, Scott Nelson, Sam Norris, Crystal Yang, and numerous seminar participants for helpful comments and suggestions. Emily Battaglia, Nicole Gandre, Jared Grogan, Ashley Litwin, Alexia Olaizola, Bailey Palmer, Elise Parrish, Emma Rackstraw, and James Reeves provided excellent research assistance. The data we analyze are provided by the New York State Division of Criminal Justice Services (DCJS), and the Office of Court Administration (OCA). The opinions, findings, and conclusions expressed in this publication are those of the authors and not those of DCJS, OCA, or the National Bureau of Economic Research. New York State, DCJS, OCA or NBER do not assume liability for its contents or use thereof.

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

© 2020 by David Arnold, Will S. Dobbie, and Peter Hull. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Measuring Racial Discrimination in Bail Decisions David Arnold, Will S. Dobbie, and Peter Hull NBER Working Paper No. 26999 April 2020 JEL No. C26,J15,K42

## ABSTRACT

We develop new quasi-experimental tools to measure racial discrimination in the context of bail decisions. Observational comparisons of white and black pretrial release rates suffer from omitted variables bias when there are unobserved racial differences in pretrial misconduct potential. We show that the bias in these observational comparisons is a function of average white and black misconduct risk, which can be estimated from the quasi-random assignment of bail judges. Estimates from New York City show that less than one-third of the release rate disparity between white and black defendants is explained by unobserved differences in misconduct potential, with more than two-thirds explained by racial discrimination. We then develop a hierarchical marginal treatment effects model that imposes additional structure on the quasi-experimental variation to investigate the drivers of this discrimination. Model estimates show that discrimination in bail decisions is driven by both racial bias and statistical discrimination, with the latter coming from a higher level of average risk and less precise risk signals for black defendants.

David Arnold Industrial Relations Section Louis A. Simpson International Bldg. Princeton University Princeton, NJ 08544-2098 dharnold@princeton.edu

Will S. Dobbie Harvard Kennedy School 79 John F. Kennedy St. Cambridge, MA 02138 and NBER will\_dobbie@hks.harvard.edu Peter Hull University of Chicago 5757 South University Avenue Chicago, IL 60637 and NBER hull@uchicago.edu

## 1 Introduction

Racial disparities are pervasive throughout much of the U.S. criminal justice system. Compared to observably similar white individuals, black individuals are more likely to be searched by the police, charged with a serious crime, detained before trial, convicted of an offense, and ultimately incarcerated.<sup>1</sup> These racial disparities are often taken as prima facie evidence of racial discrimination, driven by some form of racial bias. But there are two alternative explanations. The first is that observed disparities reflect legally relevant differences in criminal behavior that are partially observed by police officers, prosecutors, and judges but unobserved by the econometrician. The second is that these disparities do reflect racial discrimination, but that they are driven by some form of statistical discrimination and not racial bias. Empirically distinguishing between these various explanations remains difficult, hampering efforts to formulate the appropriate policy response.

This paper develops new tools to measure racial discrimination when decision-makers are as-goodas-randomly assigned. We study bail decisions, where the sole legal objective of judges is to allow most defendants to be released before trial while minimizing the risk of pretrial misconduct (such as failing to appear in court or being arrested for a new crime). Bail judges therefore risk violating U.S. law if they release white and black defendants with the same objective misconduct potential at different rates. Correspondingly, we measure discrimination by a judge's release rate disparity between white and black individuals with identical misconduct potential. This measure is consistent with mainstream legal views on what constitutes discrimination in the criminal justice system (Yang and Dobbie, 2019), as well as economic notions of discrimination that compare the treatment of white and black individuals with the same productivity (Aigner and Cain, 1977). Our measure can be understood as isolating each judge's legally unwarranted release rate disparity, leading us to use the terms racial discrimination and unwarranted disparity interchangeably throughout this paper.

Estimating discrimination by isolating unwarranted disparities among individuals with identical misconduct potential is difficult, even when bail judges are as-good-as-randomly assigned. Observational comparisons of white and black release rates cannot control for unobserved misconduct potential and therefore suffer from omitted variables bias (OVB) when there are unobserved racial differences in misconduct risk. Randomized audit studies can purge such OVB by experimentally manipulating the observed race of fictitious individuals while holding all other observable characteristics constant (e.g., Bertrand and Mullainathan, 2004; Ewens et al., 2014), but here the high-stakes and face-to-face nature of bail decisions make such randomization infeasible. Standard instrumental variables (IV) methods can test for discrimination driven by racial bias (e.g., Arnold et al., 2018; Marx, 2018), but do not speak to the presence of statistical discrimination. Standard IV methods also require an assumption of first-stage monotonicity (Imbens and Angrist, 1994; Heckman and Vytlacil, 2005), which here imposes a strong restriction on how judges choose which defendants to release before trial.

Our primary methodological contribution is to show that racial discrimination in bail decisions, regardless of its source, can be measured with observational release rate comparisons that are rescaled using quasi-experimental estimates of average white and black misconduct risk. The bias in observational release rate comparisons for individual judges comes from the correlation between race and unobserved misconduct potential in a given judge's pool of defendants. Under quasi-random assign-

<sup>&</sup>lt;sup>1</sup>There is a large literature documenting racial disparities in the criminal justice system. See, for example, work by Gelman et al. (2007), Antonovics and Knight (2009), Anwar et al. (2012), Abrams et al. (2012), McIntyre and Baradaran (2013), and Rehavi and Starr (2014), among many others.

ment, this correlation is a common function of average misconduct risk by race. We can therefore use quasi-experimental estimates of average misconduct risk by race to rescale observational release rate comparisons for individual judges in such a way that released white and black defendants are directly comparable in terms of their unobserved misconduct potential. The rescaled observational comparisons reveal the extent to which each judge releases white and black defendants with the same objective misconduct potential at different rates, even though misconduct potential is unobserved and cannot be directly conditioned on. The key econometric challenge is thus to estimate the required average misconduct risk parameters, which is difficult since misconduct outcomes are only observed for defendants who are endogenously released before trial.

We show how the average misconduct risk inputs required for our discrimination measure can be estimated by extending recent approaches to estimating average treatment effects (ATEs) with multiple discrete instruments (Brinch et al., 2017; Hull, 2020). To build intuition for our approach, consider a setting with a supremely lenient bail judge who releases nearly all defendants regardless of their race or misconduct potential. The supremely lenient judge's release rates among white and black defendants are close to one, meaning that the misconduct rates among white and black released defendants are close to the average misconduct risk parameters needed for our discrimination measure (by quasi-random assignment). In the absence of such a supremely lenient judge, the required average misconduct risk inputs can be estimated using model-based or non-parametric extrapolations of pretrial release and misconduct rates across quasi-randomly assigned judges. Importantly, this quasiexperimental approach to estimating discrimination does not require a model of judge behavior or a first-stage monotonicity assumption, only that the extrapolations of pretrial release and misconduct rates and the judges' legal objective are well-specified by the econometrician.

We use our new quasi-experimental approach to measure racial discrimination in bail decisions made in New York City (NYC), one of the largest pretrial systems in the country. We find that more than two-thirds of the average release rate disparity between observationally similar white and black defendants is explained by racial discrimination (68 percent, or 3.6 percentage points out of 5.3 percentage points), with less than one-third explained by unobserved racial differences in misconduct risk. This estimate of system-wide discrimination is robust to different extrapolation methods, definitions of a defendant's race, and definitions of a judge's legal objective. Judge-specific estimates show that the vast majority of bail judges discriminate against black defendants (88 percent), with higher levels of discrimination among more stringent judges, judges assigned to a lower share of cases with black defendants, and judges who are newly appointed in our sample period.<sup>2</sup> The judge-specific estimates are also highly correlated over time, raising the possibility that discrimination can be effectively targeted across individual bail judges.

Our second methodological contribution is to develop a hierarchical marginal treatment effects (MTE) model that imposes additional structure on the quasi-experimental variation to investigate whether unwarranted disparities in bail decisions are driven by racial bias or statistical discrimination. The model allows for judge- and race-specific risk preferences and signal quality, with the latter allowing for heterogenous race-specific predictive skill across judges (in violation of the conventional first-stage monotonicity assumption). This variation implies a distribution of judge- and race-specific MTE curves that can be used to test for racial bias at the margin of release, as well as to measure

 $<sup>^{2}</sup>$ We define a judge as newly appointed if he or she enters the data after our sample period begins and works three consecutive months of regular caseloads.

racial differences in average risk or signal quality that can generate statistical discrimination. The distribution of MTE curves is identified by the quasi-experimental variation in pretrial release and misconduct rates across as-good-as-randomly assigned judges, and can be estimated by a tractable simulated minimum distance (SMD) procedure that matches moments of this quasi-experimental variation. We find evidence of both racial bias and statistical discrimination, with the latter coming from a higher level of average risk (that exacerbates discrimination) and less precise risk signals (that alleviates discrimination) for black defendants. Estimates of racial bias alone (e.g., Arnold et al., 2018; Marx, 2018) would therefore omit an important source of racial discrimination in our setting.

We conclude by using the MTE model to investigate whether discrimination can be effectively targeted (and potentially reduced) with existing data. We suppose that judges can be subjected to race-specific release rate quotas that eliminate racial disparities as estimated by a policymaker. We find that targeting the most discriminatory judges with a quota based on our quasi-experimental estimates can reduce the average level of discrimination by 35 percent and targeting all judges with such a quota can essentially eliminate discrimination, despite the noise in our estimation procedure. By comparison, targeting judges with a quota based on observational release rate disparities can lead to a reduced but non-zero level of discrimination against white defendants, due to OVB.

Our analysis builds on recent work using quasi-experimental variation to test for different forms of discrimination in the criminal justice system, bridging the gap between internally valid (but narrowly applicable) randomized audit studies and widely applicable (but potentially biased) observational analyses. Arnold et al. (2018) use the release tendencies of quasi-randomly assigned bail judges to test for racial bias at the margin of release under the strong first-stage monotonicity assumption, while Marx (2018) uses a similar approach to test for racial bias at the margin of police stops. We show how quasi-experimental variation can be used to measure all forms of racial discrimination, not just racial bias, without any such behavioral assumptions. We further show how the drivers of discrimination can be investigated by imposing more structure on the quasi-experimental variation.<sup>3</sup>

This paper also extends recent methodological advances in ATE and MTE estimation with multiple discrete instruments (Kowalski, 2016; Brinch et al., 2017; Mogstad et al., 2018; Hull, 2020). An important feature of our approach is that we do not impose the usual first-stage monotonicity assumption, which has received recent increased scrutiny both in general (Mogstad et al., 2019) and in the specific context of judge IV designs (Mueller-Smith, 2015; Frandsen et al., 2019; Norris, 2019).<sup>4</sup> Our extrapolation-based solution to estimating ATEs without monotonicity is most closely related to Hull (2020), who considers non-parametric extrapolations of quasi-experimental data in a similar setting. Our hierarchical solution to estimating a distribution of MTE curves without monotonicity is related to the contemporaneous approach of Chan et al. (2020), who estimate a structural model of doctor preferences and skill in making pneumonia diagnoses.<sup>5</sup>

<sup>&</sup>lt;sup>3</sup>In related work, Rose (2020) shows that a policy reform reducing imprisonment punishments for technical probation violations nearly eliminated a racial disparity in incarceration without significantly increasing differences in reoffending rates. These results suggest technical rule violations convey less precise risk signals for black individuals on probation.

<sup>&</sup>lt;sup>4</sup>Skepticism of conventional monotonicity in judge-IV designs is as old as the assumption itself. In their initial paper on the identification of local average treatment effects, Imbens and Angrist (1994) note that in the context of administrative screening "[monotonicity] requires that if official A accepts applicants with probability P(0), and official B accepts people with probability P(1) > P(0), official B must accept *any* applicant who would have been accepted by official A. This is unlikely to hold if admission is based on a number of criteria" (Example 2; p. 472).

 $<sup>{}^{5}</sup>$ Chan et al. (2020) model doctor decisions as following a hierarchical bivariate probit with variation in the latent index correlation across doctors. By comparison, we model judges as acting on posteriors from noisy risk signals with variation in signal quality across judges. We also show how this model can be used to form posterior MTE frontiers for each judge and race, and link these MTE frontiers to the different drivers of racial discrimination.

The remainder of the paper is organized as follows. Section 2 provides an overview of the NYC pretrial system. Section 3 outlines the conceptual framework underlying our analysis. Section 4 describes our data and documents release rate differences for observationally similar black and white defendants. Section 5 develops and implements our quasi-experimental approach to measuring racial discrimination in bail decisions. Section 6 develops and estimates our hierarchical MTE model to explore the drivers of this discrimination. Section 7 conducts policy counterfactuals. Section 8 concludes.

## 2 Setting

We study racial discrimination in the New York City pretrial system, one of the largest in the country. The U.S. pretrial system is meant to allow most criminal defendants to be released from legal custody while minimizing the risk of pretrial misconduct. Bail judges in both New York and the country as a whole are granted considerable discretion in determining which defendants should be released before trial, but they cannot discriminate against minorities and other protected classes even when membership in a protected class contains information about the underlying risk of criminal misconduct (Yang and Dobbie, 2019). Judges are also not meant to assess guilt or punishment when determining which individuals should be released from custody, nor are they meant to consider the political consequences of their bail decisions. Bail judges therefore risk violating U.S. law if they release white and black individuals with the same objective pretrial misconduct potential at different rates.<sup>6</sup>

In NYC, bail conditions are set by a judge at an arraignment hearing held shortly after an arrest. These hearings usually last just a few minutes and are held through a videoconference to the detention center. The judge typically receives detailed information on the defendant's current offense and prior criminal record, as well as a release recommendation based on a six-item checklist developed by a local nonprofit (New York City Criminal Justice Agency Inc., 2016). The judge has several options when setting bail conditions given this information. First, she can release defendants who show minimal risk on a promise to return for all court appearances, known broadly as release on recognizance (ROR) or release without conditions. Second, she can require defendants to post some sort of bail to be released. The judge can also send higher-risk defendants to a supervised release program as an alternative to cash bail. Finally, the judge can detain defendants pending trial by denying bail altogether.<sup>7</sup>

We exploit three features of the pretrial system in our analysis. First, the legal objective of bail judges is both narrow and measurable among the set of released defendants where pretrial misconduct outcomes are observed (although not among detained defendants, where such outcomes are unobserved). Second, bail judges can be effectively viewed as making binary "treatment" decisions, releasing low-risk defendants (generally by releasing without conditions or setting a low cash bail

<sup>&</sup>lt;sup>6</sup>Legally unwarranted racial disparities are not sufficient to establish unconstitutional behavior, but are a critical condition for such a determination. The Equal Protection Clause of the U.S. Constitution guarantees that all citizens have the right to equal justice under the law, including bail decisions. However, the Supreme Court has clarified that "official action will not be held unconstitutional [under the Equal Protection Clause] solely because it results in a racially disproportionate impact... Proof of racially discriminatory intent or purpose is required." (Arlington Heights v. Metropolitan Housing Development Corp., 429 U.S. 252, 264-65, 1977). In McCleskey v. Kemp, for example, the Supreme Court rejected a challenge to Georgia's capital punishment system despite statistical evidence of large racial disparities in death penalty decisions because the evidence was "clearly insufficient to support an inference that any of the decisionmakers in [the defendant's] case acted with discriminatory purpose." 481 U.S. 279, 281-82 (1987).

<sup>&</sup>lt;sup>7</sup>Cases such as murder, kidnapping, arson, and high-level drug possession and sale almost always result in a denial of bail, though these cases make up only about 0.8 percent of our sample. By comparison, about 70 percent of defendants in NYC are released ROR each year, nearly 30 percent are assigned cash bail or one less commonly used bail options such as insurance company bail bonds, and about 1.5 percent are sent to a supervised release program.

amount) and detaining high-risk defendants (generally by setting a high cash bail amount). We also explore different definitions of bail decisions in our analysis, such as viewing judges as deciding between release without conditions and any cash bail amount. Third, the case assignment procedures used in most jurisdictions, including NYC, generate quasi-random variation in judge assignment for defendants arrested at the same time and place. The quasi-random variation in judge assignment, in turn, generates exogenous variation in the probability of a defendant being released before trial.

There are also two differences between the NYC pretrial system and other pretrial systems around the country that are potentially relevant for our analysis. First, New York state instructs judges to only consider the risk that defendants will not appear for a required court appearance when setting bail conditions (a so-called failure to appear, or FTA), not the risk of new criminal activity as in most states. We explore robustness to this New York specific definition of pretrial misconduct in our analysis. Second, many defendants in NYC will never have bail set, either because the police gave them a desk appearance ticket that does not require an arraignment hearing or because the case was dismissed or otherwise disposed at the arraignment hearing before bail was set. The decision of whether or not to issue a desk appearance ticket is made before the bail judge is assigned, however, and cases should only be dismissed or otherwise disposed at arraignment if there is a clear legal defect in the case (Leslie and Pope, 2017). We show below that there is no relationship between the assigned bail judge and the probability that a case exits our sample due to case disposal or dismissal at arraignment, and exclude these cases from our analysis.

## **3** Conceptual Framework

### 3.1 Defining Racial Discrimination

We study racial discrimination in a setting where a set of decision-makers j make binary decisions  $D_{ij} \in \{0,1\}$  for an *iid* set of individuals i. Each decision-maker's goal is to align  $D_{ij}$  with an unobserved binary state  $Y_i^* \in \{0,1\}$ . In the context of bail decisions,  $D_{ij} = 1$  indicates the decision of judge j to release defendant i (with  $D_{ij} = 0$  otherwise) while  $Y_i^* = 1$  indicates that the defendant would subsequently fail to appear in court or be rearrested for a new crime if released (with  $Y_i^* = 0$  otherwise). Each judge's objective is to release individuals without misconduct potential ( $Y_i^* = 1$ ).<sup>8</sup>

We measure racial discrimination, both overall and for each judge, by the release rate disparity between white and black individuals with identical misconduct potential  $Y_i^*$ .<sup>9</sup> Letting  $R_i \in \{w, b\}$ index the race of white and black individuals, the level of discrimination for each judge j is given by:

$$\Delta_j = E[E[D_{ij} \mid Y_i^*, R_i = w] - E[D_{ij} \mid Y_i^*, R_i = b]]$$
(1)

The system-wide level of discrimination is given by the case-weighted average  $\Delta_j$  across all judges.

<sup>&</sup>lt;sup>8</sup>Appendix B.1 discusses how our approach can be extended to multi-valued or continuous  $Y_i^*$ .

<sup>&</sup>lt;sup>9</sup>Our measure is consistent with economic notions of discrimination that compare the treatment of white and black individuals with the same productivity such as Aigner and Cain (1977). By comparison, Phelps (1972) suggests measuring discrimination by comparing the treatment of white and black individuals with the same subjective signal of labor market productivity. Discrimination measures based on objective potential outcomes (as in both our paper and Aigner and Cain (1977)) and subjective signals of potential outcomes (as in Phelps (1972)) are generally identical when the quality of the subjective signals is identical by race, but can differ when individuals of different races tend to generate more or less informative subjective signals. We return this issue in Section 6, where we estimate a structural model that allows for more or less informative risk signals defendants of different races.

The inner difference of Equation (1) compares the potential release rates of white  $(R_i = w)$  and black  $(R_i = b)$  defendants assigned to judge j with the same misconduct potential  $Y_i^*$ . The outer expectation averages this conditional release rate comparison over the distribution of  $Y_i^*$ . We say that judge j discriminates against black defendants when  $\Delta_j > 0$ , that she discriminates against white defendants when  $\Delta_j < 0$ , and that she does not discriminate against either black or white defendants when  $\Delta_j = 0$ . By holding the potential defendant population fixed, estimates of  $\Delta_j$  can be used to calculate both the average level of racial discrimination in a bail system, as well as any variation in the level of discrimination across judges. As mentioned above, we interchangeably refer to  $\Delta_j$  as the level of racial discrimination for judge j and her unwarranted disparity.

With binary  $D_{ij}$  and  $Y_i^*$ , the  $\Delta_j$  parameters can also be understood as capturing racial differences in the tendency to correctly and incorrectly classify individuals with the same misconduct potential. Let  $\delta_{jr}^T = Pr(D_{ij} = 1 | Y_i^* = 0, R_i = r)$  denote the probability that judge j correctly releases defendants of race r without misconduct potential (her "true negative rate" for this race) and  $\delta_{jr}^F =$  $Pr(D_{ij} = 1 | Y_i^* = 1, R_i = r)$  denote the probability that judge j incorrectly releases defendants of race r with misconduct potential (her "false negative rate"). Equation (1) can then be written:

$$\Delta_j = \left(\delta_{jw}^T - \delta_{jb}^T\right) \left(1 - \bar{\mu}\right) + \left(\delta_{jw}^F - \delta_{jb}^F\right) \bar{\mu} \tag{2}$$

where  $\bar{\mu} = E[Y_i^*]$  denotes the overall risk of pretrial misconduct in the population of white and black defendants. Equation (2) shows that  $\Delta_j$  is a weighted average of racial differences in true and false negative rates for judge j. Since  $1 - \delta_{jr}^T = Pr(D_{ij} = 0 | Y_i^* = 0, R_i = r)$  denotes the probability that judge j incorrectly detains defendants of race r without misconduct potential (her "false positive rate" for this race), Equation (2) also shows that  $\Delta_j$  captures racial differences in type-I and type-II error rates. The system-wide level of discrimination similarly captures the case-weighted average racial difference in these error rates across all judges.

### 3.2 Theoretical Drivers of Discrimination

Racial discrimination in the sense of  $\Delta_j \neq 0$  can be driven by two distinct theoretical channels. The first is racial bias, in which judges discriminate against black defendants at the margin of pretrial release due to either taste-based discrimination (Becker, 1957) or inaccurate racial stereotyping (Bordalo et al., 2016). The second is statistical discrimination, in which judges act on accurate risk predictions but discriminate due to racial differences in average risk or the precision of received risk signals (Phelps, 1972; Arrow, 1973; Aigner and Cain, 1977). Racial discrimination can be defined and measured without a model for judicial decision-making, but understanding these theoretical channels requires more structure to be imposed on the binary release decisions  $D_{ij}$ .

We formalize the relationship between racial discrimination, racial bias, and statistical discrimination by considering a population of white and black defendants assigned to a single risk-neutral bail judge. Following the classic analysis of Aigner and Cain (1977), we suppose the judge observes each defendant's race  $R_i$  and a noisy signal of pretrial misconduct  $\nu_i = Y_i^* + \eta_i$  with normally distributed noise:  $\eta_i \mid Y_i^*, (R_i = r) \sim N(0, \sigma_r^2)$ . We allow both average misconduct risk  $\mu_r = E[Y_i^* \mid R_i = r]$ and the quality of risk signals  $\tau_r = 1/\sigma_r$  to vary by defendant race  $r \in \{w, b\}$ . We first assume the judge forms accurate posterior beliefs  $p(\nu_i, R_i) = Pr(Y_i^* = 1 \mid \nu_i, R_i)$  given the defendant's signal and race. We also assume the judge has a subjective benefit of releasing individuals of race r, given by  $\pi_r \in (0,1)$ . The risk-neutral judge then releases all defendants whose benefit exceeds the posterior risk cost, yielding the decision rule:

$$D_i = \mathbf{1}[\pi_{R_i} \ge p(\nu_i, R_i)] \tag{3}$$

Appendix B.2 derives the specific form of the posterior function  $p(\cdot)$ , completing the model.<sup>10</sup>

Racial bias in the sense of Becker (1957) arises when the judge perceives a lower benefit from releasing black defendants relative to white defendants with the same posterior risk, so that  $\pi_b < \pi_w$ . All else equal, such bias will lead to racial discrimination. By applying different thresholds to posterior misconduct risk, the judge generally makes different decisions for white and black defendants with the same misconduct potential  $Y_i^*$ . If, for example,  $\pi_b < \pi_w$  but both mean risk  $\mu_r$  and signal quality  $\tau_r$  are the same across race (implying a common distribution of  $p(\nu_i, R_i)$  given  $Y_i^*$ ), the judge will release fewer black defendants conditional on  $Y_i^*$ , implying  $\Delta_j > 0$ . Inaccurate racial stereotyping in the sense of Bordalo et al. (2016) can similarly result in discrimination and tends to be observationally equivalent to such racial animus (Arnold et al., 2018). In this case, even though judges believe they are applying the same threshold ( $\pi_b = \pi_w$ ), inaccurate posterior beliefs will lead judges to effectively set different release standards by race. Since inaccurate stereotyping and racial animus tend to be observationally equivalent, we use the term racial bias for both drivers of discrimination.

Statistical discrimination in the sense of Aigner and Cain (1977) arises when judges act on accurate race-specific predictions of defendant risk but discriminate because these predictions are affected by racial differences in either the average misconduct risk  $\mu_r$  or signal quality  $\tau_r$ . Differences in average misconduct risk  $\mu_r$  will tend to lead to lower release rates for defendants in the group with higher average misconduct risk, thereby resulting in discrimination against that group. Suppose, for example, that signal quality and release benefits are the same across race ( $\tau_b = \tau_w$  and  $\pi_b = \pi_w$ ) but the average level of risk is higher for black defendants ( $\mu_b > \mu_w$ ). The judge uses both the risk signal  $\nu_i$  and the defendant's race to accurately predict misconduct, so the judge's posterior  $p(\nu_i, R_i)$  will be higher among black defendants for every  $\nu_i$ . Black defendants will thus be less likely to be released conditional on  $Y_i^*$ , such that  $\Delta_j > 0$ , even though the judge's posterior threshold  $\pi_r$  and the distribution of risk signals  $\nu_i$  do not depend on race given  $Y_i^*$ . Statistical discrimination due to differences in signal quality  $\tau_r$  will instead have an ambiguous effect on release rates disparities. If, for example, a judge's release threshold  $\pi_r$  is higher than the average level of misconduct risk in the population  $\mu_r$ , then noisier risk signals will lead to fewer defendants of that race being detained given true misconduct potential, as judges place more weight on the mean risk  $\mu_r$  which falls below the threshold.<sup>11</sup>

Racial bias and statistical discrimination can both generate discrimination in the sense of  $\Delta_j > 0$ but yield different predictions for misconduct outcomes at the margin of release. When risk posteriors are accurate, misconduct outcomes at the margin of release capture the race-specific benefits of release:

$$E[Y_i^* \mid p(\nu_i; R_i) = \pi_{R_i}, R_i] = E[Y_i^* \mid E[Y_i^* \mid \nu_i, R_i] = \pi_{R_i}, R_i] = \pi_{R_i}$$
(4)

<sup>&</sup>lt;sup>10</sup>An alternative model of judicial decision-making specifies race-specific costs of misconduct classification errors. Appendix B.3 shows how such a model also leads to a threshold decision rule, with  $\pi_r$  denoting the relative cost of releasing defendants with misconduct potential.

<sup>&</sup>lt;sup>11</sup>The theoretical literature typically considers racial bias and statistical discrimination in isolation, while our empirical analysis allows racial differences in risk thresholds  $\pi_r$ , signal quality  $\tau_r$ , and mean risk  $\mu_r$  to each affect unwarranted disparity  $\Delta_j$ . We continue to refer to the case of  $\pi_w \neq \pi_b$  as racial bias in the model, while referring to  $\tau_w \neq \tau_b$  or  $\mu_w \neq \mu_b$  as statistical discrimination.

The model therefore predicts that marginal white and marginal black defendants should have the same misconduct rate at the margin of release if the judge is racially unbiased ( $\pi_w = \pi_b$ ), but that marginal white defendants should have a higher probability of misconduct if the judge is racially biased against black defendants ( $\pi_w > \pi_b$ ). Finding  $\pi_w \neq \pi_b$  thus rejects accurate statistical discrimination as the sole reason for finding  $\Delta_i \neq 0$ .

### 3.3 Empirical Challenges

Estimating racial discrimination is difficult because observational comparisons of white and black release rates cannot control for unobserved misconduct potential and are therefore likely to suffer from omitted variables bias (OVB). Testing for drivers of discrimination such as racial bias is also difficult unless judges have a common ordering of defendants by their appropriateness for release, satisfying a conventional but strong first-stage monotonicity assumption.

To formalize these empirical challenges, we introduce new notation for the data observed by an econometrician. Let  $Z_{ij} = 1$  if defendant *i* is assigned to judge *j*, let  $D_i = \sum_j Z_{ij}D_{ij}$  indicate the defendant's release status, and let  $Y_i = D_i Y_i^*$  indicate the observed pretrial misconduct outcome. The expression for observed pretrial misconduct reflects the fact that an individual who is detained  $(D_i = 0)$  cannot fail to appear in court or be rearrested for a new crime, such that  $Y_i = 0$  when  $D_i = 0$  despite individual *i*'s pretrial misconduct potential  $Y_i^*$ . The econometrician observes  $(R_i, (Z_{ij})_{j=1}^J, D_i, Y_i)$  for each defendant, and records whether the defendant is white in an indicator,  $W_i = \mathbf{1}[R_i = w]$ .

Rotational assignment of arraignment shifts can generate quasi-random assignment of individuals to different bail judges. To show how such quasi-experimental variation can and cannot help with measuring racial discrimination and testing its drivers, we assume here that judges are simply randomly assigned such that  $Z_{ij}$  is independent of  $(R_i, D_{ij}, Y_i^*)$  for each j. In practice, we relax this assumption to allow for the conditional quasi-random assignment often found in the bail system.

#### **Omitted Variables Bias in Observational Comparisons**

Observational comparisons of white and black release rates generally yield biased measures of racial discrimination. Such observational comparisons, whether in bail decisions or another area of the criminal justice system, usually come from "benchmarking" regressions of outcomes such as pretrial release on an indicator for an individual's race and controls for the observed characteristics of those individuals (e.g., Gelman et al., 2007; Abrams et al., 2012). Since pretrial misconduct potential  $Y_i^*$  is both unobserved and likely to affect release decisions, these benchmarking regressions are likely to produce biased estimates of the parameters of interest,  $\Delta_j$ .

To illustrate the OVB challenge, consider a simple judge-specific benchmarking regression:

$$D_i = \sum_j \alpha_j W_i Z_{ij} + \sum_j \phi_j Z_{ij} + \epsilon_i \tag{5}$$

where  $D_i$  is again an indicator for pretrial release,  $W_i Z_{ij}$  is the interaction of the indicator for the defendant being white and a fixed effect of each judge, and  $Z_{ij}$  are non-interacted judge fixed effects. We omit the constant term so that all judge fixed effects are included, and abstract away from other defendant observables for simplicity. The interaction coefficients thus measure the difference in judge j's release rates for white defendants relative to black defendants:

$$\alpha_j = E[D_i \mid R_i = w, Z_{ij} = 1] - E[D_i \mid R_i = b, Z_{ij} = 1]$$
(6)

While we focus here on a judge-specific benchmarking regression, the same conclusions emerge from an analysis of a simpler system-wide benchmarking regression of  $D_i = \phi + \alpha W_i + \epsilon_i$ .

Even with random judge assignment, the release rate disparities  $\alpha_j$  will tend to differ from the legally unwarranted disparities  $\Delta_j$ . When  $Z_{ij}$  is independent of  $(R_i, D_{ij}, Y_i^*)$ :

$$\alpha_j = E[D_{ij} \mid R_i = w] - E[D_{ij} \mid R_i = b]$$
(7)

Defining, as above,  $\mu_r = E[Y_i^* \mid R_i = r]$  as the average misconduct risk among individuals of race r and  $(\delta_{jr}^T, \delta_{jr}^F)$  as the judge's true and false negative rates for individuals of race r, these release rate disparities can be written:

$$\alpha_j = \left(\delta_{jw}^T (1 - \mu_w) + \delta_{jw}^F \mu_w\right) - \left(\delta_{jb}^T (1 - \mu_b) + \delta_{jb}^F \mu_b\right) \tag{8}$$

In contrast, judge j's unwarranted release rate disparity given by Equation (2) can be written:

$$\Delta_j = \left(\delta_{jw}^T (1-\bar{\mu}) + \delta_{jw}^F \bar{\mu}\right) - \left(\delta_{jb}^T (1-\bar{\mu}) + \delta_{jb}^F \bar{\mu}\right) \tag{9}$$

where  $\bar{\mu} = E[Y_i^*] = p_w \mu_w + p_b \mu_b$  is the average misconduct risk in the population of defendants, with  $p_r = Pr(R_i = r)$  denoting racial shares.

The difference between the benchmarking regression coefficient  $\alpha_j$  in Equation (8) and the judge discrimination measure  $\Delta_j$  in Equation (9) measures OVB in the simple benchmarking regression given by Equation (5). This difference is:

$$\xi_j \equiv \alpha_j - \Delta_j = \left(\delta_{jw}^T (\bar{\mu} - \mu_w) + \delta_{jw}^F (\mu_w - \bar{\mu})\right) - \left(\delta_{jb}^T (\bar{\mu} - \mu_b) + \delta_{jb}^F (\mu_b - \bar{\mu})\right)$$
$$= (\mu_b - \mu_w) \times \left[\left(\delta_{jw}^T - \delta_{jw}^F\right) p_b + \left(\delta_{jb}^T - \delta_{jb}^F\right) p_w\right]$$
(10)

where the second line follows by definition of the population risk  $\bar{\mu}$ . The regression coefficient  $\alpha_j$  will be biased upward for  $\Delta_j$  when  $\xi_j > 0$  and biased downward when  $\xi_j < 0$ .

Three insights follow from the OVB formula (10). First, conventional benchmarking regressions will generally yield biased estimates of the absolute level of discrimination  $\Delta_j$ , even with quasi-random judge assignment, unless either judge release decisions are independent of potential misconduct for each race (i.e.,  $E[D_{ij} | Y_i^*, R_i]$  does not depend on  $Y_i^*$ , so  $\delta_{jr}^T = \delta_{jr}^F$  for each r) or mean misconduct risk is identical across race (i.e.,  $\mu_w = \mu_b$ ). Randomized audit studies recover unbiased measures of discrimination by ensuring that race is itself as-good-as-randomly assigned across fictitious individuals, thereby ensuring that  $\mu_r$  does not depend on r within this population. But such randomization is infeasible in settings with high-stakes and face-to-face interactions like bail decisions.

Second, conventional benchmarking regressions will also yield biased estimates of the differences in the extent of racial discrimination across judges, even when judges are as-good-as-randomly assigned. The extent of OVB can also vary across judges in Equation (10), such that difference in benchmarking coefficients between judge j and k identifies  $\alpha_j - \alpha_k = \Delta_j - \Delta_k + \xi_j - \xi_k$  and not  $\Delta_j - \Delta_k$ . In general, OVB will vary across judges whenever there are differential responses to misconduct potential differences, such that  $\delta_{jr}^T - \delta_{jr}^F$  varies across j for either race r.<sup>12</sup>

Third, Equation (10) suggests a potential avenue for estimating racial discrimination when bail judges are as-good-as-randomly assigned, using familiar econometric objects. One of the terms driving the OVB of each  $\alpha_j$  is the difference in race-specific misconduct risk in the population:  $\mu_b - \mu_w$ , which is common to all judges. With  $Y_i^*$  capturing defendant *i*'s potential for pretrial misconduct when released and  $Y_i = 0$  for all detained individuals, the  $\mu_r = E[Y_i^* \mid R_i = r]$  can be understood as average treatment effects (ATEs) of pretrial release on pretrial misconduct among individuals of race *r*. We show in Section 5 how such ATEs can be estimated from quasi-experimental judge assignment and used to purge OVB from conventional benchmarking estimates, recovering valid estimates of  $\Delta_j$ .

#### Monotonicity Violations in Standard IV Estimates

Testing between the drivers of racial discrimination is also difficult unless judges have a common ordering of defendants by their appropriateness for release, satisfying a conventional first-stage monotonicity assumption. For example, standard IV methods can be used to test for racial bias (whether due to taste-based discrimination or inaccurate stereotyping) given the quasi-random assignment of judges and such first-stage monotonicity (Arnold et al., 2018; Marx, 2018). Monotonicity is, however, an especially strong assumption in this setting: it implies that judges are equally skilled in predicting an individual's propensity for pretrial misconduct and only differ in terms of the thresholds they set on a common posterior risk ordering.

To illustrate this potential limitation of the standard IV-based test for racial bias, we consider a multiple-judge generalization of the decision model in Section 3.2. The release rule for each judge j is given by  $D_{ij} = \mathbf{1}[\pi_{jR_i} \ge p_j(\nu_{ij}, R_i)]$ , where  $\pi_{jr}$  is the race-specific release benefit of judge j, and  $p_j(v, r)$  is the judge's posterior for the misconduct risk of a defendant of race r who sends a signal of v. The most general version of this model allows risk posteriors to differ across judges because of heterogeneous beliefs and risk signal qualities. Correspondingly, we index the signals  $\nu_{ij}$  of heterogeneous quality  $\tau_{jr}$  by j as well as by r. Judges with higher  $\tau_{jr}$  can be thought of as being more skilled, in that they base release decisions on more predictive signals of misconduct potential.

Conventional first-stage monotonicity identifies marginal misconduct outcomes for white and black defendants, which can be used to test for racial bias, by assuming judges form common risk posteriors. Per Imbens and Angrist (1994), when  $p_j(\nu_{ij}, R_i) = p(\nu_i, R_i)$  does not vary by j a linear IV regression of misconduct outcomes  $Y_i$  on release  $D_i$  instrumented by two quasi-randomly assigned judge indicators  $Z_{ij}$ , in a sample of either white or black individuals, identifies a local average treatment effect:

$$\frac{E[Y_i \mid Z_{ij} = 1, R_i] - E[Y_i \mid Z_{ik} = 1, R_i]}{E[D_i \mid Z_{ij} = 1, R_i] - E[D_i \mid Z_{ik} = 1, R_i]} = E[Y_i^* \mid \pi_{jR_i} \ge p(\nu_i, R_i) > \pi_{kR_i}]$$
(11)

where here the effect of "treating" individual *i* with release is simply her misconduct potential  $Y_i^*$ . Equation (11) thus gives the average misconduct risk for "compliers" of race  $R_i$ , whose posterior risk  $p(\nu_i, R_i)$  lies between the two judge benefit thresholds  $\pi_{jR_i}$  and  $\pi_{kR_i}$  (where  $\pi_{jR_i} \ge \pi_{kR_i}$  without loss).

<sup>&</sup>lt;sup>12</sup>To see this simply, suppose all judges are non-discriminatory:  $\delta_{jr}^T = \delta_j^T$  and  $\delta_{jr}^F = \delta_j^F$  for each j and r, such that  $\Delta_j = 0$  for each j. Suppose further that judges release all defendants without misconduct potential, such that  $\delta_j^T = 1$ . Differences in judge leniency are then solely due to differences in their rate of releasing defendants with misconduct potential,  $\delta_j^F$ . Equation (10) shows that these differences drive differences in OVB, since  $\xi_j = (\mu_b - \mu_w)(1 - \delta_j^F)$  in this case. Consequently, a benchmarking analysis would tend to incorrectly suggest not only racial discrimination ( $\xi_j > 0$ ) but also differential discrimination across judges ( $\xi_j \neq \xi_k$ ) when the average risk differs by race.

As these two thresholds become closer, the IV estimand in Equation (11) approaches the marginal released outcomes of each judge in Equation (4) and can therefore be used to test whether  $\pi_{jw} = \pi_{jb}$ . Arnold et al. (2018) show how standard linear and local IV procedures yield such tests in settings with many quasi-randomly assigned bail judges.

When judge skill varies, however, first-stage monotonicity is violated and standard IV procedures no longer capture average misconduct risk for marginal defendants. If  $\tau_{jr} \neq \tau_{kr}$ , then  $p_j(\nu_{ij}, R_i) \neq$  $p_k(\nu_{ik}, R_i)$  and the same linear IV regression instead identifies a non-convex linear combination of treatment effects for "complier" and "defier" populations:

$$\frac{E[Y_i \mid Z_{ij} = 1, R_i] - E[Y_i \mid Z_{ik} = 1, R_i]}{E[D_i \mid Z_{ij} = 1, R_i] - E[D_i \mid Z_{ik} = 1, R_i]} = p_{cR_i} E[Y_i^* \mid \pi_{jR_i} \ge p_j(\nu_{ij}, R_i), p_k(\nu_{ik}, R_i) > \pi_{kR_i}] \quad (12)$$
$$- p_{dR_i} E[Y_i^* \mid \pi_{kR_i} \ge p_k(\nu_{ik}, R_i), p_j(\nu_{ij}, R_i) > \pi_{jR_i}]$$

where  $p_{cr}$  is proportional to the complier share of the population of race r who is newly released when switching assignment from judge k to judge j, and  $p_{dr}$  is proportional to the defier share who is newly detained (with  $p_{cr} - p_{dr} = 1$ ). Unlike Equation (11), Equation (12) generally cannot be used to isolate marginal released outcomes and test whether  $\pi_{jw} = \pi_{jb}$ . Consequently, the IV-based tests for racial bias proposed by Arnold et al. (2018) are generally invalid when judge skill varies. In Section 6, we develop an alternative approach to test for racial bias in a model that allows for variation in judge skill. We also show how statistical discrimination due to racial differences in average risk or signal quality across race can be measured in this more realistic model with heterogeneous judge skill.

## 4 Data and Observational Comparisons

## 4.1 Sample and Summary Statistics

We observe the universe of 1,458,056 arraignments made in NYC between November 1, 2008 and November 1, 2013. Our data contain information on a defendant's gender, race, date of birth, and county of arrest, as well as the (anonymized) identity of the assigned bail judge. In our primary analysis, we categorize defendants as white (including both non-Hispanic and Hispanic white individuals), black (including both non-Hispanic and Hispanic black individuals), or neither. We explore alternative categorizations of race in robustness checks below.

In addition to detailed demographics, our data contain information on each defendant's current offense, history of prior criminal convictions, and history of past pretrial misconduct (both rearrests and FTA). We also observe whether the defendant was released at the time of arraignment and whether this release was due to release without conditions or some form of money bail. We categorize defendants as either released (including both release without conditions and with paid cash bail) or detained (including cash bail that is not paid) at the first arraignment, though we again explore robustness to other categorizations of the initial pretrial release decision below. Finally, we observe whether a defendant subsequently failed to appear for a required court appearance or was subsequently arrested for a new crime before case disposition. We take either form of pretrial misconduct as the primary outcome of our analysis, but again explore robustness to other measures below.

We make four key restrictions to arrive at our estimation sample. First, we drop cases where the defendant is not charged with a felony or misdemeanor (N=26,057). Second, we drop cases that were

disposed at arraignment (N=364,051) or adjourned in contemplation of dismissal (N=230,517). This set of restrictions drops cases that are likely to be dismissed by virtually every judge: Appendix Table A1 confirms that judge assignment is not systematically related to case disposal or case dismissal. Third, we drop cases in which the defendant is assigned a cash bail of \$1 (N=1,284). This assignment occurs in cases in which the defendant is already serving time in jail on an unrelated charge; the \$1 cash bail is set so that the defendant receives credit for jail time served, and does not reflect a new judge decision. Fourth, we drop defendants who are non-white and non-black (N=45,529). Finally, we drop defendants assigned to judges with fewer than 100 cases (N=3,785) and court-by-time cells with fewer than 100 total cases or only one unique judge (N=191,647), where a court-by-time cell is defined using the assigned courtroom, shift, day-of-week, month and year (e.g., the Wednesday night shift in Courtroom A of the Kings County courthouse in January 2012). The final sample contains 595,186 cases from 367,434 defendants assigned to 268 bail judges.<sup>13</sup>

Table 1 summarizes our estimation sample, both overall and by race. Panel A shows that 73.0 percent of defendants are released before trial. A defendant is defined as released before trial if either the defendant is released without conditions (ROR) or the defendant posts the required bail amount before disposition. The vast majority of these releases are without conditions, with only 14.4 percent of defendants being released after being assigned money bail. White defendants are more likely to be released before trial than black defendants, with a 76.7 percent release rate relative to a 69.5 percent release rate, respectively. Among released defendants, however, the distribution of release conditions is virtually identical across race.

Judges may release white defendants at a higher rate than black defendants because of relevant differences in observed defendant or charge characteristics. Consistent with this idea, Panel B of Table 1 shows that black defendants are 4.9 percentage points more likely to have been arrested for a new crime before trial in the past year compared to white defendants, as well as 3.0 percentage points more likely to have a prior FTA in the past year. Panel C further shows that black defendants are 1.3 percentage points more likely to have been charged with a felony compared to white defendants, as well as 3.6 percentage points more likely to have been charged with a violent crime. Finally, Panel D shows that black defendants who are released are 6.6 percentage points more likely to be rearrested or have an FTA than white defendants who are released (though the composition of such misconduct is similar). Importantly, and in contrast to the other statistics in Table 1, the risk statistics in Panel D are only measured among released defendants. An individual's potential for pretrial misconduct if released is, by definition, unobserved among detained individuals despite being the key legal objective for bail judges.

### 4.2 Quasi-Experimental Judge Assignment

Our empirical strategy exploits variation in pretrial release from the quasi-random assignment of judges who vary in the leniency of their bail decisions. There are three features of the NYC pretrial system that make it an appropriate setting for this research design.

First, NYC uses a rotation calendar system to assign judges to arraignment shifts in each of the

<sup>&</sup>lt;sup>13</sup>Appendix Table A2 compares the full sample of NYC bail cases to our estimation sample. By construction, our estimation sample has a somewhat higher release rate, although the ratio of release rates by race is similar. Our estimation sample is also broadly representative in terms of defendant and charge characteristics, with a slightly higher share of defendants with prior FTAs and rearrests, and a somewhat lower share of defendants charged with drug and property crimes. Pretrial misconduct rates are also elevated in our sample, though again with similar ratios by race.

five county courthouses in the city, generating quasi-random variation in bail judge assignment for defendants arrested at the same time and in the same place. Each county courthouse employs a supervising judge to determine the schedule that assigns bail judges to the day (9 a.m. to 5 p.m.) and night arraignment shift (5 p.m. to 1 a.m.) in one or more courtrooms within each courthouse. Individual judges can request to work certain days or shifts but, in practice, there is considerable variation in judge assignments within a given arraignment shift, day-of-week, month, and year cell.

Second, there is limited scope for influencing which bail judge will hear any given case, as most individuals are brought for arraignment shortly after their arrest. Each defendant's arraignment is also scheduled by a coordinator, who seeks to evenly distribute the workload to each open courtroom at an arraignment shift. Combined with the rotating calendar system described above and the processing time required before the arraignment, it is unlikely that police officers, prosecutors, defense attorneys, or defendants could accurately predict which judge is presiding over any given arraignment.

Finally, the rotation schedule used to assign bail judges to cases does not align with the schedule of any other actors in the criminal justice system. For example, different prosecutors and public defenders handle matters at each stage of criminal proceedings and are not assigned to particular bail judges, while both trial and sentencing judges are assigned to cases via different processes. As a result, we can study the effects of being assigned to a given bail judge as opposed to, for example, the effects of being assigned to a given set of bail, trial, and sentencing judges.

Appendix Table A3 verifies the quasi-random assignment of judges to bail cases in the estimation sample. Each column reports coefficient estimates from an ordinary least squares (OLS) regression of judge leniency on various defendant and case characteristics, with court-by-time fixed effects that control for the level of quasi-experimental bail judge assignment. We measure leniency using the leaveone-out average release rate among all other defendants assigned to a defendant's judge, following the standard approach in the literature (e.g., Arnold et al., 2018; Dobbie et al., 2018). Most coefficients in this balance table are small and not statistically significantly different from zero, both overall and by defendant race. A joint F-test fails to reject the null of quasi-random assignment at conventional levels of statistical significance.<sup>14</sup>

Appendix Table A4 further verifies that the assignment of different judges meaningfully affects the probability an individual is released before trial. Each column of this table reports coefficient estimates from an OLS regression of an indicator for pretrial release on judge leniency, court-by-time fixed effects, and, to boost precision, the baseline controls from Table 1. A one percentage point increase in the predicted leniency of an individual's judge leads to a 0.95 percentage point increase in the probability of release, with a somewhat smaller first-stage effect for white defendants and a somewhat larger effect for black defendants.

<sup>&</sup>lt;sup>14</sup>Even with the quasi-random assignment of bail judges, the exclusion restriction in our framework could be violated if judge assignment impacts the probability of pretrial misconduct through channels other than pretrial release. While the assumption that judges only systematically affect defendant outcomes through pretrial release is fundamentally untestable, we join Arnold et al. (2018) in viewing it as reasonable here. Bail judges only handle one decision, limiting the potential channels through which they could affect defendants. Pretrial misconduct is also a relatively short-run outcome, further limiting the role of alternative channels. In a similar setting, Dobbie et al. (2018) find that there are no independent effects of the assigned money bail amount on defendant outcomes. We explore the robustness of our findings to such effects below.

### 4.3 Observational Comparisons

Table 2 investigates system-wide racial disparity in NYC pretrial release rates. We estimate OLS regressions of the form:

$$D_i = \phi + \alpha W_i + X'_i \beta + \epsilon_i \tag{13}$$

where  $D_i$  is an indicator equal to one if defendant *i* is released,  $W_i$  is an indicator for the defendant being white, and  $X_i$  is a vector of baseline controls. Column 1 of Table 2 omits any controls in  $X_i$ , column 2 adds the defendant and case observables from Table 1 to  $X_i$ , and column 3 further adds court-by-time fixed effects to adjust for unobservable differences at the level of quasi-experimental bail judge assignment to  $X_i$ . Such regressions generally follow the conventional benchmarking approach from the literature (e.g., Gelman et al., 2007; Abrams et al., 2012).

Table 2 documents both statistically and economically significant release rate disparities between white and black defendants in NYC. The unadjusted white-black release rate difference  $\alpha$  is estimated in column 1 at 7.2 percentage points, with a standard error (SE) of 0.5 percentage points. This release rate gap is around 10 percent of the mean release rate of 73 percent. The release rate gap falls by 26 percent, to 5.3 percentage points (SE: 0.4), when we control for defendant and case observables, and by an additional 2 percent, to 5.2 percentage points (SE: 0.4), when we include court-by-time fixed effects. These estimates are similar in magnitude to the association, reported in column 3, between the probability of release and having an additional drug charge (-5.7 percentage points) or pretrial arrest (-6.8 percentage points) in the past year.

Figure 1 summarizes the distribution of judge-specific release rate disparities across the 268 bail judges in our sample. We estimate judge-specific disparities from OLS regressions of the form:

$$D_i = \sum_j \alpha_j W_i Z_{ij} + \sum_j \phi_j Z_{ij} + X'_i \beta + \epsilon_i$$
(14)

where  $D_i$  is again an indicator equal to one if defendant *i* is released,  $W_i Z_{ij}$  is the interaction between an indicator for the defendant being white and the fixed effects for each judge,  $Z_{ij}$  are the noninteracted fixed effects for each judge, and  $X_i$  is again a baseline control vector. We first estimate Equation (14) with  $X_i$  demeaned, such that the  $\alpha_j$  captures regression-adjusted difference in release rates for white and black individuals assigned to judge *j*. We then compute empirical Bayes posteriors of  $\alpha_j$  using standard shrinkage procedures (Morris, 1983). Figure 1 shows the distribution of the racial disparity posteriors that adjust only for the main judge fixed effects, following column 1 of Table 2, as well as the distribution of posteriors when we add both baseline controls and court-by-time fixed effects, following column 3 of Table 2. We also report in Figure 1 an estimate of the empirical Bayes prior mean and standard deviation of  $\alpha_j$  across judges, as well as the fraction of judges with positive  $\alpha_i$  by applying the posterior average effect approach of Bonhomme and Weidner (2020).<sup>15</sup>

The distributions of release rate disparity posteriors in Figure 1 are located well above zero, suggesting that nearly all judges in our sample release white defendants at a higher rate. We estimate that only 3.4 percent (SE: 1.3) of judges in our sample release a larger share of black defendants in the unadjusted specification, while only 6.0 percent (SE: 1.5) are estimated to release a larger share when we adjust for defendant and case observables and court-by-time fixed effects. Figure 1 nevertheless

 $<sup>^{15}\</sup>mathrm{See}$  Appendix  $\mathrm{B.4}$  for the details of the conventional empirical Bayes procedures we apply in this section.

shows considerable variation in the magnitude of the release rate disparities across judges. The standard deviation of  $\alpha_j$  is estimated at 3.9 percentage points (SE: 0.3) in the unadjusted specification, and 3.3 percentage points (SE: 0.3) when we adjust for baseline controls and court-by-time fixed effects. The average judge-specific disparities, which differ from the system-wide averages in Table 2 due to differences in weightings, are 6.9 percentage points (SE: 0.2) for the unadjusted specification and 5.0 percentage points (SE: 0.2) in the covariate-adjusted specification. When we weight the covariate-adjusted disparities by judge caseloads we obtain a system-wide disparity estimate of 5.3 percentage points (SE: 0.2), similar to the regression-weighted 5.2 percentage point disparity in Table 2.

Taken together, this observational analysis confirms large and pervasive racial disparities in NYC bail decisions, both in the raw data and after accounting for observable differences between white and black defendants. These results are, of course, consistent with bail judges discriminating against black defendants. But we cannot rule out the possibility that these disparities are driven by legally relevant differences between white and black defendants that are observed by bail judges but unobserved by the econometrician. This OVB concern is heightened by the fact that there is considerable variation in both release rates and race after conditioning on the defendant and case observables included in  $X_i$ . For example, by applying the method of Oster (2017), we find that defendant unobservables can completely explain the release rate gap between white and black defendants in column 2 of Table 2 if these unobservables are only half as predictive of race as the observables in  $X_i$ .

## 5 Quasi-Experimental Estimates of Racial Discrimination

### 5.1 Methods

We estimate racial discrimination in release decisions with observational release rate comparisons that are rescaled using quasi-experimental estimates of average white and black misconduct risk. This quasi-experimental approach leverages first-stage variation in judge leniency but, unlike standard IV methods, does not require a first-stage monotonicity assumption. We only require that average misconduct risk among white and black defendants can be extrapolated from the quasi-experimental data and that the judges' legal objective is well-specified by the econometrician.

The first key insight underlying our approach is that when judges are as-good-as-randomly assigned, the problem of measuring legally unwarranted release rate disparities for individual judges is equivalent to the problem of estimating the average misconduct risk among the full population of black and white defendants. The source of OVB in an observational benchmarking comparison is the correlation between race and unobserved misconduct potential among a given judge's pool of white and black defendants. With quasi-random judge assignment, this correlation is common to all judges and captured by race-specific population misconduct risk. Given estimates of these race-specific risk parameters, observed release outcomes can be appropriately rescaled to make released white and black defendants comparable in terms of their unobserved misconduct potential.

The rescaling that purges OVB from observational comparisons is given by expanding the true

and false negative rates from our definition of racial discrimination in Equation (2):

$$\delta_{jr}^{T} = E[D_{ij} \mid Y_{i}^{*} = 0, R_{i} = r] = \frac{E[D_{ij}(1 - Y_{i}^{*}) \mid R_{i} = r]}{E[1 - Y_{i}^{*} \mid R_{i} = r]} = \frac{E[D_{i}(1 - Y_{i}) \mid R_{i} = r, Z_{ij} = 1]}{1 - \mu_{r}}$$
(15)

$$\delta_{jr}^F = E[D_{ij} \mid Y_i^* = 1, R_i = r] = \frac{E[D_{ij}Y_i^* \mid R_i = r]}{E[Y_i^* \mid R_i = r]} = \frac{E[D_iY_i \mid R_i = r, Z_{ij} = 1]}{\mu_r}$$
(16)

where the third equalities in both lines follow from quasi-random judge assignment and the definition of mean risk  $\mu_r = E[Y_i^* \mid R_i = r]$ . Substituting these expanded true and false negative rates into Equation (2) yields:

$$\Delta_{j} = E[D_{i}(1 - Y_{i}) \mid R_{i} = w, Z_{ij} = 1] \frac{1 - \bar{\mu}}{1 - \mu_{w}} + E[D_{i}Y_{i} \mid R_{i} = w, Z_{ij} = 1] \frac{\bar{\mu}}{\mu_{w}} - E[D_{i}(1 - Y_{i}) \mid R_{i} = b, Z_{ij} = 1] \frac{1 - \bar{\mu}}{1 - \mu_{b}} - E[D_{i}Y_{i} \mid R_{i} = b, Z_{ij} = 1] \frac{\bar{\mu}}{\mu_{b}} = E[\omega_{i}D_{i} \mid R_{i} = w, Z_{ij} = 1] - E[\omega_{i}D_{i} \mid R_{i} = b, Z_{ij} = 1]$$
(17)

where

$$\omega_i = (1 - Y_i) \frac{1 - \bar{\mu}}{1 - \mu_{R_i}} + Y_i \frac{\bar{\mu}}{\mu_{R_i}}$$
(18)

The rewritten definition of discrimination in Equation (17) shows that judge j's level of discrimination  $\Delta_j$  is given by the  $\alpha_j$  coefficients in a simple benchmarking regression, where the release decisions  $D_i$  of each individual are rescaled by a positive factor  $\omega_i$ . This  $\omega_i$  reweights the sample to make released white and black defendants comparable in terms of their unobserved misconduct potential, thereby revealing the extent to which each judge discriminates against white and black defendants with identical misconduct potential (even though misconduct potential is unobserved and cannot be directly conditioned on). Equation (18) further shows that  $\omega_i$  is a function of observed misconduct outcomes  $Y_i$  and the unobserved average race-specific misconduct risk parameters  $\mu_r$ , where again  $\bar{\mu} = \mu_w p_w + \mu_b p_b$ . The key econometric challenge is therefore to estimate average misconduct risk among the full population of white and black defendants.

Appendix Table A5 uses a simple numerical example to illustrate how our rescaling approach allows us to measure racial discrimination in bail decisions, even though misconduct potential is unobserved and cannot be directly conditioned on. This example supposes that there are 100 defendants of each race and a single race-neutral judge who can perfectly predict misconduct potential, such that she releases all defendants with  $Y_i^* = 0$  (so  $\delta_{jr}^T = \delta_j^T = 1$ ) and detains all defendants with  $Y_i^* = 1$  (so  $\delta_{jr}^F = \delta_j^F = 0$ ). We also assume that 75 of the 100 hypothetical black defendants have misconduct potential ( $Y_i^* = 1$ ) but only 25 of the 100 hypothetical white defendants have misconduct potential, such that  $\mu_b = 0.75$  and  $\mu_w = 0.25$ . Panel A shows that the perfectly predictive judge therefore has a white release rate of 0.75 but a black release rate of 0.25, meaning that a conventional benchmarking regression would find that white defendants have a 50 percentage point higher release rate than black defendants ( $\alpha_j = 0.5$ ) despite the judge being race-neutral.

Panel B of Appendix Table A5 shows how discrimination can be measured in this simple numerical example with observational release rate comparisons that are rescaled using average white and black misconduct risk. Following Equations (17) and (18), we compute  $\omega_i = \frac{0.50}{0.75} = 2/3$  for released white

defendants with  $Y_i = 0$  and released black defendants with  $Y_i = 1$ , and  $\omega_i = \frac{0.50}{0.25} = 2$  for released white defendants with  $Y_i = 1$  and released black defendants with  $Y_i = 0$ . The rescaling factor thus up-weights the release rates of individuals who are relatively less common in each race (risky white defendants and non-risky black defendants), while down-weighting the release rates of individuals who are relatively more common (non-risky white defendants and risky black defendants).<sup>16</sup> In this way, the rescaling factor equalizes the proportion of risky and non-risky defendants by race, meaning that a rescaled benchmarking regression would correctly find that white and black defendants with the same misconduct potential have identical release rates ( $\Delta_i = 0$ ). This would continue to hold when the judge is not perfectly predictive, so long as she has consistent true and false negative rates across the races. Conventional and rescaled benchmarking regressions are identical when average misconduct risk does not vary by race; in this special case of  $\mu_w = \mu_b = \bar{\mu}$ ,  $\omega_i = 1$  for all defendants.

The second key insight underlying our approach is that the average race-specific misconduct risk parameters that enter Equation (17) can be estimated by extending recent advances in ATE estimation with multiple discrete instruments (Brinch et al., 2017; Hull, 2020), which build on a long literature on sample selection models (e.g. Heckman, 1990; Andrews and Schafgans, 1998). To build intuition for this approach, consider a setting with as-good-as-random judge assignment and a supremely lenient bail judge  $j^*$  who releases nearly all defendants regardless of their race or potential for pretrial misconduct. This supremely lenient judge's race-specific release rate among both black and white defendants is close to one:

$$E[D_i \mid Z_{ij^*} = 1, R_i = r] = E[D_{ij^*} \mid R_i = r] \approx 1$$
(19)

and the race-specific misconduct rate among defendants she releases is close to the race-specific average misconduct risk in the full population:

$$E[Y_i \mid D_i = 1, Z_{ij} = 1, R_i = r] = E[Y_i^* \mid D_{ij^*} = 1, R_i = r] \approx E[Y_i^* \mid R_i = r] = \mu_r$$
(20)

where the first equality in both expressions follows by quasi-random assignment. Without further assumptions, the decisions of a supremely lenient and quasi-randomly assigned judge can therefore be used to estimate the average misconduct risk parameters needed for our discrimination measure.

In the absence of such a supremely lenient judge, the required average misconduct risk parameters can be estimated using model-based or non-parametric extrapolations of release and misconduct rate variation across quasi-randomly assigned judges. These extrapolations use local IV variation to estimate the race-specific ATEs of pretrial release  $D_i$  on misconduct  $Y_i$ , which are equivalent to the race-specific average misconduct risk parameters  $\mu_r$ . Mean risk estimates may, for example, come from the vertical intercept, at one, of linear, quadratic, or local linear regressions of estimated released misconduct rates  $E[Y_i^* \mid D_{ij} = 1, R_i = r]$  on estimated release rates  $E[D_{ij} \mid R_i = r]$  across judges j within each race r. This extrapolation-based approach does not require conventional monotonicity, in contrast to the related discrete instrument methods of Kowalski (2016) and Brinch et al. (2017).<sup>17</sup>

<sup>&</sup>lt;sup>16</sup>This pattern of up- and down-weighting generally arises when black defendants have higher misconduct risk: i.e., when  $\mu_b > \bar{\mu} > \mu_w$ . In such cases, observations of released white defendants who subsequently offend are up-weighted  $(Y_i - \mu_w > 0 \text{ and } \bar{\mu} - \mu_w > 0 \text{ so } \omega_i > 1)$ , as are observations of released black defendants who do not subsequently offend  $(Y_i - \mu_b < 0 \text{ and } \bar{\mu} - \mu_b < 0, \text{ so again } \omega_i > 1)$ . Equation (17) also shows that  $\Delta_j = \alpha_j - (E[(1 - \omega_i)D_i \mid R_i = w, Z_{ij} = 1] - E[(1 - \omega_i)D_i \mid R_i = b, Z_{ij} = 1])$ , so that our rescaling can be understood as subtracting OVB from the observational comparisons with OVB given by a  $(1 - \omega_i)$ -scaled release rate disparity. <sup>17</sup>Formally, suppose in a population of individuals each judge's release decisions are given by  $D_{ij} = \mathbf{1}[\kappa_j \geq v_{ij}]$ ,

Our approach is instead analogous to the standard regression discontinuity approach of extrapolating average potential outcomes to a treatment cutoff from nearby observations. Here, variation in misconduct rates is extrapolated from quasi-randomly assigned judges with high release rates to the maximal release rate of a hypothetical supremely lenient judge.

## 5.2 Results

#### Mean Risk by Race

Figure 2 illustrates the quasi-experimental variation in judge release rates and released misconduct rates in NYC. The horizontal axis plots estimates of release rates  $E[D_{ij} | R_i = r]$  for each judge jand each race r, obtained from the earlier benchmarking regression in Equation (14) that adjusts for baseline observables and the court-by-time fixed effects that account for the level at which bail judges are quasi-randomly assigned to cases. The vertical axis plots corresponding estimates of released misconduct rates  $E[Y_i^* | D_{ij} = 1, R_i = r]$ , obtained from the analogous OLS regression:

$$Y_i = \sum_j \rho_j W_i Z_{ij} + \sum_j \zeta_j Z_{ij} + X'_i \gamma + u_i \tag{21}$$

estimated among released individuals ( $D_i = 1$ ), where again  $X_i$  contains baseline observables and court-by-time fixed effects and is demeaned to include all judge indicators.<sup>18</sup>

Figure 2 shows significant variation in race-specific release rates across judges, with several judges releasing a high fraction of their defendants for each race. Released misconduct rates tend to increase with judge leniency for both races, as would be predicted by a behavioral model in which more lenient judges release riskier defendants at the margin. This pattern is shown by the two solid lines in Figure 2, representing the race-specific lines-of-best-fit through the quasi-experimental data. The lines-of-best-fit are obtained by OLS regressions of judge-specific released misconduct rate estimates on judge-specific release rate estimates, with the judge-level regressions weighted inversely by the variance of misconduct rate estimation error. We also plot curves-of-best-fit from judge-level quadratic and local linear specifications as dotted and dashed lines, respectively, with both specifications again weighted inversely by the variance of misconduct rate estimation error. The simple linear specification fits the local IV variation well, with quadratic and local linear specifications yielding similar fits across most of the leniency distribution.

The vertical intercepts of the different curves-of-best-fit, at one, provide different estimates of the race-specific mean risk  $\mu_r$ . These estimates are reported in Panel A of Table 3. The simplest linear extrapolation, summarized in column 1, yields precise mean risk estimates of 0.352 (SE: 0.007) for white defendants and 0.395 (SE: 0.006) for black defendants. This extrapolation suggests that the

where  $v_{ij} \mid \kappa_j, \lambda_j \sim U(0, 1)$  without loss and  $E[Y_i^* \mid v_{ij}, \kappa_j, \lambda_j] = \mu + \lambda_j (v_{ij} - \frac{1}{2})$ . This model violates monotonicity, since judges differ both in their orderings of individuals by the appropriateness of release  $(v_{ij})$  and their relative skill at predicting misconduct outcomes  $(\lambda_j)$ . Nevertheless, when  $E[\lambda_j \mid \kappa_j]$  is constant (linear) in  $\kappa_j$ , average released misconduct rates  $E[Y_i^* \mid D_{ij} = 1, \kappa_j] = E[\mu + \frac{1}{2}\lambda_j(\kappa_j - 1) \mid \kappa_j]$  are linear (quadratic) in release rates  $E[D_{ij}] = \kappa_j$ , so that these extrapolations identify the ATE  $\mu$ . More flexible extrapolations generally accommodate a broader range of judge decision-making models by leveraging richer quasi-experimental variation. In the limit, local linear regressions can yield non-parametric estimates of mean misconduct risk provided there are many lenient judges (Hull, 2020).

can yield non-parametric estimates of mean misconduct risk provided there are many lement judges (run, 2020). <sup>18</sup>These specifications leverage an auxiliary assumption of linear conditional expectations of  $D_{ij}$  and  $Y_i^*$  to tractably accommodate the conditional random assignment of bail judges given the court-by-time fixed effects in  $X_i$ . If  $Z_i$  is independent of  $(Y_i^*, D_{i1}, \ldots, D_{iJ}, R_i)$  given  $X_i$  and  $E[Y_i^* | D_{ij} = 1, R_i = r, X_i] = \psi_{jr} + X'_i \gamma$ , then  $E[Y_i | R_i, Z_i, X_i, D_i = 1]$ is linear in  $(W_i Z_{i1}, \ldots, W_i Z_{iJ}, Z_{i1}, \ldots, Z_{iJ}, X'_i)'$ , as in Equation (21). Analogously, if  $E[D_{ij} | R_i = r, X_i] = \phi_{jr} + X'_i \beta$ , under conditional random assignment  $E[D_i | R_i, Z_i, X_i]$  is linear as in Equation (14).

average misconduct risk within the population of potential black defendants is 4.3 percentage points higher than among the population of potential white defendants in this setting. Equivalently, viewing mean risk as an ATE, it suggests that a quasi-randomly assigned judge in NYC would, by releasing all defendants, see a 4.3 percentage point higher rate of pretrial misconduct among black defendants than among white defendants. Per the discussion in Section 3.3, such a difference in misconduct risk is likely to generate OVB in observational release rate comparisons.

The quadratic and local linear extrapolations of quasi-experimental variation yield similar racespecific mean risk estimates, as can be seen from Figure 2. The quadratic fit suggests a slight nonlinearity in the relationship between judge leniency and released misconduct rates, with a slightly convex dashed line for white defendants and slightly concave dashed line for black defendants. Column 2 of Table 3 shows that these nonlinearities translate to a somewhat lower estimate of white mean risk, at 0.333 (SE: 0.019), and a higher estimate of black mean risk, at 0.415 (SE: 0.021). Near one, the non-parametric fit of Figure 2 coincides with the linear fit for white defendants and the quadratic fit for black defendants, yielding mean risk estimates in column 3 of 0.352 (SE: 0.014) and 0.424 (SE: 0.016), respectively. The implied racial gap in risk – and thus the potential for OVB – rises with these more flexible extrapolations, to 8.2 percentage points in column 2 and 7.2 percentage points in column 3. We take the most flexible local linear extrapolation as our baseline specification for analyzing racial discrimination in NYC, which we show below gives the most conservative estimate of average racial discrimination. We also explore the robustness of our results to a wide range of alternative mean risk estimates below.

The extrapolations in Figure 2 yield accurate mean risk estimates when judge release rules are accurately parameterized or when there are many highly lenient judges. Appendix Figure A1 validates our extrapolations by plotting race-specific extrapolations of average predicted misconduct outcomes, among released defendants, in place of actual released misconduct averages in Figure 2. We first construct predicted misconduct outcomes  $\hat{Y}_i^*$  using the fitted values from an OLS regression of actual pretrial misconduct  $Y_i^*$  on the baseline observables in column 3 of Table 2 in the subsample of released defendants. We then plot estimates of  $E[\hat{Y}_i^* \mid D_{ij} = 1, R_i = r]$  and  $E[D_{ij} = 1 \mid R_i = r]$ , constructed as in Figure 2, in Appendix Figure A1. Since  $\hat{Y}_i^*$  can be computed for the entire sample, we also include the overall averages  $E[\hat{Y}_i^* \mid, R_i = r]$  that are analogous to the race-specific ATEs of interest. Figure A1 shows that each of the linear, quadratic, and local linear extrapolations of predicted misconduct rates yields similar and accurate estimates of the overall actual averages. The 95 percent confidence intervals of the local linear extrapolations, for example, include the actual black average and only narrowly exclude the actual white average. These results build confidence for the extrapolations of actual pretrial misconduct outcomes in this setting.

#### **Racial Discrimination**

Panels B and C of Table 3 summarize the estimates of legally unwarranted racial disparities  $\Delta_j$  given the corresponding ATE estimates in Panel A. These estimates are obtained from the sample analogue of Equation (9), noting that a judge's true negative rates can be written:

$$\delta_{jr}^{T} = E[D_{ij} \mid Y_{i}^{*} = 0, R_{i} = r] = (1 - E[Y_{i}^{*} \mid D_{ij} = 1, R_{i} = r]) \frac{E[D_{ij} \mid R_{i} = r]}{1 - \mu_{r}}$$
(22)

and similarly for her false negative rate  $\delta_{jr}^F$ , while  $\bar{\mu} = \mu_w p_w + \mu_b p_b$ . We use the regression-adjusted estimates of  $E[D_{ij} | R_i = r]$  and  $E[Y_i^* | D_{ij} = 1, R_i = r]$  from Figure 2 and the sample share of black defendants to complete this formula. Case-weighted averages of the resulting  $\Delta_j$  estimates, reported in Panel B, estimate system-wide discrimination. We also compute empirical Bayes posteriors for individual  $\Delta_j$  again via standard shrinkage procedures (Morris, 1983). Summary statistics for the judge-level prior distribution (estimated as in Figure 1) are reported in Panel C.

We find that more than two-thirds of the system-wide release rate disparity between observably similar white and black defendants in NYC is explained by racial discrimination, with less than onethird explained by unobserved differences in misconduct risk. The most conservative estimate of system-wide discrimination in Table 3, which uses local linear extrapolations to estimate race-specific mean risk, is 68 percent (3.6 percentage points) of the case-weighted average disparity of 5.3 percentage points. By comparison, the least conservative estimate of the case-weighted average  $\Delta_j$ , which uses the linear extrapolations to estimate race-specific mean risk, implies that 83 percent (4.4 percentage points) of the average benchmarking disparity in Table 2 can be explained by discrimination. We thus find that unobservable differences in defendant risk can explain 17 to 32 percent (0.9 to 1.7 percentage points) of the average benchmarking disparity that remains after adjusting for baseline observables, similar to the 30 percent (1.9 percentage points) of the unadjusted average disparity explained by baseline observables in Table 2.<sup>19</sup>

Appendix Table A6 illustrates how our rescaling approach yields this finding of significant racial discrimination in NYC bail decisions, following the simple numerical example in Appendix Table A5. We use the benchmark local linear estimates of mean risk to estimate the number of white and black defendants with and without misconduct potential in column 2 of Panel A. In column 3, we combine these estimates with covariate-adjusted estimates of release and released misconduct rates to compute the number of released defendants in each race and misconduct category, as in Equation (22). This calculation yields the case-weighted average observational disparity of 5.3 percentage points in column 5. In Panel B, we use the local linear estimates of mean risk to compute and apply the appropriate rescaling factor  $\omega_i$ . Our baseline estimates of average misconduct risk are  $\mu_w = 0.352$ for white defendants and  $\mu_b = 0.424$  for black defendants. Combining these estimates with the share of white and black defendants in our sample yields an overall average misconduct risk of  $\bar{\mu} = 0.390$ . Following Equations (17) and (18), these estimates yield a rescaling factor of  $\omega_i = \frac{1-0.390}{1-0.353} = 0.942$  for released white defendants with  $Y_i = 0$ ,  $\omega_i = \frac{0.390}{0.424} = 0.919$  for released black defendants with  $Y_i = 1$ ,  $\omega_i = \frac{0.390}{0.353} = 1.107$  for released white defendants with  $Y_i = 1$ , and  $\omega_i = \frac{1-0.390}{1-0.424} = 1.060$  for released black defendants with  $Y_i = 0$ . Thus, the rescaling factor up-weights the release rates of risky white defendants and non-risky black defendants (who are relatively less common) while down-weighting the release rates of non-risky white defendants and risky black defendants (who are relatively more common). Applying these rescaling factors to the observational release rates yields a system-wide discrimination estimate of 3.6 percentage points in column 5, as also reported in Table 3.

Figure 3 plots the full distribution of discrimination posteriors across individual bail judges using the most conservative local linear mean risk estimates. For comparison, we also include the distribution of observed racial disparities from our most complete benchmarking model. The former distribution

<sup>&</sup>lt;sup>19</sup>We can also use the decomposition (2) to compute the case-weighted disparity in true and false negative rates generating the overall 3.6 percentage point release rate disparity. From our baseline local linear extrapolation we obtain an average  $\delta_{jw}^T - \delta_{jb}^T$  of 0.017 (SE: 0.029) and an average  $\delta_{jw}^F - \delta_{jb}^F$  of 0.064 (SE: 0.037). While noisy, these estimates suggest judges favor white defendants over black defendants in both the  $Y_i^* = 0$  and  $Y_i^* = 1$  subpopulations.

is shifted evenly to the left of the latter distribution, consistent with nontrivial OVB across the judge-specific estimates. Around 68 percent of the judge-weighted average benchmarking disparity (3.4 percentage points, out of 5.0 percentage points) is found to be due to discrimination, similar to the case-weighted decomposition from Panel B of Table 3. The standard deviation of judge-specific unwarranted disparities remains large, at 3.1 percentage points, though it shrinks somewhat from the 3.3 percentage point standard deviation of observed release rate disparities. The clear majority of NYC judges have positive  $\Delta_j$ , at 87.5 percent, though this share is also somewhat smaller than the 94.0 percent predicted by the conventional benchmarking model. Panel C of Table 3 shows that these statistics are again precisely estimated and similar across different mean risk estimates.

Our estimates show that there are both statistically and economically significant inequalities in the release rate decisions of black and white defendants with identical potential for pretrial misconduct. The most conservative estimate in Table 3, for example, implies that the unwarranted release rate gap could be closed if NYC judges released roughly 2,240 more black defendants each year (or detained roughly 2,240 more white defendants). Using the estimate of Dobbie et al. (2018), releasing this many defendants would lead to around \$67 million in recouped earnings and government benefits annually. We can also compare the average unwarranted disparity to other observed determinants of pretrial release. Table 2 shows, for example, that the most conservative 3.6 percentage point unwarranted disparity estimate corresponds to more than half of the decreased probability in release associated with a defendant having an additional pretrial arrest in the past year (-6.8 percentage points).

#### Robustness

Appendix Figure A2 examines the sensitivity of our system-wide discrimination estimate to different estimates of average white and black misconduct risk. We plot the range of unwarranted disparity estimates that we would obtain from different values of these risk inputs, with our linear, quadratic, and local linear estimates of average white and black risk indicated by solid, dashed, and dotted lines, respectively. The estimated level of discrimination against black defendants generally decreases as the assumed value of black misconduct risk increases, holding the value of white misconduct risk constant. Racial differences in misconduct potential would have to be extremely large, however, before we could conclude there is no discrimination against black defendants. At our baseline estimate of white mean risk, for example, the white-black difference in misconduct risk would need to be more than 15 percentage points (108 percent) larger than our most conservative estimates to conclude there is no discrimination.

Appendix Tables A7–A9 explore the robustness of our findings to alternative definitions of the judge's legal objective, the judge's decision variable, and the defendant's race. We find similar results when using a measure of pretrial misconduct that only includes FTA (column 2 of Appendix Table A7) or only includes new arrests (column 3 of Appendix Table A7). We also find a slightly higher case-weighted average unwarranted disparity, at 5.5 percentage points, when using a measure of pretrial misconduct that only includes new arrests for a violent crime (column 4 of Appendix Table A7), though this estimate is extremely imprecise due to the rareness of the outcome. We also find similar results when we specify the judge's binary decision as between release without conditions and setting any cash bail, with racial discrimination explaining at least 63 percent (2.6 percentage points) of the covariate-adjusted white-black ROR rate difference of 4.1 percentage points (column 3 of Appendix Table A8). Finally, we obtain similar results when categorizing defendants as non-Hispanic white or

any racial minority (including Hispanic white individuals and both non-Hispanic and Hispanic black individuals), with racial discrimination explaining at least 81 percent (5.9 percentage points) of the non-Hispanic white-minority differences in release decisions of 7.3 percentage points (column 3 of Appendix Table A9).<sup>20</sup>

#### Judge Heterogeneity

Table 4 explores variation in the level of discrimination across judges in our sample. Columns 1-5 report OLS estimates of the unwarranted disparity posteriors on indicators for whether a judge is newly appointed during our sample period, has above-average leniency, and has an above-median share of black defendants (as measured before adjusting for court-by-time fixed effects). We also include indicators for what county courtroom the judge hears most cases in. Columns 6-7 investigate the persistence of our discrimination measure over time by computing separate unwarranted disparity posteriors in the first and second half of cases that each judge sees in our sample period, recomputing the race-specific mean risk estimates in each half, and estimating OLS regressions of current unwarranted disparity posteriors on lagged unwarranted disparity posteriors and judge observables. In both sets of analyses, regressions of discrimination posteriors on judge observables can be interpreted through the posterior average effect framework of Bonhomme and Weidner (2020). We weight these regressions by estimates of the inverse posterior variance of the unwarranted disparities, with very similar results when weighting by judge caseload.

We find that there are significantly lower levels of discrimination among newly appointed judges, more lenient judges, and judges with a higher share of black defendants. Judges who are newly appointed in our sample period have 1.2 percentage point lower unwarranted disparities on average, while judges with above-average leniency have 0.8 percentage point lower unwarranted disparities. Judges assigned an above-median share of black defendants have 0.7 percentage point lower unwarranted disparities. We also find that judges who primarily see cases in the Manhattan, Queens, and Richmond county courtrooms tend to exhibit higher levels of discrimination, while those who primarily see cases in Brooklyn (the omitted reference category) and the Bronx have lower levels of discrimination. We find, for example, that unwarranted disparities are 2.3 percentage points higher for Manhattan judges compared to Brooklyn judges. Together, the observable judge characteristics available in our data explain about 31 percent of the variation in the unwarranted disparity posteriors, with the courtroom indicators alone explaining about 22 percent of the variation in unwarranted disparities.

We also find that the judge-specific discrimination estimates are highly correlated over time, with an autoregression coefficient of 0.52. Lagged unwarranted disparities alone explain about 28 percent of the variation in the current unwarranted disparities, with the lagged disparity and observable judge characteristics explaining about 34 percent of the variation in the unwarranted disparities. We also note that the average unwarranted disparity in the second half of judge cases is somewhat larger, at 4.7 percentage points, suggesting that discrimination may increase with judge experience.

<sup>&</sup>lt;sup>20</sup>Another potential concern is that measurement error in the judge's legal objective is systematically correlated with race. This could be an issue if, for example, judges minimize all new crime, not just new crime that results in an arrest, and the police are more likely to rearrest black defendants conditional on having committed a new crime. In this scenario, we will tend to overestimate the misconduct risk for black defendants compared to white defendants and, as a result, underestimate the true amount of racial discrimination in bail decisions (Knox et al., Forthcoming). It is therefore possible that our estimates reflect a lower bound on the true amount of racial discrimination in NYC, at least under the plausible assumption that the police are more likely to rearrest black defendants conditional on having committed a new crime. Reassuringly, column 2 of Appendix Table A7 shows a similar level of discrimination when we measure pretrial misconduct using just FTA, which is better measured and less subject to this concern.

Taken together, our results show that there is substantial racial discrimination in NYC bail decisions, both on average and for most judges, and that judge-specific estimates of discrimination are predicted by observable characteristics and highly-correlated over time. But our results do not speak to whether such discrimination is driven by racial bias or statistical discrimination, or whether we can effectively target and potentially reduce racial discrimination using existing data. We next consider an empirical framework to answer these questions.

# 6 MTE Estimates of Bias and Statistical Discrimination

### 6.1 Methods

We develop and estimate a novel hierarchical marginal treatment effects (MTE) model that imposes additional structure on the quasi-experimental variation to investigate whether discrimination in bail decisions is driven by racial bias or statistical discrimination and to conduct policy simulations. Building on the illustrative model in Section 3.2, we suppose that judges base release decisions on noisy signals of true misconduct potential. We allow for judge- and race-specific risk preferences and signal quality, with the latter allowing heterogeneous race-specific predictive skill across judges (in violation of the conventional first-stage monotonicity condition). The model implies a distribution of judge- and race-specific MTE curves that can be used to test for racial bias at the margin of release, as well as to measure racial differences in average risk or signal quality that can generate statistical discrimination.

We model judge risk signals as  $\nu_{ij} = Y_i^* + \eta_{ij}$ , where  $\eta_{ij} | Y_i^*, (R_i = r) \sim N(0, \sigma_{jr}^2)$  denotes the noise in judge j's risk signals for defendants of race r. Signal quality is given by the inverse standard deviation of noise,  $\tau_{jr} = 1/\sigma_{jr}$ , such that higher  $\tau_{jr}$  corresponds to more precise risk signals. Judges with higher  $\tau_{jr}$  can be thought to have a richer information set or being more skilled at inferring true misconduct potential from a common information set. Judges combine these race-specific signals  $\tau_{jr}$  with potentially biased prior beliefs  $\tilde{\mu}_{jr}$  of mean misconduct risk  $\mu_r$  for each race r and an understanding of the signal-generating process. The judges' risk posteriors  $p_j(\nu_{ij}; R_i)$  are therefore potentially biased solutions to the binary classification problem of whether defendant *i* would fail to appear or be rearrested for a new crime if released  $(Y_i^* = 1)$ , given the individual's race r and noisy misconduct signal  $\nu_{ij}$ . Appendix B.2 derives these posterior functions and shows that they are strictly increasing in the risk signal. Given release benefits  $\pi_{jr}$ , the release decisions of each risk-neutral judge therefore follow a signal-threshold rule of:

$$D_{ij} = \mathbf{1}[\pi_{jR_i} \ge p_j(\nu_{ij}; R_i)] = \mathbf{1}[\kappa_{jR_i} \ge Y_i^* + \eta_{ij}]$$
(23)

where  $\kappa_{jr} = p_j^{-1}(\pi_{jr}; r)$  is an implicit function of judge j's release benefit  $\pi_{jr}$ , subjective risk belief  $\tilde{\mu}_{jr}$ , and risk signal quality  $\tau_{jr}$  for defendants of race r. Appendix B.5 shows that when judges respond to misconduct risk, such that  $\delta_{jr}^T > \delta_{jr}^F$ , there exists a signal threshold  $\kappa_{jr}$  and signal quality  $\tau_{jr} > 0$  which rationalize the reduced-form true and false negative rates. Absent further restrictions, this model is thus without observational loss provided judge release decisions are better-than-random.

When known for each race, a judge's risk threshold  $\kappa_{jr}$  and signal quality  $\tau_{jr}$  can be used to characterize the extent of racial bias in release decisions. As discussed in Section 3.2, average misconduct outcomes at the margin of pretrial release capture the race-specific release benefits  $\pi_{jr} = E[Y_i^* \mid p_j(\nu_{ij}; r) = \pi_{jr}] = E[Y_i^* \mid Y_i^* + \eta_{ij} = \kappa_{jr}]$ , which can be used to compute racial bias for judge j.<sup>21</sup> These marginal released outcomes are known functions of  $\kappa_{jr}$  and  $\tau_{jr}$ , and represent marginal treatment effects (of release on pretrial misconduct) for defendants at the margin of release. Arnold et al. (2018) use marginal released outcomes to test for racial bias among quasi-randomly assigned bail judges under an assumption of first-stage monotonicity, which here requires  $\tau_{jr} = \tau_r$ to be common to all judges such that judges act as though there is a common ordering of defendants (of each race) with regards to their appropriateness for release. Under this restriction, the race-specific marginal released outcomes needed to test for bias can be estimated with conventional MTE estimation methods.

Here, our first insight is that when  $\kappa_{jr}$  and  $\tau_{jr}$  are known, we can also measure the extent of statistical discrimination. As discussed in Section 3.2, statistical discrimination arises when judges act on risk predictions that are affected by racial differences in either mean misconduct risk  $\mu_r$  or signal quality  $\tau_{jr}$ . Mean risk for each race r is given by integrating the marginal released outcome (or MTE) curve  $\mu_{jr}(\kappa) = E[Y_i^* | Y_i^* + \eta_{ij} = \kappa]$  of each judge j over the distribution of her risk signals. The slopes of these curves capture the quality of a judge's risk signals. Relatively more precise signals for white defendants relative to black defendants will, for example, lead to a steeper-sloping  $\mu_{jw}(\kappa)$  relative to  $\mu_{jb}(\kappa)$ . More generally, the judge- and race-specific MTE curves  $\mu_{jr}(\kappa)$  can be used to calculate the extent of racial discrimination in counterfactuals calculations where a judge's release rates are set to equalize marginal released outcomes and eliminate racial bias.

Our second insight is that we can estimate the key  $\kappa_{jr}$  and  $\tau_{jr}$  parameters without imposing a strong assumption of first-stage monotonicity. By requiring  $\tau_{jr} = \tau_r$ , monotonicity can be understood to restrict the MTE curves  $\mu_{jr}(\cdot)$  to be common across judges for each race r, such that variation in judge release rates only reflects differences in risk thresholds  $\kappa_{jr}$ . An implication of this monotonicity restriction is that, absent estimation error, the race-specific release rates  $E[D_{ij} | R_i = r]$  and released misconduct rates  $E[Y_i^* | D_{ij} = 1, R_i = r]$  plotted in Figure 2 will lie on a single curve determined by the common signal quality  $\tau_r$  and mean risk  $\mu_r$ . Tests of monotonicity based on this and similar implications have been developed in the context of quasi-randomly assigned judges (Mueller-Smith, 2015; Frandsen et al., 2019; Norris, 2019) and elsewhere (Kitagawa, 2015). These tests reject in our setting (see Appendix Table A10), suggesting that conventional monotonicity is unlikely to hold.<sup>22</sup>

We therefore substitute the conventional monotonicity restriction with an alternative parameterization of heterogeneity in judge skill, permitting a distribution of MTE curves  $\mu_{jr}(\cdot)$  across judges rather than restricting  $\mu_{jr}(\cdot) = \mu_r(\cdot)$  for each judge j. We specify the signal quality parameters  $\tau_{jr}$ as being log-normally distributed (imposing the domain restriction of  $\tau_{jr} > 0$ ), jointly with the signal thresholds  $\kappa_{jr}$ :  $\ln \tau_{jr} \sim N(\alpha_r, \psi_r^2)$  and  $\kappa_{jr} \sim N(\gamma_r, \delta_r^2)$  with nonzero correlations allowed across j and r. Appendix B.5 shows how this hierarchical approach can be viewed as parameterizing differences in how judges weigh different defendant characteristics, such as demeanor or prior arrest record.

The hyperparameters governing the distributions of judge-specific MTE curves are identified by quasi-experimental variation in pretrial release and misconduct rates, and can be estimated by a simulated minimum distance (SMD) procedure that matches moments of such quasi-experimental

<sup>&</sup>lt;sup>21</sup>Appendix B.2 shows how differences in release benefits and prior risk beliefs are observationally equivalent in this model: both enter the  $\kappa_{jr}$  multiplicatively, such that for any  $\kappa \in \mathbb{R}$  and  $\tau_{jr} > 0$  there exists a set of  $\pi_{jr}$  and  $\tilde{\mu}_{jr}$  (each ranging from 0 to 1) with  $\kappa_{jr} = \kappa$ . This equivalence reflects the general difficulty of disentangling racial bias due to biased beliefs (as in Bordalo et al., 2016) from racial bias due to taste-based discrimination (as in Becker, 1957).

 $<sup>^{22}</sup>$ Appendix Table A10 applies the generalized Sargan test of Frandsen et al. (2019) to samples of white and black defendants with increasingly flexible b-spline approximations to the function linking outcomes to judge release propensities. The chi-squared test statistics are consistently larger than the corresponding test degrees of freedom, suggesting violations of conventional first-stage monotonicity.

variation. This procedure, described in full detail in Appendix B.6, first estimates race-specific curvesof-best-fit through race-specific release and released misconduct rates (as in Section 5.2). We then match the estimated intercept, slope, and curvature of these curves-of-best-fit, as well as the residual variation in first-step estimates, to the corresponding moments of simulated quasi-experimental data drawn from different parameterizations of the hierarchical MTE model. Finally, we use the SMD estimates to compute empirical Bayes posteriors of the marginal released outcomes and signal quality of each judge and race given the hyperparameter estimates and observed quasi-experimental data.

Figure 4 builds intuition for the model's identification and SMD estimation by showing how differences in key hyperparameters manifest in the quasi-experimental variation. We construct this figure by first simulating draws of  $\kappa_{jr}$  and  $\tau_{jr}$  for a given race r across a large population of judges j with arbitrarily varying leniency. We then plot the implied distribution of judge release rates  $E[D_{ij} | R_i = r]$ and released misconduct rates  $E[Y_i^* \mid D_{ij} = 1, R_i = r]$ , abstracting away from first-step estimation error. Panels A and B set the variance of signal quality across judges to zero, satisfying the usual first-stage monotonicity restriction and ensuring that the judge moments fall on a common frontier. Panel A shows how differences in mean misconduct risk  $\mu_r$  lead to differences in the vertical intercept of these curves at one, or (per the discussion in Section 5.1) the release rate of a hypothetical supremely lenient judge. Panel B shows how differences in mean signal quality instead lead to different slopes of the curves, with higher  $\tau_r$  resulting in a steeper relationship between the share of defendants that a judge releases and the extent of pretrial misconduct among the released. Panels C and D of Figure 4 then relax first-stage monotonicity by allowing signal quality to vary across judges. In this case, the quasi-experimental variation no longer falls on a common frontier even without estimation error. Panel C shows that a higher variance in signal quality manifests as more dispersion in released misconduct rates among judges with similar release rates. Panel D shows that the trend in this distribution of points becomes more nonlinear when judge signal quality is more highly correlated with judge leniency.

### 6.2 Results

#### **Racial Bias and Statistical Discrimination**

Table 5 reports SMD estimates of mean misconduct risk  $\mu_r$ , the average of misconduct outcomes for marginally released defendants  $\mu_{jr}(\kappa_{jr})$ , and the average of judge signal quality  $\tau_{jr}$ , with the underlying hierarchical MTE model hyperparameter estimates reported in Appendix Table A11. The average difference in marginal misconduct outcomes between white and black defendants captures the overall extent of racial bias, while differences in either mean risk or signal quality by race capture statistical discrimination. Columns 1-3 of Table 5 report estimates under the conventional firststage monotonicity restriction that signal quality for defendants of a given race is constant across judges. Columns 4-6 relax this restriction, allowing judges to have different rankings of defendant appropriateness for pretrial release.<sup>23</sup>

 $<sup>^{23}</sup>$ The estimates in columns 1-3 of Table 5 are derived from the hyperparameter estimates in columns 1 and 4 of Appendix Table A11, while columns 4-6 of Table 5 come from columns 2 and 5 of Appendix Table A11. The latter assumes log signal quality and release thresholds are uncorrelated. A richer model that allows for such correlation is estimated in columns 3 and 6 of Appendix Table A11. This richer model produces estimates that are very similar to columns 2 and 5 but also considerably less precise. We therefore take the uncorrelated model as our baseline in Table 5. We note that our baseline model still allows for correlation between judge signal quality and marginal released outcomes, which we find to be large in Table 5. Appendix Figure A4 shows how these model hyperparameters fit the quasi-experimental variation by plotting the model-implied average released misconduct rate across races and judges of

In both sets of model estimates, we find evidence of both racial bias and statistical discrimination. with the latter coming from a higher level of average risk (that exacerbates discrimination) and less precise risk signals (that alleviates discrimination) for black defendants. Columns 4-6 of Table 5 show, for example, that the expected misconduct rate of typical white defendants at the margin of pretrial release is 0.609 (SE: 0.026), compared to 0.543 (SE: 0.018) for black defendants. The difference in these mean marginally released outcomes is a statistically significant 6.6 percentage points (SE: 0.030), indicating the existence of racial bias at the margin of release. Table 5 further shows considerable scope for statistical discrimination. First, the model estimates confirm the finding in Section 5.2 that mean risk is lower among white defendants than black defendants. In columns 4 and 5, this difference in mean misconduct risk is 2.9 percentage points, less than half the size of the 7.2 percentage point difference from our most conservative local linear extrapolation in Table 3 but similar to the 4.2 percentage point gap from our simple linear extrapolation. Second, we find that the typical judge acts on higher-quality risk signals for white defendants than for black defendants. Columns 4 and 5 report an average signal quality of 1.18 (SE: 0.08) for white defendants and 0.83 (SE: 0.07) for black defendants, implying that the typical noise in black risk signals is roughly 40 percent more dispersed. Per the discussion of Figure 4, this result is consistent with the white line-of-best-fit in Figure 2 being somewhat steeper than the black line-of-best-fit.<sup>24</sup> With a majority of white and black defendants released, higher white signal quality is likely to offset racial discrimination against black defendants arising from other channels (see Section 3.2). Together, the racial differences in mean risk and signal quality imply that analyses of racial bias alone (as in Arnold et al. (2018) and Marx (2018)) would omit an important source of discrimination in this setting.

Table 5 further suggests that the conventional first-stage monotonicity restriction is inconsistent with judge behavior in this setting. We find significant variation in judge signal quality when we relax this restriction and allow judges to have different rankings of defendant appropriateness for pretrial release in columns 4-6, with standard deviations of 0.14 (SE: 0.02) for white defendant signal quality and 0.12 (SE: 0.01) for black defendant signal quality. This variation in judge skill is highly correlated with variation in judge release preferences (which we also find to be sizable), with covariances between judge signal quality and marginal released outcomes of 0.006 for white defendants and 0.004 for black defendants (implying respective correlation coefficients of 0.72 and 0.68). While point estimates of the mean parameters with and without conventional monotonicity are qualitatively similar, the precision is higher without. The standard error on average racial bias, for example, falls by 29 percent, from 0.042 to 0.030, from column 3 to column 6. These precision gains also suggest that the model without monotonicity provides a better fit to quasi-experimental data, consistent with a visual analysis of Figure 2 and the formal tests in Appendix Table A10.

Appendix Table A14 uses the unrestricted model to quantify the joint role of racial bias and statistical discrimination in driving racial discrimination in NYC bail decisions. Column 1 summarizes the baseline degree of discrimination, racial bias, and differences in signal quality. The model-based

different leniencies, along with the estimates of release rates and released misconduct rates from Figure 2. Both modelimplied curves-of-best-fit are approximately linear, with slight upward curvature and a more steeply sloping curve for white defendants.

<sup>&</sup>lt;sup>24</sup>The mean signal quality estimates in Table 5 suggest that the typical NYC judge predicts misconduct risk with considerable accuracy for both races. In terms of the model, a  $\tau_{jr}$  of 1.18 (0.83) yields a receiver operating characteristics curve with an Area Under the Curve (AUC) statistic of 0.801 (0.724) for white (black) defendants. By comparison, Kleinberg et al. (2017b) obtain an AUC of 0.707 with a machine learning algorithm trained on FTA outcomes among released NYC defendants. Simpler logit models which use the observables in column 3 of Table 2 to predict  $Y_i^*$  among released defendants in our sample have AUCs of around 0.65 for both white and black individuals.

estimate of average unwarranted disparity, at 4.0 percentage points, is somewhat higher than our most conservative estimate in Table 3 but similar to the estimate we obtain from the simple linear extrapolation.<sup>25</sup> Column 2 shows that average racial discrimination significantly declines when judge leniency is counterfactually raised or lowered to equalize marginal released outcomes across white and black defendants (with Panel A generally raising black release rates and Panel B generally lowering white release rates). The average unwarranted disparity falls from 4.0 percentage points to -5.1 percentage points in Panel A and -2.0 percentage points in Panel B. This result shows that absent racial bias the average unwarranted disparity is reversed, with white defendants becoming less likely to be released than black defendants of identical misconduct potential. As expected, columns 3 and 4 show that this reversal is driven by the relatively higher signal quality for white defendants. Equalizing signal quality across races for each judge, with and without racial bias, again results in average racial discrimination against black defendants. The remaining statistical discrimination solely due to mean risk differences in column 4 yields a mean unwarranted disparity of 2.8 percentage points when black leniency and signal quality are counterfactually set, and an average unwarranted disparity of 4.3 percentage points when adjusting the corresponding white parameters.

#### Judge Heterogeneity

Appendix Tables A15–A16 explore variation in empirical Bayes posteriors of racial bias and signal quality differences, following our analysis of the unwarranted disparity posteriors in Section 5.2. We again report OLS estimates of the indicated posteriors on indicators for whether a judge is newly appointed during our sample period, has above-average leniency, has an above-median share of black defendants, and for what county courtroom the judge hears most cases in. We again weight these regressions by estimates of the inverse posterior variance of the outcome variables, with very similar results again obtained when weighting by judge caseload.

In Table A15, we find significantly higher levels of racial bias among newly appointed judges, more lenient judges, and judges with a below-median share of black defendants. Courtroom indicators are also highly predictive: together, the observable judge characteristics explain about 91 percent of the variation in racial bias across judges, with the courtroom indicators alone explaining 73 percent. We also find a moderately strong relationship between racial bias and overall discrimination, with our discrimination measure explaining 12 percent of the variation in the judge-specific bias.

In Table A16, we find a relatively smaller racial gap in signal quality among newly appointed judges and judges with an above-median share of black defendants. Here, judge leniency is not a significant predictor of judge-specific signal quality by race. Courtroom indicators and other observable characteristics of the judges again explain much of the variation in signal quality differences, with 79 percent of the variation explained when we include all judge observables. We find an even stronger relationship between signal quality differences and overall discrimination than between racial bias and discrimination, with our discrimination measure explaining 71 percent of the variation in the judge-specific signal quality.

<sup>&</sup>lt;sup>25</sup>All conclusions in Section 5.2, including the fraction of discriminatory NYC judges and heterogeneity results, continue to hold with the MTE model estimates of  $\mu_r$  (see Appendix Figure A3 and Appendix Table A12).

# 7 Policy Simulations

Lastly, we use our hierarchical MTE model estimates to investigate whether racial discrimination can be effectively targeted and potentially reduced with existing data. The reduced-form analysis in Section 5 shows that judge-specific unwarranted disparities are relatively stable over time, suggesting that identifying and targeting highly discriminatory judges for an appropriate intervention could help reduce future discrimination. This analysis also shows that no more than one-third of the observed release rate disparity between white and black defendants is explained by unobserved differences in misconduct risk, suggesting that observational regressions may also be useful for targeting judgespecific discrimination even in the absence of our quasi-experimental analysis. By linking unobserved differences in misconduct risk, racial bias, and statistical discrimination in the release decisions of each judge, the hierarchical MTE model provides the necessary structure to simulate the effects of reducing racial discrimination using existing observational and quasi-experimental data.

Table 6 summarizes simulations that target both unwarranted disparity posteriors (columns 2 and 3) and observational disparities (columns 4 and 5). The simulations suppose that individual bail judges can be subjected to race-specific release rate quotas that eliminate racial disparities, as estimated by a policymaker using either an observational or quasi-experimental analysis. The simulation based on the unwarranted disparity posteriors gauges the reliability of the individual predictions given the noise in our estimation procedure. The simulation based on observational disparities further tests whether conventional benchmarking regressions may be useful for targeting discrimination despite OVB. To simulate both sets of policies, we redraw all judge-specific parameters for each race from the estimated hierarchical MTE model 250 times, along with draws of appropriate estimation error. We use these to simulate 250 draws of the quasi-experimental variation plotted in Figure 2. We then re-estimate the MTE model in each draw and compute empirical Bayes posteriors, as in our analysis of the true data. Finally, we force all or a subset of simulated judges to adjust their race-specific leniencies to the point where their racial disparities are expected to be eliminated given the simulated model estimates and posteriors. Panel A simulates closing the targeted disparities for all judges, while Panel B simulates closing the targeted disparities for all judges, while Panel B simulates closing the targeted disparities for all judges to the point where their race disparities only for judges in the top quintile of the estimated disparities.<sup>26</sup>

The simulations suggest that racial discrimination can be reliably targeted using our estimated unwarranted disparity posteriors, despite estimation error. Targeting the disparities of all judges using the unwarranted disparity posteriors results in the virtual elimination of racial discrimination (columns 2 and 3 of Panel A), while targeting only judges in the top quintile results in a 35 percent reduction in the average level discrimination (columns 2 and 3 of Panel B). These simulated reductions are essentially unchanged when the targeted judges are forced to increase their leniency (typically for black defendants) in column 2 or decrease their leniency (typically for white defendants) in column 3. The average standard deviation of unwarranted disparity across judges, reported in brackets, is also reduced from around 3 percentage points to around 2 percentage points in all cases. Observational release rate disparities still remain when eliminating discrimination, however, as the higher level of mean risk for black defendants leads to OVB in the policy target.

 $<sup>^{26}</sup>$ Column 1 of Table 6 displays the baseline simulated average of the unwarranted disparities, observational disparities, and racial bias. Column 1 reports an average unwarranted disparity of 4.0 percentage points, which is a bit larger than the 3.6 percentage point average unwarranted disparity found with our most conservative local linear extrapolation in Section 5 due to the difference in model mean risk estimates. The estimated gap in mean risk implies a larger average observational disparity of 4.9 percentage points, which roughly matches the benchmarking regression estimate of 5.0 percentage points in Section 4.

Targeting judges with observational comparisons can also reduce discrimination, as the observed release rate disparities are highly correlated with the unwarranted disparity posteriors. Appendix Figure A5 shows, for example, that we obtain a high "forecast" coefficient of 0.874 (SE: 0.01) from regressing estimated judge-specific unwarranted disparity posteriors on observational disparity posteriors, along with a very high R-squared of 0.963. Consequently, we find in Table 6 that targeting all judges with simulated observational disparity posteriors reduces average unwarranted disparity by 5.0 percentage points (columns 4 and 5 of Panel A). The resulting average unwarranted disparity of -1.0 percentage point reflects the fact that the level of observed disparities is too high because of OVB. When targeting just the observational disparity posteriors in the top quintile of judges, the average unwarranted disparity is reduced by 43 percent but not reversed (columns 4 and 5 of Panel B). This finding, that observational benchmarking regressions can be useful for monitoring and targeting racial discrimination despite OVB, mirrors a result in the education setting on the utility of biased observational value-added measures (e.g., Angrist et al., 2017). There, as here, observational rankings prove to be highly predictive of policy-relevant parameters.<sup>27</sup>

# 8 Conclusion

There are large racial disparities at every stage of the criminal justice system, but it is unclear whether these disparities reflect racial discrimination or omitted variables bias. This paper shows that racial discrimination in bail decisions can be measured using observational comparisons of white and black release rates that are rescaled with quasi-experimental estimates of average white and black misconduct risk. Estimates from NYC show that more than two-thirds of the observed racial disparity in release decisions is due to racial discrimination, with less than one-third due to unobserved racial differences in misconduct risk. Leveraging a novel hierarchical MTE model, we show that this discrimination is driven by both racial bias and statistical discrimination, with the latter due to a higher level of average risk (that exacerbates discrimination) and less precise risk signals (that offsets discrimination) for black defendants. Policy simulations suggest that discrimination for individual judges can be reliably monitored and targeted with existing data.

The methods we use to study racial discrimination in bail decisions may prove useful for measuring discrimination in several other high-stakes settings, both within and outside the criminal justice system. One key requirement is the quasi-random assignment of decision-makers, such as judges, police officers, employers, government benefits examiners, or medical providers. A second requirement is that the objective of these decision-makers is both known and well-measured among the subset of individuals that the decision-maker endogenously selects. Mapping these settings to the quasi-experimental approach of this paper can help bridge the gap between internally valid (but narrowly applicable) experimental audit studies and often-deployed (but potentially biased) observational measures.

<sup>&</sup>lt;sup>27</sup>Our simulations also highlight the impossibility of simultaneously eliminating racial discrimination (on average) and racial bias (at the margin) when either mean misconduct risk or the risk signal quality differ for white and black defendants (Kleinberg et al., 2017a). The simulation based on the unwarranted disparity posteriors, for example, results in nonzero racial bias against black defendants of between 2.3 and 3.9 percentage points at the margin of release.

## References

- ABRAMS, D. S., M. BERTRAND, AND S. MULLAINATHAN (2012): "Do Judges Vary in Their Treatment of Race?" Journal of Legal Studies, 41, 347–383.
- AIGNER, D. J. AND G. G. CAIN (1977): "Statistical Theories of Discrimination in Labor Markets," Industrial and Labor Relations Review, 30, 157–187.
- ANDREWS, D. W. K. AND M. M. A. SCHAFGANS (1998): "Semiparametric Estimation of the Intercept of a Sample Selection Model," *Review of Economic Studies*, 65, 497–517.
- ANGRIST, J. D., P. D. HULL, P. A. PATHAK, AND C. R. WALTERS (2017): "Leveraging Lotteries for School Value-Added: Testing and Estimation," *The Quarterly Journal of Economics*, 132, 871–919.
- ANTONOVICS, K. AND B. KNIGHT (2009): "A New Look at Racial Profiling: Evidence from the Boston Police Department," *Review of Economics and Statistics*, 91, 163–177.
- ANWAR, S., P. BAYER, AND R. HJALMARSSON (2012): "The Impact of Jury Race in Criminal Trials," *Quarterly Journal of Economics*, 127, 1017–1055.
- ARNOLD, D., W. DOBBIE, AND C. S. YANG (2018): "Racial Bias in Bail Decisions," Quarterly Journal of Economics, 133, 1885–1932.
- ARROW, K. J. (1973): "The Theory of Discrimination," in *Discrimination in Labor Markets*, ed. by O. Ashenfelter and A. Rees, Princeton University Press, 3–33.
- BECKER, G. S. (1957): The Economics of Discrimination, University of Chicago Press.
- BERTRAND, M. AND S. MULLAINATHAN (2004): "Are Emily and Greg More Employable than Lakisha and Jamal? A Field Experiment on Labor Market Discrimination," *American Economic Review*, 94, 991–1013.
- BONHOMME, S. AND M. WEIDNER (2020): "Posterior Average Effects," Unpublished Working Paper.
- BORDALO, P., K. COFFMAN, N. GENNAIOLI, AND A. SHLEIFER (2016): "Stereotypes," *The Quarterly Journal of Economics*, 131, 1753–1794.
- BRINCH, C. N., M. MOGSTAD, AND M. WISWALL (2017): "Beyond LATE with a Discrete Instrument," *Journal of Political Economy*, 125, 985–1039.
- CHAN, D., M. GENTZKOW, AND C. YU (2020): "Selection with Variation in Diagnostic Skill: Evidence from Radiologists," *NBER Working Paper No. 26467.*
- DOBBIE, W., J. GOLDIN, AND C. YANG (2018): "The Effects of Pre-Trial Detention on Conviction, Future Crime, and Employment: Evidence from Randomly Assigned Judges," *American Economic Review*, 108, 201–240.
- EWENS, M., B. TOMLIN, AND L. CHOON WANG (2014): "Statistical Discrimination or Prejudice? A Large Sample Field Experiment," *Review of Economics and Statistics*, 96, 119–134.
- FRANDSEN, B. R., L. J. LEFGREN, AND E. C. LESLIE (2019): "Judging Judge Fixed Effects," NBER Working Paper No. 25528.

- GELMAN, A., J. FAGAN, AND A. KISS (2007): "An Analysis of the New York City Police Department's "Stop-and-Frisk" Policy in the Context of Claims of Racial Bias," *Journal of the American Statistical Association*, 102, 813–823.
- HECKMAN, J. J. (1990): "Varieties of Selection Bias," American Economic Review Papers and Proceedings, 80, 313–318.
- HECKMAN, J. J. AND E. VYTLACIL (2005): "Structural Equations, Treatment Effects, and Econometric Policy Evaluation," *Econometrica*, 73, 669–738.
- HULL, P. (2020): "Estimating Hospital Quality with Quasi-Experimental Data," Unpublished Working Paper.
- IMBENS, G. AND J. ANGRIST (1994): "A Least Squares Correction for Selectivity Bias," *Econometrica*, 62, 467–475.
- KITAGAWA, T. (2015): "A Test for Instrument Validity," Econometrica, 83, 2043–2063.
- KLEINBERG, J., S. MULLAINATHAN, AND M. RAGHAVAN (2017a): "Inherent Trade-Offs in Algorithmic Fairness," *Proceedings of Innovations in Theoretical Computer Science*, 43:1–43:23.
- KLEINBERG, J., , H. LAKKARAJU, J. LESKOVEC, J. LUDWIG, AND S. MULLAINATHAN (2017b): "Human Decisions and Machine Predictions," *The Quarterly Journal of Economics*, 133, 237–293.
- KNOX, D., W. LOWE, AND J. MUMMOLO (Forthcoming): "The Bias Is Built In: How Administrative Records Mask Racially Biased Policing," *American Political Science Review*.
- KOWALSKI, A. (2016): "Doing More When You're Running LATE: Applying Marginal Treatment Effect Methods to Examine Treatment Effect Heterogeneity in Experiments," *NBER Working Paper No. 22363.*
- LESLIE, E. AND N. G. POPE (2017): "The Unintended Impact of Pretrial Detention on Case Outcomes: Evidence from NYC Arraignments," *Journal of Law and Economics*, 60, 529–557.
- MARX, P. (2018): "An Absolute Test of Racial Prejudice," Unpublished Working Paper.
- McINTYRE, F. AND S. BARADARAN (2013): "Race, Prediction, and Pretrial Detention," *Journal of Empirical Legal Studies*, 10, 741–770.
- MOGSTAD, M., A. SANTOS, AND A. TORGOVITSKY (2018): "Using Instrumental Variables for Inference About Policy-Relevant Treatment Parameters," *Econometrica*, 86, 1589–1619.
- MOGSTAD, M., A. TORGOVITSKY, AND C. R. WALTERS (2019): "Identification of Causal Effects with Multiple Instruments: Problems and Some Solutions," *NBER Working Paper No. 25691*.
- MORRIS, C. N. (1983): "Parametric Empirical Bayes Inference: Theory and Applications," *Journal* of the American Statistical Association, 78, 47–55.
- MUELLER-SMITH, M. (2015): "The Criminal and Labor Market Impacts of Incarceration," Unpublished Working Paper.
- NEW YORK CITY CRIMINAL JUSTICE AGENCY INC. (2016): "Annual Report 2015," Tech. rep.

- NORRIS, S. (2019): "Examiner Inconsistency: Evidence from Refugee Appeals," Unpublished Working Paper.
- OSTER, E. (2017): "Unobservable Selection and Coefficient Stability: Theory and Evidence," *Journal of Business and Economic Statistics*, 37, 187–204.
- PAKES, A. AND D. POLLARD (1989): "Simulation and the Asymptotics of Optimization Estimators," Econometrica: Journal of the Econometric Society, 1027–1057.
- PHELPS, E. S. (1972): "The Statistical Theory of Racism and Sexism," *American Economic Review*, 62, 659–661.
- REHAVI, M. M. AND S. B. STARR (2014): "Racial Disparity in Federal Criminal Sentences," *Journal* of *Political Economy*, 122, 1320–1354.
- ROSE, E. (2020): "Who Gets a Second Chance? Effectiveness and Equity in Supervision of Criminal Offenders," *Unpublished Working Paper*.
- YANG, C. AND W. DOBBIE (2019): "Equal Protection Under Algorithms: A New Statistical and Legal Framework," Unpublished Working Paper.



Figure 1: Observational Release Rate Disparities

*Notes.* This figure plots the distribution of observational release rate disparity posteriors for the 268 judges in our sample. Estimates are from the coefficients of an OLS regression of an indicator for pretrial release on white×judge fixed effects, controlling for judge main effects. Empirical Bayes posteriors are computed using a standard shrinkage procedure, as described in Appendix B.4. The unadjusted line shows the distribution of posteriors when controlling only for the main judge fixed effects. The covariate-adjusted posterior distribution adds the baseline controls from Table 2 and court-by-time fixed effects. Means and standard deviations refer to the estimated prior distribution. The fractions of positive disparities are computed as posterior average effects, as described in Appendix B.4.



Figure 2: Judge-Specific Release Rates and Conditional Misconduct Rates

*Notes.* This figure plots race-specific release rates for the 268 judges in our sample against rates of pretrial misconduct among the set of released defendants. All estimates adjust for the baseline controls in Table 2 and court-by-time fixed effects. The figure also plots race-specific linear, quadratic, and local linear curves of best fit, obtained from judge-level regressions that inverse-weight by the variance of the estimated misconduct rate among released defendants. The local linear regressions use a Gaussian kernel with a race-specific rule-of-thumb bandwidth.


Figure 3: Observational and Unwarranted Release Rate Disparities

Notes. This figure plots the distribution of observational and unwarranted release rate disparity posteriors for the 268 judges in our sample. Covariate-adjusted disparities are estimated by the coefficients of an OLS regression of an indicator for pretrial release on white×judge fixed effects, controling for judge main effects, the baseline controls from Table 2, and court-by-time fixed effects. Unwarranted disparities are estimated as described in Section 5, using the local linear extrapolations from Figure 2 to estimate the mean risk of each race. Empirical Bayes posteriors are computed using a standard shrinkage procedure, as described in Appendix B.4. Means and standard deviations refer to the estimated prior distribution. The fractions of positive disparities are computed as posterior average effects, as described in Appendix B.4.



Figure 4: Identification of Hierarchical MTE Model Parameters

Notes. This figure plots simulated race- and judge-specific release rates against rates of pretrial misconduct among the set of released defendants under different parameterizations of the hierarchical MTE model described in the text. Panel A plots differences in mean misconduct risk ( $\mu = 0.4$  vs.  $\mu = 0.3$ ) when conventional MTE monotonicity holds ( $\psi = 0$ ). Panel B plots differences in mean signal quality ( $\alpha = 1$  vs.  $\alpha = 0$ ) when conventional MTE monotonicity holds ( $\psi = 0$ ). Panel C plots differences in signal quality variance ( $\psi = 0.4$  vs.  $\psi = 0.1$ ). Panel D plots differences in the covariance between judge signal quality and judge leniency ( $\beta = 2$  vs.  $\beta = 0.1$ ). The default parameterization is  $\mu = 0.4$ ,  $\alpha = 0.2$ ,  $\psi = 0.1$ ,  $\beta = 0$ ,  $\gamma = 1.3$ , and  $\delta = 1$ .

	All	White	Black
	Defendants	Defendants	Defendants
Panel A: Pretrial Release	(1)	(2)	(3)
Released Before Trial	0.730	0.767	0.695
Share ROR	0.852	0.852	0.851
Share Money Bail	0.144	0.144	0.145
Share Other Bail Type	0.004	0.004	0.004
Share Remanded	0.000	0.000	0.000
Panel B: Defendant Characteris	tics		
White	0.478	1.000	0.000
Male	0.821	0.839	0.804
Age at Arrest	31.97	32.06	31.89
Prior Rearrest	0.229	0.204	0.253
Prior FTA	0.103	0.087	0.117
Panel C: Charge Characteristics			
Number of Charges	1.150	1.184	1.118
Felony Charge	0.362	0.355	0.368
Misdemeanor Charge	0.638	0.645	0.632
Any Drug Charge	0.256	0.257	0.256
Any DUI Charge	0.046	0.067	0.027
Any Violent Charge	0.143	0.124	0.160
Any Property Charge	0.136	0.127	0.144
Panel D: Pretrial Misconduct, W	Vhen Released		
Pretrial Misconduct	0.299	0.266	0.332
Share Rearrest Only	0.499	0.498	0.499
Share FTA Only	0.281	0.296	0.269
Share Rearrest and FTA	0.220	0.205	0.232
Total Cases	595,186	284,598	310,588
Cases with Defendant Released	434,201	$218,\!256$	$215,\!945$

Table 1: Descriptive Statistics

*Notes.* This table summarizes the NYC analysis sample. The sample consists of bail hearings that were quasirandomly assigned judges between November 1, 2008 and November 1, 2013, as described in the text. Information on demographics and criminal outcomes is derived from court records as described in the text. Pretrial release is defined as meeting the bail conditions set by the first assigned bail judge. ROR (released on recognizance) is defined as being released without any conditions. FTA (failure to appear) is defined as failing to appear at a mandated court date.

	(1)	(2)	(3)
White	0.072	0.053	0.052
	(0.005)	(0.004)	(0.004)
Male		-0.097	-0.092
		(0.005)	(0.004)
Age at Arrest		-0.005	-0.005
		(0.000)	(0.000)
Prior Rearrest		-0.066	-0.068
		(0.004)	(0.004)
Prior FTA		-0.209	-0.208
		(0.005)	(0.005)
Felony Charge		-0.192	-0.171
		(0.006)	(0.005)
Any Drug Charge		-0.055	-0.057
		(0.007)	(0.007)
Any DUI Charge		0.116	0.119
		(0.004)	(0.004)
Any Violent Charge		-0.137	-0.146
		(0.007)	(0.007)
Any Property Charge		-0.070	-0.072
		(0.005)	(0.005)
Baseline Controls	No	Yes	Yes
Court x Time FE	No	No	Yes
Mean Release Rate	0.730	0.730	0.730
Cases	595.186	595.186	595.186

Table 2: Regression Estimates of System-Wide Release Rate Disparity

*Notes.* This table reports OLS estimates of regressions of an indicator for pretrial release on defendant characteristics. The regressions are estimated on the sample described in Table 1. Robust standard errors, two-way clustered at the individual and the judge level, are reported in parentheses.

	Linear	Quadratic	Local Linear
	Extrapolation	Extrapolation	Extrapolation
Panel A: Mean Risk by Race	(1)	(2)	(3)
White Defendants	0.352	0.333	0.352
	(0.007)	(0.019)	(0.014)
Black Defendants	0.395	0.415	0.424
	(0.006)	(0.021)	(0.016)
Panel B: System-Wide Discrim	nination		
Mean Across Cases	0.044	0.037	0.036
	(0.002)	(0.006)	(0.005)
Panel C: Judge-Level Discrime	ination		
Mean Across Judges	0.043	0.035	0.034
	(0.002)	(0.006)	(0.005)
Std. Dev. Across Judges	0.031	0.030	0.031
	(0.003)	(0.003)	(0.003)
Fraction Positive	0.922	0.884	0.875
	(0.017)	(0.041)	(0.035)
Judges	268	268	268

Table 3: Mean Risk and Unwarranted Disparity Estimates

*Notes.* This table summarizes estimates of mean risk and unwarranted racial disparities from different extrapolations of the variation in Figure 2. Panel A reports estimates of race-specific average misconduct risk, Panel B reports estimates of system-wide (case-weighted) unwarranted disparity, and Panel C reports empirical Bayes estimates of summary statistics for the judge-level unwarranted disparity prior distribution. To estimate mean risk, column 1 uses a linear extrapolation of the variation in Figure 2, while column 2 uses a quadratic extrapolation and column 3 uses a local linear extrapolation with a Gaussian kernel and a rule-of-thumb bandwidth. Robust standard errors, two-way clustered at the individual and judge level, are obtained by a bootstrapping procedure and appear in parentheses.

						Split-S	Sample
		Full-Sample Disparities					arities
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
New Judge	-0.012				-0.011		-0.004
	(0.004)				(0.003)		(0.004)
Lenient Judge		-0.008			-0.010		-0.005
		(0.003)			(0.003)		(0.002)
Above-Median Black Share			-0.007		-0.006		0.002
			(0.003)		(0.004)		(0.003)
Manhattan Courtroom				0.023	0.021		0.014
				(0.004)	(0.004)		(0.003)
Bronx Courtroom				-0.003	-0.006		0.007
				(0.003)	(0.004)		(0.004)
Queens Courtroom				0.014	0.008		0.009
				(0.004)	(0.005)		(0.004)
Richmond Courtroom				0.010	0.005		0.016
				(0.004)	(0.006)		(0.004)
Lagged Disparity						0.518	0.416
						(0.062)	(0.071)
Mean Disparity	0.034	0.034	0.034	0.034	0.034	0.047	0.047
R2	0.043	0.035	0.027	0.223	0.312	0.280	0.342
Judges	268	268	268	268	268	252	252

Table 4: Unwarranted Disparities and Judge Characteristics

*Notes.* This table reports OLS estimates of regressions of unwarranted disparity posteriors on judge characteristics. Unwarranted disparities are estimated as described in Section 5, using the benchmark local linear estimate of mean risk. New judges are defined as judges appointed during our estimation period. Lenient judges are defined as judges with above-average leniency, controlling for court-by-time fixed effects. Courtroom locations are defined using the location of the modal case heard by each judge. Split-sample disparities are computed by splitting each judge's sample of cases at the median case and constructing two samples, a before-median case sample and an after-median case sample. Unwarranted disparities are then re-estimated within each subsample. The estimation procedure conditions on court-by-time effects, which causes a small number of judge effects to become collinear with the court-by-time effects and dropped. All specifications are weighted by the inverse variance of the unwarranted disparity posteriors. Robust standard errors are reported in parentheses.

	With	n Monotor	nicity	Without Monotonicity		
	White	Black	Diff.	White	Black	Diff.
	(1)	(2)	(3)	(4)	(5)	(6)
Mean Misconduct Risk	0.374	0.429	-0.055	0.400	0.429	-0.029
	(0.014)	(0.012)	(0.019)	(0.007)	(0.007)	(0.009)
Mean Marginal Released Outcome	0.585	0.526	0.059	0.609	0.543	0.066
	(0.026)	(0.042)	(0.042)	(0.026)	(0.018)	(0.030)
Mean Signal Quality	1.866	1.321	0.546	1.184	0.828	0.356
	(0.074)	(0.211)	(0.223)	(0.084)	(0.065)	(0.101)
Marginal Outcome Std. Dev.	0.192	0.105	0.087	0.060	0.050	0.010
	(0.020)	(0.028)	(0.030)	(0.006)	(0.004)	(0.007)
Signal Quality Std. Dev.				0.138	0.118	0.020
				(0.020)	(0.013)	(0.023)
Covariance of Signal Quality and				0.006	0.004	0.003
Marginal Released Outcomes				(0.002)	(0.001)	(0.002)
Judges	268	268	268	268	268	268

Table 5: Hierarchical MTE Model Estimates

*Notes.* This table reports simulated minimum distance estimates of moments of the MTE model described in Section 6. See Table A11 for the underlying hyperparameter estimates. Columns 4-6 estimate the baseline model, while columns 1-3 impose conventional monotonicity. Robust standard errors, two-way clustered at the individual and judge level, are obtained by a bootstrapping procedure and appear in parentheses.

	Target		warranted	Target Observational	
	Dessline	Disparity	Posteriors	Disparity Posteriors	
	Dasenne	Increase	Decrease	Increase	Decrease
		Leniency	Leniency	Leniency	Leniency
Panel A: Close All Disparities	(1)	(2)	(3)	(4)	(5)
Mean Unwarranted Disparity	0.040	0.001	0.001	-0.009	-0.010
	[0.032]	[0.020]	[0.025]	[0.019]	[0.025]
Mean Observational Disparity	0.049	0.010	0.010	0.000	-0.001
	[0.032]	[0.020]	[0.025]	[0.019]	[0.025]
Racial Bias	0.065	0.039	0.023	0.032	0.012
	[0.049]	[0.045]	[0.038]	[0.045]	[0.037]
Panel B: Close Top-Quintile Disp	arities				
Mean Unwarranted Disparity		0.026	0.026	0.023	0.023
		[0.030]	[0.032]	[0.032]	[0.034]
Mean Observational Disparity		0.035	0.035	0.032	0.032
		[0.031]	[0.032]	[0.033]	[0.034]
Racial Bias		0.056	0.050	0.054	0.047
		[0.049]	[0.051]	[0.050]	[0.053]
Judges	268	268	268	268	268

Table 6: Policy Simulations

*Notes.* This table reports the results from a series of policy simulations. Column 1 reports the mean unwarranted disparity, observational disparity, and racial bias across judges and 250 simulations of the hierarchical MTE model. Average standard deviations across judges are included in brackets. Simulations are based on the estimates from columns 2 and 4 of Appendix Table A11. Column 2 of Panel A recomputes the statistics for a counterfactual in which the lower of the black or white release rate of each judge is raised to equalize unwarranted disparity posteriors, while column 3 of Panel A does the same by lowering one of the two release rates. Columns 4 and 5 of Panel A instead adjust release rates to equalize observational disparity posteriors. Panel B conducts the counterfactual exercises only on judges ranked in the top quintile of unwarranted (columns 2 and 3) or observational (columns 4 and 5) disparity posteriors. Estimates of the model hyperparameters and empirical Bayes posteriors of all judge-specific parameters are recomputed in each simulation draw via the SMD procedure outlined in the text, using moments simulated according to the estimated distribution of reduced-form estimates in Figure 2.

# A Appendix Figures and Tables



Appendix Figure A1: Placebo Mean Risk Extrapolation

*Notes.* This figure plots race-specific release rates for the 268 judges in our sample against rates of predicted pretrial misconduct among the set of released defendants. Predicted misconduct is given by the fitted values of an OLS regression of misconduct on the regressors in column 3 of Table 2, estimated in the set of released defendants. Average predicted misconduct rates in the full sample of white and black defendants are indicated with solid markers at the maximal release rate of one. All estimates adjust for the baseline controls in Table 2 and court-by-time fixed effects. The figure also plots race-specific linear, quadratic, and local linear curves of best fit, obtained from judge-level regressions that inverse-weight by the variance of the estimated predicted misconduct rate among released defendants. The local linear regression uses a Gaussian kernel with a race-specific rule-of-thumb bandwidth. 95 percent confidence intervals for the local linear extrapolations' intercept estimates at one, obtained from robust standard errors two-way clustered at the individual and judge level, are indicated with brackets.



Appendix Figure A2: Sensitivity Analysis

*Notes.* This figure shows how our estimate of system-wide discrimination changes under different estimates of white and black mean risk. The mean risk estimates obtained from the linear, quadratic, and local linear extrapolations in Figure 2 are indicated by solid, dashed, and dotted lines.



Appendix Figure A3: Unwarranted Release Rate Disparities, Model-Based Mean Risk Estimates

*Notes.* This figure plots the distribution of observational and unwarranted release rate disparity posteriors for the 268 judges in our sample. Covariate-adjusted disparities are estimated by the coefficients of an OLS regression of an indicator for pretrial release on white×judge fixed effects, controling for judge main effects, the baseline controls from Table 2, and court-by-time fixed effects. Unwarranted disparities are estimated as described in Section 5, using the hierarchical MTE model estimates of mean risk for each race. Empirical Bayes posteriors are computed using a standard shrinkage procedure, as described in Appendix B.4. Means and standard deviations refer to the estimated prior distribution. The fractions of positive disparities are computed as posterior average effects, as described in Appendix B.4.



Appendix Figure A4: Hierarchical MTE Model Fit

*Notes.* This figure plots race-specific release rates for the 268 judges in our sample against rates of pretrial misconduct among the set of released defendants. All estimates adjust for the baseline controls in Table 2 and court-by-time fixed effects. The figure also plots race-specific curves of best fit implied by our baseline hierarchical MTE model hyperparameter estimates.



Appendix Figure A5: Predictiveness of Observational Release Rate Disparities

Notes. This figure plots unwarranted white-black release rate disparity posteriors against the corresponding covariate adjusted release rate disparity posteriors for the 268 judges in our sample. Observational disparities are estimated by the coefficients of an OLS regression of an indicator for pretrial release on white×judge fixed effects, controling for judge main effects, the baseline controls from Table 2, and court-by-time fixed effects. Unwarranted disparities are estimated as described in Section 5, using the local linear extrapolation from Figure 2 to estimate the mean risk of each race. Empirical Bayes posteriors are computed using a standard shrinkage procedure, as described in Appendix B.4. The slope of the solid line indicates the forecast coefficient.

	All	White	Black
	Defendants	Defendants	Defendants
	(1)	(2)	(3)
Dropped from Sample	0.00015	0.00010	0.00020
	(0.00014)	(0.00014)	(0.00017)
Baseline Controls	Yes	Yes	Yes
Court x Time FE	Yes	Yes	Yes
Mean Sample Attrition	0.416	0.409	0.424
Cases	$1,\!425,\!652$	$726,\!284$	697, 597

Appendix Table A1: Judge Leniency and Sample Attrition

*Notes.* This table reports OLS estimates of regressions of judge leniency on an indicator for leaving the sample due to case adjournment or case disposal, the baseline controls in Table 2, and court-by-time fixed effects. The regressions are estimated on the sample described in Table 1. Judge leniency is estimated using data from other cases assigned to a given bail judge, following the procedure described in Section 4.1. Robust standard errors, two-way clustered at the individual and the judge level, are reported in parentheses.

	All Defendants		White I	Defendants	Black I	Black Defendants	
	Full	Estimation	Full	Estimation	Full	Estimation	
	Sample	Sample	Sample	Sample	Sample	Sample	
Panel A: Pretrial Release	(1)	(2)	(3)	(4)	(5)	(6)	
Released Before Trial	0.856	0.730	0.879	0.767	0.832	0.695	
Share ROR	0.603	0.852	0.620	0.852	0.586	0.851	
Share Disposed	0.295	0.000	0.266	0.000	0.327	0.000	
Share Adjourned	0.192	0.000	0.201	0.000	0.183	0.000	
Share Money Bail	0.068	0.144	0.069	0.144	0.066	0.145	
Share Other Bail Type	0.329	0.004	0.311	0.004	0.348	0.004	
Share Remanded	0.000	0.000	0.000	0.000	0.000	0.000	
Panel B: Defendant Characte	eristics						
Black	0.495	0.522	0.000	0.000	1.000	1.000	
Male	0.820	0.821	0.826	0.839	0.813	0.804	
Age at Arrest	31.871	31.969	31.667	32.055	32.080	31.890	
Prior Rearrest	0.189	0.229	0.164	0.204	0.214	0.253	
Prior FTA	0.083	0.103	0.068	0.087	0.099	0.117	
Panel C: Charge Characteris	tics						
Number of Charges	1.100	1.150	1.122	1.184	1.078	1.118	
Felony Charge	0.183	0.362	0.177	0.355	0.188	0.368	
Misdemeanor Charge	0.817	0.638	0.823	0.645	0.812	0.632	
Any Drug Charge	0.340	0.256	0.327	0.257	0.352	0.256	
Any DUI Charge	0.033	0.046	0.048	0.067	0.017	0.027	
Any Violent Charge	0.071	0.143	0.062	0.124	0.081	0.160	
Any Property Charge	0.217	0.136	0.209	0.127	0.226	0.144	
Cases	$1,\!417,\!434$	595,186	715,867	284,598	701,567	$310,\!588$	

Appendix Table A2: Descriptive Statistics by Sample

*Notes.* This table summarizes the difference between the NYC analysis sample and the full sample of NYC arraignments. The full sample consists of all bail hearings between November 1, 2008 and November 1, 2013. The analysis sample consists of bail hearings that were quasi-randomly assigned to judges between November 1, 2008 and November 1, 2008 and November 1, 2013, as described in the text. Information on demographics and criminal outcomes is derived from court records as described in the text. Pretrial release is defined as meeting the bail conditions set by the first assigned bail judge. ROR (released on Recognizance) is defined as being released without any conditions. FTA (failure to appear) is defined as failing to appear at a mandated court date.

	All Defendente	White Defendents	Black
	Derendants	Derendants	Defendants
White	(1)	(2)	(3)
white	(0.00013)		
Mala	(0.00009)	0.00009	0.00004
male	(0.00003)	(0.00003)	(0.00004)
A	(0.00014)	(0.00019)	(0.00018)
Age at Arrest	-0.00011	-0.00015	-0.00008
	(0.00004)	(0.00006)	(0.00005)
Prior Rearrest	-0.00021	0.00007	-0.00044
	(0.00011)	(0.00018)	(0.00015)
Prior F'IA	0.00016	-0.00014	0.00039
	(0.00016)	(0.00024)	(0.00023)
Number of Charges	-0.00001	-0.00001	-0.00001
	(0.00001)	(0.00001)	(0.00003)
Felony Charge	0.00025	0.00011	0.00039
	(0.00020)	(0.00023)	(0.00025)
Any Drug Charge	-0.00022	-0.00017	-0.00027
	(0.00016)	(0.00021)	(0.00018)
Any DUI Charge	0.00045	0.00051	0.00008
	(0.00027)	(0.00032)	(0.00045)
Any Violent Charge	-0.00008	-0.00023	0.00001
	(0.00023)	(0.00033)	(0.00025)
Any Property Charge	-0.00033	-0.00028	-0.00036
	(0.00018)	(0.00019)	(0.00027)
oint p-value	[0.10521]	[0.30945]	[0.07931]
Court x Time FE	Yes	Yes	Yes
Cases	$595,\!186$	284,598	310,588

Appendix Table A3: Tests of Quasi-Random Judge Assignment

*Notes.* This table reports OLS estimates of regressions of judge leniency on defendant characteristics. The regressions are estimated on the sample described in Table 1. Judge leniency is estimated using data from other cases assigned to a given bail judge, following the procedure described in Section 4.1. All regressions control for court-by-time fixed effects. The p-values reported at the bottom of each column are from F-tests of the joint significance of the variables listed in the rows. Robust standard errors, two-way clustered at the individual and the judge level, are reported in parentheses.

	All Defendants	White Defendants	Black Defendants
	(1)	(2)	(3)
Judge Leniency	0.953	0.774	1.112
	(0.024)	(0.029)	(0.031)
Baseline Controls	Yes	Yes	Yes
Court x Time FE	Yes	Yes	Yes
Mean Release Rate	0.730	0.767	0.695
R2	0.178	0.172	0.187
Cases	$595,\!186$	$284,\!598$	$310,\!588$

Appendix Table A4: First Stage Effects of Judge Leniency

*Notes.* This table reports OLS estimates of regressions of an indicator for pretrial release on judge leniency. The regressions are estimated on the sample described in Table 1. Judge leniency is estimated using data from other cases assigned to a bail judge, following the procedure described in Section 4.1. All regressions control for the baseline controls in Table 2 and court-by-time fixed effects. Robust standard errors, two-way clustered at the individual and the judge level, are reported in parentheses.

		Scaling Factor	Number of Defendants	Number Released	Release Rate	Release Disparity
Panel A: Observationa	al Estimates	(1)	(2)	(3)	(4)	(5)
White Defendents	$Y_i^* = 0$	1	75	75	0.75	
white Defendants	$Y_{i}^{*} = 1$	1	25	0	0.75	0 50
Black Defendants	$Y_{i}^{*} = 0$	1	25	25	0.25	0.50
	$Y_i^* = 1$	1	75	0		
Panel B: Rescaled Esta	imates					
White Defendents	$Y_{i}^{*} = 0$	2/3	50	50	0 50	
white Defendants	$Y_{i}^{*} = 1$	2	50	0	0.50	0.00
Plaak Defendante	$Y_i^* = 0$	2	50	50	0 50	0.00
Black Defendants	$Y_{i}^{*} = 1$	2/3	50	0	0.50	

Appendix Table A5: Simple Numerical Example of Unwarranted Disparity Estimation

Notes: This table uses a simple numerical example to illustrate how unwarranted disparities can be measured with observational release rate comparisons that are rescaled using average white and black misconduct risk. We assume there is one race-neutral judge who can perfectly predict potential misconduct  $Y_i^*$ , 100 black defendants where 75 have  $Y_i^* = 1$ , and 100 white defendants where 25 have  $Y_i^* = 1$ . Panel A shows that the perfectly predictive judge has a white release rate of 0.75 but a black release rate of 0.25, meaning that an observational comparison would find that white defendants have a 50 percentage point higher release rate than black defendants despite the judge being race-neutral. Panel B shows that the true unwarranted disparity of zero can be measured by rescaling this observational release rate comparison with the scaling factor described in the text. Column 1 of Panel B shows the scaling factor ( $\omega_i$ ) in this example, and column 5 shows the resulting unwarranted disparity estimate.

		Scaling Factor	Number of Defendants	Number Released	Release Rate	Release Disparity
Panel A: Observation	al Estimates	(1)	(2)	(3)	(4)	(5)
White Defendents	$Y_{i}^{*} = 0$	1.000	189,551	$157,\!550$	0.756	
white Defendants	$Y_{i}^{*} = 1$	1.000	102,963	$63,\!543$	0.750	0.052
Black Defendants	$Y_i^* = 0$	1.000	184,403	150,232	0 709	0.055
	$Y_{i}^{*} = 1$	1.000	135,756	$74,\!633$	0.702	
Panel B: Rescaled Est	imates					
White Defendants	$Y_i^* = 0$	0.942	$178,\!540$	$148,\!398$	0.748	
white Defendants	$Y_{i}^{*} = 1$	1.107	$113,\!974$	70,338	0.140	0.036
Plack Defendants	$Y_{i}^{*} = 0$	1.060	$195,\!413$	159,202	0 711	0.050
Black Defendants	$Y_{i}^{*} = 1$	0.919	124,746	$68,\!580$	0.711	

Appendix Table A6: Unwarranted Disparity Estimation for NYC Release Decisions

Notes: This table calculates system-wide unwarranted disparity in NYC by rescaling observational release rate comparisons using estimates of average white and black misconduct risk. In Panel A we use the local linear estimates of mean risk in Table 3 to estimate the number of defendants with and without misconduct potential (column 2) as well as the number of such defendants that are released (column 3). These estimates imply that an observational comparison would find that white defendants have a 5.3 percentage point higher release rate than black defendants. In Panel B we use the same mean risk estimates to rescale this observational release rate comparison with the scaling factor described in the text. Column 1 of Panel B shows the scaling factor ( $\omega_i$ ) given by these estimates, and column 5 shows the resulting unwarranted disparity estimate.

	Any	Case	Any	Violent
	Misconduct	FTA	Rearrest	Rearrest
Panel A: Mean Risk by Race	(1)	(2)	(3)	(4)
White Defendants	0.352	0.181	0.247	0.009
	(0.014)	(0.013)	(0.017)	(0.004)
Black Defendants	0.424	0.231	0.307	0.012
	(0.016)	(0.012)	(0.017)	(0.005)
Panel B: System-Wide Discrin	nination			
Mean Across Cases	0.036	0.042	0.041	0.055
	(0.005)	(0.004)	(0.004)	(1.351)
Panel C: Judge-Level Discrime	ination			
Mean Across Judges	0.034	0.041	0.040	0.054
	(0.005)	(0.004)	(0.004)	(1.202)
Std. Dev. Across Judges	0.031	0.033	0.032	0.038
	(0.003)	(0.003)	(0.003)	(1.012)
Fraction Positive	0.875	0.903	0.902	0.935
	(0.035)	(0.026)	(0.027)	(0.089)
Judges	268	268	268	268

Appendix Table A7: Robustness to Pretrial Misconduct Outcome

*Notes.* This table summarizes estimates of mean risk and unwarranted racial disparities for different outcome variables. Panel A reports estimates of race-specific average misconduct risk, Panel B reports estimates of system-wide (case-weighted) unwarranted disparity, and Panel C reports empirical Bayes estimates of summary statistics for the judge-level unwarranted disparity prior distribution. Column 1 adjusts for differences by race in the mean risk of any misconduct (either rearrest or FTA). Column 2 adjusts for differences by race in the mean risk of FTA. Column 3 adjusts for differences by race in the mean risk of rearrest for a violent crime. Robust standard errors, two-way clustered at the individual and judge level, are obtained by a bootstrapping procedure and appear in parentheses.

	Linear	Quadratic	Local Linear
	Extrapolation	Extrapolation	Extrapolation
Panel A: Mean Risk by Race	(1)	(2)	(3)
White Defendants	0.359	0.351	0.354
	(0.007)	(0.024)	(0.030)
Black Defendants	0.401	0.434	0.430
	(0.006)	(0.023)	(0.037)
Panel B: System-Wide Discrir	nination		
Mean Across Cases	0.035	0.025	0.026
	(0.002)	(0.007)	(0.011)
Panel C: Judge-Level Discrime	ination		
Mean Across Judges	0.033	0.023	0.025
	(0.002)	(0.007)	(0.011)
Std. Dev. Across Judges	0.034	0.034	0.034
	(0.003)	(0.003)	(0.003)
Fraction Positive	0.839	0.756	0.770
	(0.019)	(0.055)	(0.084)
Judges	268	268	268

Appendix Table A8: Robustness to Judge Decision Variable

*Notes.* This table summarizes estimates of mean risk and unwarranted racial disparities from different extrapolations of the variation in Figure 2. The judge's decision variable in this table is release on recognizance (ROR) versus the assignment of any monetary bail, where there is a 4.1 percentage point release rate disparity after adjusting for covariates. Panel A reports estimates of race-specific average misconduct risk, Panel B reports estimates of system-wide (case-weighted) unwarranted disparity, and Panel C reports empirical Bayes estimates of summary statistics for the judge-level unwarranted disparity prior distribution. To estimate mean risk, column 1 uses a linear extrapolation of the variation in Figure 2, while column 2 uses a quadratic extrapolation and column 3 uses a local linear extrapolation with a Gaussian kernel and a rule-of-thumb bandwidth. Robust standard errors, two-way clustered at the individual and judge level, are obtained by a bootstrapping procedure and appear in parentheses.

	Linear	Quadratic	Local Linear
	Extrapolation	Extrapolation	Extrapolation
Panel A: Mean Risk by Race	(1)	(2)	(3)
White Defendants	0.283	0.206	0.273
	(0.010)	(0.028)	(0.018)
Black or Hispanic Defendants	0.386	0.401	0.401
	(0.005)	(0.018)	(0.012)
Panel B: System-Wide Discrimin	ation		
Mean Across Cases	0.058	0.108	0.059
	(0.003)	(0.027)	(0.008)
Panel C: Judge-Level Discrimina	tion		
Mean Across Judges	0.058	0.108	0.059
	(0.004)	(0.025)	(0.008)
Std. Dev. Across Judges	0.021	0.000	0.018
	(0.004)	(0.015)	(0.005)
Fraction Positive	0.997	1.000	1.000
	(0.022)	(0.008)	(0.025)
Judges	250	250	250

Appendix Table A9: Robustness to Definition of Defendant Race

*Notes.* This table summarizes estimates of mean risk and unwarranted racial disparities from different extrapolations of the variation in Figure 2. The racial comparison in this table is between black or Hispanic defendants to non-Hispanic white defendants, where there is a 7.3 percentage point release rate disparity after adjusting for covariates. Panel A reports estimates of race-specific average misconduct risk, Panel B reports estimates of system-wide (case-weighted) unwarranted disparity, and Panel C reports empirical Bayes estimates of summary statistics for the judge-level unwarranted disparity prior distribution. To estimate mean risk, column 1 uses a linear extrapolation of the variation in Figure 2, while column 2 uses a quadratic extrapolation and column 3 uses a local linear extrapolation with a Gaussian kernel and a rule-of-thumb bandwidth. Robust standard errors, two-way clustered at the individual and judge level, are obtained by a bootstrapping procedure and appear in parentheses.

	Number of Spline Knots						
	5	10	15	20			
Panel A: White Defendants	(1)	(2)	(3)	(4)			
Test Statistic	303.5	293.9	279.8	272.1			
Deg. of Freedom	260	255	250	245			
p-value	0.033	0.047	0.094	0.113			
Cases	$284,\!598$	$284,\!598$	$284,\!598$	284,598			
Panel B: Black Defendants							
Test Statistic	392.7	389.0	379.7	348.9			
Deg. of Freedom	260	255	250	245			
p-value	[<0.001]	[<0.001]	[<0.001]	[< 0.001]			
Cases	310,588	$310,\!588$	310,588	310,588			

Appendix Table A10: Tests of Conventional Monotonicity

*Notes.* This table reports the results of the tests of conventional MTE monotonicity proposed by Frandsen et al. (2019), computed separately by defendant race. Test statistics are based on quadratic b-spline estimates of the relationship between misconduct outcomes and judge leniency, with the number of knots specified in each column, controlling for court-by-time fixed effects. The regressions are estimated on the sample described in Table 1.

Appendix Table A11: Hierarchical MTE Model Hyperparameter Estimates

	White Defendants			Bla	Black Defendants		
	(1)	(2)	(3)	(4)	(5)	(6)	
Mean Misconduct Risk $(\mu)$	0.374	0.400	0.380	0.429	0.429	0.447	
	(0.014)	(0.007)	(0.018)	(0.012)	(0.007)	(0.014)	
Mean ln(Signal Quality) ( $\alpha$ )	0.624	0.163	0.271	0.278	-0.198	-0.249	
	(0.039)	(0.071)	(0.102)	(0.160)	(0.079)	(0.140)	
Mean Release Threshold $(\gamma)$	0.761	1.101	1.211	0.731	1.164	1.072	
	(0.035)	(0.029)	(0.123)	(0.028)	(0.045)	(0.059)	
Release Threshold Std. Dev. $(\delta)$	0.266	0.115	0.159	0.253	0.214	0.168	
	(0.031)	(0.012)	(0.052)	(0.028)	(0.024)	(0.033)	
$\ln(\text{Signal Quality})$ Std. Dev. $(\psi)$		0.115	0.112		0.141	0.146	
		(0.012)	(0.011)		(0.014)	(0.013)	
Regression of ln(Signal Quality)			-0.273			0.227	
on Release Threshold $(\beta)$			(0.173)			(0.249)	
Judges	268	268	268	268	268	268	

Notes. This table reports simulated minimum distance estimates of the MTE model described in the text. 500 simulation draws are used. Columns 3 and 6 estimate the full model with all hyperparameters. Columns 2 and 5 restrict  $\beta = 0$  (omitting the quadratic regression coefficient moment), while columns 1 and 4 also restrict  $\psi = 0$  (omitting the residual variance moment). The baseline model used in the text and summarized in Table 5 comes from columns 2 and 5 of this table. Robust standard errors, two-way clustered at the individual and the judge level, are reported in parentheses.

						Split-S	Sample
		Full-Sa	ample Disp	parities		Dispa	rities
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
New Judge	-0.013				-0.011		-0.004
	(0.004)				(0.003)		(0.004)
Lenient Judge		-0.009			-0.011		-0.009
		(0.003)			(0.003)		(0.003)
Above-Median Black Share			-0.007		-0.006		-0.000
			(0.003)		(0.004)		(0.004)
Manhattan Courtroom				0.022	0.020		0.016
				(0.004)	(0.004)		(0.004)
Bronx Courtroom				-0.004	-0.007		-0.001
				(0.003)	(0.004)		(0.005)
Queens Courtroom				0.013	0.008		0.009
				(0.004)	(0.005)		(0.005)
Richmond Courtroom				0.010	0.005		0.011
				(0.004)	(0.006)		(0.006)
Lagged Disparities						0.305	0.201
						(0.050)	(0.053)
Mean Disparity	0.042	0.042	0.042	0.042	0.042	0.042	0.042
R2	0.044	0.043	0.022	0.212	0.309	0.235	0.332
Judges	268	268	268	268	268	252	252

#### Appendix Table A12: Unwarranted Disparities and Judge Characteristics, Model-Based Mean Risk

*Notes.* This table reports OLS estimates of regressions of unwarranted disparity posteriors on judge characteristics. Unwarranted disparities are estimated as described in Section 5, using the hierarchical MTE model estimate of mean risk. New judges are defined as judges appointed during our estimation period. Lenient judges are defined as judges with above-average leniency, controlling for court-by-time fixed effects. Courtroom locations are defined using the location of the modal case heard by each judge. Split-sample disparities are computed by splitting each judge's sample of cases at the median case and constructing two samples, a before-median case sample and an after-median case sample. Unwarranted disparities are then re-estimated within each subsample. The estimation procedure conditions on court-by-time effects, which causes a small number of judge effects to become collinear with the court-by-time effects and dropped. All specifications are weighted by the inverse variance of the unwarranted disparity posteriors. Robust standard errors are reported in parentheses.

	White Defendants	Black Defendants	Diff.
Panel A: MTE Estimates	(1)	(2)	(3)
Marginal Released Outcome	0.484	0.483	0.002
	(0.026)	(0.022)	(0.029)
Panel B: IV Estimates			
Marginal Released Outcome	0.385	0.396	-0.011
	(0.067)	(0.047)	(0.054)
Baseline Controls	Yes	Yes	_
Court x Time FE	Yes	Yes	_
Mean Misconduct	0.266	0.332	_
Cases	284,598	310,588	_

Appendix Table A13: MTE Estimates of Racial Bias

*Notes.* This table reports conventional MTE estimates and IV estimates of marginal released outcomes and racial bias. The IV estimate of mean marginal released outcomes instrument for pretrial release in a regression of pretrial misconduct using a leave-one-out judge leniency measure while controlling for the baseline controls in Table 2 and courtby-time fixed effects. To estimate the MTE results, we first compute judge-specific release and misconduct rates that control for the baseline controls in Table 2 and court-by-time fixed effects. We then fit a quadratic relationship between misconduct rates and release rates. The MTE estimate of mean marginal released outcomes is the average derivative of this quadratic function. Both estimation procedures require a conventional first-stage monotonicity assumption. The difference in mean marginal released outcomes across the races estimates mean racial bias. The regressions are estimated on the sample described in Table 1. Robust standard errors, two-way clustered at the individual and judge level, are obtained by a bootstrapping procedure and appear in parentheses.

	Baseline	No Racial Bias	Equal Signal Quality	Both
Panel A: Change Black Parameters	(1)	(2)	(3)	(4)
Unwarranted Disparity	0.040	-0.051	0.090	0.028
Release Rates $(W/B)$	$0.758 \ / \ 0.709$	$0.758 \ / \ 0.801$	$0.758 \ / \ 0.657$	$0.758 \ / \ 0.720$
Racial Bias	0.065	0.000	0.065	0.000
Marginal Outcomes (W/B)	$0.609 \ / \ 0.544$	$0.609 \ / \ 0.609$	$0.609 \ / \ 0.544$	$0.609 \ / \ 0.609$
Signal Quality (W/B)	1.184 / 0.828	1.184 / 0.828	1.184 / 1.184	1.184 / 1.184
Panel B: Change White Parameters				
Unwarranted Disparity		-0.020	0.124	0.043
Release Rates $(W/B)$		$0.699 \ / \ 0.709$	$0.840 \ / \ 0.709$	$0.760 \ / \ 0.709$
Racial Bias		0.000	0.065	0.000
Marginal Outcomes (W/B)		$0.544 \ / \ 0.544$	$0.609 \ / \ 0.544$	$0.544 \ / \ 0.544$
Signal Quality (W/B)		$1.184 \ / \ 0.828$	$0.828 \ / \ 0.828$	$0.828 \ / \ 0.828$
Judges	268	268	268	268

Appendix Table A14: Unwarranted Disparity Decompositions

*Notes.* Column 1 of this table reports average unwarranted disparity and racial bias across judges and 250 simulations of the hierarchical MTE model, along with average release rates, marginal released outcomes, and signal quality of black and white defendants. Simulations are based on the estimates from columns 2 and 4 of Appendix Table A11. Column 2 recomputes the statistics for a counterfactual in which black (Panel A) or white (Panel B) release rates are set to eliminate racial bias, while column 3 adjusts black (Panel A) or white (Panel B) signal quality to equalize signal quality across race. Column 4 applies both counterfactuals simultaneously.

Appendix 18	IDIC 1115.	Itaciai D	las and J	uuge ona	1400011501	60	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
New Judge	0.051				-0.010		-0.002
	(0.004)				(0.004)		(0.004)
Lenient Judge		0.066			0.016		0.025
		(0.003)			(0.003)		(0.002)
Above-Median Black Share			-0.060		-0.038		-0.028
			(0.003)		(0.005)		(0.003)
Manhattan Courtroom				0.025	0.008		-0.015
				(0.004)	(0.004)		(0.003)
Bronx Courtroom				-0.007	-0.039		-0.029
				(0.004)	(0.006)		(0.004)
Queens Courtroom				-0.052	-0.048		-0.052
-				(0.003)	(0.003)		(0.002)
Richmond Courtroom				0.018	-0.014		-0.029
				(0.007)	(0.005)		(0.008)
Unwarranted Disparities						1.001	1.016
						(0.160)	(0.088)
Mean Bias	0.065	0.065	0.065	0.065	0.065	0.065	0.065
R2	0.087	0.384	0.425	0.725	0.832	0.119	0.907
Judges	268	268	268	268	268	268	268

Appendix Table A15: Racial Bias and Judge Characteristics

*Notes.* This table reports OLS estimates of regressions of racial bias posteriors on judge characteristics. Posteriors are obtained from the heirarchical MTE model as described in Section 6. New judges are defined as judges appointed during our estimation period. Lenient judges are defined as judges with above-average leniency, controlling for courtby-time fixed effects. Courtroom locations are defined using the location of the modal case heard by each judge. All specifications are weighted by the inverse variance of the racial bias posteriors. Robust standard errors are reported in parentheses.

прреник тавие инс	. Digital	Quanty 1	merence	s and Juc	ige Onara		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
New Judge	-0.038				-0.033		-0.001
	(0.016)				(0.013)		(0.007)
Lenient Judge		0.014			0.011		0.045
		(0.010)			(0.009)		(0.006)
Above-Median Black Share			-0.017		-0.025		-0.007
			(0.010)		(0.014)		(0.008)
Manhattan Courtroom				0.085	0.071		0.002
				(0.013)	(0.013)		(0.009)
Bronx Courtroom				-0.016	-0.035		-0.018
				(0.012)	(0.016)		(0.009)
Queens Courtroom				0.043	0.020		-0.010
				(0.016)	(0.020)		(0.013)
Richmond Courtroom				0.064	0.044		-0.027
				(0.027)	(0.019)		(0.011)
Unwarranted Disparities						3.203	3.324
						(0.150)	(0.160)
Mean Difference	0.379	0.379	0.379	0.379	0.379	0.379	0.379
R2	0.028	0.008	0.010	0.240	0.276	0.708	0.790
Judges	268	268	268	268	268	268	268

Appendix Table A16: Signal Quality Differences and Judge Characteristics

*Notes.* This table reports OLS estimates of regressions of differences in signal quality on judge characteristics. Posteriors are obtained from the heirarchical MTE model as described in Section 6. New judges are defined as judges appointed during our estimation period. Lenient judges are defined as judges with above-average leniency, controlling for court-by-time fixed effects. Courtroom locations are defined using the location of the modal case heard by each judge. All specifications are weighted by the inverse variance of the signal quality difference posteriors. Robust standard errors are reported in parentheses.

## **B** Econometric Appendix

### **B.1** Defining and Measuring Discrimination with Multi-Valued $Y_i^*$

This appendix first generalizes our definition of racial discrimination and derivation of OVB in observational comparisons to settings where the decision-maker's objective is non-binary. We then discuss how our quasi-experimental framework for measuring racial discrimination extends to this case.

Our initial definition of racial discrimination,  $\Delta_j = E[E[D_{ij} | Y_i^*, R_i = w] - E[D_{ij} | Y_i^*, R_i = b]]$ , remains sensible in the case of non-binary  $Y_i^*$ , provided the support of  $Y_i^*$  is the same in the white  $(R_i = w)$  and black  $(R_i = b)$  subpopulations. Natural generalizations of Equation (2) are given by:

$$\Delta_j = \sum_{y \in Supp(Y_i^*)} \left( \delta_{jw}^y - \delta_{jb}^y \right) p_y \tag{B1}$$

in the multi-valued  $Y_i^*$  case, where  $p_y = Pr(Y_i^* = y)$ , and:

$$\Delta_j = \int_{Supp(Y_i^*)} \left(\delta_{jw}^y - \delta_{jb}^y\right) dF(y) \tag{B2}$$

in the case of continuous  $Y_i^*$ , where  $F(\cdot)$  is the cumulative distribution function of  $Y_i^*$ . In both cases,  $\delta_{ir}^y = E[D_{ij} \mid Y_i^* = y, R_i = r]$  gives conditional release rates for each race r and each  $y \in Supp(Y_i^*)$ .

As in Section 3.3, the bias of observational benchmarking regressions relative to these parameters, when judges are as-good-as-randomly assigned, is given by:

$$\xi_{j} = \sum_{y \in Supp(Y_{i}^{*})} \delta_{jw}^{y} p_{yw} - \sum_{y \in Supp(Y_{i}^{*})} \delta_{jb}^{y} p_{yb} - \sum_{y \in Supp(Y_{i}^{*})} \left( \delta_{jw}^{y} - \delta_{jb}^{y} \right) (p_{yw} p_{w} + p_{yb} p_{b})$$
$$= \sum_{y \in Supp(Y_{i}^{*})} \left( \delta_{jw}^{y} p_{b} + \delta_{jb}^{y} p_{w} \right) (p_{yw} - p_{yb})$$
(B3)

in the multi-valued  $Y_i^*$  case, where  $p_{yr} = Pr(Y_i^* = y \mid R_i = r)$  and again  $p_r = Pr(R_i = r)$ , and:

$$\xi_{j} = \int_{Supp(Y_{i}^{*})} \delta_{jw}^{y} dF_{w}(y) - \int_{Supp(Y_{i}^{*})} \delta_{jb}^{y} dF_{b}(y) - \int_{Supp(Y_{i}^{*})} \left( \delta_{jw}^{y} - \delta_{jb}^{y} \right) d(F_{w}(y)p_{w} + F_{b}(y)p_{b})$$

$$= \int_{Supp(Y_{i}^{*})} \left( \delta_{jw}^{y}p_{b} + \delta_{jb}^{y}p_{w} \right) d(F_{w}(y) - F_{b}(y))$$
(B4)

in the case of continuous  $Y_i^*$ , where  $F_r(\cdot)$  is the cumulative distribution function of  $Y_i^*$  given  $R_i = r$ .

As in Section 5, discrimination is identified by the distribution of misconduct outcomes  $Y_i^*$  within each race when judges are quasi-randomly assigned. By Bayes' law:

$$\delta_{jr}^{y} = Pr(Y_{i}^{*} = y \mid D_{ij} = 1, R_{i} = r) \frac{E[D_{ij} \mid R_{i} = r]}{Pr(Y_{i}^{*} = y \mid R_{i} = r)}$$
(B5)

for multi-valued  $Y_i^*$  and similarly for continuous  $Y_i^*$ . The first two terms,  $Pr(Y_i^* = y \mid D_{ij} = 1, R_i = r)$ and  $E[D_{ij} \mid R_i = r]$ , are identified by  $Pr(Y_i = y \mid D_i = 1, Z_{ij} = 1, R_i = r)$  and  $E[D_i \mid Z_{ij} = 1, R_i = r]$ under quasi-random judge assignment as before. In the continuous  $Y_i^*$  case, the first term is given by the conditional density of  $Y_i^*$  given  $D_i = 1, Z_{ij} = 1$ , and  $R_i = r$ . Estimates of the race-specific misconduct distribution corresponding to the third  $Pr(Y_i^* = y \mid R_i = r)$  term (which might be obtained from similar extrapolations of quasi-experimental data as in the binary  $Y_i^*$  case) thus yield a plug-in estimator of each  $\delta_{ir}^y$ , which can be combined to estimate  $\Delta_j$  according to the initial definitions.

#### B.2 Discrimination and Bias with Normally Distributed Signal Noise

This appendix derives the decision-making model discussed in Section 2. A judge observes noisy risk signals  $\nu_i = Y_i^* + \eta_i$  with normally distributed noise:  $\eta_i \mid Y_i^*, R_i \sim N(0, 1/\tau_{R_i}^2)$ . The judge has potentially incorrect beliefs  $\tilde{\mu}_r$  on race-specific average misconduct risk  $\mu_r = E[Y_i^* \mid R_i = r]$  and knows the potentially race-specific quality of risk signals  $\tau_r$ .

The judge's subjective posterior of misconduct risk, given a signal of  $\nu_i = v$  for a defendant of race  $R_i = r$ , is derived from Bayes' rule:

$$p(\nu; r) = \frac{\widetilde{Pr}(\nu_i = v \mid Y_i^* = 1, R_i = r) \widetilde{Pr}(Y_i^* = 1, R_i = r)}{\widetilde{Pr}(\nu_i = v, R_i = r)}$$
$$= \frac{\phi(\tau_r(v-1))\tau_r \tilde{\mu}_r}{\phi(\tau_r(v-1))\tau_r \tilde{\mu}_r + \phi(\tau_r v)\tau_r(1-\tilde{\mu}_r)}$$
(B6)

where  $\widetilde{Pr}(\cdot)$  denotes subjective probabilities and  $\phi(x) \propto \exp(-x^2/2)$  is the standard normal density. Simplifying, we have:

$$p(\nu; r) = \left(1 + \exp(\tau_r^2 (1 - 2\nu)/2) \frac{1 - \tilde{\mu}_r}{\tilde{\mu}_r}\right)^{-1}$$
(B7)

This specifies a risk-neutral judge's release rule,  $D_i = \mathbf{1}[\pi_{R_i} \ge p(\nu_i; R_i)]$ .

Equation (B7) shows that risk posteriors are strictly increasing in v, such that they can be inverted to write the judge's release decision as a cutoff rule for her observed signals  $\nu_i$ :

$$D_{i} = \mathbf{1} \left[ \frac{1}{2} - \ln \left( \frac{\tilde{\mu}_{R_{i}} (1 - \pi_{R_{i}})}{\pi_{R_{i}} (1 - \tilde{\mu}_{R_{i}})} \right) / \tau_{R_{i}}^{2} \ge \nu_{i} \right]$$
(B8)

Equation (B8) shows that variation in risk beliefs  $\tilde{\mu}_r$  and risk tolerances  $\pi_r$  are observationally equivalent in this model, in the sense that as one of these parameters varies in (0, 1) the other can be set to keep the index  $I_r = \frac{\tilde{\mu}_r(1-\pi_r)}{\pi_r(1-\tilde{\mu}_r)}$ , and thus the decision rule, constant.

A consequence of Equation (B8) is that the average misconduct rate of white and black defendants at the margin of release,  $E[Y_i^* | p(\nu_i; R_i) = \pi_r, R_i = r]$ , is a function of the judges risk tolerance  $\pi_r$ and prior risk belief  $\tilde{\mu}_r$ . By Equation (4), the marginal outcomes under correct beliefs  $\mu_r$  equals the judge's risk tolerance. More generally:

$$E[Y_i^* \mid p(\nu_i; R_i) = \pi_r, R_i = r] = \left(1 + I_r \left(\frac{1 - \mu_r}{\mu_r}\right)\right)^{-1}$$
(B9)

by the observational equivalence of Equation (B8). Racial bias is found when this expression varies by race r, which could be due to racial animus ( $\pi_w \neq \pi_b$ ) or inaccurate beliefs ( $\tilde{\mu}_r \neq \mu_r$ ).

To characterize discrimination in this model, note that Equation (B8) and the conditional normal-

ity of  $\nu_i$  implies that the judge's true and false negative rates can be written, respectively:

$$\delta_r^T = Pr(D_i = 1 \mid Y_i^* = 0, R_i = r) = \Phi\left(\frac{1}{2}\tau_r - \frac{1}{\tau_r}\ln I_r\right)$$
(B10)

$$\delta_r^F = Pr(D_i = 1 \mid Y_i^* = 1, R_i = r) = 1 - \Phi\left(\frac{1}{2}\tau_r + \frac{1}{\tau_r}\ln I_r\right)$$
(B11)

and the extent of racial discrimination is given by the extent to which  $\Delta = (\delta_w^T - \delta_b^T)(1 - \bar{\mu}) + (\delta_w^F - \delta_b^F)\bar{\mu}$ varies by race, for  $\bar{\mu} = E[Y_i^*]$ . With common signal quality,  $\tau_w = \tau_b$ , a lack of racial discrimination requires  $I_w = I_b$ . By comparison with Equation (B9), this scenario will generally lead to racial bias unless white and black average misconduct risk are also equal ( $\mu_w = \mu_b$ ). More generally, the fact that  $\Delta$  is strictly decreasing (to zero) in the white index  $I_w$  and strictly increasing (to one) in the black index  $I_b$  implies that there exist a set of thresholds ( $I_w, I_b$ ) resulting in no racial discrimination on average, even when signal quality differs. Again, this will typically yield racial bias, per Equation (B9), to the extent either mean risk or signal quality differs by race.

## **B.3** Bail Release and Classification Error

This appendix shows how a judge minimizing the cost of type-I and type-II error in the bail setting implicitly uses a posterior risk threshold-crossing rule, as in Section 2. Suppose the cost of a type-I "false positive" decision (detaining an individual with no pretrial misconduct risk) is given by  $c^I > 0$ and the cost of a type-II "false negative" decision (releasing an individual with pretrial misconduct risk) is given by  $c^{II} > 0$ . A judge's utility given a release decision  $D_i \in \{0, 1\}$  is then:

$$U_i = -c^{II} D_i Y_i^* - c^I (1 - D_i)(1 - Y_i^*)$$
(B12)

Let D(v) be a decision rule mapping risk signals  $\nu_i$  to binary release decisions  $D_i$ . Suppose D(v) is set to maximize the judge's expected utility (or minimize her expected disutility):

$$D(v) = \arg\min_{d(v)} c^{II} d(v) p(v) + c^{I} (1 - d(v)) (1 - p(v))$$
(B13)

where  $p(\nu)$  denotes the judge's subjective expectation of pretrial misconduct given a signal of  $\nu_i = \nu$ . It is clear that this solution is a cutoff rule:

$$D(v) = \mathbf{1}[\pi \ge p(\nu_{ij})] \tag{B14}$$

where  $\pi = \frac{c^{II}}{c^{I}+c^{II}} \in (0,1)$  gives the judge's relative cost of type-II error. Per Equation (4), this also shows that when judge beliefs are accurate, the expected outcome of a marginally released defendant identifies this relative cost parameter.

## **B.4** Conventional Empirical Bayes Methods

This appendix summarizes the two conventional empirical Bayes approaches used in this paper: the posterior mean calculation of Morris (1983) and the posterior average effect calculation of Bonhomme and Weidner (2020). We use the former to plot the distribution of disparity posteriors in Figures 1, 3, and A3, and also to compute the prior means and standard deviations in these exhibits. We use the

latter to compute the fraction of judges with positive disparities in these figures, and also to interpret the coefficient estimates in Tables 4, A12, A15, and A16.

Let  $\hat{\theta}_j$  be an estimate of an unknown judge-specific parameter  $\theta_j$ , such as an observational benchmarking coefficient or our rescaled unwarranted disparity measure. Applying to the usual asymptotic approximation, we write  $\hat{\theta}_j = \theta_j + \varepsilon_j$  where  $\varepsilon_j \sim N(0, \Sigma_j)$  for known  $\Sigma_j$ . Conventional empirical Bayes methods further assume  $\theta_j \sim N(\mu, \Omega)$ , where  $\mu$  and  $\Omega$  are unknown hyperparameters. Given this prior distribution, the posterior mean of  $\theta_j$  after observing the estimate  $\hat{\theta}_j$  is given by:

$$E[\theta_j \mid \hat{\theta}_j] = \frac{\Sigma_j}{\Omega + \Sigma_j} \mu + \frac{\Omega}{\Omega + \Sigma_j} \hat{\theta}_j$$
(B15)

More generally, Equation (B15) gives the minimum mean-squared error prediction of  $\theta_j$  given  $\hat{\theta}_j$  when the normality of  $\theta_j$  is relaxed, provided  $\mu$  and  $\Omega$  continue to parameterize the mean and variance of the prior distribution.

Empirical Bayes posteriors estimate  $\mu$  and  $\Omega$  and plug these hyperparameter estimates into Equation (B15). We estimate  $\mu$  and  $\Omega$  by the weighted iterative procedure studied by (Morris, 1983), which is equivalent to a maximum likelihood procedure. At iteration k the hyperparameter estimates are:

$$\hat{\mu}_k = \sum_j \frac{\omega_{jk}}{\sum_{j'} \omega_{j'k}} \hat{\theta}_j \tag{B16}$$

$$\hat{\Omega}_k = \sum_j \frac{\omega_{jk}}{\sum_{j'} \omega_{j'k}} \left( (\hat{\theta}_j - \hat{\mu}_k)^2 - \Sigma_j \right)$$
(B17)

with inverse-variance weights that are proportional to  $\omega_{jk} = (\hat{\Omega}_{k-1} + \Sigma_j)^{-1}$  and where  $\omega_{j0} = 1$ . We iterate this procedure to convergence.

Bonhomme and Weidner (2020) discuss posterior average effect estimators of the cumulative distribution function for  $\theta_j$ , given by:

$$\hat{F}_{\theta}(t) = \frac{1}{J} \sum_{j} E[\mathbf{1}[\theta_{j} \le t] \mid \hat{\theta}_{j}]$$
(B18)

for each t in the support of  $\theta_j$ . Note that  $1 - \hat{F}_{\theta}(0)$  is a posterior average effect estimate of the fraction of  $\theta_j$  in the population that is positive. Under the normality assumption:

$$E[\mathbf{1}[\theta_j \le t] \mid \hat{\theta}_j] = \Phi\left(-\frac{E[\hat{\theta}_j \mid \hat{\theta}_j]}{\sqrt{\frac{\Omega\Sigma_j}{\Omega + \Sigma_j}}}\right)$$
(B19)

which can, as with Equation (B15), be estimated by plugging in the estimates of the mean and variance hyperparameters. Just as with the empirical Bayes posterior estimator, Bonhomme and Weidner (2020) show that this posterior average effect estimator has certain robustness properties: it is optimal in terms of local worst-case bias, and its global bias is bounded by the minimum worst-case bias within a large class of estimators. They further show how regressions of the empirical Bayes posterior means on judge characteristics also have a posterior average effect interpretation and thus the same robustness properties for estimating conditional mean functions.

## B.5 Conventional Monotonicity Violations and Judge Signal Quality

This appendix shows how differences in the way judges consider defendant and case characteristics, which lead to violations of conventional MTE monotonicity, can be viewed as differences in judge signal quality within models like the one we develop in Section 3.2. In doing so we show that such models are without observational loss, provided judge release decisions are better-than-random.

Consider a setting with a binary potential misconduct outcome  $Y_i^*$  and a set of binary judge release decisions  $D_{ij}$ . The distribution of these random variables is fully specified by the mean risk  $\mu = E[Y_i^*]$ and the true and false negative rates  $\delta_j^T = E[D_{ij} \mid Y_i^* = 0]$  and  $\delta_j^F = E[D_{ij} \mid Y_i^* = 0]$ . With mean risk fixed, any restriction on judicial decision-making – such as conventional MTE monotonicity or alternative parameterizations – can thus be understood as restricting the set of  $(\delta_i^T, \delta_j^F)$ .

We first show that when judges are making better-than-random release decisions, in the sense of  $0 < \delta_j^T < \delta_j^F < 1$  for each j, it is without observational loss to assume a decision-making model of  $D_{ij} = \mathbf{1}[\kappa_j \ge Y_i^* + \eta_i/\tau_j]$ , with  $\eta_i \mid Y_i^*$  following a known continuous distribution and  $\tau_j > 0$ . This follows since then  $\tau_j = G_\eta^{-1}(\delta_j^T) - G_\eta^{-1}(\delta_{1j}^F) > 0$  and  $\kappa_j = G_\eta^{-1}(\delta_j^T)/\tau_j$  rationalize each  $(\delta_j^T, \delta_j^F)$ , where  $G_\eta(\cdot)$  specifies the cumulative distribution of  $\eta_i \mid Y_i^*$ :

$$E[D_{ij} | Y_i^* = y] = Pr(\kappa_j \ge y + \eta_i / \tau_j)$$
  
=  $G_{\eta}((\kappa_j - y)\tau_j)$   
=  $G_{\eta}(G_{\eta}^{-1}(\delta_j^T)) + y(G_{\eta}^{-1}(\delta_j^F) - G_{\eta}^{-1}(\delta_j^T))$   
=  $\delta_j^T + y(\delta_j^F - \delta_j^T)$  (B20)

In particular, Equation (B20) shows that our risk signal threshold decision rule (23), in which  $\eta_i \mid Y_i^* \sim N(0, 1)$ , is without loss in this case. In general, we may think of  $\tau_j$  as capturing judge j's signal quality: how less likely she is to release defendants with  $Y_i^* = 1$  than those with  $Y_i^* = 0$ .

We next relate differences in such signal quality to conventional monotonicity violations in a simple behavioral model of judicial decision-making. Suppose judges observe a vector of defendant and case characteristics  $X_i^*$  which are, without loss, mean zero and positively correlated with misconduct potential:  $\mu_X(1) \equiv E[X_i^* | Y_i^* = 1] > E[X_i^* | Y_i^* = 0] \equiv \mu_X(0)$ . Judges place different weights  $\beta_j$  on the elements of this vector and also vary in their overall leniency  $\pi_j$ , such that:

$$D_{ij} = \mathbf{1}[\pi_j \ge X_i^{*\prime}\beta_j + U_i] \tag{B21}$$

where we assume  $U_i \mid X_i^*, Y_i^*$  is uniformly distributed. In this model  $E[D_{ij} \mid Y_i^* = y] = \pi_j - \mu_X(y)'\beta_j$ , assuming the parameters are such that these are all between zero and one.

Conventional monotonicity in this model requires  $Pr(D_{ij} \ge D_{ik} = 1)$  or  $Pr(D_{ik} \ge D_{ij} = 1)$  for each (j, k), which generally restricts the weights  $\beta_j$  to be the same across judges. If some elements of  $X_i^*$ were observed to the econometrician, one could relax this assumption by a conditional analysis within sets of defendants with identical observables (e.g., Mueller-Smith, 2015). Conditional monotonicity would then generally constrain the weights corresponding to unobserved characteristics to be constant.

Judicial decision-making is here better-than-random when  $\delta_j^T - \delta_j^F = (\mu_X(1) - \mu_X(0))'\beta_j > 0$  or when the weights in each  $\beta_j$  are non-negative with at least one element strictly positive. In this case

we have from the above result an equivalent representation of:

$$D_{ij} = \mathbf{1}[\kappa_j \ge Y_i^* + V_i/\tau_j] \tag{B22}$$

where  $V_i \mid Y_i^* \sim U(0,1)$ . Here judge signal quality is given by  $\tau_j = (\mu_X(1) - \mu_X(0))'\beta_j$  and has an straightforward interpretation: with only one element in  $X_i^*$ , for example, differences in  $\tau_j$  are proportional to differences in the behavioral weights  $\beta_j$ . More generally, this discussion shows how parameterizations of the distribution of signal quality across judges can be thought to structure differences in how judges weigh defendant and case characteristics when making release decisions.

#### **B.6** SMD Estimation of the Hierarchical MTE Model

We estimate the hierarchical model described in Section 6.1 and Appendix B.2 by a simulated minimum distance (SMD) procedure that targets moments of the distribution of race-specific judge release rates  $\rho_{jr} = E[D_{ij} \mid R_i = r]$  and released misconduct rates  $\lambda_{jr} = E[Y_i^* \mid D_{ij} = 1, R_i = r]$ , estimated from quasi-experimental judge assignments. This appendix formally specifies this procedure.

We first obtain estimates of  $\rho_{jr}$  and  $\lambda_{jr}$  from OLS regressions of pretrial release  $D_i$  and pretrial misconduct  $Y_i$  on judge-by-race interactions, adjusting for the quasi-experimental strata (courtroomby-time effects) and baseline controls as discussed in Section 5.2. Subject to the usual asymptotic approximation, the resulting estimates  $\hat{\rho}_{jr}$  and  $\hat{\lambda}_{jr}$  can be modeled as noisy measures of the true parameters, with a known distribution of sampling error. Specifically:

$$\hat{\rho}_{jr} = \rho_{jr} + \varepsilon_{jr}^{\rho} \tag{B23}$$

$$\hat{\lambda}_{jr} = \lambda_{jr} + \varepsilon_{jr}^{\lambda} \tag{B24}$$

where  $\varepsilon \mid \rho, \lambda \sim N(0, \Sigma)$  for a variance-covariance matrix  $\Sigma$  that is given by conventional asymptotics. Let  $\mathcal{X} = ((\hat{\rho}_{jr}, \hat{\lambda}_{jr})_{j=1,\dots,268, r \in \{w,b\}})$  collect these estimates across the 268 judges in our sample and both races w and b.

The model in Appendix B.2 specifies  $\rho_{jr}$  and  $\lambda_{jr}$  as functions of mean misconduct risk  $\mu_r$ , judge signal quality  $\tau_{jr}$ , and risk thresholds  $\pi_{jr}$ :

$$\rho_{jr} = \Phi((f(\pi_{jr}, \mu_r, \tau_{jr}) - 1)\tau_{jr}))\mu_r + \Phi(f(\pi_{jr}, \mu_r, \tau_{jr})\tau_{jr}))(1 - \mu_r)$$
(B25)

$$\lambda_{jr} = \Phi((f(\pi_{jr}, \mu_r, \tau_{jr}) - 1)\tau_{jr}))\mu_r / \rho_{jr}$$
(B26)

where  $\Phi(\cdot)$  denotes the standard normal cumulative distribution function and  $f(\cdot)$  is as defined in Section B.2. We further model signal thresholds  $\kappa_{jr} = f(\pi_{jr}, \mu_r, \tau_{jr})$  and log signal quality  $\ln \tau_{jr}$ as being joint-normally distributed across judges, with reisdual correlation across races. That is, we specify:

$$\ln \tau_{jr} = \alpha_r + \beta_r \kappa_{jr} + \epsilon_{jr} \tag{B27}$$

for each race r, with  $(\kappa_{jw}, \kappa_{jb})' \sim N(\mu_{\kappa}, \Omega_{\kappa})$  and  $(\epsilon_{jw}, \epsilon_{jb})' \mid \kappa \sim N(0, \Omega_{\tau})$ .

Equations (B23)–(B27) specify a complete distribution for the observed quasi-experimental estimates  $\mathcal{X}$  in terms of a hyperparameter vector  $\Theta = (\mu_w, \mu_b, \alpha_w, \alpha_b, \beta_w, \beta_b, \mu'_{\kappa}, vec(\Omega_{\kappa}^{1/2})', vec(\Omega_{\tilde{\tau}}^{1/2})')'$ . In practice, there is no simple closed form expression for this likelihood, complicating maximum likelihood estimation. Instead, we estimate  $\Theta$  by SMD, targeting moments of  $\mathcal{X}$  as motivated by the discussion in Section 6.1. Specifically, let  $\hat{M}$  be a vector with the first two race-specific elements of:

$$\hat{M}_{1r} = \sum_{j=1}^{268} \omega_{jr}^{\rho} \hat{\rho}_{jr}$$
(B28)

$$\hat{M}_{2r} = \sum_{j=1}^{268} \omega_{jr}^{\rho} (\hat{\rho}_{jr} - \hat{M}_{1r})^2$$
(B29)

the next three race-specific elements corresponding to coefficient estimates from the  $\omega_{jr}^{\lambda}$ -weighted quadratic OLS regression of:

$$\hat{\lambda}_{jr} = \hat{M}_{3r} + \hat{M}_{4r}\hat{\rho}_{jr} + \hat{M}_{5r}\hat{\rho}_{jr}^2 + \hat{\upsilon}_{jr}$$
(B30)

and the sixth race-specific element corresponding to the  $\omega_{ir}^{\lambda}$ -weighted residual variance estimate:

$$\hat{M}_{6r} = \sum_{j=1}^{268} \omega_{jr}^{\lambda} \hat{v}_{jr}^2$$
(B31)

The weights are derived from the estimation error matrix  $\Sigma$ :  $\omega_{jr}^{\rho}$  is proportional to the inverse variance of  $\hat{\rho}_{jr} - \rho_{jr}$  while  $\omega_{jr}^{\lambda}$  is proportional to the inverse variance of  $\hat{\lambda}_{jr} - \lambda_{jr}$ , with both weights rescaled to sum to one in the population of judges. We further include in  $\hat{M}$  the  $\sqrt{\omega_{jw}^{\rho}\omega_{jb}^{\rho}}$ -weighted covariance of  $\hat{\rho}_{jw}$  and  $\hat{\rho}_{jw}$  as well as the  $\sqrt{\omega_{jw}^{\lambda}\omega_{jb}^{\lambda}}$ -weighted covariance of  $\hat{\lambda}_{jw}$  and  $\hat{\lambda}_{jw}$ . Together this gives 14 elements in  $\hat{M}$ , the same number of hyperparameters in  $\Theta$ .

To estimate  $\Theta$  we use a just-identified SMD procedure that matches the empirical moments in  $\hat{M}$  with the corresponding model-implied moments averaged across 500 simulated draws of the above data-generating process. That is, we estimate:

$$\hat{\Theta} = \arg\min_{\Theta} \sum_{m=1}^{14} \left( \hat{M}_m - \frac{1}{500} \sum_{s=1}^{500} M_{ms}(\Theta) \right)^2$$
(B32)

where the functions  $M_{ms}(\cdot)$  of candidate hyperparameters  $\Theta$  are given by applying the previous moment calculations to data generated from 500 fixed simulation draws s. Conventional asymptotic theory for  $\hat{\Theta}$  applies under appropriate regularity conditions (e.g., Pakes and Pollard, 1989).

Columns 3 and 6 of Appendix Table A11 report SMD estimates and standard errors for the full model. As discussed in the main text, our baseline model estimates set  $\beta_r = 0$ . Per the intuition in Section 6.1 and to keep the model just-identified, we correspondingly drop the quadratic term from the moment regression in Equation (B30). The resulting estimates are reported in columns 2 and 5 of Appendix Table A11. To impose conventional MTE monotonicity, we further set the variance of  $\tau_{jr}$  to zero (again, per the intuition in Section 6.1), and drop the residual variance moment given in Equation (B31). The resulting estimates are reported in columns 1 and 4 of Appendix Table A11.

Lastly, given  $\hat{\Theta}$ , we compute maximum *a posteriori* probability estimates (also known as posterior modes) of the judge-specific parameters  $\theta_j = (\kappa_{jw}, \ln \tau_{jw}, \kappa_{jb}, \ln \tau_{jb})'$ , following an approach similar to that which Angrist et al. (2017) apply for a similar hierarchical model. Note that the log-likelihood
of  $\theta = (\theta'_1 \dots, \theta'_{268})'$  and quasi-experimental estimates  $\mathcal{X}$  can be written:

$$\mathcal{L}(\theta, \mathcal{X}) = \ln \phi_m \left( \mathcal{X} - \bar{X}(\theta); \Sigma \right) + \ln \phi_m \left( \theta - \mu_\theta; \Omega_\theta \right)$$
(B33)

where  $\phi_m(\cdot; V)$  gives the density of a mean-zero multivariate normal vector with variance-covariance matrix  $V; \bar{X}(\cdot)$  collects the formulas from Equations (B25) and (B26), for  $\rho_{jr}$  and  $\lambda_{jr}$  in terms of  $\mu_w$ ,  $\mu_b$ , and  $\theta$ ; and both  $\mu_{\theta}$  and  $\Omega_{\theta}$  are derived from the  $\alpha_r$  and  $\beta_r$ ,  $\mu_{\kappa}$ ,  $\Omega_{\kappa}$ , and  $\Omega_{\tau}$ . Our estimates of  $\theta$ are given by maximizing this likelihood, plugging in our baseline hyperparameter estimates  $\hat{\Theta}$ .