

University of Chicago Law School

## Chicago Unbound

---

Public Law and Legal Theory Working Papers

Working Papers

---

2018

### Racial Equity in Algorithmic Criminal Justice

Aziz Z. Huq

Follow this and additional works at: [https://chicagounbound.uchicago.edu/public\\_law\\_and\\_legal\\_theory](https://chicagounbound.uchicago.edu/public_law_and_legal_theory)



Part of the [Law Commons](#)

Chicago Unbound includes both works in progress and final versions of articles. Please be aware that a more recent version of this article may be available on Chicago Unbound, SSRN or elsewhere.

---

#### Recommended Citation

Aziz Z. Huq, "Racial Equity in Algorithmic Criminal Justice". Public Law and Legal Theory Working Papers, no. 663 (2018).

This Working Paper is brought to you for free and open access by the Working Papers at Chicago Unbound. It has been accepted for inclusion in Public Law and Legal Theory Working Papers by an authorized administrator of Chicago Unbound. For more information, please contact [unbound@law.uchicago.edu](mailto:unbound@law.uchicago.edu).

## Racial Equity in Algorithmic Criminal Justice

Aziz Z. Huq\*

(forthcoming, 68 *Duke Law Journal* – (2019))

### Abstract

*Algorithmic tools for predicting violence and criminality are being increasingly used in policing, bail, and sentencing. Scholarly attention to date has focused on their procedural due process implications. My aim here is to consider these instruments' interaction with the enduring racial legacies of the criminal justice system. There are two competing lenses for evaluating the racial effects of algorithmic criminal justice: constitutional doctrine and emerging technical standards of "algorithmic fairness." I argue first that constitutional doctrine is poorly suited to the task. It will often fail to capture the full range of racial issues that potentially arise in the use of algorithmic tools in criminal justice. While the emerging technical standards of algorithmic fairness are at least fitted to the specifics of the relevant technology, the technical literature has failed to ask how various conceptions of fairness track (or fail to track) policy-significant consequences. Drawing on the technical literature, I propose a reformulated metric for considering racial equity concerns in algorithmic design. Rather than asking about abstract definitions of fairness, a criminal justice algorithm should be evaluated in terms of its long-term, dynamic effects on racial stratification. The metric of nondiscrimination for the algorithmic context should focus on the net burden placed on a racial minority. A precise formulation of this metric suggests that it can converge with the socially efficient decision rule under certain conditions.*

---

\* Frank and Bernice J. Greenberg Professor of Law, University of Chicago Law School. Thanks in particular to Sharad Goel, Ravi Shroff, and Sam Corbett-Davies for their help over several conversations, in which they patiently explained the interaction of machine learning and statistical concepts to me. Emily Berman, Vincent Chiao, Jessica Clarke, Nancy King, Kiel Brennan Marquez, Debbie Hellman, Andrew Selbst, Chris Slobogin all gave me helpful comments that improved my thinking and corrected my errors. I was also helped greatly by workshops at Vanderbilt Law School, Northwestern Law School, and the University of Toronto Law School. Faith Laken provided terrific research assistance. All errors remain my own.

## Table of Contents

<b>Introduction</b> .....	<b>3</b>
<b>I. Algorithmic Criminal Justice: Scope and Operation</b> .....	<b>13</b>
<b>A. A Definition of Algorithmic Criminal Justice</b> .....	<b>14</b>
<b>B. The Operation of Algorithmic Criminal Justice</b> .....	<b>15</b>
1. Machine Learning and Deep Learning.....	16
2. The Impact of Machine Learning on Criminal Justice.....	18
<b>C. Algorithmic Criminal Justice on the Ground</b> .....	<b>20</b>
1. Policing.....	20
2. Bail.....	23
3. Sentencing.....	25
<b>D. The Emerging Evidence of Race Effects</b> .....	<b>26</b>
1. Policing and the Problem of Tainted Training Data.....	26
2. Bail/Sentencing Predictions and the Problem of Distorting Feature Selection.....	29
3. Conclusion: An Incomplete Evidentiary Record.....	30
<b>II. Equal Protection and Algorithmic Criminal Justice</b> .....	<b>31</b>
<b>A. What Equal Protection Protects</b> .....	<b>31</b>
<b>B. How Equal Protection Fails to Speak in Algorithmic Terms</b> .....	<b>34</b>
1. The Trouble With Intent.....	35
2. The Trouble with Classification.....	39
3. The Lessons of Algorithmic Technology for Equal Protection Doctrine.....	45
<b>III. Racial Equity in Algorithmic Criminal Justice Beyond Constitutional Law</b> .....	<b>46</b>
<b>A. The Stakes of Racial Equity in Contemporary American Criminal Justice</b> .....	<b>47</b>
<b>B. A Racial Equity Principle for Algorithmic Criminal Justice</b> .....	<b>53</b>
<b>C. Benchmarks for Algorithmic Discrimination</b> .....	<b>56</b>
<b>D. Prioritizing Conceptions of Algorithmic Discrimination</b> .....	<b>62</b>
1. Conflicts Between Algorithmic Fairness Definitions.....	63
2. The Irrelevance of False Positive Rates.....	64
3. Evaluating the Impact of Algorithmic Criminal Justice on Racial Stratification.....	66
<b>Conclusion</b> .....	<b>70</b>

## Introduction

From the cotton gin to the camera phone, new technologies have scrambled, invigorated, and refashioned the terms on which the state coerces. Today, we are in the midst of another major reconfiguration. Police, criminal courts, and parole boards across the country are turning to sophisticated algorithmic instruments to guide decisions about the ‘where,’ ‘whom,’ and ‘when’ of law enforcement.<sup>1</sup> The new predictive algorithms trawl immense quantities of data, exploit massive computational power, and leverage new machine-learning technologies to generate predictions no human could conjure. These tools are likely to have enduring effects on the criminal justice system. Yet law remains far behind in thinking through the difficult questions that arise when machine learning substitutes for human discretion.

My aim in this Article is to isolate one important design margin for evaluating algorithmic criminal justice: the effect of algorithmic criminal justice tools on *racial equity*. I use this capacious term to capture the complex ways in which the state’s use of a technology can implicate normative and legal concerns related to racial dynamics. The Article considers a number of ways in which legal scholars and computer scientists have theorized race. It evaluates these distinct approaches in terms of the way in which criminal justice in practice interacts with racial patterning. A primary lesson concerns the parameter that best captures racial equity concerns in an algorithmic setting. A secondary lesson relates to the fit between problems of race in the algorithmic context on the one hand, and legal or technical conceptions of equality on the other.

‘Racial equity’ merits a discrete, detailed inquiry given the fraught racial history of American criminal justice institutions. Since the turn of the twentieth century, public arguments about criminality have been entangled, often invidiously, with generalizations about race and the putative criminality of racial minorities.<sup>2</sup> Today, pigmentation remains de facto a proxy for criminality; that proxy distorts everything from residential patterns to labor market opportunities.<sup>3</sup> Police respond to black and white suspects in different ways.<sup>4</sup>

---

<sup>1</sup> Reed E. Hundt, *Making No Secrets About It*, 10 IS J. L. & POL. 581, 588 (2014) (“[The G]overnment now routinely asks computers to suggest who has committed crimes.”).

<sup>2</sup> KHALIL GIBRAN MUHAMMAD, *THE CONDEMNATION OF BLACKNESS: RACE, CRIME, AND THE MAKING OF MODERN URBAN AMERICA* (2010) (exploring the ways in which at the beginning of the twentieth century, policymakers in Northern cities began linking crime to African-Americans on the basis of genetic and predispositional arguments).

<sup>3</sup> See, e.g., Robert J. Sampson & Stephen W. Raudenbush, *Seeing Disorder: Neighborhood Stigma and the Social Construction of “Broken Windows,”* 67 SOC. PSYCHOL. Q. 319, 319-20 (2004) (finding that perceptions of disorder in a neighborhood were better predicted by the racial composition of a neighborhood than by actual disorder); Lincoln Quillian & Devah Pager, *Black neighbors, higher crime? The role of racial stereotypes.* 107 AM. J. SOC. 717, 718 (2001) (finding “that the percentage of a neighborhood’s black population, particularly young black men, is significantly associated with perceptions of the severity of a neighborhood’s crime problems”).

<sup>4</sup> For evidence, see CHARLES EPP, STEVEN MAYNARD-MOODY, & DONALD HAIDER-MARKEL, *PULLED OVER: HOW POLICE STOPS DEFINE RACE AND CITIZENSHIP* 32-33 (2014); Aziz Z. Huq, *The Consequences*

So do judges and prosecutors.<sup>5</sup> Partly as a result of these dynamics, roughly one in three black men (and one in five Latino men) will be incarcerated during their lifetime.<sup>6</sup> At the same time, the criminal justice system imposes substantial socioeconomic costs on minority citizens not directly touched by policing or prosecutions. In particular, minority children of the incarcerated bear an unconscionable burden as a result of separation from their parents.<sup>7</sup> More generally, there is substantial evidence that spillover costs of producing public safety fall disproportionately on minority groups.<sup>8</sup> As a result, criminal justice shapes racial stratification.<sup>9</sup> Such downstream consequences of existing criminal-justice institutions raise weighty moral and legal questions.<sup>10</sup> Even if one demurs to the analogy commonly drawn between our criminal justice system and early twentieth-century debt peonage,<sup>11</sup> I think it is clear that the criminal justice system is an institution in which racial identity has meaningful effects, and that these in turn have influences on the role of race in larger American society.<sup>12</sup>

To sharpen this point, it is useful to have at hand two examples of how new technologies can prompt debates about racial equity. I present the first at greater length because it has become a focal point in public debates. First, the Compas software application, created by the Northpointe Institution for Public Management, is used from Florida to Wisconsin to inform bail and parole decisions. Compas is organized around an algorithm that uses the answers to some 137 questions about a criminal suspect to rank them on a scale of 1 to 10. This scale is supposed to capture the suspect's risk of reoffending

---

*of Disparate Policing: Evaluating Stop and Frisk As A Modality of Urban Policing*, 101 MINN. L. REV. 2397, 2408 (2017) [hereinafter "Huq, *Disparate Policing*"] (discussing evidence of such disparities).

<sup>5</sup> For two different perspectives, emphasizing intentional bias and disparate racial impacts, see Sonja B. Starr & Marit Rehavi, *Mandatory Sentencing and Racial Disparity: Assessing the Role of Prosecutors and the Effects of Booker*, 123 YALE L.J. 2, 25-30 (2013) (documenting racial disparities in federal prosecutorial charging decisions related to the application of mandatory minimum sentences in drug cases); Richard Frase, *What Explains Persistent Racial Disproportionality in Minnesota's Prison and Jail Populations?*, 38 CRIME & JUST. 201, 265 (2009) (finding that "seemingly legitimate sentencing factors such as criminal history scoring can have strongly disparate impacts on nonwhite defendants").

<sup>6</sup> BRUCE WESTERN, PUNISHMENT AND INEQUALITY IN AMERICA 31-39 (2006) (describing the growth of the incarcerated population over time, and describing racial inequalities); Cassia Spohn, *Race, Crime, and Punishment in the Twentieth and Twenty-First Centuries*, 44 CRIME & JUST. 49, 55 (2015) (noting that in 2001 "the chances of ever going to prison were highest among black males (32.2 percent) and Hispanic males (17.2 percent)").

<sup>7</sup> See SARA WAKEFIELD & CHRISTOPHER WILDEMAN, CHILDREN OF THE PRISON BOOM: MASS INCARCERATION AND THE FUTURE OF AMERICAN INEQUALITY 41 (2014) (discussing the racially disparate spillover effects of incarceration).

<sup>8</sup> *Id.*

<sup>9</sup> For a synoptic view of this claim that is dated, but still insightful, see RANDALL KENNEDY, RACE, CRIME, AND THE LAW (1997).

<sup>10</sup> I think it is important for legal scholars to be candid in distinguishing their normative judgments from their analytic, doctrinal, and empirical claims. The following paragraph states my normative position; it is a premise of what follows, not a conclusion I seek to defend here. See also *supra* Part III.A (further defending this position).

<sup>11</sup> See, e.g., MICHELE ALEXANDER, THE NEW JIM CROW: MASS INCARCERATION IN AN AGE OF COLORBLINDNESS (2010). For nuanced criticism of Alexander's paradigm, James Forman, Jr., *Racial Critiques of Mass Incarceration: Beyond the New Jim Crow*, 87 N.Y.U. L. REV. 21, 42-43 (2012).

<sup>12</sup> James Q. Whitman, *Equality in Criminal Law: The Two Divergent Western Roads*, 1 J. L. ANALYSIS 119, 122 (2009).

and violent recidivism.<sup>13</sup> Higher scores indicate a greater risk of recidivism. In 2016, journalists from the Pro Publica organization did a quantitative analysis of Compas scores for roughly ten thousand people arrested and evaluated in Broward County, Florida. By comparing Compas scores to a person's behavior in the following two years, Pro Publica was able to evaluate the instrument's accuracy, and in particular to investigate whether it had differential effects on different racial groups.

Pro Publica concluded that the Compas instrument correctly predicted recidivism rates 61 percent of the time, and violent recidivism rates 21 percent of the time.<sup>14</sup> Pro Publica also concluded that the algorithm "was particularly likely to falsely flag black defendants as future criminals, wrongly labeling them this way at almost twice the rate as white defendants."<sup>15</sup> To reach this conclusion, Pro Publica isolated the group of black suspects who had not reoffended in the two years following their evaluation. It found that 45 percent of that group was labeled high risk by the algorithm.<sup>16</sup> Pro Publica then looked at the group of white suspects who had not reoffended, and found that only 23 percent of that group had been labeled high risk. In other words, the ratio of false positives to true negatives within the pool of 'innocent' defendants was higher for blacks than for whites.<sup>17</sup> Correspondingly, Pro Publica also found that the ratio of false negatives to true positives was lower for whites than for blacks.<sup>18</sup>

Not surprisingly, the company responded by sharply contesting Pro Publica's analysis. Northpointe data scientists insisted that Compas was well-calibrated in the sense that a white and a black defendant assigned the same risk score were equally likely to recidivate.<sup>19</sup> This constituted evidence, the company argued, that where it mattered to the imposition of state coercion (i.e., where there was a prediction of high risk), the Compas

---

<sup>13</sup> For descriptions of the Compas algorithm, see Northpointe, *Practitioners Guide to COMPAS* 17 (2012), [http://www.northpointeinc.com/files/technical\\_documents/FieldGuide2\\_081412.pdf](http://www.northpointeinc.com/files/technical_documents/FieldGuide2_081412.pdf). [https://perma.cc/4FX T-6U9M]; see also *In re Hawthorne v. Stanford*, 135 A.D.3d 1036, 1037-38 (3d Dep't 2016) (describing the COMPAS assessment tool).

<sup>14</sup> Jeff Larsen et al., *How we Analyzed the Compas Recidivism Algorithm*, PROPUBLICA (May 23, 2016), <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>. Note that rates of violent crime tend to be so low that an 'accurate' instrument would be one that simply classified everyone as low risk.

<sup>15</sup> Julia Angwin et al., *Machine Bias: There's Software Used Across the Country to Predict Future Criminals. And It's Biased Against Blacks*, PROPUBLICA (May 23, 2016), <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing/>. Pro Publica treated 'medium' and 'high' risk rankings as higher risk.

<sup>16</sup> Larsen et al., *supra* note 14. This disparity remained once Pro Publica controlled for "prior crimes, future recidivism, age, and gender." *Id.*

<sup>17</sup> *Id.*

<sup>18</sup> *Id.*

<sup>19</sup> William Dietrich et al., *COMPAS Risk Scales: Demonstrating Accuracy Equity and Predictive Parity*, Northpointe Inc. Research Department, July 8, 2016, at 3 <https://www.documentcloud.org/documents/2998391-ProPublica-Commentary-Final-070616.html> (flagging "equal discriminative ability" of the algorithm for blacks and whites); see also Anthony W. Flores et al., *False Positives, False Negatives, and False Analyses: A Rejoinder to "Machine Bias: There's Software Used Across the Country to Predict Future Criminals. And it's Biased Against Blacks,"* 80 FED PROBATION 34, 35 (2016) (describing the Pro Publica analysis as "faulty"). For a different result using a reconstruction of the Compas algorithm, see Razieh Nabi & Ilya Shpitser, *Fair inference on outcomes*, 8 (2017), <https://arxiv.org/abs/1705.10378>.

algorithm had equal error rates across groups. In addition, Northpointe made a number of (contested) technical complaints about Pro Publica’s analysis, related to the way it accounted for base recidivism rates and how it cut its sample between low and high risk defendants.<sup>20</sup> These complaints lacked the force of Northpointe’s central claim—that its risk predictions were equally accurate where it counted regardless of race. This dialogue was not the end of the matter. Other analysts raised a cautionary flag to warn against accepting the terms of the debate as framed by Pro Publica and Northpointe: Something more complex, they worried, seemed at stake, although they did not explain fully how the debate should be settled.<sup>21</sup> As a result, debate on Compas—and in particular the question of which measure of fairness should be used to evaluate a predictive algorithm—persists as a locus for normative concern.

A second example of the race-related questions potentially raised by algorithmic criminal justice arises in the policing context, where officers are increasingly using such tools in determining where to deploy and who to apprehend.<sup>22</sup> In Chicago, police faced with a wave of deadly street violence<sup>23</sup> have deployed a “Strategic Subjects List,” or SSL. This is an algorithm developed by data scientists at the Illinois Institute of Technology using U.S. Department of Justice funds.<sup>24</sup> The SSL ranks individuals known to police for the risk of involvement in a shooting using eight data points.<sup>25</sup> Its aim, according to the chief of organizational development for the department, was “to figure out now ... how does that data inform what happens in the future.”<sup>26</sup> Yet despite the fact that the SSL algorithm explicitly accounted for neither race nor gender,<sup>27</sup> interventions based on SSL were quickly lambasted for focusing solely on African-American men.<sup>28</sup> Other algorithms

---

<sup>20</sup> Dietrich et al., *supra* note 19, at 32-33.

<sup>21</sup> See, e.g., Avi Feller et al., *A computer program used for bail and sentencing decisions was labeled biased against blacks. It’s actually not that clear*, WASH. POST, Oct 16, 2016, [https://www.washingtonpost.com/news/monkey-cage/wp/2016/10/17/can-an-algorithm-be-racist-our-analysis-is-more-cautious-than-propublicas/?utm\\_term=.f8164ea2cd2c](https://www.washingtonpost.com/news/monkey-cage/wp/2016/10/17/can-an-algorithm-be-racist-our-analysis-is-more-cautious-than-propublicas/?utm_term=.f8164ea2cd2c); Matthias Spielkamp, *Inspecting Algorithms for Bias*, MIT TECH. REV., Jun 17, 2017, <https://www.technologyreview.com/s/607955/inspecting-algorithms-for-bias/>. For my analysis of the Compas algorithm, see text accompanying *infra* notes 325 to 328.

<sup>22</sup> Mara Hvistendahl, *Can ‘predictive policing’ prevent crime before it happens?*, SCIENCE, Sept. 28, 2016, <http://www.sciencemag.org/news/2016/09/can-predictive-policing-prevent-crime-it-happens> (noting the adoption of policing tools “which incorporate everything from minor crime reports to criminals’ Facebook profiles”); see also Andrew Guthrie Ferguson, *Policing Predictive Policing*, 94 WASH. U.L. REV. 1109, 1122-44 (2017) (providing a careful catalogue of predictive policing tools).

<sup>23</sup> Monica Davey, *Chicago Tactics Put a Major Dent in Killing Trend*, N.Y. TIMES, June 11, 2013, at A1.

<sup>24</sup> City of Chicago, *Strategic Subjects List: Public Safety*, Dec. 7, 2017, <https://data.cityofchicago.org/Public-Safety/Strategic-Subject-List/4aki-r3np>.

<sup>25</sup> *Id.*

<sup>26</sup> Jeremy Gorner, *Chicago Police Use “Heat List” As Strategy to Prevent Violence*, CHI. TRIB., Aug. 21, 2013, [http://articles.chicagotribune.com/2013-08-21/news/ct-met-heat-list-20130821\\_1\\_chicago-police-commander-andrew-papachristos-heat-list](http://articles.chicagotribune.com/2013-08-21/news/ct-met-heat-list-20130821_1_chicago-police-commander-andrew-papachristos-heat-list).

<sup>27</sup> Other predictive policing instruments, however, do explicitly account for suspects’ race. David Robinson & Logan Koepke, *Upturn, Stuck in a Pattern: Early Evidence on “Predictive Policing” and Civil Rights* 4-5 (2016), [https://www.teamupturn.com/static/reports/2016/predictive-policing/files/Upturn\\_-\\_Stuck\\_In\\_a\\_Pattern\\_v.1.01.pdf](https://www.teamupturn.com/static/reports/2016/predictive-policing/files/Upturn_-_Stuck_In_a_Pattern_v.1.01.pdf).

<sup>28</sup> Matt Stoudt, *The minority report: Chicago’s new police computer predicts crimes, but is it racist?*, THE VERGE, Feb. 14, 2014, <https://www.theverge.com/2014/2/19/5419854/the-minority-report-this-computer-predicts-crime-but-is-it-racist>.

used to guide the allocation of policing resources on geographic rather than individual terms have elicited kindred concerns about racial targeting.<sup>29</sup>

Questions about algorithmic criminal justice are poised to become more complex. Compas and the SSL are both quite straightforward instruments. Each applies a fixed regression equation with a limited array of parameters to a static data set. Advances in what is called *machine learning*, however, will soon render this sort of tool passé. Machine learning is a “general purpose”<sup>30</sup> technology that, in broad terms, encompasses “algorithms and systems that improve their knowledge or performance with experience.”<sup>31</sup> A supervised machine-learning instrument—the species of machine learning likely most relevant in the criminal justice space<sup>32</sup>—begins with a so-called training set of examples that are ‘labeled’ with some parameter values. The algorithm examines relations between various aspects of those examples to develop a wholly new criterion to classify new examples.<sup>33</sup> Unlike more familiar econometric tools such as regression analysis, a supervised machine learning process classifies on the basis of rules that the algorithm itself has developed. Refining this process, deep learning tools deploy “multilayered” processes, account for billions of data points, and constantly adjust their classification rule.<sup>34</sup> Machine learning is now being deployed, for instance, in Cambridge, MA, to predict house burglaries,<sup>35</sup> and in Durham, England, to predict individual recidivism.<sup>36</sup> Deep learning is used in facial recognition and machine translation; it will likely find new uses as its

---

<sup>29</sup> Justin Jouvenal, *Police are Using Software to Predict Crime. Is it a ‘Holy Grail’ or Biased Against Minorities?*, WASH. POST (Nov. 17, 2016), [https://www.washingtonpost.com/local/public-safety/police-are-using-software-to-predict-crime-is-it-a-holy-grail-or-biased-against-minorities/2016/11/17/525a6649-0472-440a-aae1-b283aa8e5de8\\_story.html?utm\\_term=.72a9d2eb22ae](https://www.washingtonpost.com/local/public-safety/police-are-using-software-to-predict-crime-is-it-a-holy-grail-or-biased-against-minorities/2016/11/17/525a6649-0472-440a-aae1-b283aa8e5de8_story.html?utm_term=.72a9d2eb22ae).

<sup>30</sup> Erik Brynjolfsson and Tom Mitchell, *What can machine learning do? Workforce implications*, SCIENCE, Dec. 22, 2017, <http://science.sciencemag.org/content/358/6370/1530.full>.

<sup>31</sup> PETER FLACH, MACHINE LEARNING: THE ART AND SCIENCE OF ALGORITHMS THAT MAKE SENSE OF DATA 3 (2012); ETHEM ALPAYDIN, INTRODUCTION TO MACHINE LEARNING 2 (3d ed. 2014); *see also* text accompanying *infra* notes 73 to 82 for a fuller account of machine learning, and text accompanying *infra* notes 84 to 88 for a discussion of deep learning.

<sup>32</sup> Susan Athey, *Beyond Prediction: Using Big Data for Policy Problems*, SCIENCE, Feb 3, 2017, at 483, <http://science.sciencemag.org/content/355/6324/483.full> (noting the use of structured machine learning to solve prediction problems).

<sup>33</sup> M. I. Jordan & T. M. Mitchell, *Machine learning: Trends, perspectives, and prospects*, SCIENCE, Jul. 17, 2015, at 255 (defining supervised learning as a process in which “the training data take the form of a collection of (x, y) pairs and the goal is to produce a prediction y\* in response to a query x\*”); Athey, *supra* note 32, at 483 (explaining that machine learning “programs take as input training data sets and estimate or ‘learn’ parameters that can be used to make predictions on new data”); Comm. on the Analysis of Massive Data et al., *Frontiers in Massive Data Analysis* 104 (2013), [http://www.nap.edu/catalog.php?record\\_id=18374](http://www.nap.edu/catalog.php?record_id=18374) (noting that in supervised learning, the analyst must actively specify a variable of interest).

<sup>34</sup> Jordan & Mitchell, *supra* note 33, at 256. Deep learning uses a process called stochastic gradient ascent to improve predictive quality continuously. Yann Le Cuan et al., *Deep Learning*, 521 SCIENCE 436, 437 (2015).

<sup>35</sup> Cynthia Rudin, *Predictive Policing: Using Machine Learning to Detect Patterns of Crime*, WIRED, Aug. 2013, <https://www.wired.com/insights/2013/08/predictive-policing-using-machine-learning-to-detect-patterns-of-crime/>.

<sup>36</sup> Chris Baraniuk, *Durham Police AI to help with custody decisions*, BBC NEWS, May 10, 2017 <http://www.bbc.com/news/technology-39857645>; *see infra* Part I.C for a catalog of more examples of how machine learning is used in the criminal justice context.



capabilities are further explored. My use of the term “algorithmic criminal justice” is intended to capture both existing instruments such as Compas and the SSL, and also machine learning (including deep learning) tools that are likely to be deployed for prediction purposes in the future. Such synoptic consideration is warranted because all these tools leverage historical data to generate ‘predictions’ for new, out-of-sample data.<sup>37</sup>

Algorithmic tools in criminal justice warrant distinct treatment because they are a new technology that is likely to become pervasive in respect to which normative intuitions remain inchoate and hence malleable. They also represent a qualitative change from the crude evaluative tools embodied in present bail and sentencing practices. These build on imprecise measures of recidivism risk, fail to account for immediate or downstream costs, and cannot be calibrated with the precision of emerging tools. The precision enabled by the algorithmic turn pries open a substantively new domain of policy-design possibilities.

There are, very generally stated, two distinct analytic frameworks in use now for evaluating the racial effects on machine-learning tools in criminal justice. The first derives from constitutional law. The second is derived from the computer-science literature on algorithm design.<sup>38</sup> Both, in my view, fall short. The constitutional law of racial inequality directs attention to trivial or irrelevant design margins. It is at times counterproductive. In contrast, technical discussions of algorithmic fairness have yielded an array of parameters that capture different elements of an algorithm’s operation. But as the debate between *Pro Publica* and *Northpointe* shows, the computer-science literature has generated no clear consensus about which parameter *matters*. This Article fills the gap left by the irrelevance of constitutional law and the under-theorization of computer-science. It offers a novel, normatively grounded, and empirically pertinent framework for thinking about racial equity in this emerging technological context.<sup>39</sup>

Consider first the current constitutional framework for the regulation of race effects in policing. The doctrine, in rough paraphrase, has two main prongs. One concern in the jurisprudence turns on the use of “racial stereotypes or animus” held by individual actors.<sup>40</sup> A focus on animus or stereotypes, though, doesn’t easily translate into contexts in which an algorithm blends data streams to estimate unknown parameter values. At best, a concern with intent captures a subset of problematic cases in which data inputs are tainted. Second, Equal Protection law is also concerned with the use of racial classifications. But in the emergent context of algorithmic criminal justice, where decision rules are computed

---

<sup>37</sup> I use the term prediction not because all these instruments aim at the future. Rather, the term captures the possibility that one data set will be used to generate an instrument for drawing inferences about a different sample of data. It is a prediction in the sense of being an out-of-sample estimate.

<sup>38</sup> I will not work through all of the relevant computer science literature here. For a brief survey that touches on some of the questions analyzed here, see Joshua A. Kroll et al., *Accountable Algorithms*, 165 U. PA. L. REV. 633, 682-90 (2017).

<sup>39</sup> To the extent that algorithmic tools are more generally replacing diffused human discretion, my reconceptualization of equality norms may have more general application.

<sup>40</sup> *Pena-Rodriguez v. Colorado*, 137 S. Ct. 855, 869 (2017); see also Aziz Z. Huq, *What is Discriminatory Intent?*, 103 CORNELL L. REV. --, at 10-21 (forthcoming 2018) [hereinafter “Huq, *Discriminatory Intent*”], [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3033169](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3033169) (analyzing the central role of intent in the context of the Fourteenth Amendment).

endogenously from historical data and then applied without being broadcast to the public, the expressive or distortive harms of racial classifications may well not be present. Crudely stated, an algorithm's use of racial data is unlikely to stigmatize or otherwise impose any harm putatively linked to the use of suspect classifications. Eliminating such criteria, moreover, can leave actual outcomes unchanged. Worse, it can generate needless public-safety-related costs. This is because algorithmic use of a proscribed criterion, such as race, might in some instances improve the quality of predictions. Thinking about Equal Protection jurisprudence in relation to algorithmic criminal justice therefore suggests that the former is not a coherent or morally acute metric.

If constitutional law provides no creditable guidance, what about the burgeoning computer-science scholarship on “algorithmic fairness” and “algorithmic discrimination,” terms to date used to cover a number of different means of evaluating predictive tools?<sup>41</sup> At a very high level of abstraction, the technical literature usefully distinguishes between two different ways in which race effects might emerge in algorithmic criminal justice. The first is the use of racially tainted historical data to build an algorithm. For example, a policing algorithm used to predict who will be involved in crime, such as SSL, might employ data gathered by police, such as records of past street stops or past arrests. If the pattern of this historical policing activity is informed by racial considerations, then the algorithm's predictions will be accordingly skewed. Fixing this first problem of polluted training data is straightforward. As several legal scholars have noted, algorithms can simply be constructed without tainted training data.<sup>42</sup> Whatever difficulties this might present in terms of implementation, it raises no great theoretical impediment.

But the second way in which a racial problem can arise from the use of algorithmic tools does.<sup>43</sup> It turns on the possibility that an algorithm will generate patterns of error that are systematically skewed between racial groups. As the debate between Pro Publica and Northpointe illustrates, however, there is more than one way of measuring errors, and more than one way of thinking about racial skewing. Indeed, the computer science literature has generated an embarrassment of possible metrics. Simplifying this literature by stripping away redundant and irrelevant conceptual trappings, I suggest that an analysis of racial equity might focus on one of four different parameters.

One might *first* simply look at whether equal fractions of each racial group are labeled as risky—such that they will be subject to additional policing or detention. Where risk is measured as a continuous variable, this would mean looking at whether the average risk scores of different racial groups varied. *Second*, one might ask whether the same

---

<sup>41</sup> See *infra* Part II for a survey of the relevant work.

<sup>42</sup> See, e.g., Kroll et al., *supra* note 38, at 680 (“[A]lgorithms that include some type of machine learning can lead to discriminatory results if the algorithms are trained on historical examples that reflect past prejudice or implicit bias”); Anupam Chander, *The Racist Algorithm?*, 115 MICH. L. REV. 1023, 1039 (2017) (same); Kate Crawford, *Artificial Intelligence's White Guy Problem*, N.Y. TIMES, June 25, 2016, <https://www.nytimes.com/2016/06/26/opinion/sunday/artificial-intelligences-white-guy-problem.html> (same).

<sup>43</sup> Kroll et al., recognize that “machine learning models can build in discrimination through choices in how models are constructed.” Kroll et al., *supra* note 38, at 681. This is not, however, the central focus of their wide-ranging and useful analysis.

classification rule is being used to assign racial groups to the high-risk category. This condition is satisfied if the same numerical risk score is used as a cut-off for all groups. *Third*, one might separate each racial group and then look at the rate of false positives conditional on being categorized as high risk. That is, one could examine the rate of false positives conditional on being identified as risky. This is the parameter that Northpointe stressed. And *fourth*, one might separate each racial group and ask how frequently false positives are conditional on being in fact a low risk person. This is the parameter Pro Publica underscored.

Each of these metrics tracks a subtly different conception of nondiscrimination. So which fits best a normatively relevant conception of racial equity? The question is complicated by two considerations. First, there is an irreconcilable tension between the first and second criteria. If the average risk score of two racial groups diverge, it is not possible to use the same classification rule and also to ensure that an equal fraction of each group is categorized as high risk. The same risk threshold applied to different populations, that is, yields different results. Second, computer-science scholars (in collaboration with legal scholars, including myself) have developed in the past two years an impossibility result concerning the third and fourth metrics. Under most empirically plausible conditions, a risk instrument cannot satisfy both the third and the fourth criterion. That is, if the proportion of false positives as a fraction of all positives is equalized between races, then the ratio of low-risk individuals subject to coercion will diverge between the two groups. There is hence an irreconcilable tension in many states of the world between having equally accurate predictions of high risk and equalizing the rates of false positives within the pool of ‘innocent’ suspects.

To prioritize between these conceptions of racial equity, it is necessary to elaborate an account of the normative stakes of racial equity in criminal law. In the ordinary course, we might look to constitutional law to this end. But we have already seen that constitutional law does not provide a fit and tractable frame for analysis. I thus return to first principles. In my view, the primary reason for concern with racial equity in the algorithmic criminal justice context is that efforts to suppress crime entrench wider social patterns of racial stratification. In important part, stratification effects arise because of the asymmetrical spillovers from criminal justice for minority but not majority populations. A parameter for measuring racial equity, therefore, should be selected on the ground that it captures this causal effect of criminal justice on racial stratification.

An algorithm that recommends coercion for a member of the subordinated racial group at the margin when it is not justified in terms of benefits to that racial group will likely increase racial stratification. When coercion of the marginal minority group member is unjustified, it imposes a net burden on the minority group, compounding social stratification. Further, if the majority group does not benefit from the policy, or if its net gain is less than the costs imposed on the minority group, that policy is also socially inefficient.<sup>44</sup> I suspect that governments often over-estimate the crime-suppression benefits

---

<sup>44</sup> Only if the gains to a majority group exceed the costs to a minority group is there a tension between efficiency and racial equity. As I explain below, I think it is plausible to prioritize equality norms in many of these conflicts.

of coercive actions while underestimating their costs. Racial equity is therefore served in the first instance today by ratcheting coercion down to socially optimal levels,<sup>45</sup> and then by selecting for criminal justice tools that do not burden minority groups.

In designing an algorithm, this intuition must be translated into instructions for the classification protocol. This should be done differently for serious and less serious crime. For serious violent and property crime, the most important costs and benefits of crime (and crime prevention) accrue directly to the perpetrator and the victim. Spillovers are by comparison small. In these conditions, a single, socially optimal classification rule will advance racial equity and satisfy an efficiency criterion. Rates of false positives, underscored by *Pro Publica* and *Northpointe*, are less relevant. For less serious crimes and misdemeanors, however, empirical evidence studies identify large spillover costs asymmetrically imposed on minority but not majority communities. At the margin, these spillovers mean that coercion of the minority is both less likely to be efficient and more likely to generate racial stratification. Accordingly, a bifurcated classification rule employing different risk thresholds for stratified racial groups is appropriate to account for asymmetrical spillovers.

Such binary thresholds will be socially efficient and racially just, but face practical and legal hurdles. First, evaluating algorithmic tools in light of social externalities will require much more information about downstream costs than is presently available. Governments have been woefully deficient in collecting such data; present “risk assessment” instruments embody information about recidivism risk, but neither the direct nor the indirect costs of criminal justice coercion.<sup>46</sup> There is a large epistemic void here that scholars can fill, and it is possible that other big-data tools will be important in this regard. Second, the use of racially bifurcated thresholds would raise constitutional concerns akin to those engendered by affirmative action programs. To the extent current doctrine mandates an outcome that is both socially inefficient and also racially iniquitous, it is legally and morally indefensible.

Some limitations on my analysis in this Article should be demarcated up front. First, I should again underscore that the costs and benefits of algorithmic tools vary depending on where in the criminal justice process they are deployed. My aim here is to set out a general framework and not to pass judgment on any particular computational tool. Second, this Article does not address the integration of algorithmic outputs into individualized suspicion determinations under the Fourth Amendment<sup>47</sup> or the issues related to procedural

---

<sup>45</sup> In using the terms “social efficiency,” I mean to capture a static (and in my view naïve) account of welfare that looks only to proximate costs and benefits. It is my view that racial stratification is plausibly described as an ‘inefficient’ equilibrium to the extent that it dissipates large amounts of human capital while inflicting onerous psychological and stigmatic burdens. But since my view is not orthodox, I do not insist on it here and instead use “efficiency” in its more common sense.

<sup>46</sup> Christopher Slobogin, *Principles of Risk Assessment: Sentencing and Policing*, 15 OHIO J. CRIM. L. 583, 583 (2018).

<sup>47</sup> The best works include Kiel Brennan-Marquez, “*Plausible Cause*”: *Explanatory Standards in the Age of Powerful Machines*, 70 VAND. L. REV. 1249, 1254 (2017) (arguing that for predictions to be used as a basis for searches under the Fourth Amendment, they have to be “intelligible,” in the sense of being amenable to explanation), Michael L. Rich, *Machine Learning, Automated Suspicion Algorithms, and the Fourth*

due process rights from the Fifth and Fourteenth Amendments.<sup>48</sup> These constitutional rules engage different elements of algorithm design. For example, an important recent article develops a concept of “procedural regularity” to ensure that algorithmic decisions are “made using consistently applied standards and practices.”<sup>49</sup> This is an important concern. But it is distinct from racial equity. I also do not address the statutory standard supplied by Title VII of the Civil Rights Act of 1964, which has been the lens of other recent work on algorithmic justice.<sup>50</sup> Nor do I address algorithms’ use outside the criminal justice context.<sup>51</sup>

Finally, one prominent article examines the racial effects of a larger class of “evidence based” predictive instruments. It condemns those instruments in general, arguing that they elicit “overt discrimination based on demographics and socioeconomic status.”<sup>52</sup> Its legal analysis is premised on the dubious proposition that “[c]urrent” constitutional law “calls into serious question the variables related to socioeconomic status, such as employment status, education, income, dependence on government assistance, and job skills.”<sup>53</sup> I am not convinced this is an accurate statement of current law; my analysis thus

---

*Amendment*, 164 U. PA. L. REV. 871, 929 (2016) (developing a “framework” for integrating machine-learning technologies into Fourth Amendment analysis), and Andrew Guthrie Ferguson, *Big Data and Predictive Reasonable Suspicion*, 163 U. PA. L. REV. 327, 383-84 (2015). Judicial consideration of this issue has been limited. *Cf.* *Commonwealth v. Smith*, 709 S.E.2d 139, 143 (Va. 2011) (relying on constructive knowledge doctrine to allow officer use of a predictive algorithm in a Terry stop).

<sup>48</sup> Danielle Keats Citron, *Technological Due Process*, 85 WASH. U. L. REV. 1249, 1256-57 (2008) (criticizing the “crudeness” of then-extant algorithms, and urging greater opportunities for individual challenges).

<sup>49</sup> Kroll et al., *supra* note 38, at 637-38. For another argument focused on process, see Kate Crawford & Jason Schultz, *Big Data and Due Process: Toward A Framework to Redress Predictive Privacy Harms*, 55 B.C. L. REV. 93, 109 (2014) (arguing for “procedural data due process [to] regulate the fairness of Big Data’s analytical processes with regard to how they use personal data (or metadata ...) in any adjudicative process”).

<sup>50</sup> Solon Barocas & Andrew D. Selbst, *Big Data’s Disparate Impact*, 104 CALIF. L. REV. 671, 694-712 (2016) (examining “[i]ability under Title VII for discriminatory data mining [which] will depend on the particular mechanism by which the inequitable outcomes are generated.”); *accord* Kroll et al., *supra* note 38, at 692-95.

<sup>51</sup> In addition, there is a small body of insightful popular literature about the distributive effects of algorithmic instruments more generally. *See* CATHY O’NEILL, *WEAPONS OF MASS DESTRUCTION: HOW BIG DATA INCREASES INEQUALITY AND THREATENS DEMOCRACY* 202-206 (2016) (decrying the regressive tendencies of big-data technologies generally); *accord* VIRGINIA EUBANKS, *AUTOMATING INEQUALITY: HOW HIGH-TECH TOOLS PROFILE, POLICE, AND PUNISH THE POOR* (2018) (same).

<sup>52</sup> Sonja B. Starr, *Evidence-Based Sentencing and the Scientific Rationalization of Discrimination*, 66 STAN. L. REV. 803, 806 (2014). One other article contains the glib assertion that “if racial and ethnic variables significantly improved the predictive validity of risk-needs models, then including them would appear to be narrowly tailored to the government’s compelling interests.” Melissa Hamilton, *Risk-Needs Assessment: Constitutional and Ethical Challenges*, 52 AM. CRIM. L. REV. 231, 259 (2015). Hamilton equates ‘narrow tailoring’ with ‘minimal efficacy.’ She fails to meaningfully grapple with existing precedent. And she is opaque as to what kind of race effects might have legal or normative significance. Her analysis is thus quite limited. Finally, a brief 2016 article suggests that the application of certain algorithmic tools in a sentencing context might violate the Bill of Attainder Clause. Gregory Cui, *Evidence-Based Sentencing and the Taint of Dangerousness*, 125 YALE L.J. F. 315 (2016).

<sup>53</sup> Starr, *supra* note 52, at 830. Starr also argues that evidence-based methods do worse in sheer accuracy terms than readily available alternatives such as clinical assessments. *Id.* at 842-62. This is also orthogonal to my analysis here.

proceeds on the basis of different doctrinal predicates. In any case, the earlier article does not explicate carefully both the costs and benefits of algorithmic criminal justice.<sup>54</sup> A more meticulous approach is needed that disaggregates possible technological approaches and normative effects.

The Article unfolds in three steps. Part I defines algorithmic criminal justice and illustrates it by isolating discrete clusters of related instruments now employed in criminal justice or likely soon to be used. I also supply non-technical exposition of the relevant technologies. Part II explores legal criteria of racial equity with special attention to the Equal Protection Clause. It identifies deficiencies in that framework as it applies to algorithmic criminal justice. Part III then turns to the nascent computer-science literature on technical standards of fairness for algorithmic criminal justice. I begin by articulating a normative account of racial equity concerns in criminal justice. I then work through the various metrics identified in the literature to measure racial equity, as well as the tensions between those metrics. Finally, I set forth my own account of racial equity, and explain how it can be operationalized—both in theory and in practice.

## I. Algorithmic Criminal Justice: Scope and Operation

Predictive criminal justice was old when Captain Renault told his men in *Casablanca* to “round up the usual suspects.”<sup>55</sup> The meaningful use of “criminal justice determinations that do not rest simply on probabilities but on statistical correlations between group traits and group criminal-offending rates” can be traced back to the beginning of the twentieth century.<sup>56</sup> The resulting profusion of predictive instruments extends well beyond algorithmic criminal justice instruments to be considered here. For example, an array of evidence-based interventions from interviews to actuarial scoring have long been employed in the sentencing context.<sup>57</sup>

To sharpen the ensuing analysis, I think it is useful to define with some precision a discrete domain of practices as “algorithmic criminal justice.” This Part offers such a

---

<sup>54</sup> Starr notes that “[t]here appears to be a general consensus that using race would be unconstitutional,” *id.* at 812, but this assertion is not based on a comprehensive appreciation of the ways in which racial effects might be embedded in, or emerge from, algorithmic instruments.

<sup>55</sup> *CASABLANCA* (Warner Bros. 1942); *see also* *Jurek v. Texas*, 428 U.S. 262, 275 (1976) (“[P]rediction of future criminal conduct is an essential element in many of the decisions rendered throughout our criminal justice system.”).

<sup>56</sup> BERNARD E. HARCOURT, *AGAINST PREDICTION: PROFILING, POLICING, AND PUNISHING IN AN ACTUARIAL AGE* 18 (2007); *see also* Shima Baradaran, *Race, Prediction, and Discretion*, 81 *GEO. WASH. L. REV.* 157, 176-77 (2013) (“Criminal justice actors often predict which defendants are going to commit an additional crime in determining whether to arrest defendants, to release them on bail, or to release them on parole, or in determining their sentence. This prediction is often based not only on individual evaluation, but also on a group’s criminality and past behavior.”); Richard Berk, *Forecasting Methods in Crime and Justice*, 4 *ANN. REV. LAW AND SOC. SCI.* 219, 221-23 (2008) (setting out the history of formal crime prediction models).

<sup>57</sup> Cecelia Klingele, *The Promises and Perils of Evidence-Based Corrections*, 91 *NOTRE DAME L. REV.* 537, 539 (2015) (discussing “the use of actuarial risk and need assessment instruments, motivational interviewing and counseling techniques, deterrence-based sanction programs, and incentives to probationers and parolees for successful compliance with court orders” with attention to their effects on aggregate incarceration levels).

definition, and then fleshes out that concept with a series of examples from the policing, bail, and post-conviction (parole and probation) contexts. Where salient, I offer capsule accounts of relevant technologies central to my analysis.

### A. A Definition of Algorithmic Criminal Justice

Algorithmic criminal justice, as I define the term, is *the application of an automated protocol to a large volume of data to classify new subjects in terms of the probability of expected criminal activity, in relation to the application of state coercion*. This definition has three elements. Once explicated, those elements provide a justification for treating this domain as a distinct object of legal and normative inquiry.

First, my definition requires an *automated* protocol, or algorithm, that routinizes a decision, here about state coercion.<sup>58</sup> In contrast to such a structured decision-making context, American criminal justice is replete with instances in which officials such as police officers, sentencing judges, parole boards or probation officers exercise partially structured discretion to determine the legality of coercing a particular person. Even where a written protocol is used, as in the sentencing context, substantial residual discretion remains.<sup>59</sup> In a larger domain of cases, though, criminal justice actors act unbounded by either protocol or clear rules. For example, the Fourth Amendment imposes thresholds of “reasonable articulable suspicion” for certain street stops,<sup>60</sup> and “probable cause” for certain arrests.<sup>61</sup> The Supreme Court has resisted efforts to formalize these concepts into “technical”<sup>62</sup> rules, and instead has preferred “practical, common-sense judgment.”<sup>63</sup> Algorithmic criminal justice represents a categorical rejection of such ad hoc, situated judgments as an instrument of regulation.

Second, automation is required because of the sheer *volume* of data used by these tools. Law enforcement agencies increasingly have access to pools of data that are “vast, fast, disparate, and digital.”<sup>64</sup> Colloquially, the instruments at issue here rely on “big data” as that term is used in computational science.<sup>65</sup> In a two-year field study of the Los Angeles

---

<sup>58</sup> THOMAS H. CORMEN ET AL., INTRODUCTION TO ALGORITHMS 10 (2d ed. 2001) (defining an algorithm as “any well-defined computation process that takes some value, or set of values, as input and produces some value, or set of values, as outcome” (emphases omitted)); *see also* Reuben Binns, *Algorithmic Accountability and Public Reason*, PHIL & TECH. 1, 3 (2017) (describing algorithms in terms of whether a system will “take certain inputs and produce certain outputs by computational processes”); MARTIN ERWIG, ONCE UPON AN ALGORITHM: HOW STORIES EXPLAIN COMPUTING 26-27 (2017) (offering an illuminating conceptual account of algorithms).

<sup>59</sup> For an analysis of the scope of discretion in the federal context at present, see Kevin R. Reitz, “*Risk Discretion*” at *Sentencing*, 30 FED. SENT. REP. 68, 68-69 (2017).

<sup>60</sup> *Terry v. Ohio*, 392 U.S. 1, 22 (1968).

<sup>61</sup> *Brinegar v. United States*, 338 U.S. 160, 174 (1949).

<sup>62</sup> *Id.* at 175.

<sup>63</sup> *Illinois v. Gates*, 462 U.S. 213, 244 (1983). The Court has stressed police expertise rather than formal rules. *United States v. Cortez*, 449 U.S. 411, 418 (1981) (“[A] trained officer draws inferences and makes deductions ... that might well elude an untrained person.”).

<sup>64</sup> Sarah Brayne, *Big Data Surveillance: The Case of Policing*, 82 AM. SOC. REV. 977, 980 (2017).

<sup>65</sup> DAWN E. HOLMES, BIG DATA: A VERY SHORT INTRODUCTION 15-16 (2017) (characterizing big data as “huge amounts of data that has not been collected with any specific questions in mind and is often unsorted,” and that is characterized by “volume, variety, and velocity”).

Police Department (“LAPD”), for example, sociologist Sarah Brayne documented how traditional law enforcement databases of persons arrested or convicted of crimes have been supplemented with information about all contacts, of any sort, with police, social services, health services, and child welfare services.<sup>66</sup> This data is integrated with data from “dragnet surveillance tools”; close-circuit television (“CCTV”) cameras used to acquire and track license plate numbers; and “privately collected data.”<sup>67</sup> Because the ensuing massive data pools cannot be sorted by hand, they are only useful because of advances in processing power and computational software. The IC Realtime Company, for example, offers an application called “Ella,” which can recognize and execute natural language queries for CCTV footage.<sup>68</sup> Such changes in the speed and accuracy of queries effect a step change in the quality of surveillance-based evidence available to police.

Third, the algorithmic instruments at issue here make *out-of-sample predictions* about new actors’ likely criminal conduct. It is true that algorithmic instruments can be applied also to extant pools of big data in order to identify historical crimes. For example, the Securities and Exchange Commission analyzes large volumes of trading to identify investors who might be engaged in insider trading.<sup>69</sup> Pattern analysis of this kind can raise questions of racial effects, but in distinct and different ways from out-of-sample prediction methods. The instruments I’m focused on here are generally calibrated using one pool of data, and then applied to new data as a means for identifying or predicting crime that was previously unknown, and that typically has not yet occurred. For example a parole board might have information on historical patterns of reoffending. It supplies that data to a machine leaning tool, which in turn generates a test for forecasting recidivism by suspects yet to interact with the criminal justice system.<sup>70</sup>

So defined, algorithmic criminal justice tools are inductive rather than deductive. They lack opportunities for verification via the collation of other indicia of law-breaking. Algorithmic criminal justice, moreover, claims no insight into the *causes* of crime or criminality.<sup>71</sup> It is just an arrow pointing at crime’s likely next incidence.

## **B. The Operation of Algorithmic Criminal Justice**

We have already seen two instances of algorithmic criminal justice, the Compas algorithm and the SSL list. These examples, though, do not provide a good measure of the scope and effects of algorithmic criminal justice’s operation. New technologies of machine learning (and in particular the subspecies of deep learning) are likely to dominate algorithmic criminal justice in the future. As a result, both Compas and the SSL algorithm

---

<sup>66</sup> Brayne, *supra* note 64, at 995.

<sup>67</sup> *Id.* at 993-94.

<sup>68</sup> James Vincent, *Artificial Intelligence is Going to Supercharge Surveillance*, THE VERGE, Jan, 23, 2018.

<sup>69</sup> Mary Jo White, Chair, SEC, *Keynote Address at the 41st Annual Securities Regulation Institute*, Jan. 27, 2014, <http://www.sec.gov/News/Speech/Detail/Speech/1370540677500> [<https://perma.cc/M7YV-33PR>] (describing the SEC’s NEAT program, which can identify and analyze insider trading activity around times of major corporate events).

<sup>70</sup> Richard Berk, *An impact assessment of machine learning risk forecasts on parole board decisions and recidivism*, 13 J. EXP. CRIM. 193. 195 (2017).

<sup>71</sup> *Cf.* Usama Fayyad, *The Digital Physics of Data Mining*, 44 COMM. ACM, Mar. 2001, at 62.



are likely soon to be relics. Newer tools will combine powerful computational instruments with large volumes of data to enable prediction of a kind that is qualitatively distinct from historical antecedents.<sup>72</sup> A survey of the potential uses of these new instruments reveals a fuller sense of the scope and effects that algorithmic prediction tools will have on criminal justice.

### 1. *Machine Learning and Deep Learning*

A machine-learning algorithm solves a “learning problem ... of improving some measure of performance when executing some task through some type of training experience.”<sup>73</sup> The basic task a supervised machine-learning algorithm must perform can be framed as follows: The algorithm is prompted to define a function  $f(x)$  which produces an output  $y$  for any given input  $x$ . In other words, it classifies  $x$  in terms of  $y$ .<sup>74</sup> Its outputs take the form of a sorting of  $x$  onto categories of  $y$ ,<sup>75</sup> although the resulting classifications are correlational rather than causal in nature.<sup>76</sup> Finally, its performance is measured in terms of how well it captures the relation of  $x$  to  $y$ .<sup>77</sup>

To begin with, a supervised machine-learning algorithm is assigned a set of “training” data already labeled in terms of  $y$  so it can develop a model, represented by the mathematical function  $f(x)$ , that best represents the relationship between features of each observation in the training data and the known classification  $y$ . This function  $f(x)$  is then applied to a new “test set” of data.<sup>78</sup> The algorithm predicts how to classify this new data by applying  $f(x)$  to generate predictions of  $y$ .<sup>79</sup> Such supervised tools are but one kind of machine learning. There is also a species of *unsupervised* machine-learning algorithm. These begin with unlabeled training data, and tend to be tasked with the development of classifications based on the data’s immanent structure.<sup>80</sup>

---

<sup>72</sup> See JERRY KAPLAN, ARTIFICIAL INTELLIGENCE: WHAT EVERYONE NEEDS TO KNOW 39 (2016) (pointing to “improvements in computer speed and memory, the transition from physically stored data to electronically stored data, easier access (mainly due to the Internet), and low-cost high-resolution digital sensors” as the technological predicates of machine learning).

<sup>73</sup> Jordan & Mitchell, *supra* note 33, at 255.

<sup>74</sup> *Id.* This process can also be described in terms of a “classifier” rather than a function, that examines inputs with “feature values” and outputs a class variable. Pedro Domingos, *A Few Useful Things to Know About Machine Learning*, COMM. ACM, Oct. 2012, at 78-80 (“A classifier is a system that inputs (typically) a vector of discrete and/or continuous feature values and outputs a single discrete value, the class.”).

<sup>75</sup> PETER FLACH, MACHINE LEARNING: THE ART AND SCIENCE OF ALGORITHMS THAT MAKE SENSE OF DATA 14 (2012) (noting that “multi-class classification” is “a machine learning task in its own right”).

<sup>76</sup> Consider in this regard recommendation algorithms employed by consumer-facing companies such as Amazon and Netflix. *Cf.* KAPLAN, *supra* note 72, at 32 (arguing that machine learning algorithms operate like “incredibly skilled mimics, finding correlations and responding to novel inputs as if to say, ‘This reminds me of ....’”).

<sup>77</sup> Jordan & Mitchell, *supra* note 33, at 255-57 (noting that performance can be defined in terms of accuracy, with false positive and false negative rates being assigned a variety of weights).

<sup>78</sup> HOLMES, *supra* note 65, at 24 (discussing classification and distinguishing training and test sets of data); ALPAYDIN, *supra* note 31, at 40 (describing the use of training and validation data).

<sup>79</sup> STUART RUSSELL & PETER NORVIG, ARTIFICIAL INTELLIGENCE: A MODERN APPROACH 693 (3d ed. 2010).

<sup>80</sup> FLACH, *supra* note 75, at 14-15.

No machine-learning algorithm is given ex ante a functional form  $f(x)$  that defines the relationship between observations and classifications. Rather, the algorithm employs one of a wide number of procedures to ascertain  $f(x)$  through a process called “feature selection.”<sup>81</sup> The latter include decision trees, decision forests, logistic regression, support vector machines, neural networks, kernel machines, and Bayesian classifiers.<sup>82</sup> By sorting through many different possible  $f(x)$ s on the basis of its training data using one of these methods, the algorithm homes in upon an  $f(x)$  that optimizes the accuracy of its performance metric. Many people experience this in the ‘learning’ they experience on the part of Siri, Alexa, or other ‘assistants.’<sup>83</sup>

Deep learning is a subset of machine learning whereby the algorithm is made up of “multiple layers of representation,” each of which transform the raw data into a slightly more abstract form.<sup>84</sup> Given enough layers of transformation, the algorithm can perform very complex functions, such as playing the Chinese game Go or recognizing specific images from representational input. What distinguishes deep learning is the fact that its “layers of features are not designed by human engineers: they are learned from data using a general-purpose learning procedure.”<sup>85</sup> The most well-known forms of deep learning tools are based on “neural networks,” inspired by patterns observed in human brain.<sup>86</sup> Deep-learning instruments are especially apt for unsupervised tasks, with no specification of features, and little “manual interference,” such that designers “just wait and let the learning algorithm discover all that is necessary by itself.”<sup>87</sup> The utility to police of an instrument that can extract speech or visual patterns from large quantities of audio-visual inputs (e.g., CCTV footage, cellphone call content) is self-evident.<sup>88</sup>

---

<sup>81</sup> Kroll et al., *supra* note 38, at 681 (describing feature selection as concerning the “choices about which data models should consider”); *see also* David Lehr & Paul Ohm, *Playing with the Data: What Legal Scholars Should Learn About Machine Learning*, 51 U.C. DAVIS L. REV. 653, 700 (2017) (describing feature selection); Avrim L. Blum, & Pat Langley, *Selection of relevant features and examples in machine learning*, 97 ARTIFICIAL INTELLIGENCE 245, 250-53 (1997) (decomposing feature selection into a nested sequence of analytic tasks).

<sup>82</sup> David J. Hand, *Classifier Technology and the Illusion of Progress*, 21 STAT. SCI. 1, 1, 4 (2006) (documenting these instruments, and contending that in “real world” conditions, simpler instruments often perform better).

<sup>83</sup> “9 Applications of Machine Learning From Day-to-Day Life,” MEDIUM, July 30, 2017, <https://medium.com/app-affairs/9-applications-of-machine-learning-from-day-to-day-life-112a47a429d0/>.

<sup>84</sup> Le Cuan et al., *supra* note 34, at 436; *id.* at 438 (“A deep learning architecture is a multilayer stack of simple modules, all (or most) of which are subject to learning, and many of which compute non-linear input-output mappings.”); ALPAYDIN, *supra* note 31, at 85-109 (describing neural networks). For a non-technical account of back propagation, the key element of deep learning, see James Somers, *Is AI Riding a One-Trick Pony?*, 120(8) MIT TECH. REV. 29, 31 (2017).

<sup>85</sup> Le Cuan et al., *supra* note 34, at 436.

<sup>86</sup> Jurgen Schmidhuber, *Deep learning in neural networks: An overview*, 61 NEURAL NETWORKS 85, 86-87 (2015).

<sup>87</sup> ALPAYDIN, *supra* note 31, at 107-08.

<sup>88</sup> Maryam M. Najafabadi, et al., *Deep learning applications and challenges in big data analytics*, 1 J. BIG DATA 1, 11-14 (2015) (describing uses of deep learning tools). Deep learning has also been used to play “games of perfect information” such as chess and Go. David Silver et al., *Mastering the game of Go with deep neural networks and tree search*, 529 NATURE 484, 490 (2016).

## 2. *The Impact of Machine Learning on Criminal Justice*

Adoption of machine learning within the criminal justice system changes the scale, reach, and operation of state power. Consider each factor in turn.

First, these tools dramatically inflate the state's ability to acquire otherwise inaccessible information.<sup>89</sup> For instance, police in London and in South Wales now track individuals' locations and movements over days and weeks by applying machine learning tools to thousands of hours of CCTV footage.<sup>90</sup> Machine-learning tools also facilitate predictions that would be far more imprecise if based solely upon more familiar regression analyses.<sup>91</sup>

Second, machine learning instruments sever the connection between the human operator and the function  $f(x)$  used to solve the classification problem. Unstructured human discretion, which once infused the criminal justice system, is displaced by an algorithmically structured logic that is not the function of any human hand. As a result, it will often not be possible to speak of the 'intent' or the 'anticipated' consequences of a classification protocol. Rather, the algorithm will "sift through vast numbers of variables, looking for combinations that reliably predict outcomes," handling "enormous numbers of predictors—sometimes, remarkably, more predictors than observations—combining them in nonlinear and highly interactive ways,"<sup>92</sup> and hence generating utterly unexpected outcomes.

To the extent that the design of a machine learning process involves the intentional crafting and selection of training data, feature sets, or the like, moreover, there will often be no way to ascertain the role of designers' racial sentiments (if any) directly,<sup>93</sup> and no easy way to infer intentionality indirectly from the instrument's results.<sup>94</sup> There is no such thing as code that bespeaks racial animus. Design choices that might be molded by racial animus also cannot be reverse engineered to cast light on background human motivations. And it is difficult to know how to disentangle the effect of background differences in criminality and bad designer intent when evaluating the outputs of an algorithm. As a

---

<sup>89</sup> For recognition of this general point, see *United States v. Garcia*, 474 F.3d 994, 998 (7th Cir. 2007) (Posner, J.) ("Technological progress poses a threat to privacy by enabling an extent of surveillance that in earlier times would have been prohibitively expensive," thereby "giving the police access to surveillance techniques that are ever less expensive and ever more effective.").

<sup>90</sup> David Bond, *CCTV watchdog warns UK police over use of facial recognition*, FIN. TIMES, Oct. 29, 2017, <https://www.ft.com/content/ab60f9f2-bb26-11e7-8c12-5661783e5589>.

<sup>91</sup> Jon Kleinberg et al., *Prediction policy problems*, 105 AM. ECON. REV. 491, 493-94 (2015).

<sup>92</sup> Ziad Obermeyer, & Ezekiel J. Emanuel, *Predicting the future—big data, machine learning, and clinical medicine*, 13 NEW ENGLAND J. MED. 1216, 1217 (2016).

<sup>93</sup> Barocas & Selbst, *supra* note 50, at 710 ("The idea that the representation of different social groups can be brought into proportions that better match those in the real world presumes that analysts have some independent mechanism for determining these proportions.").

<sup>94</sup> Tal Z. Zarsky, *Transparent Predictions*, 2013 U. ILL. L. REV. 1503, 1519-20 (noting how predictions can be generated in processes "which [are] not explainable in human language," such that "[i]t would be difficult for the government to provide a detailed response when asked why an individual was singled out to receive differentiated treatment by an automated recommendation system").

result, the effects of, and evidence for, human intentions—a central term of legal and constitutional analysis—are likely to be elusive in practice.

Third, algorithmic tools can be as sticky or stickier than the forms of human discretion. Hence, whereas it is always a possibility that human agents will observe the unintended effects of human action, machine decision-making can be opaque and hence resistant to change. Algorithmic systems can thus be “stuck in time until engineers dive in to change them.”<sup>95</sup> Indeed, it will often not be clear to a human operator that an algorithmic criminal justice tool needs reconsideration. That human operator necessarily sees only a limited and unrepresentative tranche of case outcomes. She must also grapple with the sheer technological complexity of algorithmic tools. Hence, algorithmic errors are often liable to prove more durable than human errors.

Fourth, the consequences of switching between unstructured human discretion and algorithmically structured prediction can often be unexpected. This happens even when a semi-structured instrument is altered. For example, in the wake of the Supreme Court’s decision eliminating the mandatory character of the Federal Sentencing Guidelines, studies found “significantly increased racial disparities after controlling for extensive offender and crime characteristics.”<sup>96</sup>

Fifth, the emerging crop of algorithmic tools are potentially very different from risk assessment tools employed currently in bail and sentencing. Current instruments rely on a relatively small number of variables—two leading models use 16 and 20 parameters respectively—and fixed classification rules to generate recidivism risks.<sup>97</sup> These instruments focus solely on recidivism risk. They make no effort to estimate either the direct or the remote costs of coercive action. In contrast, tools such as Compas include recommended cut-off points that at least imply an evaluation of social costs. There is no reason, moreover, that an algorithm could not be trained with data that reflected both the costs and the benefits of coercive action, although this does not yet appear to be standard practice.

To summarize, the operation and the effects of predictions offered by algorithmic criminal justice are qualitatively distinct from the unstructured and semi-structured forms of human discretion that have until now dominated the criminal justice system. Not all such tools use machine learning or deep learning. But it is only a question of time before these powerful instruments crowd out more simple models. Indeed, it is striking that both the Compas algorithm and the SSL instrument described in the introduction have been criticized on the basis of their weak predictive power.<sup>98</sup> A likely, if not inevitable,

---

<sup>95</sup> O’NEILL, *supra* note 51, at 204.

<sup>96</sup> Crystal S. Yang, *Free at Last? Judicial Discretion and Racial Disparities in Federal Sentencing*, 44 J. LEGAL STUD. 75, 77 (2015); accord Max Schanzenbach & Emerson H. Tiller, *Strategic Judging under the United States Sentencing Guidelines: Positive Political Theory and Evidence*, 23 J. L. & ECON. & ORG. 23 (2007) (similar finding).

<sup>97</sup> Slobogin, *supra* note 46, at 584-85 (explaining the OxRec and VRAG assessment tools); see also Lauryn P. Gouldin, *Disentangling Flight Risk from Dangerousness*, 2016 BYU L. REV. 837, 869-70 (2016) (describing the PSA tool).

<sup>98</sup> On the Compas algorithm: Julia Dressel and Hany Farid, *The accuracy, fairness, and limits of predicting recidivism*, SCIENCE ADVANCES, Jan. 17, 2018, <http://advances.sciencemag.org/content/4/1/eaao5580>

consequence of such critiques is the adoption of new, more powerful computational tools to achieve the same end. In any event, a phase shift in the quality of criminal-justice action can already be observed across the spectrum of criminal justice functionalities. Even if machine-learning and deep-learning tools are not now omnipresent, they are likely to be so soon.<sup>99</sup>

### C. Algorithmic Criminal Justice on the Ground

Algorithmic tools are used now in three main criminal-justice contexts: policing, bail decisions, and post-conviction matters. This section provides a capsule summary of the ways in which predictive instruments are operationalized across those three distinct domains.

#### 1 Policing

In the policing context, algorithmic tools are employed to make predictions about both places and people.<sup>100</sup> Place-focused tools aggregate “real-time” information on the frequency and geographic location of crimes to “determine staffing needs or allocate resources” as between different regions.<sup>101</sup> Consonant with a focus on the location of crime, police departments across the country have increasingly adopted the Compstat, or Crime Control Strategy Meeting, structure first developed in New York. Under Compstat, precinct commanders are subject to biweekly questioning by senior departmental leadership in a “data-saturated environment” about how they are responding to crime trends.<sup>102</sup> While Compstat itself does not necessarily incorporate algorithmic tools, its focus on data-driven predictions of crime’s geographic dispersion invites the use of algorithmic tools. Further, a number of criminologists have identified promise in a place-based prediction approach involving the “the application of police interventions at very

---

(finding that the Compas algorithm performs no better than people with no experience of the criminal justice system in making recidivism predictions). On the SSL: Jessica Saunders et al., *Predictions put into practice: A quasi-experimental evaluation of Chicago’s predictive policing pilot*, 12 J. EXP. CRIMINOLOGY 347, 363 (2016) (finding that “while using arrestee social networks improved the identification of future homicide victims, the number was still too low in the pilot to make a meaningful impact on crime).

<sup>99</sup> One reason for this, of course, is the promotion of algorithmic implements by the companies that manufacture them, and stand to gain financially from their adoption. Elizabeth E. Joh, *The Undue Influence of Surveillance Technology Companies on Policing*, N.Y.U. L. REV. 101, 114-120 (2017) (describing mechanisms of private influence on public adoption of computational technologies in the criminal justice sector).

<sup>100</sup> For overviews, see Walter L. Perry et al., Rand Corp., *Predictive Policing: Forecasting Crime for Law Enforcement 2* (2013), [https://www.rand.org/pubs/research\\_briefs/RB9735.html](https://www.rand.org/pubs/research_briefs/RB9735.html); see also Jennifer Bachner, *Predictive Policing: Preventing Crime with Data and Analytics* 14 (2013), available at <http://www.businessofgovernment.org/sites/default/files/Predictive%20Policing.pdf> (“The fundamental notion underlying the theory and practice of predictive policing is that we can make probabilistic inferences about future criminal activity based on existing data.”).

<sup>101</sup> Andrew Guthrie Ferguson, *Crime Mapping and the Fourth Amendment: Redrawing “High-Crime Areas,”* 63 HASTINGS L.J. 179, 182 (2011).

<sup>102</sup> James J. Willis et al., *Making Sense of COMPSTAT: A Theory-Based Analysis of Organizational Change in Three Police Departments*, 41 LAW & SOC’Y REV. 147, 147 (2007); see also David L. Carter & Jeremy G. Carter, *Intelligence-Led Policing: Conceptual and Functional Considerations for Public Policy*, 20 CRIM. JUST. POL’Y REV. 310, 316-19 (2009) (describing Compstat’s operation).

small geographic units of analysis,” or hot spots.<sup>103</sup> A number of randomized, controlled experiences have found evidence that such place-focused tools are effective in suppressing crime.<sup>104</sup>

Consistent with these developments, influential jurisdictions have adopted machine-learning tools to facilitate place-based policing.<sup>105</sup> One of the earliest adaptors, starting in 2015, was the New York Police Department. This force embarked on a two-year pilot program using HunchLab, an algorithm developed by the Philadelphia-based Azavea company.<sup>106</sup> According to Azavea’s web site, HunchLab’s “ensemble machine learning” algorithm uses “temporal patterns (day of week, seasonality); weather; risk terrain modeling (locations of bars, bus stops, etc.); socioeconomic indicators; historic crime levels; and near repeat patterns” as a means of “predicting individual crime expectations across the jurisdiction.”<sup>107</sup> Other departments, such as Los Angeles, have adopted a system created by the PredPol company. PredPol produces a propriety algorithm based on a “near-repeat” machine-learning model. This assumes that if a crime occurs at a given location, the immediate surroundings are at increased risk for future crime.<sup>108</sup> The PredPol model is an extrapolation first developed by anthropologist Jeffrey Brantingham and mathematician Andrea Bertozzi, of an algorithm used to predict the distribution of earthquake shocks.<sup>109</sup> One randomized, controlled study of the use of a machine-learning tool derived from models of epidemic aftershocks to implement hot-spot policing found that the instrument predicted crime well, and led to a 7.4 percent reduction in crime volume as a function of patrol time.<sup>110</sup>

In the last five years, however, place-focused tools have started to be supplemented with person-focused tools. Chicago, for example, started to build a database of alleged gang

---

<sup>103</sup> ANTHONY A. BRAGA & DAVID L. WEISBURD, *POLICING PROBLEM PLACES* 9 (2010). More generally, proactive policing of various kinds (not necessarily involving stops) is also associated with crime-control effects. Charis E. Kubrin et al., *Proactive Policing and Robbery Rates Across U.S. Cities*, 48 *CRIMINOLOGY* 57, 62 (2010).

<sup>104</sup> See Anthony A. Braga & Brenda J. Bond, *Policing Crime and Disorder Hot Spots: A Randomized Controlled Trial*, 46 *CRIMINOLOGY* 577 (2008); Anthony Braga et al., *Problem-Oriented Policing in Violent Crime Places: A Randomized Control Experiment*, 37 *Criminology* 541 (1999).

<sup>105</sup> See, e.g., J. Brian Charles, *How Police are Using Tech to Fight Crime*, *GOVERNING*, Apr. 11, 2018, <http://www.governing.com/topics/public-justice-safety/gov-gang-violence-predictive-policing-high-point-ic.html?r> (describing the use of ONESolution predictive software by the High Point, NC, police department).

<sup>106</sup> Laura Nahmias and Miranda Newbauer, *NYPD testing crime forecast software*, *POLITICO*, July 8, 2015, <https://www.politico.com/states/new-york/city-hall/story/2015/07/nypd-testing-crime-forecast-software-090820>.

<sup>107</sup> Hunchlab, *Hunchlab under the Hood* 12 (2015), <https://cdn.azavea.com/pdfs/hunchlab/HunchLab-Under-the-Hood.pdf>.

<sup>108</sup> Brayne, *supra* note 64, at 989; see generally PredPol, *How Predictive Policing Works*, *PREPPOL.COM* (2015), <http://www.predpol.com/how-predictive-policing-works> (providing a predictably rosy overview of the algorithm’s uses).

<sup>109</sup> Aaron Shapiro, *Reform Predictive Policing*, 541 *NATURE* 458, 459 (2017); see also Elizabeth E. Joh, *Policing by Numbers: Big Data and the Fourth Amendment*, 89 *WASH. L. REV.* 35, 44 (2014) (describing PredPol’s use in Santa Cruz, California).

<sup>110</sup> George O. Mohler et al., *Randomized controlled field trials of predictive policing*, 110 *J. AM. STAT. ASS’N* 1399, 1407 (2015).

members in order to draw inferences about their propensity to commit violent crimes.<sup>111</sup> That city's SSL predicts the likelihood of an individual becoming a victim of a homicide using an analysis of that person's known social network, and in particular by counting the number of first degree co-arrest links and the number of second-degree co-arrest links with previous homicide victims.<sup>112</sup> Names generated by the SSL algorithm were disseminated to district commanders, who had discretion about what interventions to apply.<sup>113</sup> The algorithm, however, identified only one percent of the pool of eventual homicide victims, and yielded no identifiable crime-control gains.<sup>114</sup>

A related use of deep-learning tools involves facial recognition algorithms that can search for dangerous persons in a specific place at a particular time. This emerging use is not a matter of out-of-sample prediction. It is a matching exercise based on new data, and so falls at the periphery of my analysis. For instance, the Metropolitan Police of London combine dense CCTV with facial recognition instruments in monitoring public events that are thought to be attractive targets for terrorist attacks.<sup>115</sup> In June 2017, the deployment of facial recognition algorithms to real-time CCTV inputs generated a first arrest for British police.<sup>116</sup>

The situation in the United States is less clear. As of 2016, at least five metropolitan police departments in the United States—including Chicago's, Dallas's, and Los Angeles'—claimed to use a facial recognition algorithm to comb public CCTV data.<sup>117</sup> Facial images have been made available by the Federal Bureau of Investigation since 2011.<sup>118</sup> In 2017, Orlando, FL, and Washington County, OR, became some of the first governmental purchasers of Amazon's "Rekognition" tool, which uses "artificial intelligence" to scan and identify up to a hundred faces in a single CCTV shot.<sup>119</sup>

---

<sup>111</sup> John Buntin, *Social Media Transforms the Way Chicago Fights Gang Violence*, GOVERNING, Oct. 2013, at 26, 28 (describing how data was acquired for social media analysis).

<sup>112</sup> Saunders et al., *supra* note 98, at 354.

<sup>113</sup> *Id.*; see also Mark Guarnio, *Can Math Stop Murder?*, CHRISTIAN SCI. MONITOR (July 20, 2014), <http://www.csmonitor.com/USA/2014/0720/Can-math-stop-murder-video>, archived at <http://perma.cc/G3TA-9SPT> (discussing predictive policing techniques in Chicago including sending officers to the houses of suspected gang leaders).

<sup>114</sup> Saunders et al., *supra* note 98, at 363; *id.* at 365 (noting that those included on the SSL list were in fact less likely to be a victim of a shooting, although that difference was not statistically significant).

<sup>115</sup> Mark Townsend, *Police to Use Facial Recognition Cameras at Cenotaph Service*, THE GUARDIAN, Nov. 12, 2017, <https://www.theguardian.com/technology/2017/nov/12/metropolitan-police-to-use-facial-recognition-technology-remembrance-sunday-cenotaph>.

<sup>116</sup> Cara McGoogan, *British police arrest suspect spotted with facial recognition technology*, TELEGRAPH, June 7, 2017, <http://www.telegraph.co.uk/technology/2017/06/07/british-police-arrest-suspect-spotted-facial-recognition-technology/>.

<sup>117</sup> Claire Garvie et al., *The Perpetual Line-Up: Unregulated Face Recognition in America*, Oct. 16, 2016, <https://www.perpetuallineup.org/>.

<sup>118</sup> U.S. Gov't Accountability Office, GAO-16-267, *Face Recognition Technology: FBI Should Better Ensure Privacy and Accuracy* 7, 15 (May 2016)

<sup>119</sup> Matt Cagle & Nicole Ozer, *Amazon Teams Up With Government to Deploy Dangerous New Facial Recognition Technology*, ACLU, May 18, 2018, <https://www.aclu.org/blog/privacy-technology/surveillance-technologies/amazon-teams-government-deploy-dangerous-new?redirect=blog/privacy-technology/surveillance-technologies/amazon-teams-law-enforcement-deploy-dangerous-new>.

Nevertheless, real-time application of facial recognition technologies to CCTV data still appears rare, in particular because of technological barriers. It is telling that between June and September 2017, the National Institute for Science and Technology ran a prize challenge for facial recognition technology. The winner of the contest, NTechLab, created an algorithm with a rate of 0.22 false non-matches for every 0.0001 false match.<sup>120</sup> And, as noted above, in 2018, the IC Realtime company introduced a commercially available algorithm called Ella that can recognize and respond to natural language queries to search large quantities of video footage for specific images.<sup>121</sup>

## 2. *Bail*

The second use of algorithmic tools is in the pre-trial context of arraignment hearings in which judges determine whether defendants are to be detained pending criminal trial, or released having posted a money bail or otherwise. Pretrial detainees comprise roughly 60 percent of the jail population, and between 2005 and 2013 some 450,000 people were incarcerated awaiting trial on any given day.<sup>122</sup> Pre-trial detention decisions impose considerable costs on individuals, in relation to employment, health outcomes, and childcare costs.<sup>123</sup> One study, for example, estimates a lower-bound net cost of detention for the marginal individual of \$55,385 and an upper-bound net cost of \$101,223.<sup>124</sup> At the same time, “[r]elatively little is known about the charge characteristics and case dispositions” for that pretrial detention population.<sup>125</sup> But studies in a range of jurisdictions find evidence of racial disparities in bail decisions; predictably, black and Latino defendants receive systematically less favorable treatment.<sup>126</sup>

Much of the impetus of recent bail reform has hinged on the much-criticized effect of wealth upon access to pretrial release.<sup>127</sup> Algorithmic criminal justice does not

---

<sup>120</sup> Patrick Grother et al., *The 2017 IARPA Face Recognition Prize Challenge (FRCP)*, Nov. 2017, at 2, <http://nvlpubs.nist.gov/nistpubs/ir/2017/NIST.IR.8197.pdf>.

<sup>121</sup> James Vincent, *Artificial Intelligence is Going to Supercharge Surveillance*, THE VERGE, Jan. 23, 2018, <https://www.theverge.com/2018/1/23/16907238/artificial-intelligence-surveillance-cameras-security>.

<sup>122</sup> Jaeok Kim et al., *Unpacking pretrial detention: An examination of patterns and predictors of readmissions*, 29 CRIM. J. POL. REV. 663, 664 (2018); Roy Walmsley, Int'l Ctr. for Prison Studies, *World Pre-trial/Remand Imprisonment List 1* (2d ed. 2014), [http://www.prisonstudies.org/sites/default/files/resources/downloads/world\\_pre-trial\\_imprisonment\\_list\\_2nd\\_edition\\_1.pdf](http://www.prisonstudies.org/sites/default/files/resources/downloads/world_pre-trial_imprisonment_list_2nd_edition_1.pdf).

<sup>123</sup> For a discussion of the costs of pretrial bail, see Sandra G. Mayson, *Dangerous Defendants*, 127 YALE L.J. 490, 547 (2018).

<sup>124</sup> Crystal S. Yang, *Toward an Optimal Bail System*, 92 N.Y.U. L. REV. 1399, 1436 (2017).

<sup>125</sup> Kim, *supra* note 122, at 667.

<sup>126</sup> Yang, *supra* note 124, 1466–67 (finding “compelling evidence that bail judges in these jurisdictions treat defendants of different races differently in setting bail”); Traci Schlesinger, *Racial and Ethnic Disparity in Pretrial Criminal Processing*, 22 JUST. Q. 170, 187 (2005) (same result for hold rates within a county using the nationally representative State Court Processing Statistics); accord Stephen Demuth & Darrell Steffensmeier, *The Impact of Gender and Race-Ethnicity in the Pretrial Release Process*, 51 Soc. Probs. 222, 222 (2004);. *But see* Frank McIntyre & Shima Baradaran, *Race, Prediction, and Pretrial Detention*, 10 J. EMPIRICAL LEGAL STUD. 741, 769 (2013) (rejecting any finding of racial disparities after having controlled for the probability of re-arrest).

<sup>127</sup> Nick Pinto, *The Bail Trap*, N.Y. TIMES MAG. (Aug. 13, 2015), [http://www.nytimes.com/2015/08/16/magazine/the-bail-trap.html?\\_r=1](http://www.nytimes.com/2015/08/16/magazine/the-bail-trap.html?_r=1).



necessarily respond to this problem, except to the extent it enables a reduction of pretrial detention without imposing any cost on crime-related outcomes.<sup>128</sup> Rather, such tools are an obvious fit in a context where magistrates are forced to make predictive decisions about the risk of violence, criminality, or flight on the basis of relatively cursory information. Already, two simple algorithms, the Public Safety Assessment (PSA) and Canadian Level of Service Inventory Revised (LSI-R) use information from criminal history to personality patterns and age to offer recidivism predictions.<sup>129</sup> The latter instrument, however, is administered by professionals through interviews, and involves no computational element.<sup>130</sup> More sophisticated algorithmic instruments are now starting to be introduced into courtrooms to inform bond determinations in jurisdictions across the country.<sup>131</sup>

Numerous jurisdictions give judges access to the Compas system in the pre-trial arraignment context.<sup>132</sup> But there is a surprising paucity of public information about the manner of its implementation, or its effects on the rates of pretrial release or the composition of the pretrial detainee population. Two studies, one conducted in Philadelphia and the other in an unnamed large American city, compared the performance of different machine-learning algorithms with that of the existing bench. Both find that the computational method generated less mis-ranking of criminal defendants and less crime.<sup>133</sup> These studies, however, focus narrowly on the important question of gains to public safety that would result from a move from human to machine prediction. They appear to assume that jurisdictions will respond to algorithmic criminal justice instruments by using less pretrial incarceration to obtain the same levels of deterrence. It is not clear, though, why this assumption is warranted. Most (but not all) of the studies are silent as to the possibility or magnitude of racial effects, a striking omission given the large empirical literature documenting racial disparities in bail decisions.<sup>134</sup>

---

<sup>128</sup> As New Orleans has. Aviva Sen and CityLab, *New Orleans' Great Bail Reform Experiment*, THE ATLANTIC, Oct 19, 2017, <https://www.theatlantic.com/politics/archive/2017/10/new-orleans-great-bail-reform-experiment/544964/> (finding that both the pretrial detention and the pretrial crime rate had fallen using the tool).

<sup>129</sup> EPIC, *Algorithms in the Criminal Justice System*, <https://epic.org/algorithmic-transparency/crim-justice/>.

<sup>130</sup> Alexander M. Holsinger, *Implementation of actuarial risk/need assessment and its effect on community supervision revocations*, 15 J. RES. & POL. 95, 98-99 (2013).

<sup>131</sup> Ellora Thadaneey Israni, *When an Algorithm Helps Send You to Prison*, N.Y. TIMES, Oct. 26, 2017, <https://www.nytimes.com/2017/10/26/opinion/algorithm-compas-sentencing-bias.html>; see also Richard Berk & Jonathan Hyatt, *Machine Learning Forecasts of Risk to Inform Sentencing Discretion*, 27 FED. SENT. 222, 223 (2015) (explaining advantages of machine-learning tools over the LSI-R); Richard F. Lowden, *Risk Assessment Algorithms: The Answer to an Inequitable Bail System*, 19 N.C. J. L. & TECH. 221, 230 (2018) (listing jurisdictions that have adopted algorithmic tools).

<sup>132</sup> Angwin et al., *supra* note 15. It is not wholly clear how much weight judges give to the Compas scores, or whether there is even a uniform practice.

<sup>133</sup> Jon Kleinberg et al., *Human decisions and machine predictions*, 133 Q. J. ECON, 237, 237-38 (2017) [hereinafter “Kleinberg et al., *Human decision*”] (finding that large decreases in offending rates could be achieved by moving from judicial to machine predictions in the bail context for violent crimes); Richard A. Berk, Susan B. Sorenson, and Geoffrey Barnes, *Forecasting domestic violence: A machine learning approach to help inform arraignment decisions*, 13 J. EMP. L. STUD. 94, 110 (2016) (finding that the release rate of 20 percent repeat offenders in a pool of domestic violent defendants could be dropped to a 10 percent rate through a move from judicial to machine-led determinations).

<sup>134</sup> See sources cited in *supra* note 126. One study to address potential racial effects is Kleinberg et al., *Human decision*, *supra* note 133, at 237.

### 3. Sentencing

A recent survey of state sentencing practice comments that it is “improbable” that any convicted felon, whether an adult or juvenile, would be sentenced today without the aid of some sort of actuarial risk instrument, albeit not necessarily one that employs algorithmic means.<sup>135</sup> In some jurisdictions, such as Pennsylvania, New Hampshire, Arkansas, and Vermont, state law even affirmatively mandates the use of predictive instruments in the sentencing phrase.<sup>136</sup> Such instruments have emerged as part of a “full-on embrace of practices that promise to reduce the risk of reoffending by convicted persons.”<sup>137</sup>

In 2015, more than 60 risk-assessment tools were used in sentencing contexts.<sup>138</sup> These risk assessments typically evaluated where within a statutorily calibrated sentencing range an offender’s sentence should lie accounting for “utilitarian crime-control grounds.”<sup>139</sup> Some jurisdictions, such as Virginia, use a noncomputational “actuarial” instrument calibrated by age, felony record, offense type, employment, and gender, to sort nonviolent, low-risk offenders to alternative punishments such as probation, jail time, and restitution.<sup>140</sup> In other jurisdictions, computational instruments such as Compas are used for that same purpose.<sup>141</sup> Such instruments are only recently attracting judicial attention, including a high-profile constitutional challenge to Wisconsin’s algorithm.<sup>142</sup>

Finally, I have been able to locate only one well-detailed example of a machine-learning algorithm being employed in the parole context. In 2010, the Pennsylvania Board of Probation and Parole started developing a machine-learning protocol using random

---

<sup>135</sup> Zachary Hamilton et al., *Designed to fit: The development and validation of the STRONG-R recidivism risk assessment*, 43 CRIM. J. & BEHAV. 230, 231 (2016).

<sup>136</sup> 42 PA. CONSOL. STAT. § 2154.7 (2009) (mandating the creation of a “risk assessment instrument” for determining “the relative risk that an offender will reoffend and be a threat to public safety”); ARK. CODE ANN. § 16-93-615 (a)(1)(B) (2015) (“The determination... shall be made by reviewing information such as the result of the risk-needs assessment to inform the decision of whether to release a person on parole by quantifying that person’s risk to reoffend, and if parole is granted, this information shall be used to set conditions for supervision.”); N.H. REV. STAT. ANN. § 504-A:15(I) (2011) (requiring that “[e]very person placed on probation or parole... be assessed by the department of corrections, using a valid and objective risk assessment tool, to determine that person’s risk of recidivating” and that the results be used to determine the length of active supervision); Vt. STAT. ANN. tit. 13, § 7554c(a)(1) (2015) (“The objective of a pretrial risk assessment is to provide information to the Court for the purpose of determining whether a person presents a risk of nonappearance or a threat to public safety so the Court can make an appropriate order concerning bail and conditions of pretrial release.”).

<sup>137</sup> Klingele, *supra* note 57, at 551-52.

<sup>138</sup> Anna Maria Barry-Jester et al., *Should Prison Sentences Be Based on Crimes That Haven’t Been Committed Yet?*, FIVETHIRTYEIGHT POLITICS, Aug. 4, 2015, <http://fivethirtyeight.com/features/prison-reform-risk-assessment/>.

<sup>139</sup> John Monahan & Jennifer Skeem, *Risk Assessment in Criminal Sentencing*, 12 ANN. REV. CLIN. PSYCHOL. 489, 493-94 (2016).

<sup>140</sup> *Id.* at 495.

<sup>141</sup> Angwin et al., *supra* note 15.

<sup>142</sup> *State v. Loomis*, 881 N.W.2d 749, 761-62 (2016) (upholding the use of the Compas tool in sentencing in Wisconsin); *see infra* text accompanying notes 166 to 167 (discussing *Loomis*).

forests to generate forecasts of recidivism to assist members of the Board in making discrete parole decisions.<sup>143</sup> When subject to performance evaluation seven years later, the algorithm was found to have reduced re-arrests for both non-violent and violent crime.<sup>144</sup>

#### D. The Emerging Evidence of Race Effects

The interaction of criminal justice and race looms large in legal scholarship.<sup>145</sup> To date, however, consideration of the racial effects (if any) of algorithmic criminal justice has been piecemeal. This section briefly surveys existing studies of algorithmic criminal system systems that touch on questions of race. This survey suggests there are real reasons for closely analyzing the various effect of algorithmic tools on criminal justice. It also suggests that there is more than one pathway by which racial effects can emerge. Any analytic framework for capturing race effects in this context must therefore also be plural.

##### 1. Policing and the Problem of Tainted Training Data

In the policing context, some commentators have flagged the possibility that the use of algorithmic instruments will reinforce existing race-based patterns of policing.<sup>146</sup> This may occur because algorithmic predictions will vary depending on the quality of the training data used to construct the predictive function. For example, if the training data systematically omits data about certain subsets of a population—if it has what Kate Crawford calls “black holes”<sup>147</sup>—it will generate results that fail to account for some population. Such gaps can be a function of poor relations between law enforcement and certain communities. For example, imagine a jurisdiction that allocates patrol resources based on historical reports of crime. Neighborhoods characterized by poor relations with police might underreport crime, such that they receive fewer policing resources in the future. But, contra Crawford, algorithmic tools might also be used to compensate for asymmetrical data gaps. Hence, the Shotspotter system records shots fired in urban environments. It can thus reveal neighborhoods in which residents do not report shootings to police.<sup>148</sup> This has at least the potential to mitigate historical enforcement gaps. To

---

<sup>143</sup> Richard Berk, *An impact assessment of machine learning risk forecasts on parole board decisions and recidivism*, 13 J. EXP. CRIMINOLOGY 193, 195 (2017) [hereinafter “Berk, *Impact assessment*”]. For a clear and nontechnical explanation of random forests methods, see Richard Cutler, D. et al., *Random forests for classification in ecology*, 88 ECOLOGY 2783, 2884-85 (2007).

<sup>144</sup> Berk, *Impact Assessment*, *supra* note 143, at 212.

<sup>145</sup> See sources cited in *supra* notes 11 and 9.

<sup>146</sup> Brayne, *supra* note 64, at 997 (arguing that “data-driven surveillance may be implicated in the reproduction of inequality ... by deepening the surveillance of individuals already under suspicion; by widening the criminal justice dragnet unequally; and leading people to avoid ‘surveilling’ institutions that are fundamental to social integration”).

<sup>147</sup> Kate Crawford, *The Anxieties of Big Data*, NEW INQUIRY, May 30, 2014, <https://thenewinquiry.com/the-anxieties-of-big-data/>.

<sup>148</sup> Sarah Griffith, *Fighting a Losing Battle*, DAILY MAIL U.K., Apr. 19, 2016, <http://www.dailymail.co.uk/sciencetech/article-3547719/Fighting-losing-battle-AI-ShotSpotter-computer-used-track-gunfire-reveals-far-shots-fired-reported.html>. Another example is a predictive tool used by hospitals to predict readmissions missed patients with asthma from the readmission function because those patients had been triaged to intensive care units rather than being released in the training data. Rich Caruana et al., *Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day*

conclude that algorithmic instruments will either necessarily undermine, or necessarily perpetuate, historical imbalances in the allocation of criminal justice resources seems premature. They can do both.

Policy distortions might also arise if historical data of political activity, deployed as training data for an algorithmic tool, is inflected by the racial presumptions and stereotypes of the past officials. Such measurement effort “will create decision and allocation bias.”<sup>149</sup> The concern here is a variant on a worry common in medical research that “race is such a dominant category in the cognitive field that the ‘interim solution’ [of using race as a proxy for some other trait of interest] can leave its own indelible mark.”<sup>150</sup> That is, race is such a freighted category that once deployed, it cannot be taken back. Race effects can arise if data collected as a by-product of police activity does “not pertain to future instances of crime” but rather to “instances of crime that become known to police.”<sup>151</sup> If police activity is predicted by race, then subsequent policing (and hence the costs of policing) will be unevenly allocated by race. This can happen even if nonviolent crime is evenly distributed. The result is greater black exposure to arrest and incarceration.<sup>152</sup> Such distortions are not evitable. It has been demonstrated in the computer science literature that incorporating an element of randomization into the algorithm is one way of buffering this kind of distortion.<sup>153</sup>

How forceful, as an empirical matter, are these concerns? One study of PredPol’s algorithm suggests that there is reason for concern. According to that study, when the algorithm used police data to generate predictions of narcotics crimes in Oakland, the algorithm recommended that twice as much policing resources be directed to black as opposed to white areas, despite the fact that narcotics offenses were reasonably equally spread across both white and black areas.<sup>154</sup> A second study, once more focused on the PredPol algorithm, identified the possibility of “runaway feedback loops,” by which police are repeatedly sent back to the same neighborhood in a way that reinforces with growing

---

*readmission*, in 2015 PROC. 21TH ACM SIGKDD INT’L CONF. ON KNOWLEDGE DISCOVERY AND DATA MINING 1721.

<sup>149</sup> Stendhal Mullainathan & Ziad Obermeyer, *Does Machine Learning Automate Moral Hazard and Error?*, 107 AM. ECON. REV.: PAPERS & PROC. 476, 478 (2017).

<sup>150</sup> Troy Duster, *Race and reification in science*, 307 SCIENCE 1050, 1050 (2005).

<sup>151</sup> Kristian Lum & William Isaac, *To Predict and Serve?*, SIGNIFICANCE, Oct. 2106, at 14, 16.

<sup>152</sup> For findings that race, rather than criminality, predicted deployment in one city (Seattle), see Katherine Beckett et al., *Drug Use, Drug Possession Arrests, and the Question of Race: Lessons from Seattle*, 52 SOC. PROB. 419, 435 (2005); Katherine Beckett et al., *Drug Use, Drug Possession Arrests, and the Question of Race: Lessons from Seattle*, 52 SOC. PROB. 419, 435 (2005); see also Huq, *Disparate Policing*, *supra* note 4, at 2929-40 (discussing effects of such disproportionate allocations of policing resources).

<sup>153</sup> Kroll et al., *supra* note 38, at 682-83; Pauline T. Kim, *Auditing Algorithms for Discrimination*, 166 U. PA. L. REV. ONLINE 189, 195 (2017) (“[W]hen a model produces biased outcomes due to the processes generating the input values, merely tweaking the distribution of data inputs will not solve the problem.”).

<sup>154</sup> Lum & Isaac, *supra* note 151, at 18, fig. 2. The background estimate of the geographic distribution of narcotics offenses derives from separate national data. A survey of Los Angeles residents by the advocacy group Stop LAPD Spying Coalition similarly found police contact highly concentrated in less than 2 percent of the population, and suggested that this flowed from the use of an algorithmically derived “Chronic Offender Bulletin.” Maha Ahmed, *Aided by Palantir, The LAPD Uses Predictive Policing to Monitor Specific People and Neighborhoods*, THE INTERCEPT, May 11, 2018, <https://theintercept.com/2018/05/11/predictive-policing-surveillance-los-angeles/>.

force initial distortions in the training data.<sup>155</sup> Third, in a qualitative study of L.A.’s use of predictive policing, Brayne concludes that PredPol increases surveillance of low-income, minority residents who are already under surveillance; widened “the surveillance dragnet” unequal; and drove members of the aforementioned communities “to avoid ‘surveilling’ institutions.”<sup>156</sup> These studies suggest that PredPol and similar technologies do distort the optimal allocation of policing resources in part because the individuals being regulated do not hew to fixed patterns of behavior. Rather, they respect PredPol-driven interventions. Moreover, because the individuals regulated by PredPol are heterogeneous in their responses to policing—some engage in more avoidance behavior than others—the error rate across the population as a result of such variable responsiveness to new policing intervention will also be uneven.<sup>157</sup>

Companies marketing algorithmic criminal justice instruments have evinced varying levels of concern about this possibility of racial effects. On the one hand, PredPol’s manufacturer advertises its exclusion of “drug related and traffic offenses from its predictions to remove officer bias.”<sup>158</sup> Similarly, HunchLab underscores its reliance on non-crime-related data as a way of making predictions not influenced by potentially flawed past exercises of officer discretion.<sup>159</sup> On the other hand, the Sentencing Commission of Pennsylvania has incorporated arrest data into its sentencing algorithm despite the fact that there is good reason to think that police discretion as to when and whom to arrest may have racial distortions.<sup>160</sup>

Finally, concerns about the polluting effect of historical training data are not limited to predictive algorithms. Studies of facial recognition technologies also suggest racial disparities in accuracy rates. One 2012 study tested three commercial algorithms on mug shots from Pinellas County, Florida. African Americans were between five and ten percent less likely to be successfully identified—i.e., more likely to be falsely rejected—than other demographic groups. It identified a similar decline for females as compared to males and younger subjects as compared to older subjects.<sup>161</sup> A measure of caution, though, should be used in evaluating these studies today. Much has changed in the domain of machine and

---

<sup>155</sup> Danielle Ensign et al., *Runaway feedback loops in predictive policing*, 4-5 (2017), <https://arxiv.org/abs/1706.09847>.

<sup>156</sup> Brayne, *supra* note 64, at 997.

<sup>157</sup> For an account of the difference between prediction problems and causal inference problems, and the risks of confusing the two, see Athey, *supra* note 32, at 355.

<sup>158</sup> *Machine Learning and Policing*, PREDPOL BLOG, July 19, 2017, <http://blog.predpol.com/machine-learning-and-policing>.

<sup>159</sup> HunchLab, *supra* note 107, at 12.

<sup>160</sup> Barry-Jester et al., *supra* note 138.

<sup>161</sup> Brendan F. Klare et al., *Face Recognition Performance: Role of Demographic Information*, 7 IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY 1789, 1797 (2012); *see also* Jonathon Phillips et al., *An Other-Race Effect for Face Recognition Algorithms*, 8 ACM TRANSACTIONS ON APPLIED PERCEPTION 14:1, 14:5 (2011) (finding that other leading algorithms performed worse on African Americans, women, and young adults than on Caucasians, men, and older people, respectively).

deep learning since 2012,<sup>162</sup> and it cannot be assumed that limitations on computational instruments that existed then still hinder analogous tools today.

## 2. *Bail/Sentencing Predictions and the Problem of Distorting Feature Selection*

The problems with algorithmic criminal justice do not begin and end with a concern with tainted historical training data. Attention to the bail and sentencing context suggests that even when there is no allegation of tainted training data, algorithmic criminal justice can generate concerns related to racial equity as a consequence of feature selection decisions. Even if these concerns focus on arguably unanticipated results, they might nonetheless have empirically consequential magnitudes.

Perhaps the highest profile debate concerning the racial effects of algorithmic instruments in criminal justice concerns, though, has focused on the Compas algorithm. Analyzing Compas data from Broward County, Florida, Pro Publica observed that the algorithm was “likely to falsely flag black defendants as future criminals, wrongly labeling them this way almost twice the rate of white defendants,” and to mislabel white defendants as “low risk more often than black defendants.”<sup>163</sup> That is, conditional on being a non-risky type, the Compas algorithm is more likely to overstate the risk presented by a black rather than a white person. Northpointe’s response did not focus on this measure of false positives (or, correlatively, the measure of false negatives that list in favor of whites). Instead, it identified the pool of individuals assigned a certain risk score as the relevant pool of comparators, and showed that within that pool, a white and a black defendant were equally likely to recidivate.<sup>164</sup> The ratio it emphasized, that is, takes as a denominator the group identified as high risk within each racial group, and then asks how many of those identifications are erroneous. This is the rate of false positions conditional on being identified as a risky type. As one group of computer scientists has noted, the resulting debate might well be understood not in terms of whether the Compas algorithm was racially discriminated—after all, there was no dispute that the algorithm did not include race as a feature—but rather what kind of racial effects counted in a normative or legal evaluation of its performance.<sup>165</sup>

Compas’s use in criminal sentencing has been challenged on various constitutional grounds. But its race-related effects remain untested in court. The most extensive judicial treatment of Compas, offered by the Wisconsin Supreme Court in *State v. Loomis*, focused on a Due Process challenge to the extent to which criminal defendants could challenge the algorithm’s terms.<sup>166</sup> Rejecting a challenge to the way in which the algorithm accounted

---

<sup>162</sup> A particularly vivid illustration of this is the dramatic increase in the quality of machine translation tools. Gideon Lewis-Krauss, *The Great A.I. Awakening*, N.Y. TIMES, Dec, 14, 2016, at M40, <https://www.nytimes.com/2016/12/14/magazine/the-great-ai-awakening.html>.

<sup>163</sup> Angwin et al., *supra* note 15.

<sup>164</sup> Dietrich et al., *supra* note 19, at 3.

<sup>165</sup> See Feller et al., *supra* note 21.

<sup>166</sup> *State v. Loomis*, 881 N.W.2d 749, 761-62 (2016) (holding that because the algorithm employed only publicly available data, or data that a defendant has supplied, the defendant could have denied or explained any information that was employed to develop his prediction). It is worth noting that the Court’s analysis here misses the force of the defendant’s argument: The latter seemed objected not so much to the

for a suspect's gender, the Wisconsin Court noted in passing "concerns regarding how a COMPAS assessment's risk factors correlate with race."<sup>167</sup> Unfortunately, the Court did not connect that observation with either a legal theory pursuant to which such correlations might be objectionable, or, alternatively, a normative basis for concern notwithstanding legality. Otherwise, commentators have noted that actuarial sentencing tools, whether algorithmic or not, might have more or less disparate racial impact or "inequitable social consequences."<sup>168</sup> But precisely what these "consequences" might be remains unclear.

Nevertheless, the COMPAS debate suggests that concerns about racial equity can persist *even if the inputs to the algorithm are not tainted by any historical bias*. Part of my aim here, particularly in Part III, is to explain how this can be so. For now, it suffices to say that earlier commentators who have suggested that algorithmic bias can be addressed exclusively through "a transparency of inputs and outputs" may have captured only one part of a larger normative picture.<sup>169</sup>

### 3. *Conclusion: An Incomplete Evidentiary Record*

Race interacts with algorithmic criminal justice tools in one of three ways. First, racial animus or stereotypical thinking can infect and distort training data. Second, race may be a feature used for classification. Third, the classification rule may have predictable effects that seem asymmetrical between racial groups. Scholars' thinking about and responses to the racial effects of algorithmic criminal justice instruments have been ad hoc and unsystematic. We have at best fragments of a more systematic account of how such effects arise and what effects they have. Hence, understanding empirically the manner in which algorithmic tools redistribute coercive outcomes should remain an important focus of research. Still, even with limited evidence in hand, it seems there is reason for thinking about the appropriate normative framework for evaluating these instruments' racial effects—especially given the long and troubled interaction between criminal justice policy and widely held beliefs about racial differences in culture and behavior.

The conceptual tools for that investigation are plainly wanting at the moment. There is no general agreement on the ways in which racial effects might count against the adoption or continued use of an algorithm. Insufficient attention, moreover, has been paid to the difference between tainted training data and problematic feature selection. There is also no general understanding of what it means to say that feature selection is flawed. Nor is there any consideration of how different kinds of racial effects might be weighed against each other. The field is ripe, in short, for more careful theorization of what it precisely means to talk about racial equity in this context.

---

nondisclosure of information about his own circumstances, but the manner in which that information was evaluated and weighted by the Compas algorithm.

<sup>167</sup> *Id.*

<sup>168</sup> Monahan & Skeem, *supra* note 139, at 507. For an empirical study that renders these concepts with more precision, see Jennifer L. Skeem and Christopher T. Lowenkamp, *Risk, race, and recidivism: predictive bias and disparate impact*, 4 CRIMINOLOGY 680, 702-03 (2016) (analyzing the relation of the Post Conviction Risk Assessment tool and future arrests, and finding that scores tracked the same level of recidivism within each group).

<sup>169</sup> Chander, *supra* note 42, at 1039.

## II. Equal Protection and Algorithmic Criminal Justice

But is such theorization needed? The Equal Protection Clause of the Fourteenth Amendment, after all, purports to provide a general norm regulating the state's use of race. Perhaps constitutional equality jurisprudence provides the needful criterion for evaluating the race effects of algorithmic criminal justice.

Or perhaps not. I describe and apply in this Part conventional doctrinal norms articulated under the Equal Protection Clause. Its core lesson is that the dominant intent- and classification-focused calibration is ill suited to the forms and dynamics of algorithmic criminal justice tools. To be sure, one might choose to apply the litmus tests supplied in the jurisprudence. But given that these focus on qualities of state action that are irrelevant, or barely relevant, to the way that algorithms in practice work, it is hard to see why one would do so. If there is a lesson here, indeed, it is about the woeful inadequacy of our constitutional equality norms for the contemporary world.

### A. What Equal Protection Protects

Equal Protection doctrine imposes two fundamental prohibitions on governmental action touching on race.<sup>170</sup> One concerns formal racial classifications. The other pertains to racialized intentions. The Court has either rejected or ignored concerns about the illegitimate nature or delegitimizing consequences of raw racial disparities in criminal justice.

Almost since its inception, constitutional Equal Protection has been understood to prohibit most laws containing an explicit racial classification, as well as laws that assign rights or burdens based on racial classification.<sup>171</sup> The first major judicial interpretation of the Clause, *Strauder v. West Virginia*, concerned a state statute limiting jury service to “white male persons ... twenty-one years of age.”<sup>172</sup> Invalidating the conviction of an African-American man under this provision, the Court explained that the statute’s want of facial equality violated the Constitution’s “immunity from inequality of legal protection.”<sup>173</sup> Racial classifications today are not per se invalid. Rather, they now trigger searching judicial review of their tailoring and means-ends rationality, an inquiry known as “strict scrutiny.”<sup>174</sup>

---

<sup>170</sup> Concerns about racial equity in criminal law need not be expressed in terms of Equal Protection jurisprudence. Many cases formally concerning Due Process arose in the context of discriminatory law enforcement, and are plausibly understood in terms of the Court’s desire to constrain the latter’s discretion. My concern in this Part is the formal doctrinal specification of equality, not its potential jurisprudential substitutes.

<sup>171</sup> Canonical exceptions include *Plessy v. Ferguson*, 163 U.S. 537, 550 (1896).

<sup>172</sup> 100 U.S. 303, 308 (1879) (citation omitted).

<sup>173</sup> *Id.* at 310.

<sup>174</sup> *Parents Involved in Cmty. Sch. v. Seattle Sch. Dist. No. 1*, 551 U.S. 701, 720 (2007) (imposing such scrutiny whenever “the government distributes burdens or benefits on the basis of individual racial classifications”); see also *Gratz v. Bollinger*, 539 U.S. 244, 270 (2003) (describing the use of such classifications as “pernicious” (citation omitted)); *Adarand Constructors, Inc. v. Peña*, 515 U.S. 200, 235 (1995) (“Federal racial classifications, like those of a State, must serve a compelling governmental interest, and must be narrowly tailored to further that interest.”).



Notoriously, strict scrutiny is not necessarily fatal.<sup>175</sup> In *Fisher v. University of Texas at Austin*, for example, the Court upheld the University of Texas at Austin's admission program, even though it accounted for race as one element of a "Personal Achievement Index," or PAI.<sup>176</sup> This satisfied strict scrutiny because the University "articulated concrete and precise goals" in relation to educational diversity; relied on "both statistical and anecdotal evidence" of a need for affirmative action; and engaged in ongoing deliberation about admissions protocols.<sup>177</sup> Precisely how *Fisher* calibrated strict scrutiny, though, is difficult to say. Educational diversity is not easily reduced to "concrete and precise" terms. Nothing the Court said illuminated how it tested the means-ends rationality behind the University's actions.<sup>178</sup> Yet in other contexts, it has construed strict scrutiny to work a near-categorical prohibition on similarly race-conscious government action.<sup>179</sup> For instance, in an earlier capital habeas case, the Court made that errant suggestion that race is "totally irrelevant to the sentencing process."<sup>180</sup> Such evanescent dicta, however, are probably too frail to support any firm conclusion.

Second, the Equal Protection Clause's regulation of racial considerations extends to instances in which the state harms an individual because of "a racially discriminatory purpose."<sup>181</sup> This requires litigants "show both that the passive enforcement system had a discriminatory effect and that it was motivated by a discriminatory purpose."<sup>182</sup> The Court has not defined with precision what counts as a "racially discriminatory purpose."<sup>183</sup> But at a minimum, it seems to include naked, taste-based aversion to a group based exclusively

---

<sup>175</sup> *Grutter v. Bollinger*, 539 U.S. 306, 326 (2003) ("Strict scrutiny is not 'strict in theory, but fatal in fact.'").

<sup>176</sup> 136 S. Ct. 2198 (2016). On the construction of the PAI, see *Fisher v. University of Texas*, 131 S.Ct. 2411 S. Ct. 2415, 2415-15 (2013).

<sup>177</sup> *Fisher*, 136 S. Ct. at 2214-15.

<sup>178</sup> David A. Strauss, *Fisher v. University of Texas and the Conservative Case for Affirmative Action*, 2016 SUP. CT. REV. 1, 16-17 ("The central problem is that judgments about the kind and degree of diversity that a student body should have ... are simply not susceptible to precise metrics.").

<sup>179</sup> See Richard H. Fallon, Jr., *Strict Judicial Scrutiny*, 54 UCLA L. Rev. 1267, 1337 (2007) ("According to one interpretation, strict scrutiny embodies a nearly categorical prohibition against infringements of fundamental rights, regardless of the government's motivation, but subject to rare exceptions when the government can demonstrate that infringements are necessary to avoid highly serious, even catastrophic harms.").

<sup>180</sup> *Zant v. Stephens*, 462 U.S. 862, 885 (1983).

<sup>181</sup> *Washington v. Davis*, 426 U.S. 229, 240 (1976) (holding that "the basic equal protection principle that the invidious quality of a law claimed to be racially discriminatory must ultimately be traced to a racially discriminatory purpose"). Racial intent must be the but-for cause of an action. *Personnel Admr. v. Feeney*, 442 U.S. 256, 279 (1979) (proof of discriminatory purpose requires showing that government decision-maker "selected or reaffirmed a particular course of action at least in part 'because of,' not merely 'in spite of,' its adverse effects upon an identifiable group").

<sup>182</sup> *Wayte v. United States*, 470 U.S. 598, 608 (1985).

<sup>183</sup> See Huq, *Discriminatory Intent*, *supra* note 40, at 21-36 (describing five different theories of discriminatory purpose in the case law); accord David Strauss, *Discriminatory Intent and Taming of Brown*, 59 U CHI. L. REV. 935, 947 (1989) (noting that even canonical cases such as *Brown v. Board of Education* did not clarify "which conception of discrimination [the Court] embraced, or how far the principle of [Equal Protection] extended").

on race.<sup>184</sup> So the Court recently explained that evidence that a juror relied on “racial stereotypes or animus to convict a criminal defendant” would be sufficient to warrant reversal of that conviction on Sixth Amendment grounds.<sup>185</sup> Even here, the doctrine is not without ambiguity. It is not clear, for instance, whether a state actor shown to have made a decision based on racial animus could plausibly respond that their action could nonetheless be upheld because they survived strict scrutiny. Analytically, it is hard to see how a measure based on an invidious stereotype could ever be closely fitted to a legitimate state interest. So it may be that the question does not arise because it has little operational importance.<sup>186</sup>

And there are other ways that race can infiltrate the mind beyond animus. For example, a rational reliance on race as a statistically accurate proxy for some other policy-salient quality is analytically distinct from taste-based discrimination.<sup>187</sup> The Court has not been clear on whether such statistical discrimination triggers constitutional concerns. On the one hand, in the 2007 case *Johnson v. California*, a majority of the Justices held that strict scrutiny applied to an unwritten California prison policy of racially segregating prisoners for up to sixty days each time they enter a new correctional facility with the aim of mitigating violence between gangs of different races.<sup>188</sup> On the other hand, lower federal courts routinely shake off challenges to race-specific suspect descriptions. The Supreme Court has consistently and repeatedly declined to intervene.<sup>189</sup> All that can safely be said is that at least in some instances, statistical discrimination will be subject to close judicial scrutiny, and sometimes it won’t be. The cut point between those domains remains to be defined.

---

<sup>184</sup> This is what economists call taste-based discrimination. GARY BECKER, *THE ECONOMICS OF DISCRIMINATION* 14-15 (2d ed. 1971) (modeling taste-based discrimination as a “discrimination coefficient,” which “acts as a bridge between money and net costs. Suppose an employer were faced with the money wage rate  $PI$  of a particular factor; he is assumed to act as if  $PI(1 + d_i)$  were the net wage rate, with  $d_i$  as his [discrimination coefficient] against this factor”).

<sup>185</sup> *Pena-Rodriguez v. Colorado*, 137 S. Ct. 855, 869 (2017); see also *Foster v. Chatman*, 136 S. Ct. 1737, 1747-55 (2016) (that the Georgia Supreme Court had made a “clearly erroneous” decision when it declined to find that prosecution use of preemptory strikes in a capital case was *not* animated by a discriminatory purpose in the face of lurid evidence to the contrary).

<sup>186</sup> I am not sure, however, this conclusion would be warranted. Consider, for example, an action shown to be tainted by racial animus, but that could be defended as narrowly tailored given different motivational premises and additional evidentiary support. The so-called travel ban might have this character. For a discussion somewhat short on illumination, see *Trump v. International Refugee Assistance Project*, 137 S. Ct. 2080, 2086 (2017) (per curiam). For an extended discussion, see Aziz Z. Huq, *Article II and Antidiscrimination*, 117 MICH. L. REV. – (forthcoming 2019).

<sup>187</sup> Kasper Lippert-Rasmussen, “*We are all Different*”: *Statistical Discrimination and the Right to be Treated as an Individual*, 15 J. ETHICS 47, 54 (2011) (providing a formal definition of such rational discrimination). This is what economists call statistical discrimination. Kenneth J. Arrow, *The Theory of Discrimination*, in *DISCRIMINATION IN LABOR MARKETS* 3, 24-27 (Orley Ashenfelter & Albert Rees eds., 1973) (“Skin color and sex are cheap sources of information. Therefore prejudices (in the literal sense of pre-judgments, judgments made in advance of the evidence) about such differentia can be easily implemented ....”).

<sup>188</sup> *Johnson v. California*, 543 U.S. 499, 505 (2005).

<sup>189</sup> See, e.g., *Monroe v. City of Charlottesville*, 579 F.3d 380 (4th Cir. 2009), *cert. denied*, 130 S. Ct. 1740 (2010) (denying certiorari in a case exempting race-based suspect selection from equal protection scrutiny); *Brown v. City of Oneonta*, 221 F.3d 329 (2d Cir. 2000), *cert. denied*, 534 U.S. 816 (2001) (same).

Moreover, it seems likely that not all racial animus in the criminal justice system is crisply articulated in the Queen's English. Instead, we might expect the overt racial labelings to be the exception, with race more commonly embedded in "tacit," unspoken understandings.<sup>190</sup> The only evidence of the latter's operation may be downstream differential effects on suspects and defendants of different races. In contrast to its strict superintendence of overt classifications, however, the Court has rejected the argument that a constitutional violation can be made out by a showing of disparate racial impacts. In *McCleskey v. Kemp*, most importantly, the Court rejected an Equal Protection challenge to Georgia's capital punishment system based on econometric evidence of racial disparities.<sup>191</sup> Lower courts have extended that holding to the distinct context of statistical evidence about the role of race in a single decision-maker's actions over time (e.g., a single district attorney over a number of years).<sup>192</sup>

Paradoxically, both the Court's embrace of the racial intent rule and its repudiation of a disparate treatment rule have been justified by the need to maintain the criminal justice system in good working order. In *McCleskey*, Justice Powell's majority opinion expressed alarm that the defendant's challenge would "throw[] into serious question the principles that underlie our entire criminal justice system."<sup>193</sup> In Powell's view it was inconceivable that the Constitution would "require that a State eliminate any demonstrable disparity that correlates with a potentially irrelevant factor in order to operate a criminal justice system."<sup>194</sup> On the other hand, the Court has explained decisions enforcing closer invigilation of race's role in the jury deliberation context as "necessary to prevent a systemic loss of confidence in jury verdicts, a confidence that is a central premise of the Sixth Amendment trial right."<sup>195</sup> The Court thus appears to believe that the legitimacy of a criminal justice system simultaneously requires keen alertness to concerns of racial justice, and also a willful blindness to such concerns.

Stated in summary form then, current constitutional jurisprudence compels judges to maintain the stability of the criminal justice system by ignoring racial disparities, by isolating racial classifications and by extirpating (some) racial animus. It is a doctrinal status quo that fits poorly with emergent algorithmic realities.

---

<sup>190</sup> On the notion of tacit understandings, see Michael Polanyi, *The logic of tacit inference*, 41 PHILOSOPHY 1, 2-3 (1966).

<sup>191</sup> 481 U.S. 279, 292-93 (1987).

<sup>192</sup> John H. Blume et. al., *Post-McCleskey Racial Discrimination Claims in Capital Cases*, 83 CORNELL L. REV. 1771, 1794 (1998) (collecting cases); see also *Chavez v. Illinois State Police*, 251 F.3d. 612, 645, 648 (7th Cir. 2001) (finding no discriminatory purpose despite statistical showing of racial disparities in traffic stops). But cf. Reva B. Siegel, *Blind Justice: Why the Court Refused to Accept Statistical Evidence of Discriminatory Purpose in McCleskey v. Kemp-and Some Pathways for Change*, 112 NW. U. L. REV. 1269, 1288 (2018) (flagging limits to *McCleskey's* scope).

<sup>193</sup> *McCleskey*, 481 U.S. at 315.

<sup>194</sup> *Id.* at 319.

<sup>195</sup> *Pena-Rodriguez v. Colorado*, 137 S. Ct. 855, 869 (2017).

## B. How Equal Protection Fails to Speak in Algorithmic Terms

Equal Protection doctrine is sharply criticized by those who perceive it to embody a judicial failure to account for the diffusion and impact of racial effects in society, let alone our highly racially stratified criminal justice system.<sup>196</sup> I set these concerns aside here, and take the doctrine seriously on its own terms. Even then, I find reasons to doubt that the current doctrine can respond effectively to the questions of race raised by algorithmic criminal justice. The concerns of constitutional law simply do not map onto the ways in which race impinges on algorithmic criminal justice. The result is a gap between legal criteria and their objects.

Crucially, the two main doctrinal touchstones of bad intent and bad classifications provide scant traction for the analysis of algorithmic criminal justice. Both hinge on concepts that translate poorly, if at all, to the algorithmic context, and are not easily adapted for application to that end. A focus on racial animus will almost never be fruitful. A focus on classification leads to perverse and unjustified results. The replacement of unstructured discretion with algorithmic precision, therefore, thoroughly destabilizes how Equal Protection doctrine works on the ground. The resulting mismatches compel my conclusion that a new framework for thinking about the pertinent racial equity questions is needed.

### 1. *The Trouble With Intent*

Taking intent as a touchstone of Equal Protection concern directs attention to questions at best tangential to the potential role of race in algorithmic criminal justice. To be sure, problematic intent might enter into algorithmic design in different ways, one of which is easily accounted for in doctrinal terms. But, in general, intent will rarely be the crux of the matter.

To begin, I suspect that the notion of machine intentionality is sufficiently counterintuitive to find no place in constitutional law. Speculation about a future of “superintelligent” artificial intelligences aside,<sup>197</sup> the transformation of training data into new schemes of classification by machine learning or deep learning does not obviously map onto familiar forms of human intentionality. The most advanced artificial intelligences

---

<sup>196</sup> Recent critiques include Russell K. Robinson, *Unequal Protection*, 68 STAN. L. REV. 151, 154 (2016) (contending that “the Supreme Court has steadily diminished the vigor of the Equal Protection Clause in most respects”); Ian Haney-López, *Intentional Blindness*, 87 N.Y.U. L. REV. 1779, 1828 (2012) (arguing that the Court has “split equal protection into the separate domains . . . , one governing affirmative action and the other discrimination against non-Whites” in a move that has made it systematically easier for white plaintiffs to prevail”).

<sup>197</sup> Cf. NICK BOSTRAM, *SUPERINTELLIGENCE: PATHS, DANGERS, STRATEGIES* 52 (2014) (defining superintelligence as “intellects that outperform the very best human minds across many very general cognitive domains”).

can now pass the Turing test<sup>198</sup> and defeat (human) world champions at Go.<sup>199</sup> But even these machines do not obviously possess the sort of psychological interiority commonly thought to be a necessary predicate to intentionality.<sup>200</sup> Talk of machine intentionality, therefore, is either premature or a badly framed metaphor. As a result, it is better to treat the algorithm itself as irrelevant to the constitutional analysis so far as intentionality is concerned.

Bracketing the machine-learning tool as agent, however, there are two possible ways in which intention might enter the picture. First, an algorithm's designer might be motivated by either an animosity toward a racial group, or else a prior belief that race correlates with criminality, and then design the algorithm on that basis. Barocas and Selbst call this "masking."<sup>201</sup> Masking might occur through either a choice to use polluted training data, or else the deliberate selection of some features but not others. For instance, it is well-understood that when employers ignore credit score information, they tend to search for proxies that have the inadvertent effect of deepening racial disparities.<sup>202</sup> A discriminatory algorithm designer will leverage such knowledge to fashion instruments that yield the disparate racial effects they believe to be warranted a priori. Without knowing the full spectrum of features that could, conceivably, have been included in the training data—which can be "enormous"<sup>203</sup>—it will be difficult or impossible to diagnose this kind of conduct absent direct evidence of discriminatory intent.<sup>204</sup> It will, moreover, be especially difficult to show that but for race, a specific feature would have been included, as the doctrine requires.<sup>205</sup> A basic principle of "feature selection" instructs that one should "keep the important features and discard the unimportant one."<sup>206</sup> To the extent that masking occurs, therefore, it seems clear that the litigation process would rarely yield evidence of such intentional manipulation of the algorithm's design.

---

<sup>198</sup> In June 2014, an artificial intelligence passed the Turing test, arguably for the first time. Kevin Warwick & Huma Shah, *Can machines think? A report on Turing test experiments at the Royal Society*, 28 J. EXP. & THEO. ARTIFICIAL INTELLIGENCE 989, 990 (2016). The Turing test involves human judgments about natural language conversations between a computer and a machine; a machine passes the test if the human observer is unable to distinguish human from machine.

<sup>199</sup> Silver et al., *supra* note 88, at 490.

<sup>200</sup> Accounts of this interiority vary. On one influential definition, intentions are "conduct-controlling pro-attitudes, ones which we are disposed to retain without reconsideration, and which play a significant role as inputs to reasoning." MICHAEL BRATMAN, INTENTION, PLANS, AND PRACTICAL REASON 20 (1987). On another view, when *S* is doing *A* intentionally, *S* knows that she is doing *A*. G.E.M. ANSCOMBE, INTENTION 11-15 (1963). Machines lack attitudes or self-knowledge in the relevant senses.

<sup>201</sup> Barocas & Selbst, *supra* note 50, at 692 ("[D]ecision makers could knowingly and purposefully bias the collection of data to ensure that mining suggests rules that are less favorable to members of the protected class.").

<sup>202</sup> Robert Clifford and Daniel Shoag, "No More Credit Check Score": *Employer Credit Bans and Signal Substitution* 3 (Mat 2016), [http://scholar.harvard.edu/files/shoag/files/clifford\\_and\\_shoag\\_final.pdf](http://scholar.harvard.edu/files/shoag/files/clifford_and_shoag_final.pdf).

<sup>203</sup> Athey, *supra* note 32, at 483.

<sup>204</sup> Training data will often have so many potential features that inferring the reason for the inclusion of some and exclusion of others will often not be feasible. *Id.* at 483.

<sup>205</sup> *Personnel Admr. v. Feeney*, 442 U.S. 256, 279 (1979) (proof of discriminatory purpose requires showing that government decision-maker "selected or reaffirmed a particular course of action at least in part 'because of,' not merely 'in spite of,' its adverse effects upon an identifiable group").

<sup>206</sup> ALPAYDIM, *supra* note 31, at 73-74.

Perhaps the most important reason to set aside the masking phenomena, however, is the fact that it does not appear to be a significant one in practice. Part of the reason for this is that racial animus has a performative, interpersonal aspect. Racial discrimination commonly entails an effort by one group to “produce esteem for itself by lowering the status of another group,”<sup>207</sup> and correlatively producing a “set of ... privileges and benefits” of superordinate group membership.<sup>208</sup> Masking is a form of discrimination that involves no interpersonal interaction, and no esteem-affirming performance. It might therefore be no surprise that it is comparatively rare.

Intent, however, might be relevant in algorithmic criminal justice in a second more commonly salient way. As Part I explained, the training data used to create a classificatory function can be the product of biased or distorted decision-making. For example, in a jurisdiction where African-Americans were targeted for frequent and unjustified police contact, the pool of arrestees and convicted criminals may be biased by an underrepresentation of non-black individuals.<sup>209</sup> A jurisdiction in which black neighborhoods are underserved by police responses to emergency calls, in contrast, might generate data on the distribution of crime with a black (or grey) hole in respect to African-American neighborhoods.<sup>210</sup> A jurisdiction, moreover, might underserve black neighborhoods by understaffing responses to 911 calls at the same time as concentrating a disproportionate amount of street policing resources on the same neighborhoods.<sup>211</sup> An algorithm trained on police-generated data from this jurisdiction is likely to allocate resources in ways that reflect and perhaps entrench disparities in the way in which policing resources are allocated.

The relevant intent in this example, moreover, differs in two important ways from the canonical form of impermissible intent in Equal Protection case-law. First, in the absence of an express policy, the operation of racial preferences by officials in activities that produce training data will generally be highly decentralized and uncoordinated. Policing, and to a lesser extent bail determinations and sentencing, are dispersed rather than centralized forms of state action. Individual officers or magistrates have a large degree of discretion in consequence of their sheer numerosity and the difficulty of monitoring their

---

<sup>207</sup> Richard H. McAdams, *Cooperation and Conflict: The Economics of Group Status Production and Race Discrimination*, 108 HARV. L. REV. 1003, 1044 (1995). A similar idea is introduced in George A. Akerlof, *Discriminatory, Status-Based Wages Among Tradition-Oriented, Stochastically Trading Coconut Producers*, 93 J. POL. ECON. 265, 265 (1985).

<sup>208</sup> Cheryl I. Harris, *Whiteness as Property*, 106 HARV. L. REV. 1707, 1713 (1993). For a seminal account of this concept, see DAVID R. ROEDIGER, *THE WAGES OF WHITENESS: RACE AND THE MAKING OF THE AMERICAN WORKING CLASS* (1991).

<sup>209</sup> For findings of such disparities, see, e.g., *Floyd v. City of New York*, 959 F. Supp. 2d 540, 567 (S.D.N.Y. 2013) (New York); *Boston Police Commissioner Announces Field Interrogation and Observation (FIO) Study Results*, Oct 8, 2014, <http://bpdnews.com/news/2014/10/8/boston-police-commissioner-announces-field-interrogation-and-observation-fio-study-results> (Boston).

<sup>210</sup> See, e.g., *Cent. Austin Neighborhood Ass'n v. City of Chicago*, 2013 IL App (1st) 123041, ¶ 4, 1 N.E.3d 976, 979 (describing allegations of longer response times to 911 calls in minority neighborhoods in Chicago).

<sup>211</sup> As appears to be the case with Chicago. *Id.*; Aamer Madhani, *Chicago police and ACLU agree to stop-and-frisk safeguards*, USA TODAY, Aug. 7, 2015, <http://www.usatoday.com/story/news/2015/08/07/chicago-police-agree-reform-stop-and-frisk/31277041/>.

decisions. It is hardly clear how a court could or would make a determination of ‘intent’ when confronted with an extensive multitude unguided by formal decision procedures. Constitutional doctrine has not developed an intellectual tool-kit for aggregating a large number of dispersed individual motives so as to ascertain whether a but-for standard of intentionality has been met by a collectivity.

An analogous, but easier, problem arises in the legislative context, where many individuals bring to bear potentially diverse motives in order to shape singular institutional acts with the force of law. Equal Protection law has struggled with how to conceptualize the concept of “intent” in the legislative context so unsuccessfully that one influential commentator has advocated wholesale retreat from judicial accounting for legislators’ subjective intents; in his view, the task of principled aggregation is simply too hard for judges.<sup>212</sup> Unlike legislatures, a plurality of geographically and temporally diffused cohorts of officials (whether police or magistrates) lack any stable procedures or mechanisms for eliciting and formalizing a singular intent.<sup>213</sup> Their ability to form a coherent, let alone legally relevant, intent may seriously be doubted.

Even if such an intentionality could be derived from a diffuse haze of discrete policing decisions or detention-related judgments, it is not clear whether the mere incorporation by reference of such historical judgments into new, forward-looking algorithmic tools would trigger Equal Protection Clause concern. Even if historical intent can be inferred successfully, there remains a question of whether reliance on flawed historical data *counts* as a constitutionally relevant form of intent. It is certainly possible for bad intent to endure over time. Indeed, the Court has invalidated states’ laws enacted to preserve “white supremacis[m]” many decades before litigation began, and in so doing rejected the notion that “events occurring in the [intervening] years [could have] legitimated the provision.”<sup>214</sup> But there are no Equal Protection cases in which the Court has considered outcomes resulting from concededly discriminatory official action that in turn was adopted by a new and different actor as the rationale for new, forward-looking policy.<sup>215</sup> In short, there is simply no way of knowing whether a ‘relay-race’ theory of bad intent would pass muster in constitutional law.

Perhaps the closest analog to this problem of governmental reliance on flawed data arises in the Fourth Amendment context. In that domain, the Court has declined to treat the

---

<sup>212</sup> Richard H. Fallon, Jr., *Constitutionally Forbidden Legislative Intent*, 130 HARV. L. REV. 523, 533-34 (2016) (insisting that “ultimate determinations of constitutional validity should always depend on the content and effects of challenged legislation, not the subjective intentions of the enacting legislatures”).

<sup>213</sup> Cf. CHRISTIAN LIST & PHILLIP PETTIT, GROUP AGENCY THE POSSIBILITY, DESIGN, AND STATUS OF GROUP AGENTS 81 (2012) (“[A] group’s performance as an agent depends on how it is organized: on its rules and procedures for forming its propositional attitudes ....”). In most cases, the groups relevant to algorithmic criminal justice have no such rules or procedures.

<sup>214</sup> *Hunter v. Underwood*, 471 U.S. 222, 233 (1985).

<sup>215</sup> The closest analog of which I am aware arises under the Fair Housing Act, where there can be a question whether a municipal decision on, say, taxes or zoning causes a pattern of residential racial segregation. Cf. *Texas Department of Housing & Community Affairs v. Inclusive Communities Project, Inc.*, 135 S. Ct. 2507, 2523 (2015) (demanding a showing of “robust causality”). But the question here is not one of causation: It is a question of whether the intentions of the original police or magistrate ought to be imputed to the algorithm given their influence on the training data.

flaws in a first-moving official's behavior as infecting a second, subsequently acting official's decision to depend on that first officer. For example, when a police officer relies on a recalled warrant mistakenly distributed by another police force, the latter's mistake of law is not imputed to the arresting officer such that evidence must be excluded.<sup>216</sup> Although the analogy is inexact,<sup>217</sup> the Fourth Amendment's stingy treatment of imputed fault when the unlawful action of official X becomes the basis of official Y's otherwise lawful act suggests that an intent-focused Equal Protection lens will have limited traction in the algorithmic criminal justice context. For there is no reason to think that the rules of imputed motivation should vary between the Equal Protection and the Fourth Amendment contexts.

But that theoretical problem may be precisely that—theoretical. Even if flawed training data were identified, it seems unlikely that its tainted nature could suffice to establish a constitutional concern in practice. Any moderately competent municipality found using flawed data would hardly concede that it was doing so intentionally. Rather, it would be far more likely to defend its decision as the best option, *faute de mieux*, given historically shaped constraints. Because a constitutional violation cannot be shown unless the state relied on race as a ground of decision, as opposed to acting in spite of race,<sup>218</sup> this defense would likely succeed. As a practical matter, therefore, the narrow definition of intent in Equal Protection doctrine would likely insulate racially tainted training data from legal attack.

This means that none of the pathways for integrating intent into the Equal Protection analysis of algorithmic criminal justice are likely to prove fruitful. None of them are well suited for a consideration of the ways in which race in practice interacts with algorithmic criminal justice. Equal Protection doctrine was designed to police the dispersed, open-ended discretionary judgments of street-level officials. It does a poor job when applied to the very different context of algorithm design and application. It is hence necessary to consider the logic of anti-classification as an alternative lens.

## 2. *The Trouble with Classification*

The anticlassification strand of Equal Protection doctrine prohibits the government from “classify[ing] people either overtly or surreptitiously on the basis of a forbidden category” such as race.<sup>219</sup> At first blush, it seems a natural fit: Algorithms work by applying categories to training data (when defining features) and then generating novel classification rules to apply to test data. A rule to the effect that race (or, say, ethnicity) could not be used either as a feature or as an element of a classifier absent narrow tailoring to a compelling

---

<sup>216</sup> *Herring v. United States*, 555 U.S. 135, 140 (2009); *see also Arizona v. Evans*, 514 U.S. 1, 14–15 (1995) (same result for errors by a judicial administrator). The Court, however, is willing to impute another officer's knowledge of information salient to the legality of a search when doing so renders a search lawful. *Whitely v. Warden*, 401 U.S. 560, 568 (1971). Although these positions can be squared, it is striking that imputation is available only when it expands state authority.

<sup>217</sup> The availability of exclusion in Fourth Amendment cases is said to turn on the deterrent effect of that remedy. *Herring*, 555 U.S. at 144. The consequential focus on deterrence is absent when one is concerned with attributions of intentionality.

<sup>218</sup> *Personnel Adm. v. Feeney*, 442 U.S. 256, 279 (1979).

<sup>219</sup> Jack M. Balkin & Reva B. Siegel, *The American Civil Rights Tradition: Anticlassification or Antisubordination?*, 58 U. MIAMI L. REV. 9, 10 (2003).



state interest would seem to be a natural fit. Such a rule, however, would be unmoored from the justifications for an anticlassificatory rule. It would also engender results that contradict the assumed purposes of the rule.

The anticlassification account of Equal Protection is premised on two main justifications. First, it is motivated by a concern that the state's use of racial classifications will facilitate or amplify private discrimination.<sup>220</sup> This worry is in turn premised on the empirical claim that a “perception ... fostered by [government]” of differences between racial groups “can only exacerbate rather than reduce racial prejudice.”<sup>221</sup> The foundation of this empirical claim, to be sure, is hardly clear. Why would the communicative effect of state racial classifications entail a legitimation of private animus? The causal link here is not obvious.<sup>222</sup> One interpretation of the Court's argument might start with the Court's claim that race is “‘in most circumstances irrelevant’ to any constitutionally acceptable legislative purpose.”<sup>223</sup> Read sympathetically, the Court appears to be saying that because race is irrelevant to the vindication of legitimate government ends, the observation that the state is treating race nevertheless as salient has the effect of propagating a false popular belief in racial hierarchies.<sup>224</sup> The second possible interpretation of an anticlassification rule turns on a nonconsequentialist, deontological intuition. That is, according to some Justices, it is a moral axiom that the state must treat all persons as individuals, and such

---

<sup>220</sup> *Anderson v. Martin*, 375 U.S. 399, 402 (1964) (holding that a Louisiana statute, which mandated the designation of a candidate's race on election ballots, violated equal protection because it enlisted the power of the state to enforce private racial prejudices).

<sup>221</sup> *Adarand Constructors, Inc. v. Peña*, 515 U.S. 200, 229 (1995).

<sup>222</sup> It is not clear why the reasonable observer would draw an inference about *a racial group*, though, instead of an inference that the *government* was unjustified and irrational in its action. In any event, the claim that racial identity is not salient in a context where racial preferences retain a powerful hold is a deeply dubious one. For an estimate of the prevalence of racial animus using an innovative empirical method, see Seth Stephens-Davidowitz, *The cost of racial animus on a black candidate: Evidence using Google search data*, 118 J. PUB. ECON. 26, 26-28 (2014) (using Google data to estimate the prevalence and geographic variation of anti-black sentiment).

<sup>223</sup> *McLaughlin v. State of Fla.*, 379 U.S. 184, 192 (1964) (citing *Hirabayashi v. United States*, 320 U.S. 81, 100 (1943)). Notice that the Court's argument here is crucially ambiguous. It could be that an individual's *race* is irrelevant to many legitimate state ends, but that the persistence of *racism* as an ambient social phenomenon is relevant to how the state can achieve those ends. The Court's formulation elides this difference, and therefore misses the possibility that racism may be salient to the state's means-end rationality, even if race per se is not.

<sup>224</sup> Chris S. Crandall et al., *Social Norms and the Expression and Suppression of Prejudice: The Struggle for Internalization*, 82 J. PERSONALITY & SOC. PSYCHOL. 359, 359 (2002) finding that “[t]he public expression of prejudice toward 105 social groups was very highly correlated with social approval of that expression. Participants closely adhere to social norms when expressing prejudice, evaluating scenarios of discrimination, and reacting to hostile jokes”); Katie M. Duchscherer & John F. Dovidio, *When Memes are Mean: Appraisals of and Objections to Stereotypic Memes*, 2 TRANSLATIONAL ISSUES PSYCHOL. SCI. 335, 335 (2016) (online experiment involving memes about Asian stereotypes in which “seeing another person object to the meme increased the likelihood that White participants would object . . . but only when the race of the person was unstated, and not when the person was Asian”); Fletcher A. Blanchard et al., *Condemning and Condoning Racism: A Social Context Approach to Interracial Settings*, 79 J. APPLIED PSYCHOL. 993, 993 (1994) (study demonstrating that cues from other people that racial discrimination is permissible or impermissible affect whether a person will condemn a racist remark); Fletcher A. Blanchard et al., *Condemning and Condoning Racism: A Social Context Approach to Interracial Settings*, 79 J. APPLIED PSYCHOL. 993, 995-96 (1994) (studies on students showing that hearing others condemn racism led to anti-racist opinions, while hearing others condone racism weakened anti-racist opinions).

individualization precludes any taking account of their race.<sup>225</sup> This moral demand for individuation entails demanding judicial scrutiny for all racial classifications.

There are, to be sure, reasons for skepticism about these moral and theoretical premises of the anticlassification principle.<sup>226</sup> But even bracketing those hesitations, and taking those justifications on face value, there is still no reason to think that the logic of anticlassification strongly militates against the use of race either as a feature or as an element of a classifier by machine-learning tools. To the contrary, as a matter of either precedent or logic, Equal Protection law can accommodate racially sensitive algorithmic criminal justice.

Consider the first concern about the communicative effect of racial classifications. It is not clear that an algorithmic classifier is the sort of racial criterion that courts perceive to be objectionable. Rather, it is somewhat akin to the explicit use of race in criminal suspect identifications, which has to date elicited scant constitutional concern.<sup>227</sup> Suspect descriptions instead operate as ‘given’ elements of the regulatory backdrop. Courts have not been wholly clear about why such suspect descriptions do not elicit careful scrutiny. One possible explanation is that judges believe suspect descriptions to be based on extrinsic facts, rather than airy suppositions about racial types, and as such not the kind of generalizations that trigger anticlassification concerns. This logic might be extended to the algorithmic context. Race-based feature selections would then trigger no more constitutional concern than race-based suspect descriptions. The argument would be that a classifier based on training data is akin to a suspect description of a familiar sort insofar as both are predicated on historical facts about crime.<sup>228</sup> Indeed, an advocate of algorithmic

---

<sup>225</sup> *Missouri v. Jenkins*, 515 U.S. 70, 120–21 (1995) (Thomas, J., concurring) (“At the heart of this interpretation of the Equal Protection Clause lies the principle that the government must treat citizens as individuals, and not as members of racial, ethnic, or religious groups.”). The same position is articulated, with respect to gender in DAVID MILLER, *PRINCIPLES OF SOCIAL JUSTICE* 168–69 (1999) (arguing that to treat a woman on the basis of “information that relates to the whole group or class” to which she belongs is “to fail to treat her respectfully as an individual, and potentially to commit an injustice”). This argument does not rest on the stigmatizing consequences of race-based action.

<sup>226</sup> For devastating critiques of the idea of colorblindness, see ELIZABETH ANDERSON, *THE IMPERATIVE OF INTEGRATION* 155-79 (2010) (describing the concept as “confused” and “incoherent”); Reva B. Siegel, *Discrimination in the Eyes of the Law: How “Color Blindness” Discourse Disrupts and Rationalizes Social Stratification*, 88 CAL. L. REV. 77, 81-83 (2000); Reva B. Siegel, *The Racial Rhetorics of Colorblind Constitutionalism: The Case of Hopwood v. Texas*, in *RACE AND REPRESENTATION: AFFIRMATIVE ACTION* 29 (Robert Post & Michael Rogin eds., 1998). For an originalist critique of anticlassification rules as interpretations of the Equal Protection Clause, in favor of a “duty-to-protect” view, see Christopher R. Green, *The Original Sense of the (Equal) Protection Clause: Pre-Enactment History*, 19 GEO. MASON U. CIV. RTS. L.J. 1, 3 (2008).

<sup>227</sup> For a collection of cases, see R. Richard Banks, *Race-Based Suspect Selection and Colorblind Equal Protection Doctrine and Discourse*, 48 UCLA L. Rev. 1075, 1095-96 (2001). Note that this is not a function of the inclusion of other considerations. Classifications that include race as one among many elements can run afoul of the Equal Protection Clause. See Balkin & Siegel, *supra* note --, at 16-17 (noting conflicting precedent on this point).

<sup>228</sup> Are algorithms different because the historical data upon which they are based is not specifically linked to a particular crime? Consider, however, the decision in *Brown v. City of Oneonta*, for example, which declined to impose constitutional tort liability when a description of a black male suspect provoked Oneonta police to stop more than two hundred “non-white persons,” including women, encountered on the streets. Although the stops there were in a trivial sense based on a historical fact, the connection between

criminal justice might observe that human observers are more likely than a machine to err in their deployment of race as a signal of criminality than an algorithm.<sup>229</sup> They might further contend it is perverse to object to efforts to mitigate the effects of race on criminal-justice outcomes through the substitution of machine for human judgments.

A second reason to think that an anticlassificatory logic does not work well in this domain would focus upon the absence of any communicative effect from algorithmic criminal justice. Many of the algorithms discussed in Part I are sheltered from disclosure by trade-secrets law, and hence are not disclosed presently to the public.<sup>230</sup> Even if they were to be disclosed in the course of litigation, it would likely be under the auspices of a protective order. To the extent that anticlassification rules rest on a concern about the communicative effects of state action, the use of an algorithmic tool that is wholly opaque should mitigate those concerns. More generally, the Supreme Court has been more accommodating of the conceded use of race when it is somewhat obscured from public view.<sup>231</sup> A state actor that relies upon an algorithmic tool, but that muffles the precise content of that tool from the public through trade secrets law or otherwise, might mitigate the most powerful challenges on Equal Protection grounds. Stated more positively, the much maligned quality algorithmic “opacity” has the benefit of dampening troublesome communicative effects for racial classification. Advocacy of transparency has the perverse effect of courting the expressive harms that Equal Protection tries to minimize.

A related, if somewhat subtler, question arises if race is employed as a feature of the training data—i.e., for each discrete observation (individual) in the training data, race is recorded—but race plays no role in the labels used to describe the classification task, or in the tools used to identify an appropriate function. Does that approach have a constitutionally impermissible communicative effect?

Northpointe omitted race from the training data used for Compas.<sup>232</sup> But this appears to reflect corporate risk aversion rather than an effort at legal compliance. Current law does not address whether the availability of race as an input into the deliberative process that results in state action violates the Equal Protection Clause on anticlassification grounds. To be sure, there is language in earlier precedent that suggests that any racial trace

---

that fact and the subsequent police actions was very strained. *Brown v. City of Oneonta*, 235 F.3d 769, 779 (2d Cir. 2000) (Calabresi, J., dissenting from denial of rehearing en banc) (citation and quotation marks omitted). The same might be said of algorithmic tools.

<sup>229</sup> See, e.g., Kleinberg et al., *Human decision*, *supra* note 133, at 52, fig. 9 (making precisely this argument in graphical form).

<sup>230</sup> See Rebecca Wexler, *Life, Liberty, and Trade Secrets: Intellectual Property in the Criminal Justice System*, 70 STAN. L. REV. \_\_ (forthcoming 2018), [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2920883](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2920883) (collecting examples, and arguing for more transparency). More generally, algorithms used by commercial actors are also “secret.” Danielle Keats Citron & Frank Pasquale, *The Scored Society: Due Process for Automated Predictions*, 89 WASH. L. REV. 1, 10–11 (2014)

<sup>231</sup> See Strauss, *supra* note 178, at 24 (noting “the Court’s insistence on nontransparency” in affirmative action cases). An unpublished paper by Strauss on ‘do it but don’t tell me,’ makes this point even more forcefully.

<sup>232</sup> Angwin et al., *supra* note 15.

in official deliberation raises a constitutional problem.<sup>233</sup> But the weight of precedential evidence (as well as common sense) suggests that the mere fact that a decisionmaker can *observe* the race of subjects does not mean that it is therefore invalid. As a practical matter, many front-line state officials encounter suspects, defendants, and citizens and thereby directly perceive their interlocutors' race.<sup>234</sup> Similarly, the federal judiciary must—and indeed does—routinely recognize the race of litigants in order to reach judgments on statutory and constitutional discrimination claims, even when it is not strictly necessary.<sup>235</sup> Finally, recent affirmative action jurisprudence implies (without expressly stating) that the bare fact of racial awareness is not sufficient to state a constitutional violation. The University of Texas, whose admission policy was reviewed and upheld by the Court in 2016, considered race as part of its PAI, and this alone did not suffice to generate a constitutional problem.<sup>236</sup> In short, it seems tolerably likely that an algorithmic criminal justice tool can use race as a feature in training data without triggering constitutional concern.

What of the argument against the state's use of racial classifications from its putative obligation to treat individuals as individuals rather than as members of groups? The moral logic of individuation trains on “intentional uses” of racial classifications, not merely coincidental or happenstance entanglements with race.<sup>237</sup> That logic might seem to have traction here since algorithmic criminal justice entails a decision-maker relying on group membership rather than accounting for all relevant characteristics of an individual.

But this is not quite right. In the absence of masking,<sup>238</sup> there is no human decision to assign costs or benefits on the basis of a racial classification with algorithmic criminal justice. And race is generally not going to be used as a substitute for more fine-grained traits. In any case, merely withholding race information does not ensure that an algorithm will not point toward race as a salient proxy. Machine learning algorithms take training data (with or without a race parameter), and use them to generate a new classifier, which can then be applied to test data.<sup>239</sup> The fact that an algorithm is not initially supplied with an impermissible ground of decision as a feature of training data does not mean that it will not end up including that criterion in its classifier. Machine-learning tools are powerful and useful precisely because they can detect regularities in a data-set that would not manifest in the absence of computational tools. Although machine-learning tools can be designed to

---

<sup>233</sup> *Parents Involved*, 551 U.S. at 748 (“The way to stop discrimination on the basis of race is to stop discriminating on the basis of race ....”).

<sup>234</sup> Reva B. Siegel, *Equality Talk: Antisubordination and Anticlassification Values in Constitutional Struggles over Brown*, 117 HARV. L. REV. 1470, 1471 (2004) (“For a half-century now, the Constitution has prohibited state action that classifies on the basis of race, yet as Americans have debated the implications of that principle, few have thought it barred collecting racial data.”).

<sup>235</sup> Justin Driver, *Recognizing Race*, 112 COLUM. L. REV. 404, 408 (2012) (documenting courts’ “unsettled and unsettling approach” to the recognition of litigants’ racial identities).

<sup>236</sup> *Fisher v. Univ. of Texas at Austin*, 136 S. Ct. 2198, 2206 (2016) (noting that “race is given weight as a subfactor within the PAI”).

<sup>237</sup> ANDERSON, *supra* note 227, at 155. Anderson is discussing “racial preference” here, but her point applies to racial classifications too.

<sup>238</sup> See *supra* text accompanying notes 201 to 204.

<sup>239</sup> See *supra* text accompanying notes 79 to 81.

be “private,” in the sense of eschewing reliance on certain traits,<sup>240</sup> they can also “help pinpoint reliable proxies” for traits even without information about the distribution of such traits in the population.<sup>241</sup> If race emerges as part of the classifier, this is not an “intentional” action in any meaningful sense—and yet it is still a classification on the basis of race.

The official deploying the algorithm, moreover, cannot be faulted for failing to engage in sufficient individuation: She supplies granular training data, selects among different computational tools, and then applies these tools to the specific facts about the individual being classified.<sup>242</sup> Even if the training data includes race information, the official has not designated race as a salient trait in any meaningful way. A decision-making process in which no human actor has elected to employ race as a criterion of action is not fairly characterized as an instance in “the government distributes burdens or benefits *on the basis of* individual racial classifications.”<sup>243</sup> The argument against algorithmic criminal justice from the moral demand for individuation, therefore, fails.

There is a one final argument for the inapplicability of anticlassification logic here. Race is commonly thought to be already highly correlated with socioeconomic characteristics related to criminogenic and victimization distributions. It might hence be reasonably anticipated that many algorithmic tools designed to be predictive of criminality will, even absent any race feature in the training data, generate a function that either mimics, or is a good approximation of, racial distributions in the population. Given this, it is possible that “by remaining blind to sensitive attributes, a classification rule can select exactly the opposite of what is intended.”<sup>244</sup> That is, the absent of a *de facto* predictive trait from the training data can generate systematic and serious errors in prediction.

A simple example from outside the machine learning context illustrates this possibility. Imagine that wearing a particular baseball cap is used as a proxy for drug possession by police (say, because it may signal gang membership). Both blacks and whites wear this cap. For 100% of whites, and for zero percent of blacks, the cap is an accurate signal of drug possession. Let us say that police stop all those encountered wearing the cap, and this population is 75% white and 25% black. Because the cap generates a 75 percent success rate, its categorical (and colorblind) use might be deemed a meritorious criterion. But the efficacy of searches, and the avoidance of needless hassle for minorities, can be increased by limiting the instrument to white suspects.<sup>245</sup> Colorblindness here generates

---

<sup>240</sup> Cynthia Dwork & Aaron Roth, *The algorithmic foundations of differential privacy*, 9 FOUND’NS & TRENDS IN THEORETICAL COMP. SCI. 211, 216-18 (2014) (developing a related concept of differential privacy).

<sup>241</sup> Barocas & Selbst, *supra* note 50, at 693.

<sup>242</sup> Kroll et al., *supra* note 38, at 682 (noting that in machine learning, “decision rules evolve on the fly—they are not specified directly, but are inferred from the data”).

<sup>243</sup> *Parents Involved in Cmty. Sch. v. Seattle Sch. Dist. No. 1*, 551 U.S. 701, 720 (2007). What if the government actor designing the algorithm fails to prevent the algorithm from homing in on race? I read the doctrine not to problematize such culpable omissions.

<sup>244</sup> Kroll et al., *supra* note 38, at 686.

<sup>245</sup> This example is drawn from Ian Ayres, *Outcome tests of racial disparities in police practices*, 4 JUST. POL. RES. 131, 139 (2002).

substantial and avoidable social costs, and can be corrected, however, by simply accounting for race.<sup>246</sup>

In the machine-learning context, a fix entails the creation of a predictive tool that assigns individuals from different demographics to different classifications even though they exhibit the same behavioral traits.<sup>247</sup> Lest this seem obviously beyond the legal and moral pale, consider that one study of probation and parole decisions found that the decision to omit race from a machine-learning algorithm, the accuracy of recidivism predictions declined “by about 7 percentage points.”<sup>248</sup> The procedural purity demanded by an anti-classification rule, in sum, would come at a high price in terms of accuracy in algorithmic application.<sup>249</sup>

### 3. *The Lessons of Algorithmic Technology for Equal Protection Doctrine*

Current doctrinal approaches to constitutional racial equality arose after the Court had abandoned its early twentieth century interpretation of the Equal Protection Clause as “a rationality test ... invoked sporadically to strike down economic regulation.”<sup>250</sup> They were configured in a context of judicial efforts to dismantle educational segregation in the Jim Crow south, and then against a backlash to the Civil Rights Movement.<sup>251</sup> It was probably inevitable that the legal conception of racial discrimination as a matter of intention or classification would reflect the judicial concern with the discretionary choices of the police officer, school board president, or state legislator—i.e., the modal problems thrown up by mid-century civil rights law.

The institutional context of Equal Protection, however, has changed. Today, perhaps the sharpest and most controversial questions of racial justice are presented in the criminal justice domain. The emergence of algorithmic tools in that domain present questions poorly fitted to the doctrinal templates of intention and classification. This loose fit arises because the ways in which race filters into individual officials’ discretionary criminal justice decisions are very different from the ways in which it can infuse algorithmic tools. Equal Protection, as a result, poses questions that are simply not relevant to the operation of algorithmic criminal justice. It is a superseded legal technology so far

---

<sup>246</sup> See Cynthia Dwork et al., *Fairness Through Awareness*, 2012 PROC. 3RD INNOVATIONS THEORETICAL COMPUTER SCI. CONF. 214 (formally demonstrating this result); accord Pauline T. Kim, *Data-Driven Discrimination at Work*, 58 WM. & MARY L. REV. 857, 918 (2017) (arguing that the risk of omitted variable bias means that controlling for sensitive demographic variables may sometimes be necessary to avoid biased results).

<sup>247</sup> Bryce W. Goodman, *Economic Models of (Algorithmic) Discrimination 3*, in 29TH CONFERENCE ON NEURAL INFORMATION PROCESSING SYSTEMS (2016). For a parallel result using the Compas data, see Nabi & Shpitser, *supra* note 19, at 8.

<sup>248</sup> Richard Berk, *The Role of Race in Forecasts in Violent Crime*, 1 RACE & SOC. PROBS. 231, 235 (2009).

<sup>249</sup> Could such a use of race be justified as a narrowly tailored response to a compelling state interest? It is hard to say how much gain in accuracy would be required to make this claim compelling. For a discussion of how hard it is to make this judgment, see text accompanying *supra* notes 176 to 180.

<sup>250</sup> Michael J. Klarman, *An Interpretive History of Modern Equal Protection*, 90 MICH. L. REV. 213, 216 (1991).

<sup>251</sup> *Id.* at 217-18; see also BRUCE ACKERMAN, 3 WE THE PEOPLE: THE CIVIL RIGHTS REVOLUTION 328-37 (2014) (consider the judicial forms of this backlash).

as algorithmic criminal justice goes. As more state power is channeled through algorithmic channels, it will become increasingly obsolete.

On the one hand, the manner in which algorithmic criminal justice unfolds generally means that there are few opportunities for intentional discrimination of the familiar kind. The process of feature selection, to be sure, creates opportunities to use race as an input, to intentionally omit race in order to generate discriminatory patterns, or to choose an insufficient number of variables in ways that mimic the same effect.<sup>252</sup> But this sort of masking will be very hard to discover (much as prosecutorial or judicial animus is hard to identify now). It does not, at least on the basis of current evidence, appear to be a significant problem. On the other hand, the logic of anticlassification might first seem to provide a firm foundation for regulating algorithmic criminal justice. But that logic turns out again to be a bad fit. The use of race in criminal justice algorithms is akin to the use of race in suspect descriptions. It lacks both the intentionality and the expressive spillovers that render non-individuation troubling. Just as in the context of race-based suspect descriptions, moreover, it will sometimes be necessary to use race to achieve substantively accurate policy results.

In the dialog between Equal Protection and algorithmic criminal justice, therefore, I suspect that constitutional law has much to learn and less to teach. A set of tools developed for a regulatory world of dispersed state actors, occasionally motivated by naked animus, cannot be mechanically translated into a world of centralized, computational decision-making. Even after law has made its contribution, therefore, the question of racial equity in algorithmic criminal justice remains open for debate—while the relevance and moral acuity of equality jurisprudence should be viewed as in serious doubt absent more intensive rethinking.

### **III. Racial Equity in Algorithmic Criminal Justice Beyond Constitutional Law**

The failure of constitutional law to provide a meaningful benchmark of racial equality is important in its own right. Yet it leaves the study of algorithmic criminal justice unmoored. It means there is no normatively attractive, empirically tractable way of evaluating the race effects of big-data predictive tools. This Part fills that gap. In order to do so, I will start by offering my own account of the normative stakes of racial equity in criminal justice to fill the vacuum left by our deficient constitutional doctrine. My view is that the reason for concern about racial equity in criminal justice generally is that fact that our policing and adjudicative institutions play significant roles in the reproduction and entrenchment of social stratification. In a racially segmented society, when a person's life chances are defined importantly by their race, I believe this to be a moral wrong of the first order.

With that normative benchmark in hand, I turn to the extensive computer-science literature on the question. That scholarship has developed a series of definitions of what is alternatively defined as algorithmic fairness or algorithmic discrimination. The literature has focused first on precise mathematical formulations of each definition, and second on

---

<sup>252</sup> See Kroll et al., *supra* note 38, at 681; Barocas & Selbst, *supra* note 50, at 692.

the generation of impossibility theorems—i.e., formal proofs that it is not possible to maximize two or more parameters that in some fashion measure the racial effects of an algorithm. Because the computer-science literature has been “silent on the choice [between different understandings of fairness],”<sup>253</sup> mere specification of alternative conceptions of racial equity is not sufficient for any tractable conclusions about public policy. By applying my account of racial equity in criminal justice to these standards, I aim to make progress on determining which technical conception captures something of normative significance.

Two caveats are useful here. First, for the sake of clarity of exposition, I focus here on a binary between white and black defendants, even though this obscures the more complex racial dynamics of American policing today.<sup>254</sup> A focus on a black/white binary is warranted here as a way of clarifying the fundamental conceptual stakes. It is obviously inadequate as a general account of racial equity in policing, and I do not intend it as such. Moreover, I should emphasize again that my aim here is not to offer a judgment in respect to any specific algorithm, but a more general analytic approach. Much depends on the particular costs and benefits that in situ flow from a given instrument.

Second, a racial-equity analysis of algorithmic criminal justice should not be a comparative one. It is not sufficient, that is, to point to a superseded technology that relies upon flawed human discretion and that already generates large racial effects, as a justification for new, slightly less flawed technologies for allocating coercion. The mere fact that the status quo ante is characterized by racial injustice does not legitimize proposals that preserve or extend some substantial part of that injustice. No one thinks (or should think) the Jim Crow regime laudable, for example, merely because it followed slavery. Improvements in the status quo are a necessary but not sufficient condition for racial equity to be satisfied. It seems likely that the shift to algorithmic tools in criminal justice will be an enduring one. At the moment that a new policy is introduced, with potential path-dependent effects that will unfold over many iterations of policy-making, it is especially important to understand the conditions under which that policy promotes racial equity: Far better, that is, to embed that principle at a policy’s inception than to attend years of damage that cannot ever wholly be unraveled. Each technology ought to be evaluated on its merits and in light of its consequences. It is plainly inadequate to say that the technology should be adopted because and only because it makes a step change in the woeful status quo.

#### **A. The Stakes of Racial Equity in Contemporary American Criminal Justice**

Why care about racial equity in criminal law? Without an answer to that question—and we have already seen that constitutional law doesn’t give a convincing one—no

---

<sup>253</sup> Richard Berk et al., *Fairness in Criminal Justice Risk Assessments: The State of the Art* 29 (May 2017) [hereinafter “Berk et al., *Fairness in Criminal Justice*”], <https://arxiv.org/abs/1703.09207>; see also Sorelle A. Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian, *On the (im)possibility of fairness* 12 (2016) <https://arxiv.org/abs/1609.07236> (“Choice in mechanism [i.e., feature design] must thus be tied to an explicit choice in worldview,” and in particular a choice to prioritize either individual or group fairness).

<sup>254</sup> Cf. Ramiro Martinez, *Incorporating Latinos and immigrants into policing research*, 6 *CRIMINOLOGY & PUB. POL.* 57, 57 (2007) (documenting the “lack of research on Latino/as and Latino groups” in relation to the criminal justice system).



analysis of algorithmic criminal justice’s racial equity effects gets off the ground. Accordingly, I start by offering my own evaluation of the racial stakes of criminal justice. I do not intend to break new ground here, but rather aim to set forth clearly a distinct normative position respecting racial equity in the criminal justice context.

American criminal justice implicates racial equity concerns because of their dynamic effects on racial stratification. Historical and contemporary empirical evidence suggests that both in the past and now, criminal justice has been invoked in public discourse and applied in state practice so as to predictably exacerbate the subordinate status of African Americans in general. The dynamic (re-)production of iniquitous social stratification—beyond the bare facts of animus and classification—is what should grip our collective conscience.<sup>255</sup>

At a very general level of abstraction, four causal mechanisms link criminal justice institutions to racial stratification. *First*, inherent black criminality has been invoked for more than a hundred years as public justification for more punitive interventions against African-Americans, and for the withholding of social services from them on moral desert grounds. *Second*, black communities have in practice been both over-policed—in the sense of subjected to higher rates of coercive interventions—and also under-protected—in the sense of not receiving the same measure of protective legal resources that nonblack communities receive. As a result of this inefficient allocation of policing resources, state coercion has not resulted in lower levels of private coercion for African-Americans. *Third*, pivotal actors within the criminal justice system, such as police, prosecutors, and judges, have tended to treat black suspects and defendants more harshly than white ones. Hence, the per-capital cost of crime suppression has been greater for blacks than whites. *Fourth*, the spillover effects from disparate policing for black families and communities appear to be larger in magnitude than the spillover effects in white communities, even controlling for the extent of coercion. The net result of these mechanisms is that criminal justice imposes “compounding”<sup>256</sup> disadvantage upon African-Americans as a group that works as a brake on individuals’ efforts to rise in the social hierarchy. Even if not all African-Americans are impeded by this headwind, enough are that we can meaningfully talk of persisting racial stratification to which criminal justice institutions have contributed. These diverse causal pathways underpin the need for careful attention to the manner in which formal criminal justice institutions can undermine the status of African-Americans as a group.

Rather than offering normative and empirical justifications for each element of this position—a task that would require a book rather than an article—I sketch some suggestive evidence for these causal linkages between criminal justice and racial stratification. I start

---

<sup>255</sup> Racial stratification is objectionable on (at least) two grounds. First, it embodies what Tim Scanlon calls a manifest failure of equal concern on the part of the state. T.M. SCANLON, *WHY DOES INEQUALITY MATTER?* (2018). Second, stratification generates deadweight welfare losses in the form of unused human capital, psychological and social harms, and violence that flows from the latter. Of course, to the extent that such dynamic consequences have normative salience, it is because of a predicate obligation of equal concern toward the disadvantaged.

<sup>256</sup> I draw this term from Deborah Hellman, *Indirect Discrimination and the Duty to Avoid Compounding Injustice*, in *FOUNDATION OF INDIRECT DISCRIMINATION LAW* 105, 107 (Tarunadh Khaitan and High Collins, eds., 2018). Hellman’s use of the term assumes an original act of discrimination; my use does not (although discriminatory acts are woven *across* the operation of criminal justice).

with history (although I do not want to suggest that the state's obligations here rest on its historical responsibility for creating racial stratification in the first instance, rather than its role in perpetuating that condition). At the beginning of the twentieth century, national public discourse about "law and order became racialized," even as "conviction and incarceration rates for African Americans jumped disproportionately."<sup>257</sup> As the leading historical work by Khalil Gibran Mohammed vividly demonstrates, Progressive-era academics, journalists, and politicians in the north linked crime to African-Americans, at the same time as they downplayed white ethnic groups as sources of crime. By the early 1940s, Mohammed explains, "'Black' stood as the unmitigated signifier of deviation (and deviance) from the normative category of 'White.'"<sup>258</sup> Concomitant to this rhetorical shift, urban policing and carceral resources were disproportionately allocated to African-Americans who were in the process of migrating up from the rural south. In northern cities in particular, police singled out blacks for intense surveillance and coercion.<sup>259</sup> This pushed up the rate of black incarceration and the proportion of the prison population that was black.<sup>260</sup> The black share of that population never subsequently dropped.<sup>261</sup> Racialized mass incarceration, that is, was at its inception a product of a moral panic stoked by northern elites in respect to the growing presence of an African-American population that previously had been the South's 'problem.'

Today, racial disparities characterize both victimization rates and exposure to criminal-justice coercion. Black men are more likely than white men to be victims of serious (index) crimes such as murder.<sup>262</sup> They are also more likely to be arrested and incarcerated than white peers.<sup>263</sup> In many urban contexts, blacks and whites also experience widely varying chances of being stopped by police.<sup>264</sup> Moving from the policing to the adjudicative phase of the criminal justice process, common sentencing regimes impose disparate treatment on similarly situated offenders of different races by the use of different penalty structures for behavior closely associated with different racial groups.<sup>265</sup> As a result, one in eight black men in their twenties is in prison or jail on any given day, while some 69 percent of black high school dropouts are imprisoned over their lifetime,

---

<sup>257</sup> Jeffrey S. Adler, *Less Crime, More Punishment: Violence, Race, and Criminal Justice in Early Twentieth-Century America*, 102 J. AM. HIST. 34, 35 (2015).

<sup>258</sup> MUHAMMED, *supra* note 2, at 13.

<sup>259</sup> Christopher Muller, *Northward Migration and the Rise of Racial Disparity in American Incarceration, 1880-1950*, 118 AM. J. SOC. 281, 310 (2012).

<sup>260</sup> *Id.*

<sup>261</sup> *Id.*

<sup>262</sup> BRENDAN O'FLAHERTY, *THE ECONOMICS OF RACE IN THE UNITED STATES* 333-34 (2015)

<sup>263</sup> *Id.* at 335-36.

<sup>264</sup> Huq, *Disparate Policing*, *supra* note 40, at 2411-12 (summarizing data from Chicago and New York).

<sup>265</sup> See, e.g., David A. Sklansky, *Cocaine, Race, and Equal Protection*, 47 STAN. L. REV. 1283, 1303 (1995) (describing the use of racially charged language in the enactment of narcotics statutes that impose different sentences on crack and power cocaine offenses).

compared with just 15% for white high school dropouts.<sup>266</sup> For young black men, therefore, prison has thus become a predictable part of the life course.<sup>267</sup>

Note also that the intensive concentration of policing and incarceration resources along racial lines is not a rational, cost-justified response to crime. As I have argued elsewhere, there is evidence that some of the most common forms of policing black communities are inefficacious.<sup>268</sup> Black incarceration rates are also too high to be plausibly justified. One estimate suggests that reducing incarceration rates from 2004 to 1984 levels, and investing the resulting savings in an increased police presence, would lead to a net decline in violent crime nationally of about 130,000 incidents.<sup>269</sup> Even if racial minorities benefit from the public safety produced by the criminal justice system, therefore, it is also at a highly disproportionate and unnecessary direct cost.

Is part of this burden, though, justified by higher black crime rates? Even if we assume that “African Americans engage in significantly higher rates of street crime,” there is evidence that conditions of “racial segregation and concentrated disadvantage”—i.e., environmental conditions that themselves are a function of non-race-neutral policies—explain much of the difference between different racial groups’ crime rates.<sup>270</sup> That is, it is not so much that race is causally related to criminality, but that African-Americans are subject to forms of social and economic stratification and segmentation that conduce to criminality. Paradoxically, these underlying conditions are in an important respect a function of the federal government’s decision to shift resources away from building human and social capital to policing crime. The intensification of policing and incarceration since the early 1970s, the historian Elizabeth Hinton has argued, was a conscious, and racially tinged, policy substitute for Great Society programs that could have mitigated those conditions.<sup>271</sup> That substitution could be reversed. As the sociologist Patrick Sharkey has demonstrated, it is precisely the local recreation of social services, and the concomitant creation of social capital, that has been a leading contributor to recent declines in crime. In one empirical study, Sharkey and his colleagues thus estimated that “the addition of 10 community nonprofits per 100,000 residents leads to a 9 percent decline in the murder rate,

---

<sup>266</sup> See generally Bruce Western & Christopher Wildeman, *The Black Family and Mass Incarceration*, 621 ANNALS AM. ACAD. POL. & SOC. SCI. 221 (2009); Bruce Western & Christopher Muller, *Mass Incarceration, Macrosociology, and the Poor*, 647, ANNALS AM. ACADEMY OF POL. & SOC. SCI. 166 (2013).

<sup>267</sup> For an extended account of these effects, see Kristen Henning, *Boys to Men: The Role of Policing in the Socialization of Black Boys*, in POLICING THE BLACK MAN: ARREST, PROSECUTION, AND IMPRISONMENT 57 (Angela Y. Davis, ed. 2017).

<sup>268</sup> Huq, *Disparate Policing*, *supra* note 4, at 2929-40.

<sup>269</sup> Philip J. Cook & Jens Ludwig, *Economical crime control*, in CONTROLLING CRIME: STRATEGIES AND TRADEOFFS 1, 2-4 (2010).

<sup>270</sup> Callie Harbin Burt et al., *Racial discrimination, ethnic-racial socialization, and crime: A micro-sociological model of risk and resilience*, 77 AM. SOC. REV. 648, 650-52 (2012); see also Lauren J. Krivo & Ruth D. Peterson, *The structural context of homicide: Accounting for racial differences in process*, 65 *Am. Soc. Rev.* 547, 556 (2000) (examining marginal effects of social advantage of black and white communities).

<sup>271</sup> Elizabeth Hinton, ‘A War Within Our Boundaries’: Lyndon Johnson’s Great Society and the Rise of the Carceral State, 102 J. AM. HIST. 100, 101-02 (2015).

a 6 percent decline in the violent crime rate, and a 4 percent decline in the property crime rate.”<sup>272</sup>

Finally, the direct costs of black incarceration are only part of the distinctive burden imposed by the current criminal justice system on racial minorities. Current crime suppression also imposes considerable collateral costs (or externalities) asymmetrically on racial minorities. To begin with, the immediate cost of encounters with police is racially asymmetric. The black experience of a police stop is reliably correlated with “stigma, trauma, anxiety and depression,”<sup>273</sup> because of the historically fraught nature of relations between American police and racial minorities. African-Americans are, moreover, commonly subject to policing measures that are not generally employed against white citizens—such as pretextual vehicular stops—and are quite aware that they are objects of disparate treatment based on the presumption of black criminality.<sup>274</sup> They are also quite aware of the stigmatizing connection between race and criminality drawn since the beginning of the twentieth century. Even today, “demography-based suspicion is among the key social facts that define American life in the late twentieth century and early twenty-first centuries.”<sup>275</sup> Ethnographic studies paint a bleak picture of interactions between police and young black men as marked by distrust and fear, and as a source of widespread alienation and disaffection.<sup>276</sup> Against the background of this broadly shared supposition of the relationship of criminality and race, public encounters with police can, even if warranted, humiliate and rob innocent racial minorities of the “ability to present themselves to other people as the ordinary people they are.”<sup>277</sup>

These effects generate further negative spillovers. As Randall Kennedy cogently observed three decades ago, African-American men experience a “racial tax” from American criminal justice systems—even if they have no contact with it—because police and citizens are prone to perceive their race as a proxy for criminality, and hence to configure them as potential criminals rather than potential victims.<sup>278</sup> Recent empirical

---

<sup>272</sup> Patrick Sharkey, Gerard Torrats-Espinosa, and Delaram Takyar, *Community and the Crime Decline: The Causal Effect of Local Nonprofits on Violent Crime*, 82. AM. SOC. REV. 1214, 1234 (2017).

<sup>273</sup> Amanda Geller et al., *Aggressive Policing and the Mental Health of Young Urban Men*, 104 AM. J. PUB. HEALTH 2321 (2014). For a powerful account of why these costs accrue distinctly to racial minorities, see Nicholas K. Pert, *Why Is the N.Y.P.D. After Me?*, N.Y. TIMES, Dec. 17, 2011, <http://www.nytimes.com/2011/12/18/opinion/sunday/young-black-and-frisked-by-the-nypd.html>.

<sup>274</sup> CHARLES EPP, STEVEN MAYNARD-MOODY, & DONALD HAIDER-MARKEL, PULLED OVER: HOW POLICE STOPS DEFINE RACE AND CITIZENSHIP 117-18 (2014).

<sup>275</sup> WILLIAM J. STUNTZ, THE COLLAPSE OF AMERICAN CRIMINAL JUSTICE 22 (2011).

<sup>276</sup> Rod K. Brunson, “Police Don’t Like Black People”: African-American Young Men’s Accumulated Police Experiences, 6 CRIMINOLOGY & PUB. POL. 71, 85 (2007) (finding that street encounters with police perceived as unfair drive perceptions of the police as unjust among young black men); Jacinta M. Gau and Rod K. Brunson, *Procedural Justice and Order Maintenance Policing: A Study of Inner-City Young Men’s Perceptions of Police Legitimacy*, 27 JUST. Q. 255, 268 (2010) (noting that many respondents described contact with police “as demeaning and of inordinate frequency”).

<sup>277</sup> Paul Bou-Habib, *Racial Profiling and Background Injustice*, 15 J. ETHICS 33, 44 (2011). Bou-Habib is here addressing profiling, but I am extending his point.

<sup>278</sup> KENNEDY, *supra* note 9, at 158-60. The criminal justice system thus creates “shared categories of classification systems through which individuals perceive and make sense of their environment.” Michèle Lamont, Stefan Beljean, and Matthew Clair, *What is missing? Cultural processes and causal pathways to inequality*, 12 SOCIO-ECON. PERSP. 573, 583 (2014).

work has confirmed Kennedy’s account of the externalities of criminal justice for minority groups as a whole. African-American men hence continue to receive disfavored treatment in a wide array of economic and social contexts that limit an important slice of life opportunities.<sup>279</sup> The increased risk of contact with police, and hence incarceration, undermines the economic and social resources available to the larger racial cohort embedded in the same geographic community.<sup>280</sup> One in four black children also experiences parental incarceration—an experience that directly and negatively impacts their health and education outcomes.<sup>281</sup> Most notably, and dismayingly, black parental incarceration is associated with a 49% increase in infant mortality, an increase that has no parallel among white families affected by incarceration.<sup>282</sup> Not even the children, in other words, are spared. Rather, a concentration of policing and incarceration within black communities generates distinctive burdens with no parallel for majority racial groups—burdens that diffuse and concatenate across communities and generations. It is on this basis, I think, that it is plausible to characterize the contemporary American criminal justice system as “a systematic and institutional phenomenon that reproduces racial inequality and the presumption of black and brown criminality.”<sup>283</sup>

This account of racial equity in criminal justice does not hinge on the presence of discriminatory animus at any specific point in policing or the adjudicative process. Of course, disparate racial treatment happens (probably quite often).<sup>284</sup> But this account of racial equity is forward looking and consequentialist insofar as it is trained on the ways in which systems reproduce practical socioeconomic stratification over time. Moreover, this account suggests that criminal justice institutions are not presently socially efficient. Their footprint could be diminished in ways that do not create social costs from more crime. At present, however, the inefficiently large costs of criminal justice (which are not justified by sufficient offsetting social benefit) fall disproportionately on racial minorities. Many reforms that increase social efficiency will also further racial equity as a result.

A possible counter-argument is that a particular quantum of state coercion will, *ceteris paribus*, be more costly for a member of a white majority than a black minority because whites’ greater wealth and more remunerative employment outcomes mean that their economic losses from even transient coercion or incapacitation are likely to be

---

<sup>279</sup> For an effective summary of the relevant data, see Devah Pager and Hannah Shepherd, *The Sociology of Discrimination in Employment, Housing, Credit, and Consumer Markets*, 34 ANN. REV. SOC. 181 (2008); see also DEVAH PAGER, MARKED: RACE, CRIME, AND FINDING WORK IN AN ERA OF MASS INCARCERATION 93-96 (2007) (reporting effects of racialized assumptions of criminality on employment opportunities).

<sup>280</sup> See Todd Clear, *The Effect of High Imprisonment Rates on Communities*, 37 CRIM. & JUST. 97, 116 (2008) (discussing this effect); see also Amy E. Lerman & Vesla M. Weaver, *Staying Out of Sight? Concentrated Policing and Local Political Action*, 651 ANNALS AM. ACAD. POL. & SOC. SCI. 202, 204 (2014) (finding in a study of New York that “witnessing stops that occur with little justification and that feature physical force can make people feel occupied and powerless, and can incentivize disengagement with government”).

<sup>281</sup> WAKEFIELD & WILDEMAN, *supra* note 7, at 41; *id.* at 146 (noting that parental incarceration is a mechanism for the intergenerational transmission of inequality).

<sup>282</sup> *Id.* at 108.

<sup>283</sup> Naomi Murakawa & Katherine Beckett, *The Penology of Racial Innocence: The Erasure of Racism in the Study and Practice of Punishment*, 44 L & SOC. REV. 695, 701 (2010).

<sup>284</sup> See sources collected in *supra* notes 3 to 5.

greater than those of African-Americans.<sup>285</sup> I am skeptical. I find it troubling to use racial stratification by wealth and income as a lever to discount the costs imposed on African-Americans. I also do not accept that the implicit metric at work in this analysis (in effect, the capacity to pay) tracks a normatively attractive species of welfare. Finally, I have already flagged negative externalities to the status of African-Americans as a group, and to communities and families, that simply have parallel for racial majorities. I think it is more likely that black communities and families will want for the social and financial buffers that mitigate the shock of criminal-justice contacts. Hence, I think this counter-argument is both empirically and normatively flawed.

## **B. A Racial Equity Principle for (Algorithmic) Criminal Justice**

The algorithmic tools described in Part I are mechanisms to allocate coercion within the criminal justice system. As part of that system, they should be evaluated by the same criteria that are applied to other parts of the system. That is, the introduction of new computational and epistemic technologies does not alter the basic stakes of racial equity.

In this light, the key question for racial equity is whether the costs that an algorithmically driven policy impose upon a minority group outweigh the benefits accruing to that group. If an algorithmic tool generates public security by imposing greater costs (net of benefits) *for blacks as a group*, it raises a racial equity concern. That policy undermines racial equity by deepening the causal effect of the criminal justice system on race-based social stratification.<sup>286</sup> This test is consequentialist. It focuses on the effects of an algorithm's use.<sup>287</sup> It is holistic. Unlike older risk assessment tools, it accounts for both the benefits and the costs of intervention. And, to emphasize again, it is quite general: There is no reason not to apply it to criminal justice more generally. I develop the test here nevertheless because I am concerned with algorithmic tools. That test, indeed, is particularly well suited for algorithmic tools, which can develop precise cut-points for using coercion based on analyses of large volumes data.

This standard has a distant kinship to John Rawls' difference principle, which holds that "[a]ll differences in wealth and income, should work for the good of the least favored."<sup>288</sup> But the principle offered here operates within a much narrower institutional scope (criminal justice alone) and is justified on more specific grounds—i.e., to ensure that

---

<sup>285</sup> Cedric Herring and Loren Henderson, *Wealth inequality in Black and White: Cultural and structural sources of the racial wealth gap*, 8 RACE & SOC. PROB. 4, 4-5 (2016).

<sup>286</sup> This standard is analytically distinct from disparate impact as conventionally understood, not least because it does not account for benefits for a policy for those beyond the burdened group. It is an interesting question whether disparate impact, especially as applied to state action, might be reconfigured to approach the standard suggested in the text.

<sup>287</sup> Note that it is possible to take the view that there is a nonconsequentialist obligation on the state's part to show equal regard for all its citizens, and to think that my consequentialist metric is a way of honoring that obligation.

<sup>288</sup> John Rawls, *Distributive Justice: Some Addendum*, in COLLECTED PAPERS 154 (Samuel Freeman, ed., 1999). Rawls formulated the difference principle in a number of different ways. Nothing here rests on those variations, so I ignore them.

institutions purportedly operating in furtherance of public safety are not doing so in a fashion that exacerbates differences in racial strata.

What, though, of animus? Of course, individual officials do act at times with an invidious state of mind.<sup>289</sup> At present, the institutional process of adjudication and the doctrines structuring inquiries into bad intent ensure that few such instances are ever brought to light, let alone used as a basis for constitutional relief.<sup>290</sup> Still, I am skeptical that the resulting harms are of the same magnitude as the damage that comes from criminal justice's effect on racial stratification. Even if Equal Protection doctrine were more effective at identifying instances of bad motivation, a criminal-justice system purged of animus would still have substantial ramifications for racial stratification. It is the existence of racial stratification, moreover, and the channeling of anxieties about security and difference into racialized forms, that plausibly drive much animus in the first instance. Addressing stratification, on this view, is a more enduring and effective means of regulating animus than the emaciated and enfeebled investigative doctrinal instruments the Court permits.<sup>291</sup>

There are two ways of analyzing the relevant costs and benefits of an algorithmically allocated coercive measure. Application of the proposed racial equity criterion will produce a policy that enhances net social welfare on either of these two approaches. The first is to focus solely on the immediate costs and benefits of a coercive intervention, and to ignore externalities. This is a plausible approach with serious crimes, where externalities are dwarfed by immediate costs and benefits. A second approach accounts for both immediate costs and also externalities for different groups. The latter take many forms, including the effect of high incarceration rates on black communities and children, and the social signification of race as a marker of criminality. But as I argued above, the evidence summarized above suggests that these impacts are felt principally by members of racial minorities. It is, moreover, plausible to hypothesize that these spillover costs will largely be experienced by members of the same racial group as the suspect given persisting patterns of racial residential segregation.<sup>292</sup> Hence, the spillover costs of coercion of minority individuals for the minority group will be greater on a per capita basis than the costs of coercing majority group members. If the costs of coercing minorities are larger, while benefits remain static, racial justice will be satisfied by an algorithmic tool that imposed a higher threshold for black suspects than for white suspects. For less serious crimes, moreover, these spillover effects may be similar in magnitude to the direct benefits and costs of coercion. Hence, a simplified analysis that ignores spillovers would be inappropriate. Rather, a bifurcated rule with different thresholds for whites and blacks may be necessary to ensure that minority coercion does not exacerbate racial stratification for less serious offenses.

---

<sup>289</sup> For evidence of that effect, see CHARLES EPP, STEVEN MAYNARD-MOODY, & DONALD HAIDER-MARKEL, *PULLED OVER: HOW POLICE STOPS DEFINE RACE AND CITIZENSHIP* 117-18 (2014).

<sup>290</sup> For an extended argument to this effect, see Huq, *Discriminatory Intent*, *supra* note 40, at 21-36.

<sup>291</sup> *Id.*

<sup>292</sup> Matthew Hall, Kyle Crowder, and Amy Spring, *Neighborhood foreclosures, racial/ethnic transitions, and residential segregation*, *AM. SOC. REV.* 526, 527 (2015) (“[T]he modal experience for blacks (and Hispanics) in U.S. cities is high residential segregation.”).

Under either of these approaches, it will often be the case that racial equity and social efficiency (in the sense of ensuring that immediate social benefits exceed immediate social costs) will align. For example, when a majority group does not benefit from a policy, or when its net gain is less than the costs imposed on the minority group—and the latter suffers a net loss—that policy is socially inefficient. Equity and efficiency therefore align.

This approach makes certain simplifying assumptions. I believe them to be plausible. Hence, it assumes that most crime is intraracial such that costs and benefits do not cross the color line by and large. Obviously, this is not always true. But it does hold as a general matter.<sup>293</sup> Moreover, my analysis assumes away a number of unusual circumstances in which racial equity and social efficiency come apart. Because these circumstances are rare, I do not dwell on them. I mention two here briefly. First, it is possible that a policy benefits both the minority and the majority group, but the former benefit less than the latter. As a result of this gap, the extent of racial stratification increases even as the minority is benefited. The evaluation of such a policy would turn, in my view, on the magnitude of social gain and the extent to which the policy generated stratification. I do not think a general conclusion is appropriate regarding such policies.

Second, net gains from a policy for a majority group may exceed the net cost imposed on a minority group. (Imagine, for example, a national security policy that generated significant results by imposing crushing burdens on a very small ethnic or religious minority.) In this case, there is a tension between efficiency and antidiscrimination. Such conflicts have generated conflict among scholars.<sup>294</sup> In the crime control context, I suspect that this will rarely occur given the intra-group nature of much crime. Yet my own view is that gains in net social welfare should generally not be obtained by imposing burdens on minority groups subject to wider dynamics of compounding subordination.<sup>295</sup> In effect, such a policy would yield a regressive wealth transfer from blacks to whites in which the former pay for the security enjoyed by the latter.<sup>296</sup> I would hence prioritize the distribution that resulted from a policy over the sheer quantity of social welfare it yielded, at least in the absence of catastrophic general welfare losses from

---

<sup>293</sup> See Robert M. O'Brien, *The Interracial Nature of Violent Crimes: A Reexamination*, 91 AM. J. SOC. 817, 818-19 (1987) (finding evidence that crime is more intraracial than would be anticipated).

<sup>294</sup> Compare Louis Kaplow & Steven Shavell, *Should Legal Rules Favor the Poor? Clarifying the Role of Legal Rules and the Income Tax in Redistributing Income*, 29 J. LEGAL STUD. 82, 821-251 (2000) (favoring welfare maximization), with Lee Anne Fennell & Richard H. McAdams, *The Distributive Deficit in Law and Economics*, 100 MINN. L. REV. 1051, 1129, 1170 (2016) (doubting this maxim).

<sup>295</sup> This view implies that consequences are morally salient, but that welfare maximization is not the only measure of such consequences. The basic arrangements of a society are also important, and sometimes merit protection or improvement even at the cost of net social welfare. For a different view, that turns solely on purpose, and seems unconcerned with consequences, see Jed Rubenfeld, *Affirmative Action*, 107 YALE L.J. 427, 440-41 (1997) (“A law whose express purpose is racial apartheid or expulsion is unconstitutional per se, because racial purification of society is an objective that no legislature can pursue under the Fourteenth Amendment-- period.”).

<sup>296</sup> Cf. Tal Zarsky, *The Trouble with Algorithmic Decisions: An Algorithmic Road Map to Examine Efficiency and Fairness in Automated and Opaque Decision Making*, 4 SCIENCE, TECH. & HUMAN VALUES 118, 123 (2016) (noting the possibility that an algorithm can “enable[] transfers that systematically harm minorities and other protected groups).



forbearance. I do not perceive any circumstances in which that latter exception plausibly applies.<sup>297</sup>

### C. Benchmarks for Algorithmic Discrimination

A large computer-science literature on algorithmic design has generated a plethora of definitions of ‘algorithmic fairness’ and ‘algorithmic discrimination.’ One count finds twenty-one definitions.<sup>298</sup> Not all are relevant in the criminal justice context. Not every concept is analytically distinct from all others. My aim in this section is to home in upon the most relevant subset of such definitions, and to develop a quadripartite taxonomy of potential metrics for gauging racial equity. Stated otherwise, what follows is a synthesis and simplification of a much larger technical literature—a synthesis written with the aim of practical application in mind.

I begin by sketching the four most salient metrics in the literature.<sup>299</sup> These can be summarized as follows: One might *first* simply look at whether equal fractions of each racial group are labeled as risky—such that they will be subject to additional policing or detention. A similar, although not identical, analysis where risk is measured as a continuous variable without a threshold for coercive action would look for equal average risk scores across different racial groups. *Second*, one might ask whether the same classification rule is being used to assign racial groups to the high-risk category. This condition is satisfied if the same numerical risk score is used as a cut-off for all groups. *Third*, one might separate each racial group and then look at the rate of false positives conditional on being categorized as high risk. And *fourth*, one might separate each racial group and ask how frequent false positives are conditional on being in fact a non-risky person. In the literature, this has been characterized as a consideration of the population of those who in fact will not engage in subsequent criminal conduct within a racial group, and an inquiry into what proportion of that subset were erroneously categorized as warranting coercion.

---

<sup>297</sup> Minority politicians and police chiefs have at times believed that a disproportionate policing focus on African-Americans was warranted in terms of community self-preservation—a belief that the “accumulated impact” of harsh antinarcotics measures has over time shown to be erroneous. JAMES FORMAN JR., LOCKING UP OUR OWN: CRIME AND PUNISHMENT IN BLACK AMERICA 124-48, 218 (2017) (documenting these calls). If they had been correct—and Forman persuasively suggests that they were wrong on the facts—then this would have justified a less demanding risk threshold for blacks than for whites.

<sup>298</sup> See Arvind Narayan, “21 Definitions of Fairness and their Politics,” YouTube, Mar. 1, 2018, <https://www.youtube.com/watch?v=jIXIuYdnyyk>.

<sup>299</sup> There are different enumerations of competing definitions of algorithmic fairness in criminal justice in particular. Richard Berk and his co-authors identify six different definitions. Berk et al., *Fairness in Criminal Justice*, *supra* note 253, at 13. They do not include one of the definitions I consider. Another paper by Sam Corbett-Davies and colleagues (including me) identifies three definitions that are salient to criminal justice policy. Sam Corbett-Davies et al., *Algorithmic decision making and the cost of fairness*, in PROCEEDINGS OF THE 23RD ACM SIGKDD INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING 797, 798 (2017). In addition, Feldman et al. define fairness as the inability to predict a trait from the execution of an algorithmic function. Michael Feldman et al., *Certifying and Removing Disparate Impact*, in PROCEEDINGS OF THE 21TH ACM SIGKDD INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING 259, 265 (2015), <https://dl.acm.org/citation.cfm?id=2783311>.

The four concepts of fairness or nondiscrimination are summarized in Table 1, which pairs each conception to the relevant parameter (or variable) that is to be equalized.

**Table 1: Conceptions of Nondiscrimination in Algorithmic Criminal Justice**

<i>Conception of Fairness</i>	<i>Parameter that should be equalized</i>
<i>Statistical parity</i>	<i>Proportion of each group subject to coercion</i>
<i>Single threshold</i>	<i>Treatment of equally risky persons within each group</i>
<i>Equally precise coercion</i>	<i>Proportion of those ranked as risky who are erroneously classified</i>
<i>Predictive error equality</i>	<i>Proportion of innocent persons that are subject to coercion</i>

Table 1 is intended to capture the range of core conceptions of nondiscrimination that should matter in the criminal justice context. It does not, as I have already noted, capture the full range of potential conceptions. For instance, one recent survey additionally flags the idea of treatment equality,<sup>300</sup> which looks simply at the ratio of false positives to false negatives for a given racial group. To date, however, the latter concept has not played a large role in debates about racial equity. My analysis does not suggest that it should—and hence I leave it to one side for present purposes.

Figure 1 below helps clarify these four concepts. It displays the risk ranking assigned by an algorithm—represented as a continuous variable of two groups, white and black. The x-axis represents the risk value assigned to members of the population; the y-axis represents the frequency with which members of the group are assigned to a risk level. For the purposes of this analysis, I assume that the training data used to generate the risk assessments is not flawed, and in particular is not biased in ways that result in whites or blacks being subject to disproportionate coercion. I make this assumption so as to enable a narrow focus on the question whether the algorithmic classification rule *standing on its own* presents a question of racial justice.

The graphic also contains a vertical line to represent the cut-off point for the purposes of allocating coercion. Those who fall to the right of this threshold are subject to the coercive treatment (either a policing or a detention-related intervention), while those who are to the left of the threshold are not subject to any coercion. The parts of the curve that represent populations that will be coerced (assuming the algorithm’s recommendations are followed) are represented with colored blocks in the graphic. The proportion of the white and the black populations being subject to coercion is a function of the area under the respective curve to the right of the threshold.

<sup>300</sup> Berk et al., *Fairness in Criminal Justice*, *supra* note 253, at 14.

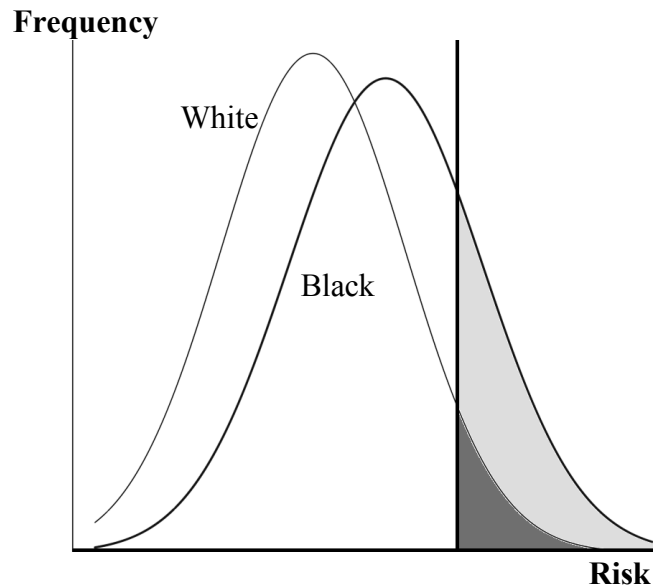
This form of graphical representation has a number of advantages. In using a continuous variable, and in capturing the way in which a threshold will distinguish between populations that are themselves quite internally varied in terms of their riskiness, this chart captures some of the key features of criminal justice algorithms in practice. In particular, it captures the fact that a decision must be taken about who the *marginal* person on the risk curve is who should be detained. It also captures the intuition that the risk curves for different racial groups might diverge.<sup>301</sup>

---

<sup>301</sup> An alternative used in the literature is a confusion table, which is a two-by-two matrix that distributes individuals in terms of whether they ultimately committed acts justifying coercion, and whether they were in fact coerced. *See, e.g.*, Tom Fawcett, *An Introduction to ROC Analysis*, 27 PATTERN RECOGNITION LETTERS 861, 862 (2006) (describing the use of confusion matrices). Confusion tables, however do not capture all the information that an algorithm generates—such as the variance in risk values—and relies on knowledge that a decision maker by construction does not know at the time the relevant decision has to be made—i.e., whether a suspect or a defendant in fact will go on commit a crime or impose a harm on others in the future. Confusion tables hence omit useful information, while including information that cannot plausibly inform the decision whether to coerce or not. They are non-ideal instruments for exploring algorithmic fairness because the latter is a standard that has to be applied at the moment the algorithm is used—and not later, once new information about potential states of the world has become available.

Moreover, confusion tables fail to distinguish the average person subject to coercion from the marginal subject of coercion. For example, imagine a single decision rule (say, a risk threshold of 10%) is applied to both a white and a black population. The white population comprises some with a 1% chance of carrying contraband, and some with a 75% chance. The black population comprises some with a 1% chance, and some with a 50% chance. A confusion table draws attention to the fact that the proportion of stops that are ‘false positives’ for the white group will be one-half that for the black group (i.e., 25% rather than 50%), but without elucidating whether this is a function of (a) a biased decision rule, or (b) a neutral and justified decision rule being applied to different distributions of rule in the population. *See* Camelia Simon, Sam Corbett-Davies, and Sharad Goel, *The Problem of Infra-Marginality in Outcome Tests for Discrimination*, 11 ANN. APP. STAT. 1193, 1194 (2017) (setting out example); *see also* Ayres, *supra* note 245, at 131. This confusion, ironically, is avoided by foregoing the use of confusion tables.

**Figure 1: Hypothetical risk distributions for white and black populations**



In Figure 1, the tails of the curve for the black population are to the right of those for the white population. This means that the algorithm assigns in general higher risk values to black than white persons in the population. If the risk distributions of both populations are equal, no interesting question of racial equity or discrimination would arise from its functioning: White and black outcomes would not be distinct. This element of the hypothetical is not meant to imply that blacks in fact are more likely to commit crimes than whites. It is rather to present a situation that is plausible, and that presents most sharply the questions of racial equity that are of interest here.

The four conceptions of algorithmic fairness, or algorithmic nondiscrimination can be elaborated as follows. First, an algorithmic classifier might exhibit **statistical parity**. This means that an equal proportion of members of each group are subject to coercion. In terms of the graphic, this means that the areas under the white and the blacks curves to the right of the threshold are equal to each other.<sup>302</sup> This can happen, it is worth noting, even if there is wide variation in the ratio of false positives to true positives for whites and for blacks. Where there is no threshold, one might instead use the average risk score for a given group. A variant on statistical parity is “conditional statistical parity,” which requires that, having controlled for a “limited set of ‘legitimate’ risk factors, an equal proportion of defendants within each race group.”<sup>303</sup> In practice, however, this definition is highly sensitive to what counts as a “legitimate” risk factor. Because my analysis does not assume an answer to the question of what counts as a legitimate risk factor, I must put aside here the possibility of conditional statistical parity.

<sup>302</sup> Corbett-Davies et al., *supra* note 299, at 2; *see also* Dwork et al., note 246, at 7 (defining statistical parity in terms of the fact that “an individual observed a particular outcome provides no information as to whether the individual is a member of S or a member of T”); *id.* at 11-15 (developing a concept of fairness that maximizes statistical parity while at the same time observing the constraint that otherwise similar people be treated alike).

<sup>303</sup> Corbett-Davies et al., *supra* note 299, at 2.

Statistical parity is a clear and simple idea. Indeed, it is employed as part of the prima facie case in disparate impact analysis in employment discrimination law.<sup>304</sup> Under longstanding administrative agency construction, a racial difference in selection rates of “greater than four-fifths” is “generally” taken as evidence of “adverse impact.”<sup>305</sup> On the other hand, there is no abstract or a priori reason why state coercion should be equally distributed among racial groups. To be sure, there is some evidence that at least for certain sorts of offenses, such as narcotics crimes, there is “no statistically significant differences” in offending rates for different racial and ethnic groups.<sup>306</sup> But on the assumption that the algorithm’s training data are not flawed, the hypothetical would simply not capture such cases.

Second, an algorithmic classifier might be viewed as fair if it treated two people who evinced the same ex ante evidence of risk, and differed by race, in the same way. The computer science literature has distinguished between a single threshold and “multiple race-specific thresholds.”<sup>307</sup> A recent paper further offers a formal proof to the effect that the “immediate utility” of a decision rule—defined in terms of the immediate benefits of crime directly suppressed and direct costs of coercion (and ignoring externalities)—is typically optimized by maintaining a single threshold rule for coercion, rather than having plural thresholds.<sup>308</sup> That is, a social planner with an algorithmic tool that is trained on unbiased data would select a single risk threshold for both whites and blacks if she wished to optimize over the costs and benefits of crime control. This analysis of social welfare, however, does not answer the question of what necessarily furthers racial equity under all conditions. In particular, it is important to observe that the formal proof of optimality is limited to the immediate effects of an algorithmic tool, rather than its indirect effects. Racial stratification is plausibly understood to be a compounding effect of the latter concept, rather than something captured by the former.

This conception of fairness in algorithmic criminal justice has not so far attracted a distinctive label. Indeed, some accounts of discrimination in the algorithmic context simply do not cite this kind of fairness, preferring to focus on the relative frequency of false (or true) positives (or negatives) in the two racial groups.<sup>309</sup> In other work, this conception has

---

<sup>304</sup> See, e.g., *Int'l Bhd. of Teamsters v. United States*, 431 U.S. 324, 339 (1977); *Hazelwood Sch. Dist. v. United States*, 433 U.S. 299, 308 n.14 (1977).

<sup>305</sup> 29 C.F.R. § 1607.4(D) (2016) (“A selection rate for any race, sex, or ethnic group which is less than four-fifths (4/5) (or eighty percent) of the rate for the group with the highest rate will generally be regarded by the Federal enforcement agencies as evidence of adverse impact, while a greater than four-fifths rate will generally not be regarded by Federal enforcement agencies as evidence of adverse impact.”); see also *Ricci v. DeStefano*, 557 U.S. 557, 587 (2009) (endorsing this four-fifths rule).

<sup>306</sup> U.S. Department of Health and Human Services, *Results from the 2013 National Survey on Drug Use and Health: Summary of National Findings 25* (2013), <http://www.samhsa.gov/data/sites/default/files/NSDUHresultsPDFWHHTML2013/Web/NSDUHresults2013.pdf>.

<sup>307</sup> Corbett-Davies et al., *supra* note 299, at 8.

<sup>308</sup> *Id.* at 3-6.

<sup>309</sup> See, e.g., Berk et al, *Fairness in Criminal Justice*, *supra* note 253, at 13 (not mentioning this kind of fairness in a six-fold taxonomy).

been characterized simply as “fairness,”<sup>310</sup> but that nomenclature is probably too vague to be helpful. I label this definition, therefore, the ‘**single threshold**’ definition of algorithmic fairness. Graphically, the single threshold definition of fairness is represented by the fact that the vertical line that marks the threshold between coercion and its absence is in the same place for both racial groups. If the vertical thresholds were placed in different locations on the x-axis, there would be a group of individuals between the two thresholds who would present the same evaluated risk, but would be treated differently solely on account of their race.

A third conception of algorithmic nondiscrimination examines only the portion of the population that lies to the right of the risk threshold. In Figure 1, this comprises the shaded areas under the curves. These capture the parts of the white and the black population subject to coercion as a consequence of the algorithm’s recommendations. Not all of these recommendations, however, will be borne out by future events. In the bail context, for example, some fraction of those subject to state coercion would not have gone on to commit crimes that justified pretrial detention. They will, in other words, be ‘false positives.’ One way of thinking about nondiscrimination is in terms of the false positive error rate conditional on being assigned state coercion by the algorithm (which can also be stated as  $P(\text{innocent}|\text{high risk})$ ). So if a greater fraction of blacks stopped or detained turn out to be innocent in the relevant sense than the same fraction of ‘innocent’ whites, then this would violate the third conception of fairness. Or (stated in yet another form) if the proportion of those ‘false positives’ under the black curve to the right of the risk threshold is greater than the proportion of ‘false positives’ under the white curve to the right of the threshold, then this conception of equality is violated.<sup>311</sup> This notion is captured by a number of different terms in the computer science literature. A leading group of analysts label it conditional use accuracy.<sup>312</sup> In my view, it is simplest to label it “**equally precise coercion**” because this conception is centrally concerned with the rate at which false positives occur conditional on the fact of being coerced.<sup>313</sup>

Equally precise coercion played a role in the debate over the Compas algorithm.<sup>314</sup> Responding to Pro Publica’s allegations of racial disparity, the Northpointe company focused on the fact that the rate of error among the black and white groups subject to coercion was the same.<sup>315</sup> In effect, the Northpointe argument was that so long as equally precise coercion obtained, there was no discrimination problem.

The fourth and final conception of fairness in the algorithmic context also focuses on false positives, but from a different angle. Rather than the subset subject to coercion, it

---

<sup>310</sup> Dwork et al., *supra* note 246, at 2.

<sup>311</sup> This conception is focused not on the absolute number of false positives but rather than percentage of those subject to coercion within a racial group that would not have gone on to engage in socially undesirable behavior. It would be perverse to define fairness in terms of a parameter that was driven primarily by the relative size of the two groups under study.

<sup>312</sup> Berk et al., *Fairness in Criminal Justice*, *supra* note 253, at 14.

<sup>313</sup> Precision is the term used by machine learning specialists, who perceive the term “accurate” to imply a normative judgment. I am grateful to Sharad Goel for discussion of this point.

<sup>314</sup> See *supra* text accompanying notes 13 to 21.

<sup>315</sup> Dietrich et al., *supra* note 19, at 3.

focuses on the subset who would not go on to commit a crime or violent act. This subset of ‘actually innocent’ persons is used as a denominator. For a numerator, it asks what fraction of that subpopulation is incorrectly subject to coercion. In the bail context, for example, this means asking whether “among those defendants who would not have gone on to commit [a violation], detention rates are equal across race groups.”<sup>316</sup> In other words, conditional on being ‘innocent’ (in whatever sense of that term is relevant), the rate of erroneous false positives across racial groups does not vary (or  $P(\text{high risk}|\text{innocent})$ ). This conception of equality is not easy to capture using Figure 1, since the baseline category of the ‘actually innocent’ are dispersed on both sides of the risk thresholds. In effect, it comprises a diffused subset of whites and blacks who in fact would not commit actions that justify coercion. This conception of fairness requires that we look for the proportion of that actually innocent subset to the right of the risk threshold. If one racial group’s ratio is larger than the other’s, there is reason for concern on this theory.

This conception has attracted a wide variety of labels, including “predictive equality,”<sup>317</sup> “conditional procedural accuracy,”<sup>318</sup> and “equalized odds.”<sup>319</sup> Another group of analysts use the label “balance for the positive class” for a related concept.<sup>320</sup> Their paper also mentions the concept of “balance for the negative class” to capture the symmetrical idea that “the assignment of scores shouldn’t be systematically more inaccurate for negative instances in one group than another.”<sup>321</sup> Deviating from my own past usage,<sup>322</sup> I will use the label **predictive error equality** here to capture the idea that what is at stake in this fourth definition of nondiscrimination is the notion that the burden placed on the innocent subset of each racial group should be the same. Predictive error equality is the focus of the Pro Publica critique of the Compas algorithm: The journalistic organization demonstrated, that is, that the proportion of ‘innocent’ black defendants recommended for detention by the Compas algorithm was substantially higher than the proportion of innocent white defendants subject to the same recommendation.<sup>323</sup> In effect, Pro Publica implicitly leveraged the intuition that what matters with an algorithm is what happens to the ‘actually innocent.’ If the treatment of actual innocents varies across racial groups, Pro Publica’s argument went, an algorithm could not be ranked as nondiscriminatory.

#### D. Prioritizing Conceptions of Algorithmic Discrimination

The range of possible ways of operationalizing the quality of nondiscrimination in the algorithmic criminal justice context raises the question of how to evaluate and rank the

---

<sup>316</sup> Corbett-Davies et al., *supra* note 299, at 2.

<sup>317</sup> *Id.*

<sup>318</sup> Berk et al, *Fairness in Criminal Justice*, *supra* note 253, at 13-14.

<sup>319</sup> Moritz Hardt et al., *Equality of Opportunity in Supervised Learning*, in *ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS* (2016), <http://papers.nips.cc/paper/6373-equality-of-opportunity-in-supervised-learning>.

<sup>320</sup> Jon Kleinberg, Sendhil Mullainathan & Manish Raghavan, *Inherent trade-offs in the fair determination of risk scores* 4, 9 (2016), <https://arxiv.org/abs/1609.05807> (Labeling this concept as “[c]alibration within groups”).

<sup>321</sup> *Id.*

<sup>322</sup> Corbett-Davies et al., *supra* note 299, at 2.

<sup>323</sup> Angwin, *supra* note 15.

four main competing conceptions. My aim in this section is twofold. First, I point to results in the technical literature that demonstrate the impossibility of pursuing all these conceptions of nondiscrimination simultaneously. Second, I offer my own normative account of which conception to prioritize. This account, detailed above, hinges on the minimization of costs net of benefits for the minority group. Contrary to both Northpointe and Pro Publica, this contends that rates of false positives (whatever denominator is used) are not compelling normative benchmarks. Instead, the analysis should focus on whether a minority risk threshold yields net costs or benefits for that group. Where there are no spillovers, it is likely that the same threshold will obtain for both minority and majority groups. Where there are large and asymmetric spillovers, both social efficiency and racial equity are served by different thresholds.

### 1. *Conflicts Between Algorithmic Fairness Definitions*

It would seem desirable to satisfy all these definitions of equality. At least at first blush, all capture colorable and important intuitions about the fair allocation of coercion. But matters are not so simple. It turns out that this is not possible in many cases—and not possible under conditions that are reasonably likely to occur in practice—for two reasons.

First, it will generally be the case that statistical parity cannot be achieved using a single threshold. This is readily apparent from Figure 1, which illustrates the case in which the risk distributions of racial groups vary. When this happens, it will always be the case that a single risk threshold will subject different proportions of each group to coercion. Hence, it is not possible—assuming differences in the distributions of risk between the two racial populations—to have both a single threshold and also statistical parity.

Second, it will generally be the case that it is also impossible to achieve both equally precise coercion and predictive error equality. This impossibility result holds under two conditions. First, it must be the case that the base rates of criminality are different for the two racial groups. Second, it also must be the case that there is no function that allows for “perfectly accurate classification” (a condition also known as “separation”).<sup>324</sup> Under these conditions, “one cannot have both conditional use accuracy equality and equality in the false negative and false positive rates,” where the latter term is simply conditional procedural accuracy.<sup>325</sup> It is for this reason that assessments of the Compas algorithm have diverged. On the one hand, the original criticism of the algorithm focused on the difference in the rate of conditional procedural errors for blacks and whites.<sup>326</sup> On the other hand, the defenses of Northpointe’s instrument focused on the fact that it was calibrated within the categories of risk—i.e., the conditional use error rate was equal for both whites and

---

<sup>324</sup> *Id.* at 18-19. For derivations of the same result, see Kleinberg, Mullainathan, & Raghavan, *supra* note 320, at 5-6; Alexandra Chouldechova, *Fair prediction with disparate impact: A study of bias in recidivism prediction instruments*, 5 *BIG DATA* 153, 154 (2017). Corbett-Davies et al. develop a related point, which is that under certain conditions it is not possible to equalize conditional procedural accuracy between groups without establishing different thresholds for black and white classifications. Corbett-Davies et al., *supra* note 299, at 6-8.

<sup>325</sup> Berk et al., *Fairness in Criminal Justice*, *supra* note 253, at 14; Kleinberg et al., *supra* note 320, at 5.

<sup>326</sup> See Angwin, *supra* note 15.



blacks.<sup>327</sup> Neither side recognized that given the possibility of underlying differences in the empirical characteristics of racial groups, and absent separation, these two metrics of algorithmic fairness were bound, mathematically, to diverge under plausible conditions.<sup>328</sup>

A choice therefore must be made about which conception of nondiscrimination to pursue. The computer-science literature, while helpful in defining the range of possible conceptions of algorithmic nondiscrimination, is less helpful in evaluating and ranking those definitions.

## 2. *The Irrelevance of False Positive Rates*

Two of the four definitions of algorithmic nondiscrimination developed above—equally precise coercion and predictive error equality—focus on the rate of false positives. These two definitions differ, however, in terms of their denominator, which is alternatively (1) being coerced, or (2) being ‘actually innocent.’ False-positive focused definitions not only played a central role in the debate between Northpointe and Pro Publica,<sup>329</sup> they have also infiltrated public debate more broadly.<sup>330</sup> A concern with false positives is not without normative appeal. But definitions of nondiscrimination that hinge on false positive rates do not index in any obvious fashion the extent to which an algorithmic instrument exacerbates racial stratification. This section is hence directed at ruling out two of the four possible metrics of racial equity that have attracted the most public attention to date.

For four interrelated reasons, the temptation to focus on false positives should be resisted. *First*, the criminal-justice decisions subject to algorithmic resolution are all made in advance of potential adverse actions. A stop is conducted, bail is denied, or a sentence is extended, that is, before the state knows, or can know, whether a suspect or defendant will in fact commit a criminal act. Officials using an algorithm, therefore, cannot know who is a true positive and who is a false positive among the pool of persons to the right of the vertical threshold illustrated in Figure 1. Even if we assume that an official responsible for applying the algorithm knows the general shape of the distribution (for example, as illustrated in Figure 1), she does not and cannot know whether a particular suspect is in fact going to inflict harm; all she knows is how the algorithm has ranked that person. A test for nondiscrimination that distinguishes false positives from true positives relies on information that is not available to that official. And it is not at all clear why the failure to account for information that the official or algorithm cannot access should be treated as a failure. Provided that the decision rule otherwise achieves valued public goals at the lowest collateral cost, that is, it is not clear why the (ordinarily unknown) distribution of false positives should matter.

---

<sup>327</sup> See sources cited in *supra* note 19.

<sup>328</sup> There are a number of computational fixes, which fall into the categories of pre-, in-, and post-processing. But none is a complete fix. Berk et al., *Fairness in Criminal Justice*, *supra* note 253, at 25-29.

<sup>329</sup> See *supra* text accompanying notes 13 to 21.

<sup>330</sup> *Perfecting in China, a threat in the West*, THE ECONOMIST, June 2, 2018, at 11 (worrying that “[s]ome sentencing algorithms are more likely to label black defendants than white ones as being at high risk of reoffending.”). This reads as a concern with predictive error equality, although this is not wholly free from doubt.

Second, the law itself has a very high tolerance for false positives. In the policing and the pretrial detention contexts in particular, the law is willing to tolerate a very high level of false positives on the ground that the gains to crime suppression offset the costs of that high rate of false positives. Hence, in the policing context a mere showing of “reasonable articulable suspicion,” which is far less than probable cause, is enough to warrant a street stop.<sup>331</sup> In the bail context, the standard for detention under federal law is framed in terms of reasonableness and envisages substantial room for error.<sup>332</sup> But disparities in the allocation of state-created goods (or bads) are generally thought to be worrisome if those goods are important. This, I think, explains the coverage of housing and employment opportunities by disparate impact regimes.<sup>333</sup> Moreover, if the law takes the view that there is no reason for concern at the prospect of absolutely high levels of stops or pretrial bail detentions, it is not clear that the law contains the normative resources to establish concern when those resources are allocated in subtly disparate ways—especially if the overall pattern of stops redounds to the net benefit of the society, and also of the subordinated group.<sup>334</sup>

Third, a failure of equally precise coercion or of predictive error equality is a mathematical function of the use of a single threshold for risk to two racial groups with different risk distributions.<sup>335</sup> Given that relationship, it is necessary to choose between unequal rates of false positives and different risk thresholds. Merely pointing to one form of inequality is question-begging. If the risk threshold is set at the socially efficient level, moreover, such that it optimizes over immediate costs and benefits for blacks as well as whites,<sup>336</sup> equalizing false positives risks the imposition of unnecessary costs on the minority group. Although not dispositive, it is worth noting that the disparate impact law does not treat unavoidable disparities generated by the pursuit of a valid governmental interest as cause for concern.<sup>337</sup> At least where the state has no other means of suppressing crime without a violation of equally precise coercion or of predictive error equality, it is not obvious why the ensuing disparities should be treated as fatally problematic.

---

<sup>331</sup> The initial delineation of rules for a street stop is contained in *Terry v. Ohio*, 392 U.S. 1, 22 (1968). The phrase “reasonable, articulable suspicion” was used first in *Brown v. Texas*, 443 U.S. 47, 51 (1979); see also *Florida v. Royer*, 460 U.S. 491, 502 (1983).

<sup>332</sup> 18 U.S.C.A. § 3142(e) (mandating pretrial conditions unless conditions can be imposed that “reasonably assure the appearance of the person as required and the safety of any other person and the community, such judicial officer shall order the detention of the person before trial”); see *Stack v. Boyle*, 342 U.S. 1, 4 (1951) (defendant charged with noncapital offense shall be released on bail if defendant gives adequate assurance that he will appear at trial and submit to sentence if convicted).

<sup>333</sup> *Tex. Dep’t of Hous. & Cmty. Affairs v. Inclusive Cmty. Project, Inc.*, 135 S. Ct. 2507, 2510 (2015) (affirming availability of disparate impact liability under the Fair Housing Act).

<sup>334</sup> Note that the mere fact of a violation of equally precise coercion or predictive error equality is not evidence that the net effect of a criminal-justice measure is to exacerbate overall racial disparities. There is no empirical equivalence between these terms.

<sup>335</sup> Corbett-Davies et al., *supra* note 299, at 3-6

<sup>336</sup> *Id.*; see *supra* text accompanying note 308.

<sup>337</sup> For a discussion of the current doctrinal position of this element of disparate impact law, see Samuel R. Bagenstos, *Disparate Impact and the Role of Classification and Motivation in Equal Protection Law After Inclusive Communities*, 101 CORNELL L. REV. 1115, 1140 (2016).

Finally, and most importantly, if one is concerned with the impact of algorithmic criminal justice on a stratified racial minority, it is not at all clear why the focus should solely be on false positives. The negative expressive effects and social harms imposed by criminal justice institutions upon African-American communities that are not merely triggered by false positives. Directing coercion toward black suspects and defendants even when such coercion is warranted can have an expressive effect on public beliefs about black criminality and the same debilitating effects on communities, families, and children. Indeed, there is no particular reason to believe that any of these spillover costs are less if the person subject to the coercion is in fact a true rather than false positive. Put another way, if you care about racial stratification, what should matter is the *absolute* cost of a coercive tactic for minority, net of benefit, for *all* members of that racial group—whether or not they ultimately would have acted in ways that justified coercion. Both kinds of stops have costs; both count for the purposes of racial equity. True, those costs are offset when an algorithm makes a correct prediction—but that is captured better by a focus on the benefits of the coercive measure being allocated.<sup>338</sup>

For these four reasons, I do not think that either equally precise coercion or predictive error equality provides an appropriate metric for thinking about racial equality in this context. Rather, it is desirable in the end to know whether crime control is inflicting more costs than benefits for the minority group as a whole—and not just those who would otherwise not go on to inflict any social harm.

### 3. *Evaluating the Impact of Algorithmic Criminal Justice on Racial Stratification*

So what does matter? The opening two movements of this Part mapped the effect of criminal justice institutions on racial stratification, and charted a general principle of racial equity. Existing criminal justice systems influence the extent of racialized social stratification in society as a whole.<sup>339</sup> Racial equity in criminal justice generally—and in particular in the algorithmic context—should be primarily concerned with mitigating these pernicious effects. It should repudiate the tight linkages that have bound criminal justice to the reproduction of racial hierarchy since the beginning of the twentieth century. Even if the present-day operation of criminal justice institutions cannot undue past harms, at a minimum they should not compound those harms.

The question therefore is which of the available technical benchmarks best captures this pathway between criminal justice and racial stratification. As argued above, an appropriate benchmark would home in upon the net cost (benefit) of an algorithmic criminal justice instrument for the racial minority in the socially subordinate position. A measure of costs net of benefits for the racial minority is relevant morally because it captures the extent to which a criminal justice measure depresses the social standing of the minority group. In the context of black/white comparisons in the American context at least,

---

<sup>338</sup> I can imagine one more reason for taking normative account of false positives only: One might posit that the ratio of false positives to true positives is a measure of intragroup transfers. The greater the proportion of false positives, that is, the more the burden of crime suppression falls on those members of the minority who are innocent. This may be a morally relevant quality, but I am not convinced it is a measure of racial equity.

<sup>339</sup> See *supra* Part III.A.

this analysis is simplified by the fact that much violent crime is intraracial. That is, the benefits of a crime suppression measure imposed on blacks are likely to accrue largely to blacks (while the same is true for whites). The analysis would be more complex if we assumed that the racial minority did not capture all or most of the benefits of crime suppression targeting members of that minority.

In my view, there is no one metric developed in the computer science literature or otherwise that captures this concern with racial stratification. Benchmarks that concern the rate of false positives capture in a very loose way the magnitude of unjustified state coercion. But they fail to acknowledge the state's inability to distinguish justified from unjustified exercises of coercion *ex ante*. Statistical parity does account for the aggregate cost of coercion on a racial minority. But it does so only through a comparative lens. It asks whether the minority is burdened more or less than a majority group. It also fails to consider offsetting benefits to the minority group. Because most crime is intraracial, it fails to account for the possibility that the benefits of crime suppression for blacks outweigh its costs. A comparative measure such as statistical parity is at best considered an evidentiary tool, therefore, rather than a direct measure of racial equity.

An inquiry into racial equity can usefully focus instead on the question whether the *marginal* decision to impose coercion within the black population can be justified. I present first a simple version of this inquiry that assumes that all costs and benefits are immediate, and that there are no spillovers. Consider again Figure 1. Imagine sliding the threshold for coercion for the minority population right from the y-axis. At first, the threshold would assign for coercion many people for whom the immediate costs of such coercion outweighed any benefit for the simple reason that their risk of causing harm was so low. At some point in the rightward movement of the threshold, however, the immediate costs of coercion would be balanced by its benefits. When the costs of this marginal decision to coerce are outweighed by its benefits, the threshold has been calibrated such that no net burden is being placed on the minority population, and all coercion generates a net gain for that group. Assuming that most relevant crime is intraracial, this means that the marginal benefits of coercion (for the black community) are greater than the costs of coercion (for the black community). Such a policy leaves that racial group no worse off than it would otherwise be.

For interventions that prevent serious crimes, there is no reason to think that the immediate costs of coercion, or the immediate benefits of crime control, vary between racial groups. Moreover, spillovers can be ignored because such costs are likely to be rounding errors in relation to the costs of murder, sexual assault, armed robbery, and the like. Such a tightly focused analysis might, for example, be appropriate in the analysis of bail decisions where a suspect may go on to commit a serious violent crime. Under these conditions, a single risk threshold calibrated to be socially optimal (in the sense of eliminating cost-unjustified coercion) will satisfy racial equity. It will also be socially efficient.

This goal has likely not been reached in practice. Even assuming that criminal-justice decision-makers are applying a single threshold rule (rather than being influenced

by animus or racial stereotypes), it is very likely that many present uses of police coercion and detention are unjustified. The benefits of state coercion are likely over-estimated, while its costs are under-estimated. Consistent with this prediction, current risk-assessment tools estimate the benefits of coercion but do not measure costs.<sup>340</sup> Still, the present lack of empirical data on the costs and benefits of many familiar criminal justice institutions, such as street stops and bail denials, means that this intuition is hard to substantiate. But the available data suggests an excess of coercion beyond the socially optimal.<sup>341</sup> When the supernumerary costs of such coercion fall on racial minorities, they intensify racial stratification. Ratcheting back the sheer volume of coercion, therefore, may be a first order task in reform projects that have racial equity in mind.

This simple analysis of racial equity accounts only for the immediate costs and benefits of coercion. It does not account, though, for the externalities set forth in Part III.A. A more complex model of racial equity would account for all negative spillovers from algorithmically allocated coercion. These externalities are substantially greater for racial minorities than for the racial majority. They are also nontrivial in scale. Where less serious crime is concerned (e.g., public order offenses), it is likely that these externalities are of the same magnitude as the immediate benefits and costs of crime-control. Second-order, downstream costs of coercion therefore cannot be safely ignored as rounding errors in an analysis of the criminal justice system's dynamic effects. The analysis for less serious crime, or for interventions that do not impede serious harms, is hence different from the analysis when serious social harm is directly at stake.

Accounting for the racially asymmetrical distribution of externalities alters the racial-equity analysis. It means that the marginal costs of coercion are likely to be greater for the racial minority. Accordingly, the point on the x-axis at which costs are equal to benefits for the minority is to the right of the same break-even point for the majority group. Because the operation of criminal-justice coercion, that is, generates asymmetrical harms to black families and black communities, as well as exacerbating Kennedy's racial tax, there will be a class of crimes for which a greater benefit will be required to achieve net positive effects for black suspects. And because the costs and benefits of crime are largely intraracial, the same higher risk threshold will be required to achieve social efficacy. Whether the focus is social efficiency or racial equity, this implies that the risk threshold for blacks should be set at a higher level (i.e., farther to the right in Figure 1) than the threshold for whites. Therefore, accounting for both the immediate and spillover costs of crime control when its immediate benefits are small conduces to a bifurcated risk threshold—one rule for the majority, and one for minority. The single vertical line in Figure 1 would bifurcate. The line for blacks would move rightward.

---

<sup>340</sup> Slobogin, *supra* note 46, at 584-85.

<sup>341</sup> See Huq, *Disparate Policing*, *supra* note 4, at 2929; Note, *Bail Reform and Risk Assessment: The Cautionary Tale of Federal Sentencing*, 131 HARV. L. REV. 1125, 1127-28 (2018) ("The pretrial imprisonment rate in the United States is among the highest in the world—more than four times the world's median pretrial imprisonment rate."); see also Mayson, *supra* note 123, at 545-48 (spelling out costs and benefits of bail in a way that clarifies its complexity).

This is akin to common affirmative action schemes, in which otherwise similar black and white persons are treated differently because of the different spillover consequences of their treatment. In the affirmative action context, the existence of a positive diversity benefit (which is another kind of spillover) warrants a less stringent threshold rule for assigning a benefit to the racial minority.<sup>342</sup> In the criminal justice context, similarly, the existence of negative spillovers for black families and communities warrants a more stringent risk threshold for the racial minority. The argument for a bifurcated classification rule is arguably stronger here than the argument for affirmative action: The alleviation of racial stratification, in my view, is a more acute interest than diversity because it directly benefits the most marginalized (which affirmative action may not) and immediately relieves stigmatic and material harms. Alleviating the effect of accumulated disadvantage caused by the historical operation of criminal justice institutions, in other words, is a more compelling goal than crafting a well-rounded university.

Unlike affirmative action, however, the case for multiple risk thresholds can be made independently on either racial equity or social efficiency grounds. So long as a policy's costs (benefits) are largely internalized by racial groups, and so long as costs are greater at the margin for the minority group, a socially optimal rule would require different risk thresholds. Where the state adopts a cost-benefit approach to criminal justice policy,<sup>343</sup> an exacting approach to cost-benefit trades offs in crime control may in some cases generate dual thresholds.<sup>344</sup> In the algorithmic context, it is worth noting that a machine-learning tool given the necessary data and asked to vindicate social efficiency (understood in a capacious sense that reached both static and dynamic effects) could converge on a bifurcated rule absent race-conscious human decision-making.

However that goal is approached, its achievement imposes large new epistemic burdens on the state. Whereas risk assessment in criminal justice to date has focused narrowly on the costs of crime, a rigorously executed algorithmic method demands data on the costs of crime control. This is not merely a matter of counting state expenditures, but also measuring spillovers. This is a massive task. But its size and difficulty ought not to be a justification for avoidance. The current dearth of information about the spillover costs of criminal justice institutions, particularly for minority communities, is causally related to their stratifying effects. Ignorance of spillovers, coupled to a myopic focus on a small number of high-profile crimes, creates the epistemic background against which actually existing state institutions compound racial stratification. That ignorance is thus a form of "hermeneutical injustice," in which "some significant area of one's social experience [is] obscured from collective understanding owing to persistent and wide-ranging hermeneutical marginalization."<sup>345</sup> Racial inequity cannot be justified by hermeneutic

---

<sup>342</sup> See, e.g., *Fisher v. Univ. of Tex.*, 136 S. Ct. 2198, 2208 (2016).

<sup>343</sup> A version of cost-benefit analysis is endorsed in Barry Friedman & Maria Ponomarenko, *Democratic Policing*, 90 N.Y.U. L. REV. 1827, 1907 (2015) (encouraging even "small steps" in that direction).

<sup>344</sup> It is also possible that a jurisdiction could pursue social efficiency by deploying a nonracial bifurcation in the risk threshold. For instance, it may in some instances be possible to employ socioeconomic stratification to much the same end.

<sup>345</sup> MIRANDA FRICKER, *EPISTEMIC INJUSTICE: POWER & THE ETHICS OF KNOWING* 154 (2007) (emphasis omitted).

injustice. Precisely how the epistemic gap will be closer is a large question, and I do not take it up here. But it is worth noting that the algorithmic tools mapped here may have a role. Determining how ‘big data’ tools can contribute to this epistemic enterprise, indeed, is perhaps the next technological frontier in criminal justice.

At the same time, a multiple threshold rule for different racial groups—runs headlong into the anti-classification rule of Equal Protection doctrine.<sup>346</sup> At a minimum, it would receive strict scrutiny.<sup>347</sup> As a result, a multiple threshold regime would be in serious constitutional jeopardy. Under these conditions, which are hardly empirically implausible, the regime imperiled by our constitutionality equality doctrine is the only one that *both* mitigates racial stratification *and also* maximizes social welfare. Why would we want to place that regime beyond reach? I can think of no good answer. Such a result, in my view, hence tells us more about our wrongheaded equality doctrine than it does about the substance of algorithmic criminal justice.

### Conclusion

This revolution, when it comes, will be digitized. Algorithmic criminal justice relying on first machine learning and then on deep learning, is only now beginning to impinge on actual, existing criminal justice institutions. The latter have been enduring sites for the production of racial stratification, not only in the form of a policing and carceral apparatus that weighs most heavily on African-Americans but also in terms of a racial tax that extends to all members of the group, whether or not they have any connection to criminality.

Given this history, it seems to me important to get algorithmic criminal justice right. Such tools, if fashioned wisely, might be useful in restoring equilibrium and mitigating the burden of racial externalities. Wrongly configured, they may prove subtle levers for preserving and exacerbating those burdens. My aim in this Article has been to demonstrate that constitutional law does not contain effectual tools to meet these problems. It is a mistake, therefore, to contort constitutional doctrine in the hope that it will do service in a context where it is so substantially ill-fitted. Far better, in my view, to recognize that constitutional law has almost nothing useful to say about what counts as a racially just algorithm—and might instead achieve the remarkable double-header of impeding both racial equity and social welfare maximization. The doctrine is thus a moral vacuity.

Reformulation of the doctrine, in my view, is desirable but unlikely. In the interim, algorithm designers, local officials, and state legislators should instead ask directly how best to achieve racial equity given the shape of existing criminal justice institutions and the technical tools at their disposal. I have offered an answer to that question that draws on, without quite tracking, existing technical definitions of algorithmic nondiscrimination. I have further stressed that my approach has the distinctive feature of aligning racial equity with social efficiency. My project has been demarcated in terms of algorithmic criminal justice. But it should not escape notice that there is no particular reason to confine the scope

---

<sup>346</sup> See text accompanying *supra* notes 219 to 225.

<sup>347</sup> *Johnson v. California*, 543 U.S. 499, 505 (2005).

of the analysis to algorithmic tools, or even to criminal justice. But those extensions are for another day. For now, a recognition of the potential convergence of equity and efficiency might move us closer to a remedy for the difficult, enduring, and damaging legacy of our racialized criminal justice past.