

# Algorithms and the Individual in Criminal Law

Renée Jorgensen, University of Michigan ([rjorgen@umich.edu](mailto:rjorgen@umich.edu))

Forthcoming in the *Canadian Journal of Philosophy*. Please see published version for citation purposes.

## 1 Background

Legal license to treat an individual in certain ways—subjecting them to special scrutiny, detaining them, or ruling that they are liable to penalties—often depends on our having a sufficiently high degree of confidence, given the available evidence, that such treatment would be appropriate. When making these determinations individual police officers or judges sometimes rely exclusively on their own observations and evaluation of the available information; sometimes they consult or defer to the judgment of an expert. Someone attentive to how vulnerable human decision-making is to cognitive biases of all sorts—not to mention the distortions introduced by prejudice—might hope to improve these decisions by basing them on algorithmic predictions.<sup>1</sup> Rather than relying on a single agent’s personal assessments of whether a particular suspect is likely to reoffend, for example, we might hope to leverage historical trends in re-arrest or re-conviction data to yield some objective measures, insulated from individual irrationality or animus.

This, in its most optimistic frame, was the driving promise and aim of the risk assessment tools first developed to aid parole decisions in Chicago in 1933. The original model leveraged factors like marital status, behavioral infractions within the detention facility, prior arrest record, and work history to sort inmates into nine rough cohorts, and furnished parole boards with the relative frequency of re-arrest among past members of an offender’s cohort as an indicator of the probability that the inmate would reoffend (Burgess 1936-7, 499). Since then statistical ‘Risk and Needs Assessment’ (RNA) tools have been refined and proliferated; they are now used in the majority of jurisdictions in the United States to guide pre-trial decisions relating to whether (and how high) to set bail, as well as post-trial determinations concerning whether to divert a defendant from incarceration and when to consider them for parole.<sup>2</sup>

---

<sup>1</sup> Suggestions of this kind are made in J. Kleinberg, S. Ludwig, et al., ‘Algorithms as Discrimination Detectors’ (2020), and Taslitz, ‘Police are People too: Cognitive Obstacles to, and Opportunities for, Police Getting Individualized Suspicion Judgment Right’ (2010).

<sup>2</sup> The most commonly used tools include the Arnold Public Safety Assessment (PSA), the Virginia Pretrial Risk Assessment Instrument (VPRAI), the Ohio Risk Assessment System (ORAS), the Correctional Assessment and Intervention System (CAIS), and the Level of Service/Case Management Inventory (LS/CMI).

Most of these tools employ straightforward statistical analysis on historical arrest databases, seeking to isolate the strongest correlations between a relatively sparse set of recorded variables and a property representing the target outcome (e.g. failure to appear, another arrest, or arrest for *violent* offense). There is some variation in the variables used: ‘third generation’ risk assessment measures improve on the original ‘second generation’ models<sup>3</sup> by using not only *static variables*—properties that do not change over time, like age at first arrest, having a prior conviction, sex, etc.—but also *dynamic variables* (e.g. years since last offense, employment status, present substance abuse) which are responsive to the subject’s current behavior, and can reflect reduced (or increased) risk over time. Simplifying a bit, these tools are ultimately algorithms taking the variable values as inputs, assigning them weights, and outputting an estimate of how often someone with those features in the database ends up with the target outcome.

There is now a new wave of tools with a somewhat different structure, more aptly described as applications of artificial intelligence. They leverage a wide array of information in vast databases to *train* a model to recognize patterns in an existing dataset, in order to predict the outcome value for a new entry. This method allows the trained model to discover previously unnoticed correlations between the target outcomes and properties in the dataset. The hazard is that the correlations might be artifacts of the particular dataset, rather than robust connections in the underlying phenomena.<sup>4</sup> As these machine-learning techniques improve, programs using them have become increasingly effective at a wide variety of recognition and classification tasks, and have been pressed into service for a range of prediction tasks, too.<sup>5</sup> The allure of these tools is that they offer a chance not just to make an informed guess about how often an outcome will occur in some set, but to identify and intervene in a case *before the predicted outcome occurs or becomes acute*: removing pre-cancerous tumors, connecting struggling students with extra resources, or, in the case of criminal justice, *preventing* a predicted victimization.

Buoyed by enthusiasm for data-driven policing and sentencing, both sorts of tools have made their way into law-enforcement at several points. To name just a few examples: PredPol and HunchLab are used by police departments across the United States to identify hotspots for property crime, assault, and auto theft. Several states use Palantir’s data analysis program *Gotham* to leverage information aggregated from service and arrest databases in order to guide the allocation of police resources and aid in suspect

---

<sup>3</sup> ‘First generation’ risk assessment refers to the unstructured clinical assessment, often based on an interview with the subject, that pre-dated the widespread use of the actuarial tools.

<sup>4</sup> There is another hazard in the context of criminal justice, which I will discuss later—that there will be robust connections *which are themselves unjust*.

<sup>5</sup> Many of which are discussed and criticized in O’Neil *Weapons of Math Destruction* (2016), Eubanks *Automating Inequality: How High-Tech Tools Profile, Police, and Punish the Poor* (2018), and Brayne *Predict and Surveil: Data, Discretion, and the Future of Policing* (2021).

identification and profiling. These include at least the Chicago Police Department's 'Strategic Subjects Initiative' (SSI), LAPD's 'Los Angeles Strategic Extraction and Restoration' (LASER) program, and the Northern California Regional Intelligence Center. Police departments in New Orleans and New York city had similar contracts.<sup>6</sup> On the post-arrest side of things, the Correctional Offender Management for Profiling Alternative Sanctions (COMPAS) program, developed in 2002, is used both to make pre-trial determinations about bail and post-trial determinations about sentencing and parole throughout Michigan, Wyoming, Wisconsin, California, and in an ever-increasing number of counties in other states.<sup>7</sup>

Not everyone welcomes the increased use of algorithmic prediction tools. Attorney General Eric Holder, for instance, cautioned that

“Although these measures were crafted with the best of intentions, I am concerned that they may inadvertently undermine our efforts to ensure *individualized and equal justice*. By basing sentencing decisions on static factors and immutable characteristics – like the defendant's education level, socioeconomic background, or neighborhood—they may exacerbate unwarranted and unjust disparities that are already far too common in our criminal justice system and in our society.” (Holder 2014), emphasis added.

Some of the obvious ethical concerns about using algorithms in criminal justice---bias in error rates, data looping, redundant encoding, etc.---have already been the subject of significant academic and media attention. COMPAS in particular has drawn substantial criticism for having racially biased error rates.<sup>8</sup> There are also a handful of epistemic concerns that have been raised about the probative value of risk scores, including worries that these systems base predictions on spurious correlations, are lacking in explanatory value, and give fixed datapoints from a person's past too much weight to be epistemically reliable in predicting whether they will reoffend.<sup>9</sup> In this paper, I will set all of these aside

---

<sup>6</sup> LAPD suspended their use of LASER after significant public protest. NOPD suspended their contract with Palantir in early 2018, after public backlash at the secrecy of the initial arrangement and terms. Palantir had confidential contracts with a number of city police departments, including NYPD, and it is unclear how many are ongoing. Other prominent clients in the United States include the Central Intelligence Agency, the Department of Homeland Security, Immigration and Customs Enforcement, Department of Health and Human Services, and the Center for Disease Control.

<sup>7</sup> See Herrschaft, *Evaluating the Reliability and Validity of the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) Tool: Implications for Community Corrections Policy* (2014); Kehl, Guo and Kessler, 'Algorithms in the Criminal Justice System: Assessing the Use of Risk Assessments in Sentencing' (2017). COMPAS was developed and is managed by a private company (Northpointe), which re-named itself 'Equivant' in 2017.

<sup>8</sup> While COMPAS is equally likely to misclassify defendants of any race *in one way or another*, it is disproportionately likely to misclassify a Black defendant as high risk, and disproportionately likely to misclassify a white defendant as low risk. See Angwin, Larson, Mattu, and Kirchner, 'Machine Bias: There's software used across the country to predict future criminals. And it's biased against Blacks' (2016). Angwin *et al.* focused their analysis on predictions for people arrested in Broward County, Florida, between 2013-2014.

<sup>9</sup> PredPol, SSI, and LASER employ a mix of arrest, crime reporting, and conviction data, raising worries that enforcement bias and differing levels of confidence in the police distort the dataset in ways that compromise

in order to focus on a third type of concern visible in Holder’s remarks: whether pursuing criminal justice with these tools is consistent with treating the defendant as an individual.

The reason for this narrow focus is straightforward: while they are serious, the problems mentioned above are mostly problems with using algorithms *badly*. But no matter how we clean, debias, or supplement, all the tools in question trade in *actuarial inference*. The score assigned to an individual does not reflect any deep insight into his “true nature”; it reports the frequency of a type of outcome among people in the database who are similar to him with respect to the values of the predictor variables. Setting aside programs focused on predicting *locations* of crimes rather than individuals, the basis for risk predictions made by these algorithmic tools is, ultimately, observations about the behavior of *other* people. The reasoning structure is actuarial in that it moves from the conjunction of *the subject has feature G* and *the relative frequency of feature F among others with G is x* to confidence of approximately *x* in *the subject has feature F*. In a slightly different context—addressing the use of statistical or probabilistic evidence to establish liability in civil trials, or settle sentencing questions in criminal trials—several legal theorists, philosophers, and judges have objected that inferences of this form functionally make it a ‘crime to belong to a reference class,’<sup>10</sup> violating the right to be ‘treated as an individual’. We might reasonably ask whether this right also forbids the use of any of the algorithmic tools mentioned above. If so, this would be a problem not just with using the tools badly, but with using them *at all*. But it isn’t immediately obvious what the right to be treated as an individual forbids, because it isn’t clear what it is a right *to*, exactly.

Rather than trace the constitutional grounds or legal interpretation of this right, my project in this paper is to explore its core: what moral interests might it protect, and are those interests threatened by relying on the outputs of algorithmic methods in determinations of probable cause, guilt, or sentencing? After exploring a few different interpretations of the right (in §2), I ultimately propose understanding it as protecting

---

the fairness of the algorithms. See Richardson, Schultz and Crawford, ‘Dirty Data, Bad Predictions: How Civil Rights Violations Impact Police Data, Predictive Policing Systems, and Justice’ (2019). See also Bolinger, ‘Explaining Justificatory Asymmetries between Statistical and Individualized Evidence’ (2021), arguing that statistical evidence in general makes at most a marginal contribution to justified credences, while introducing a risk of error concentrated on particular demographic groups, and so over the long run, the evidential value of relying on it is outweighed by the moral cost of so doing.

<sup>10</sup> See, e.g., Colyvan, Regan and Person, ‘Is it a Crime to Belong to a Reference Class?’ (2001); Enoch, Spectre and Fisher, ‘Statistical evidence, Sensitivity, and the Legal Value of Knowledge’ (2012); Risinger, ‘Unsafe verdicts: The need for reformed standards for the trial and review of factual innocence claims’ (2004); Thomson, ‘Liability and Individualized Evidence’ (1986), and Tribe, ‘Trial by Mathematics: Precision and Ritual in the Legal Process’ (1971). A quite different claim, also referred to as the ‘right to an individualized decision’, is concerned with the use of algorithms not as applied to groups of people in order to predict the *fittingness* of some specific treatment or verdict, but to actually *decide cases* on the total evidence, where the learning data is all pre-existing caselaw (for instance), and the inputs are the facts of a given new case. These applications present very different problems, and for simplicity I will set them aside. For an informative discussion, see Binns, ‘Human Judgement in Algorithmic Loops’ (forthcoming).

agents' claims to a fair distribution of the burdens and benefits of the rule of law (in §3). What it forbids is not the use of probabilistic information or statistical methods, but treating wrongdoing by some as justification for imposing extra costs on others; it demands that we respect the separateness of persons. This has significant implications for the administration of criminal justice (explored in §4): it permits the use of predictive tools in principle, but only if the predictors are transparent, subject to agents' deliberate control, and unavoidable burdens are fairly distributed or outweighed by benefits to the very same individuals. In §5 I explain that since this bars relying on indices of distributive injustice or unchosen properties to determine whom to subject to extra costs associated with criminal justice, it precludes the most common applications of algorithmic tools for bail determinations or the regulation of street crime. However, they might be justifiable and effective in combatting white-collar crime.

I will where possible treat all the applications together, but the implications vary depending on the context in which an algorithmic prediction tool is used, so it will be helpful to have a sense of the variety of uses to which they are put. So, very roughly: some of the tools mentioned are used to guide the application of pre-arrest scrutiny, establishing *probable cause* to subject someone to additional search or surveillance, arrest, or detention. At the next stage, a risk assessment might be offered to support bringing charges, or as evidence about whether the subject will either reoffend or fail to appear if released on bail before trial. In theory (if not in practice) profile evidence could be offered at trial as evidence helping to establish beyond a reasonable doubt that the defendant is guilty of the crimes with which he is charged.<sup>11</sup> Post-conviction, it could be offered at the time of sentencing to indicate the probability of re-offending if given a short sentence, or fitness for diversion into non-carceral forms of punishment (e.g. probation, supportive housing or mental health assistance programs). Finally, it can be used at any stage within the sentence duration to determine eligibility for parole.

## 2 Interpreting the Right

We can start unpacking the right to individualized treatment by reviewing how judges have discussed the purpose and value of the requirement. Justice Stevens highlighted the rule's protective function for establishing probable cause in his dissent in *Samson v. California*: "The requirement of individualized suspicion, in all its iterations, is the shield the Framers selected to guard against the evils of arbitrary action, caprice, and harassment."<sup>12</sup> In an opinion rejecting the use of actuarial evidence for sentencing in

---

<sup>11</sup> For an overview of present uses of profile evidence for probable cause and at trial, see Harris, 'Particularized Suspicion, Categorical Judgments: Supreme Court Rhetoric Versus Lower Court Reality Under *Terry v. Ohio*' (1998).

<sup>12</sup> In *Samson v. California*, 547 U.S. 843, 860 (2006), the US Supreme Court (voting 6-3) affirmed the decision of the California Court of Appeal, that it does not violate the fourth amendment protections against unreasonable search to subject parolees to suspicionless search because it is a condition of parole that one consents to search by an officer with or without cause or search warrant. Justice Stevens authored the dissent, joined by Souter and Breyer.

*United States v. Shonubi*, Judge Newman emphasized that the ‘specific evidence’ requirement is only satisfied by “evidence that points specifically to [behavior] for which the defendant is responsible.”<sup>13</sup> A flat-footed reading of these comments might yield an interpretation contrasting individualization with generalized treatment, leading us to interpret the right as something like:

- *A claim that high-stakes legal decisions be personalized, rather than being subjected to ‘one-size-fits-all justice.’*

But this is too simplistic. As Harcourt (*Against Prediction*, 2007) stresses, relying on actuarial data actually allows our determinations to be highly *tailored* to the individual. For instance, rather than having broad sentencing categories, these methods enable us to fit the sentence to the strength of the correlation between the subject’s specific features and re-arrest or re-conviction. We could in principle make similarly personalized judgments about probable cause or reasonable suspicion using algorithmic tools, given the very large databases and high number of personalizing variables these methods allow us to take into consideration. But more personalized treatment isn’t necessarily better. As Lippert-Rasmussen (2011) points out, personalization may lead to our being treated *worse* than otherwise, and is in some tension with other weighty principles of justice, such as the generality and equal application of law, and the fair social distribution of various burdens. More urgently, this form of individualization does not capture the connection to the individual’s *responsibility* stressed in Judge Newman’s comments. The problem is not that statistical determinations are inadequately *personalized*, but that they are not appropriately *responsive*: they treat the individual according to how we expect him to act based on our experience with others like him, rather than how he himself has acted. Personalization is not the point.

Moral theorists prefer to characterize the relevant obligation as a duty to be responsive to the individual’s *responsible agency*, grounded in the values of autonomy or respect. For instance, Dworkin (1977) contends that detaining a person based on actuarial prediction, however accurate, is unjust “because that denies his claim to equal respect as an individual.” Duff (1998, 155-6) also anchors the claim in respect, holding that

“[t]o respect the defendant as a responsible citizen, we must treat him and judge him as an autonomous agent, who determines his own actions in the light of his own values or commitments. His membership of this actuarial group is part of the context of that self-

---

<sup>13</sup> *United States v. Shonubi* (103 F.3d 1085 2d Cir. 1997), at 1089-1090. The defendant (Charles Shonubi) was convicted of having smuggled heroin into the United States on eight separate occasions. At the original sentencing, the court multiplied the volume of heroin he was found carrying by 8, to estimate the total quantity smuggled across all his trips. Shonubi appealed on the grounds that the total amount hadn’t been proved with specific evidence, and the case was sent back for resentencing. The prosecution then used statistics about average drug seizures using the same method, arrested at the same airport, to estimate the total amount; Shonubi appealed again, and Judge Newman again returned the case for resentencing, explaining that “The statistical and economic analyses relate to drug trafficking generally and not to Shonubi specifically.” at p. 1091.

determination; and as observers, we might think it very likely that he will have determined himself as a criminal.”

Nevertheless

“respect for autonomy, and the ‘presumption of harmlessness’ which follows from it, forbids us to ascribe criminal dangerousness to anyone, unless and until by his own criminal conduct he constitutes himself as having such a character.”

Walen (2011) articulates the content of the state’s duty to respect the autonomy of its citizens in much the same way: “A state must normally accord its autonomous and accountable citizens this presumption [that they are law-abiding] as a matter of basic respect for their autonomous moral agency.” This is consonant with suggestions by Amour (1994), Duff (1998), and Moss (2018) that in general taking statistical generalizations as reason to conclude that an individual is *probably* dangerous runs afoul of the individual’s moral claims.<sup>14</sup>

But does treating individuals with appropriate respect really require approaching them as a completely novel case, without an expectation that our knowledge of other cases will give us reliable guidance concerning them? One might be skeptical whether merely predicting that an individual is likely to offend is a failure of respect or affront to their autonomy.<sup>15</sup> Some do argue that viewing someone as predictable in this way fails to regard them appropriately as an agent, rather than a thing determined by external pressures.<sup>16</sup> But we needn’t even embrace anything this strong to identify a moral problem with using actuarial predictions to warrant *harmfully interfering* with a person. Plausibly, it is morally negligent or reckless to intentionally harm someone unless we have not only reasonably high credence (e.g. above some threshold) that the action is morally appropriate, but this credence is *resilient*. Very roughly: our present evidence must be such that little if any new information consistent with it would cause our credence to drop below the threshold.<sup>17</sup> The more harmful the interference, the more resilient the credence must be to justify it. Even if *making* a prediction based on statistics is not a failure of respect for the individual’s agential freedom, using that prediction as

---

<sup>14</sup> Many have argued that something similar holds more generally: respecting others’ moral autonomy prohibits basing our appraisals of their character on statistical evidence. See, e.g., (Walen, *A Unified Theory of Detention, with Application to Preventative Detention for Suspected Terrorists* 2011); Buchak, ‘Belief, Credence, and Norms’, (2014); Moss, *Probabilistic Knowledge* (2018). Some also maintain that we wrong others specifically when we use statistics to draw inferences that diminish the subject or would lead us to act against their interest. See e.g. Basu, ‘The Wrongs of Racist Belief’ (2019); Schroeder, ‘When Beliefs Wrong’ (2018); Wasserman, ‘The Morality of Statistical Proof and the Risk of Mistaken Liability’ (1992).

<sup>15</sup> My thanks to Patrick Tomlin for raising this concern.

<sup>16</sup> See, e.g. Marušić and White ‘How Can Beliefs Wrong? -- A Strawsonian Epistemology’ (2018); Basu, ‘What We Epistemically Owe to Each Other’ (2019); Duff, ‘Dangerousness and Citizenship’ (1998).

<sup>17</sup> For a better discussion of resilience, see Joyce, ‘How Evidence Reflects Probabilities’ (2005); Buchak (2014); and Moss *Probabilistic Knowledge* (2018). For an argument that the resilience requirement explains the intuitive justificatory limits of statistical generalizations, see Bolinger, ‘The Rational Impermissibility of Accepting (some) Racial Generalizations’ (2020), and Bolinger, ‘Explaining Justificatory Asymmetries’ (2021).

grounds for harming them is a failure of respect for their agential *status*, because unless supplemented, statistical evidence cannot be adequately resilient.

At a bare minimum, civic respect and equality still requires that the *default* orientation of law-enforcement to any member of the political community not be one expressive of suspicion or disrespect. Considering a person to be *probably law-abiding* orients police to respect and protect them; considering them to be *probably lawbreaking* activates a very different script. There are reasons to doubt that in practice this difference in default orientation is primarily responsive to *evidence*, rather than stereotypes or group-based prejudice.<sup>18</sup> But even if it tracked group-level rates of arrest, default suspicion would fail to treat citizens as they are entitled. Generalizations about ‘types of people’ or trends in broad demographic categories aren’t sufficient to justify suspending the civic respect owed to a particular person, that demands treating them as *probably law-abiding*. Borrowing heavily from Duff’s language, we might articulate this as:

- *A claim to be respected as a presumptively law-abiding citizen, unless and until one defeats this presumption through one’s own action and behavior.*

The central role this gives to respect and autonomy seems on the right track, and to capture much of the intuitive moral core of the demand that treatment be individualized. But it doesn’t explain what is objectionable about actuarial inferences after guilt has been established; once the presumption has been defeated by admissible, individualized evidence. If the right requires nothing more than that we treat agents as law-abiding until we have adequate particularized evidence that they aren’t, then there is no conflict at all between it and the use of algorithmic risk scores in making sentencing determinations. While one could simply accept this conclusion, it seems to me that the ‘presumption of law-abidingness’ does not exhaust the obligations grounded in civic respect.

What else might it entail? Perhaps

- *A claim to not be subject to extra burdens simply on account of one’s social identity.*

This is the interpretation naturally suggested by Colyvan, Regan and Person (2001), and rejected as unrealistically idealistic by Tillers (2005). There are two ways to develop the thought that equality of standing or respect entitles individuals to be free from *extra* burden, and I find both plausible. On the one hand, we might be concerned about being subject to *disproportionate* burdens associated with law-enforcement, relative to other groups; this is the central animating idea behind Bambauer (2015)’s explication of why statistical evidence should not be used to establish probable cause. On the other hand, we might worry about being subject to burden that they would not be subject to were we to

---

<sup>18</sup> Analysis of transcripts of traffic stops in Oakland, CA found that police officers speak significantly less respectfully to black than to white community members, even after controlling for officer race, infraction severity, stop location, and stop outcome. Voigt, Camp, Prabhakaran, Hamilton, Hetey, Griffiths, Jurgens, Jurafsky, and Eberhardt, ‘Racial Disparities in Police Language’ (2017)

hold all else *except their social identity* fixed; something like this the centerpiece of Underwood (1979)'s explanation of why racial membership and other protected categories are an inappropriate base for statistical prediction. She grounds protection against the use of these and other unalterable features in the value of autonomy: "Of all the factors that might be used for predictive purposes, those beyond the individual's control present the greatest threat to individual autonomy. Use of such factors in a statistical prediction device is particularly undesirable if the device is to be used in a context in which autonomy is highly valued."<sup>19</sup> This emphasis on preserving the individual's *control* gives us reason to also object to adding burdens to social identities that aren't themselves protected, but are either unchosen (e.g. socio-economic status), or reflect important personal choices (e.g. marital status).

Some offer a more procedural gloss of the right, arguing that it is actually a proxy for the right to a certain kind of *explanation* for the state's decisions in her particular case:<sup>20</sup>

- *A claim to an explanation for the State's exercise of coercive powers.*

Vredenberg (working paper) compellingly argues that the value of explanations of this kind is instrumental. Access to such explanations is a prerequisite for agents' ability to act on the political system, to hold it accountable, and form the rules which characterize the basic structure of society. The right so understood requires both more and less than that the subject be given a true account of why a legal decision concerning her has been made; the explanation offered must equip her to act intentionally to hold the decision-making body accountable. It therefore must both be intelligible to her, and bear some relation to the *actual* decision-making procedure employed. In understanding the moral core of the right to individualization as a right to the information necessary to form and reform the legal policies, this gloss aligns closely with Justice Stevens' comment that the right is a shield against arbitrary uses of state power.

Each of these interpretations highlights something of value in the content of a right to be treated as an individual. Rather than offer a competing interpretation of 'individualization', I suggest that in fact the right doesn't protect a *unique* interest—rather, the entitlement to 'be treated as an individual' simply demands that the procedures of the criminal law be justifiable not merely in the aggregate but to each individual subject to them. Rephrased, it is:

- *A claim to fair distribution of the benefits and burdens of public law.*

So interpreted, the right entails that a person must not face disproportionate burden or suspicion except as a consequence of their own *responsible action*, and that agents of the

---

<sup>19</sup> Underwood, 'Law and the Crystal Ball: Predicting Behavior with Statistical Inference and Individualized Judgment' (1979), at p. 1436.

<sup>20</sup> This interpretation is implicit in the 'explanationist' strand of the legal literature on statistical evidence, which contend that statistical evidence should be inadmissible in trials because it is inadequately *explanatory*, or supplies probabilistic support without raising the *plausibility* of the hypothesis that the defendant is guilty.

state default to respectful engagement. It is grounded ultimately in the preconditions for laws to be *fair*, both in their content and administration. The crowning virtue of the rule of public law is its ability to shape citizens' practical reason and ground reliable expectations, enabling them to hold each other to standards which they had fair opportunity to meet. Its benefits include many of the values articulated by the interpretations we've surveyed: expressing respect for autonomy, constraining the exercise of coercive power, ensuring that sanctions are responsive to responsible agency, and ensuring that those subject to law are in a position to challenge or reform it. The burdens, meanwhile, are the various costs associated with the scrutiny and punitive sanctions applied in the course of enforcing the laws. What fair distribution of burdens and benefits demands depends on context: pre-conviction, every individual must have fair opportunity to avoid hostile encounters with law enforcement; at trial they must not face disproportionate likelihood of false conviction; post-conviction they must not be subject to disproportionate punishment.

To afford all individuals a fair opportunity to avoid hostile encounters with law enforcement the laws must be public, clear, and prospective.<sup>21</sup> These are necessary conditions on its ability to structure citizens' relationships to each other and the state in a way that expresses respect for their autonomy and equality as agents. Citizens can't know *what* the law forbids if its requirements are secret or inscrutable; if it is retroactive, they cannot act intentionally to avoid violating it.<sup>22</sup> Importantly, these requirements take lexical priority over considerations of administrative efficiency: the *value* of the rule of law is lost if the laws are applied in ways that do not facilitate the mutual accountability of citizens and state. Choices about the administration of public law have drastic implications for individuals' freedom, ability to pursue their life projects, and participation in the political community. Legal transgressions expose a person to the coercive power of the state in various forms, ranging from asset forfeiture to deprivation of liberty and loss of civic rights. General appeals to the aggregate value of efficient crime reduction cannot justify or compensate a particular person for their loss of crucial protections against suffering the State's coercive imposition of these harms. It is good to lower the average risk of suffering victimization, but if an overall reduction is achieved by dramatically increasing the burdens borne by particular members of the population in a way that neither tracks their responsible action nor is offset by equally weighty benefits to *them*, the burdens and benefits of the rule of law are not fairly distributed.

---

<sup>21</sup> I am primarily focused on a notion of 'fair opportunity' that is non-comparative, demanding simply a normatively sufficient chance of avoidance. But a comparative conception of fairness is also relevant here, requiring that a subject have *not substantially worse chances* of avoidance than others in the political community (I am indebted to comments from Chad Lee-Stronach for this point).

<sup>22</sup> My analysis in this section strongly echoes Fuller's articulation of the value of the rule of law, particularly as developed and defended by Murphy, 'Lon Fuller and the Moral Value of the Rule of Law' (2005). Fuller, *The Morality of Law* (1969, p. 106) gives eight requirements for the rule of law: law must be (1) general, (2) publicly accessible, (3) prospective rather than retrospective, (4) clear, (5) non-contradictory, (6) possible to satisfy, (7) stable, and (8) there must be congruence between what the law requires and what is enforced.

I suggest, then, that the right issues an injunction not against the use of probabilistic information or generalizations, but against a familiar form of moral aggregation. Just as invocations of the separateness of persons are in other contexts made to assert that benefits to some cannot offset harms to others, the demand that we treat people *as individuals* here asserts that wrongdoing by some does not weaken others' moral claim against the imposition of extra costs, even if they belong to the same demographic group. The right to individualized treatment is the right of those subject to a criminal law that its procedures be justifiable to them individually, not only by appeal to average outcomes. This lens is both unifying and clarifying: it explains why each of the earlier glosses feels partly—but only partly—right. And unlike the other interpretations, which identify relatively all-purpose goods or moral interests, this reading of the right gives it a content specific to criminal law.<sup>23</sup>

Accepting this interpretation has relatively revisionary implications for the administration of criminal law. While I am mainly exploring the limitations the right imposes on algorithmic prediction tools, it is worth noting that it constrains non-predictive administrative decisions, too. Consider the practice of cash bail (allowing individuals to avoid pre-trial detention conditional on paying a sizable fee, typically \$10,000, refunded if they appear on their scheduled court date). The immediate effect is to ensure that some of the most severe burdens of an encounter with the law—lengthy pre-trial detention, during which the defendant incurs a variety of costs often including job-loss due to forced absence—fall disproportionately on the very poor.<sup>24</sup> The median yearly income for people who are detained pre-trial because they cannot post a bail bond is just \$15,109; 37% make less than \$9,489 per year. This practice leaves those below the poverty line exposed to disproportionately severe burdens, without affording them offsetting benefits adequate to justify the imposition. Whether bail determinations are made in highly personalized ways or in deference to an algorithmically produced risk score, insofar as they distribute the burdens and benefits of the rule of law unfairly, they violate the moral interest that animates the right to be treated as an individual.<sup>25</sup>

---

<sup>23</sup> My thanks to Tom Parr for pointing this out. Importantly, I do not mean to imply that we have a moral interest in being treated as an individual *only* in the domain of legal decisions. The relationships of respect and answerability that law formalizes may extend to informal, interpersonal interactions, and so plausibly the interest protected by a formal right to an individualized decision may persist in as a moral claim in informal contexts. My thanks to Deborah Hellman for discussion on this point.

<sup>24</sup> Rabuy and Kopf (*Detaining the Poor: How money bail perpetuates an endless cycle of poverty and jail time*, 2016)'s analysis of data released by the Bureau of Justice Statistics found that "the median bail amount [\$10,000] represents eight months of income for the typical jailed defendant."

<sup>25</sup> Reforms that waive bail if a defendant receives a low risk score reduce the number of poor who are subjected to pre-trial detention, but concentrate its effects more heavily on residents of poor communities of color. Rather than using risk scores to filter its effects, fairly distributing the burdens and benefits of law would require that we do away with cash bail as a general practice. My thanks to Tali Mendelberg for bringing this case to my attention, and to Sarah Stroud for pointing out the range of implications the right to individualized treatment has (if I am right) for the administration of criminal law beyond questions about the use of algorithms.

### 3 Individual Treatment and Algorithmic Tools

On my analysis, the right to be treated as an individual is in principle consistent with the use of algorithmic tools in criminal law. But rather than make a particular policy proposal, I think it more fruitful to fill in the general contours of the constraint, articulating more specifically what it means for a policy to be justifiable to those subject to it. As I have unpacked it, what the right requires is not that the application of legal sanctions be *personalized*, but that they be justifiable to each individual subject to them. This is necessary for the administration of criminal law to express appropriate respect for each individual's autonomy and preserve the mutual accountability of citizens and state. Applying this specifically to predictive tools, I claim that if a factor  $f$  is used as a predictor in the administration of criminal law, at least three conditions must be met:

- [CONTROL] --  $f$  must be subject to agents' deliberate control
- [TRANSPARENCY] -- It must be *transparent* that  $f$  is used as a predictor, such that the basis for decisions is sufficiently clear to facilitate civilian criticism or reform.
- [BURDENS] -- The unavoidable extra burdens imposed by using  $f$  as a predictor (increased hassle, risk of false conviction, or severity of punishment) must be outweighed by the benefits it yields *to the individuals who must bear these burdens*.

I suspect that in practice, there are few applications within the administration of criminal law where predictive algorithms can be deployed while satisfying all of these conditions.

Let's start with the *Control* condition. We said earlier that to avoid subjecting anyone to more than their fair share of burden, legal sanctions must be tied to responsible agency. It follows that it must at least be *in-principle* possible for an individual of any permissible social identity to avoid suffering a downside cost that is not born by everyone. So if a high risk-score suffices to justify extra scrutiny, conviction, or a lengthier sentence, the predictor variables can't be tied to unchangeable identity-tracking properties like race or gender. But mere in-principle avoidability is not sufficient to ensure that extra costs track *agency*. So the properties used to indicate criminality—and thus to determine the distribution of costly legal sanctions—must also be ones that at least the actually law-abiding agents could act to avoid. They can't be things that subjects have little real chance of escaping, like residence in high-crime neighborhoods, poor educational background, or an unstable family environment.<sup>26</sup> We will likely find robust correlations between these features and criminal offending rates, especially in historical databases—but it would violate the right to individual treatment to leverage such correlations to justify the

---

<sup>26</sup> This intersects with a dilemma arising from antecedent distributive injustice: children who grow up in concentrated urban poverty do not have prospects of avoiding criminality comparable (or even close) to those with different social starting positions. For a discussion of this dilemma, see especially Ewing, 'Recent Work on Punishment and Criminogenic Disadvantage' (2018); Howard, 'Moral Subversion and Structural Entrapment' (2016); Kim, 'Entrapment, Culpability, and Legitimacy' (2020); Shelby, 'Justice, Deviance, and the Dark Ghetto' (2007); and Watson 'A Moral Predicament in the Criminal Law' (2015).

predictive application of criminal sanctions, because doing so concentrates extra hassle and risk on individuals on the basis of factors over which they lack agential control.

The *Transparency* condition articulates a precondition for the rule of law. Law expresses respect for subjects' autonomy only when it enables citizens to anticipate what compliance requires from them. So, in addition to the brute ability to avoid properties that would lead to having a high risk-score, subjects must also be able to *act intentionally to avoid* them—which means they need to be able to know which variables are used, and roughly how. Finally, the *Burdens* condition acknowledges that some costs are unavoidable, and cannot always be distributed perfectly evenly across the population. It allows the imposition of these costs, but *only if* they are offset by benefits to the individuals who bear them.

#### 4 The Space for Prediction

One might worry that my interpretation of these constraints is too strict; that the population-level gains to efficiency or accuracy justify violating at least one of them in some contexts. For instance, can't the deterrence-benefits of using algorithmic predictions to guide reasonable suspicion outweigh individual subjects' moral complaints against failures of transparency, provided that the algorithms are sufficiently accurate?

##### 4.1 *Secrecy and Strategic Gaming*

One immediate argument for this kind of tradeoff appeals to the importance of keeping prediction factors secret in order to avoid "strategic gaming": subjects deliberately manipulating or avoiding the predictors while continuing to engage in the targeted behavior (e.g. criminal offenses). First let's get clear on the assumptions behind this objection.<sup>27</sup> Strategic gaming is problematic only under very specific conditions:

- (i) the proxy criteria (the predictors) only weakly or contingently correlate with the target criteria,
- (ii) the proxy properties are within subjects' deliberate control (they are alterable),
- (iii) the tradeoff costs of gaming the proxy criteria are low, and
- (iv) moreover this can be done without affecting the subject's true eligibility with respect to the target criteria.<sup>28</sup>

If any one of these conditions is not met, then either a subject's attempt to game the proxy will *also* change how they fare with respect to the target, or the difficulty involved in strategic gaming will offset the incentive. To illustrate: LSAT scores are an oft-used proxy for the facility of reasoning needed for success in law school (the target criteria).

---

<sup>27</sup> For much of the following discussion, I am indebted to immensely helpful conversations with Kathleen Creel.

<sup>28</sup> I've drawn these conditions for problematic strategic gaming from Cofone and Strandberg, 'Strategic Games and Algorithmic Transparency' (ms).

But they are also robustly connected to that target: students who strategically aim only to improve their LSATs—enrolling in test-prep courses and practicing critical reasoning skills—thereby also make themselves better candidates with respect to the target criteria. So while schools’ transparent reliance on LSATs incentivizes students to focus on improving their test scores, this facilitates, rather than undermines, the end goal of admitting well-prepared students.<sup>29</sup> The possibility of strategic gaming fails to provide even a *pro tanto* justification for keeping proxy criteria secret unless gaming would undercut the aims.

Similarly, if the proxy for criminal wrongdoing is robustly connected to wrongdoing—e.g. if affiliation with a violent organization like the *Proud Boys*, or performance of preparatory acts like buying a high-capacity magazine for a firearm or purchasing large quantities of ammonium nitrate fertilizer are the chosen proxies—publicity can be net-beneficial. By incentivizing avoidance of the proxy, transparency discourages the linked criminal behavior. It also equips those for whom the proxy was misleading to avoid or politically contest decisions that rely on it, thus reducing the false-positive error rate. Especially when the costs of a false-positive are comparatively high, these error-correcting tendencies of transparency can be expected to outweigh the costs of strategic gaming. But precisely because using predictive proxies attaches costs to behaviors that are not themselves wrongdoing, and so shapes the behavior of those subject to it, Underwood (1979, at p. 1438) cautions that “[r]espect for autonomy thus counsels not only against the use of uncontrollable factors, but also against the use of those controllable factors that involve behavior generally regarded as private and protected against official interference.” Even where predictive, the range of properties used as proxies to guide the application of criminal sanctions will need to be tightly constrained to ensure that it does not intrude too far on autonomy.

The need to keep a proxy secret arises when all four of (i)-(iv) above are met. Given the way I have characterized the right to individualized treatment, plausibly anything consistent with it will satisfy conditions (ii) and (iii): that the proxy be within subjects’ deliberate control, and low-cost to alter or avoid. So strategic gaming could be a genuine concern if there are compelling reasons to use a highly contingent proxy (satisfying [i]) that is strongly independent of criminal wrongdoing (satisfying [iv]). There may be many administrative decisions for which it is permissible to use secret proxies, but I contend that, with few exceptions, the administration of criminal law is not one of them. Using a property as a proxy for criminality imposes significant costs on those who have it—at *least* high risk of ‘hassle factor’ (the costs associated with being subjected to extra scrutiny), at worst high risk of suffering unwarranted punishment or

---

<sup>29</sup> There is a different worry about using LSAT scores as a proxy for lawschool readiness: students without access to testprep resources, but otherwise equally promising, will be excluded by this proxy. Since access to a top lawschool is one of the means of social mobility, there is a legitimate concern that using the proxy unjustly skews access to those with higher disposable family income. While important, this is not ultimately a concern about strategic gaming. My thanks to Geoff Sayre-McCord for discussion on this point.

assault by agents of the state. When the proxy is only weakly connected with criminal wrongdoing, the state can neither justify its decision to secretly use it by appeal to the harm principle, nor necessity, nor to the decision's having been ratified by a democratic decision-making process. And when reliance on the proxy concentrates the highest costs of false-positives disproportionately on an already disadvantaged subpopulation, members of that subgroup have a dual complaint against secrecy: one against the ways that attaching costs to the proxy property undermines their autonomy, and one against the way that the choice of proxy fails to treat their subgroup as political equals. When there are adequate alternative means of enforcing the law, the presumptive weight of either of these complaints defeats the marginal administrative efficiency that could be achieved by a secret proxy.

#### 4.2 *Opacity*

A thoroughly different argument against transparency holds not that it is *undesirable*, but that it is *impossible*: the ways a sophisticated algorithm arrives at its predictions can be too complex to understand, let alone explain. No matter how much we might want to be transparent about the reasons why these algorithms make the predictions they do, the best we can do is describe how the algorithm was trained. But though the thought that predictive algorithms are essentially a “black box” has captured the popular imagination, it is something of a red herring in this context. Not all predictive algorithms are uninterpretable; only those arising from unconstrained or unsupervised ‘deep’ learning using high-dimensional models present this particular challenge. When they are comparably accurate, more transparent algorithms are preferable since opacity can obscure errors and makes it difficult to troubleshoot. And it is unlikely that either high-dimensional models or deep learning methods will be necessary—or much help—for optimizing the predictive accuracy of algorithms specifically in the context of criminal law. Though great advances have been made in recognition and automated judgment tasks, machine-intelligence has yet to stably out-perform simple rules at predicting social outcomes (like arrests), consistently plateauing around 65-70% accuracy overall.<sup>30</sup>

COMPAS is no exception: though it leverages a complex model, using up to 137 features of an individual's file to predict risk of being arrested for any new offense within two years of release, it only achieves about 68% overall accuracy.<sup>31</sup> What this means is that roughly two-thirds of the time either the person was classified as low-risk and in fact was not rearrested within two years, or they were classified as medium or high risk and were rearrested. An independent audit of COMPAS's predictions by Angwin, et al. (2016) found that the program had a slightly lower accuracy rate for those it classified as high-

---

<sup>30</sup> Narayanan, How to Recognize AI Snake Oil (2019 working paper); Yang, Wong and Coid, ‘The Efficacy of Violence prediction: a meta-analytic comparison of nine risk assessment tools’ (2010).

<sup>31</sup> Northpointe invoked trade secrets to avoid disclosing the details of their model, but their in-house evaluation of their software put overall accuracy at 68%. See Dieterich, Mendoza and Brennan, *COMPAS risk scales: Demonstrating Accuracy, Equity, and Predictive Parity* (2016).

risk (61%), but *much* lower accuracy when predicting *violent* reoffending specifically: only 20% of those classified as highly likely to be rearrested for violent crimes actually were.

A predictive accuracy rate of 65-70% is roughly on par with the *far* simpler models used by the second-generation risk assessment tools developed in the 1970s. Dressel and Farid (2018) found that a standard linear predictor using just 7 static features (age, sex, number of juvenile misdemeanors, number of juvenile felonies, number of prior crimes, crime degree, and charge) yields results comparable to COMPAS's predictor.<sup>32</sup> In fact, they found that untrained subjects who were given just these datapoints about each case and asked to make a prediction (without receiving any particular instruction as to *how*) also outperformed COMPAS in overall accuracy, and displayed slightly *less* racial bias.<sup>33</sup> Perhaps most startlingly—and underscoring just how far our prediction tools are from the imagined pre-crime oracles of *Minority Report*—all of these predictive methods were out-performed by a crude predictor with just two static factors: birthdate and number of prior convictions. While these facts should raise serious moral concerns about relying on the predictions yielded by these algorithmic tools when making high-stakes decisions, they provide one point of reassurance: uninterpretable models pose no special hurdle to transparency for our purposes, because they aren't all that useful for the applications of interest to us.

#### 4.3 *Moral Hazards of Training Predictive Models*

There is a deeper reason to generally avoid using deep machine learning to develop prediction tools for criminal law. These methods require substantial training data in order to learn predictive patterns, but it is treacherous to use the extant databases (requests for service, crime reports, arrests, or convictions) for this purpose. Information recorded in these datasets is invisibly shaped both by administrative discretion and by upstream structural injustices that artificially forced overlap between communities of color and criminogenic conditions--most especially underfunded schools and depressed economic conditions.

Some hope that we can correct for this with more or bigger datasets: given rich enough data, factors that are unrelated to the outcome of interest won't correlate closely

---

<sup>32</sup> At p.3. It's worth noting that since offending is measured by *arrest* (or in some cases *conviction*), rather than directly observed, some proportion of these tools' accuracy is just their ability to predict arrest patterns, which are subject to enforcement bias.

<sup>33</sup> Dressel and Farid ('The Accuracy, Fairness, and Limits of Predicting Recidivism', 2018) ran two studies with untrained subjects. In the first condition, participants were given *just* the seven features listed; in the second, they were also told the defendant's race. COMPAS has a recorded overall accuracy of 64.9% for Black defendants, 65.7% for white. It has a 40.4% false-positive error rate for Black defendants, 25.4% for white; and false-negative error rates of 30.9% and 47.9%, respectively. By comparison, Dressel & Farid's subjects had an overall accuracy of 68.2% for Black defendants, 67.6% for white (in condition II this dropped to 66.2% and 67.6%, respectively); false-positive error rates of 37.1% (this rose to 40% in condition II) for Black defendants, and 27.2% (26.2% in condition II) for whites; and false-negative error rates of 29.2% (rose to 30% in condition II) for Black defendants and 40.3% (42.1% in condition II) for whites.

enough with it to be reliable predictors, and so will not be learned.<sup>34</sup> But this optimism is misplaced when the information in available datasets is relatively sparse, or the overlap between properties is not accidental but *artificial*, or the outcome can only be measured or represented indirectly through measures (like ‘arrests’) that are themselves shaped by unrelated factors (like the probability of detection, political influence, trust in the police, or familiarity with legal protections). As Johnson (2020) demonstrates, even explicitly coding a model *not* to use properties like race or gender as predictors will not prevent it from learning to make predictions that track these features: "Where there are robust correlations between socially sensitive attributes, proxy attributes, and target features, and we've ruled out using the socially sensitive attributes, the next best thing for the program to use will be the proxy attributes." Put simply, algorithms trained on datasets in which decisions to arrest, conviction, sentence, and re-arrest have been subject to racial bias can be expected to learn correlations that, though *causally* spurious, are genuinely “there” in the data, projecting these traces of past injustice forward.<sup>35</sup>

For street crime in particular (including robbery, vehicle theft, arson, homicide, and assault) many of the strongly correlated properties are straightforward measures of socioeconomic disadvantage: employment status, income, education level, prior contacts with police, and relative security of housing. So it is doubtful that an algorithm trained on the available datasets would be able to respect the constraint that predictors be limited to factors within subjects’ deliberate control. But even bracketing these concerns about available training data, and *even if* the predictions made were highly accurate, the right to individualization as I have interpreted it may more directly preclude using deep machine-learning in developing the algorithms. A learning method which bases the risk prediction on correlations that *emerge* between very large numbers of variables and the outcomes is necessarily backward-looking and opaque. Insofar as it finds unexpected or surprising relationships, and bases new verdicts on these, it tends toward retroactivity, imbuing properties that had been considered harmless with criminal significance after the fact. If we cannot anticipate which properties will yield a high risk-score, then we cannot satisfy the requirement to be *prospective*. So the right as I have glossed it precludes the use of unsupervised deep machine learning, not *only* because it is unexplainable, but because it cannot articulate expectations adequately transparent and avoidable to perform the functions crucial to public law.

## 5 Some upshots

We began with a simple question—is the use of algorithmic prediction tools in criminal law consistent with the right to be treated as individual?—and have arrived at a *highly* qualified ‘maybe’. On the interpretation I have offered, this right does not preclude the

---

<sup>34</sup> My thanks to Simon Goldstein and Stephen Finlay for this suggestion; it is developed in more detail in J. Kleinberg, J. Ludwig, et al. (2018), at p. 136.

<sup>35</sup> For more thorough articulation and three detailed case studies of dirty data being used to train the models for predictive policing software, see Richardson, Schultz and Crawford (2019).

use of statistical methods in principle, but does significantly constrain their design and application. Law enforcement is fundamentally different in its orientation than some other applications of predictive algorithms: the law does not—*must* not—aim to detect ‘social cancers’ before they manifest. It rather must function to announce expectations for behavior, using the coercive apparatus only to hold agents accountable to those very expectations. When legal decisions are made in ways that do not afford subjects a fair opportunity to avoid hostile encounters with law enforcement, or that impose costs disproportionately, this constitutes an unfair distribution of the burdens and benefits of the rule of law. The impulse toward secrecy must be resisted; where predictions are made, they must be based only on factors that are within agents’ deliberate control, and not core to valuable exercise of autonomy.

Requiring a fair distribution strongly constrains which variables can be used as a basis for applying extra scrutiny or criminal sanction. It rules out reliance on a great many static factors (age, gender, race), as well as a number of indexes of disadvantage (zip code or neighborhood, income level, previous police contact, number of acquaintances with police contacts or arrest records, education level). The former because they are unavoidable; the latter because using them to justify the imposition of yet more costs on the victims of upstream distributive injustice—this time in the form of increased risk of suffering unjustified state coercion—is patently unfair. But while it is clearly unjust to base the distribution of *burdens* on unavoidable factors, you might think the same cannot be said of distributing *benefits*. Suppose that rather than use high risk scores to apply sanctions, we were to simply use low risk scores to exonerate, shorten sentences, or waive bail? There is cause for concern here too. A policy of this kind channels goods towards those who *lack* the markers of disadvantage that yield a high risk score, and so can still be expected to entrench racial and economic inequalities and compound disadvantage. It may be an improvement even so, but there is a dark side to reforms which succeed in alleviating injustice for many and concentrate the remaining costs on people who are comparatively vulnerable or politically powerless. Once only marginalized groups face the worst costs, it is much more difficult to muster the political will to enact the reforms necessary to correct the remaining injustice. A partial fix may well be worse than doing nothing, then, because it allows the majority to simply look away.

Where this leaves us depends on the application. For street crime—particularly property offenses like autotheft, burglary, or mugging—the social value of predicting any particular future offense is low, especially as compared to the cost of a false positive prediction to each individual who is misclassified. This is because any given offense in this category is either quite difficult to predict with accuracy greater than chance (e.g. homicides), or imposes only relatively minor compensable harm (property damage or loss) on a small number of victims. Insofar as these sorts of crimes are also driven by inelastic social causes, a predictive proxy is unlikely to have strong deterrent effects, and *is* likely to track socioeconomic disadvantage. Even if it is possible *in principle* to design risk-assessment or crime-prediction algorithms independent of these variables, it is at best unclear what evidential or predictive value a truly unbiased tool would have. Of the

extant tools, those that conditionalize on static variables *alone* presently outperform those that *also* incorporate dynamic variables.<sup>36</sup> We can expect that both would outperform prediction based only on the subset of dynamic variables that are not ruled out by the considerations just raised. So, while it may be possible to constrain the data used so that an algorithmic risk projection is consistent with the moral interests protected by the right to individualized treatment, it is unclear whether such predictors will have evidential value sufficient to justify their use.

However, white-collar crime, wage theft, and financial fraud more generally may well be appropriate arenas for the use of predictive tools. These tend to have a higher victims-per-offense ratio, and consequently there is higher social value to predicting or identifying any single instance. They are also most commonly perpetrated by a relatively well-resourced portion of the population, for whom additional scrutiny presents little more than a hassle. The subpopulations subjected to extra scrutiny, higher risk of false conviction, or longer sentences due to reliance on algorithms in financial crimes are also less likely to overlap with either a stable subgroup (like racial or SES category) or with populations already subject to intersectional disadvantage and distributive injustice, so the presumptive reasons against using a secret proxy are far less weighty for this application. So of the possible applications for algorithmic predictions in criminal law, white collar crime enforcement looks most promising.

In closing, let me revisit the optimistic aim of using algorithmic tools to improve the high-stakes decisions of criminal law. It is laudable to try to make determinations less biased, and to reduce the number of people subjected to unjustified or disproportionate costs in the course of law enforcement. Maybe we could make some progress toward this aim by supplementing the judgment of police officers, judges, juries, and parole boards with algorithmic assessments across the board. But this says more about how badly distorted our unassisted decisions are than about the accuracy or fairness of the algorithmic tools. Whether it is wise to embrace these tools as an incremental improvement depends on several factors we haven't had space to work through in this paper, including what the alternative is, and how decisionmakers would be instructed to incorporate the risk predictions into their deliberations. Without going into detail now, it's worth emphasizing that the most natural instruction to give—that a high risk score may be sufficient for an adverse judgment, but isn't necessary—will yield the worst of both worlds. If adverse judgments are still permitted in the absence of a high risk score, then the algorithmic tool does not constrain any extant bias the decisionmakers may have toward (e.g.) giving disproportionately long sentences to defendants of color. But if a high score *is* sufficient, then any bias in the false-positive error rates of the algorithm simply

---

<sup>36</sup> Herrschaft (2014); Dressel and Farid (2018). But see Degiorgio and DiDonato (2013) for findings that adding dynamic factors to static demographic models in fact improves the fit of a model predicting probation revocation specifically for substance abuse.

combines with the extant distortions---and worse, the whole decision process has a veneer of being even-handed and objective.

## 6 Acknowledgements

I am especially grateful to Jeff Behrends, Stephen Finlay, Simon Goldstein, Ian Hamilton, Gabriel Karger, Seth Lazar, Chad Lee-Stronach, Tali Mendelberg, Tom Parr, Chelsea Rosenthal, Kate Vredenburg, and Annette Zimmerman, for written comments on earlier drafts of this paper. I am also indebted to Rima Basu, Luc Bovens, Katie Creel, Marcello Di'Bello, Tom Dougherty, Anna Edmonds, Maegan Fairchild, Sarah Stroud, Robin Guong, Bernard Harcourt, Deborah Hellman, Sarah Hirschfield, Gabrielle Johnson, Geoff Sayre-McCord, Victor Tadros, Patrick Tomlin, Alex Worsnip, and Brian Zaharatos, for helpful comments and discussion of the ideas presented here. Thanks also to other participants at the 2019 workshop on the *Democratic Implications of Algorithmic Decision-Making* at Princeton University, the 2020 *Radcliffe Institute Workshop on The Ethics of Technology at Work and in Public Institutions* at Harvard University, the Dianoia Institute of Philosophy *Ethics Working Group* at Australian Catholic University, the 2020 *Philosophy, Artificial Intelligence, and Society* workshop, the 2021 works-in-progress group at University of North Carolina Chapel Hill, as well as to seminar participants at the *Center for Ethics, Law, and Public Affairs* at the University of Warwick, Rima Basu's 2020 seminar on *Belief, Evidence, and Agency* at Claremont McKenna, the Princeton University 2020 Woodrow Wilson Fellows, and Gideon Rosen's 2021 graduate seminar on *Non-Ideal Criminal Law* at Princeton University for helpful discussion of the ideas presented here.

## 7 Works Cited

- Amour, Jody. 1994. "Race Ipsa Loquitur: Of Reasonable Racists, Intelligent Bayesians, and Involuntary Negrophobes." *Stanford Law Review* 46 (4): 781-816.
- Angwin, J., J. Larson, S. Mattu, and L. Kirchner. 2016. "Machine Bias: There's software used across the country to predict future criminals. And it's biased against Blacks." *ProPublica* [www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing](http://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing).
- Bambauer, Jane. 2015. "Hassle." *Michigan Law Review* 113 (4): 461-511.
- Basu, Rima. 2019. "The Wrongs of Racist Belief." *Philosophical Studies* (Routledge) 176: 2497-2515.
- Basu, Rima. 2019. "What We Epistemically Owe to Each Other." *Philosophical Studies* 176: 915-931.
- Binns, Reuben. forthcoming. "Human Judgement in Algorithmic Loops." *Regulation & Governance* ((on file with author)).
- Bolinger, Renee Jorgensen. 2021. "Explaining Justificatory Asymmetries between Statistical and Individualized Evidence." In *The Social Epistemology of Legal Trials*, by Zachary Hoskins and Jon Robson. Routledge.
- Bolinger, Renee Jorgensen. 2020. "The Rational Impermissibility of Accepting (some) Racial Generalizations." *Synthese* 197: 2415-2431.

- Brayne, Sarah. 2021. *Predict and Surveil: Data, Discretion, and the Future of Policing*. New York: Oxford University Press.
- Buchak, Lara. 2014. "Belief, Credence, and Norms." *Philosophical Studies* 169: 285–311.
- Burgess, Earnest. 1936-7. "Protecting the Public by Parole and by Parole Prediction." *American Institute of Criminal Law & Criminology* 27: 491.
- Cofone, and Strandberg. ms. "Strategic Games and Algorithmic Transparency." *Conference Draft, on file with author*.
- Colyvan, M, H Regan, and S Person. 2001. "Is it a Crime to Belong to a Reference Class?" *The Journal fo Political Philosophy*.
- Degiorgio, and DiDonato. 2013. "Predicting probationer rates of recarceration using dynamic factors from the Substance Abuse Questionnaire-Adult Probation III (SAQ-Adult Probation III)." *American Journal of Criminal Justice* 39: 94-108.
- Dieterich, W, T Mendoza, and T Brennan. 2016. *COMPAS risk scales: Demonstrating Accuracy, Equity, and Predictive Parity*. Technical Report, Northpointe Inc.,.
- Dressel, and Farid. 2018. "The Accuracy, Fairness, and Limits of Predicting Recidivism." *Science Advances* 4 (eaao5580): 3.
- Duff, R.A. 1998. "Dangerousness and Citizenship." In *Fundamentals of Sentencing Theory*, by A. J. Ashworth and M. Wasik, 155-6. Oxford: Oxford University Press.
- Dworkin, Ronald. 1977. *Taking Rights Seriously*.
- Enoch, David, L Spectre, and T Fisher. 2012. "Statistical evidence, Sensitivity, and the Legal Value of Knowledge." *Philosophy and Public Affairs* 40 (3): 197-224.
- Eubanks, Virginia. 2018. *Automating Inequality: How High-Tech Tools Profile, Police, and Punish the Poor*. New York: St. Martin's Press.
- Ewing, Benjamin. 2018. "Recent Work on Punishment and Criminogenic Disadvantage." *Law and Philosophy* 37: 29-68.
- Fuller, Lon. 1969. *The Morality of Law*. New Haven: Yale University Press.
- Harcourt, Bernard. 2007. *Against Prediction: Profiling, Policing, and Punishment in an Actuarial Age*. Chicago, IL: University of Chicago Press.
- Harris, D. 1998. "Particularized Suspicion, Categorical Judgments: Supreme Court Rhetoric Versus Lower Court Reality Under Terry v. Ohio." *St. John's Law Review* 72 (3): 975.
- Herrschaft. 2014. *Evaluating the Reliability and Validity of the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) Tool: Implications for Community Corrections Policy*. Dissertation, Rutgers University.
- Holder, Attorney General Eric. 2014. "Remarks to the National Association of Criminal Defense Lawyers 57th Annual Meeting, and 13th State Criminal Justice Network Conference." Philadelphia, PA, August 1.
- Howard, Jeffrey. 2016. "Moral Subversion and Structural Entrapment." *Journal of Political Philosophy* 24 (1): 24-46.

- Johnson, Gabrielle. 2020. "Algorithmic bias: on the implicit biases of social technology." *Synthese* s11229-020-02696-y.
- Joyce, James. 2005. "How Evidence Reflects Probabilities." *Philosophical Perspectives* 19 (1): 153-178.
- Kehl, Guo, and Kessler. 2017. *Algorithms in the Criminal Justice System: Assessing the Use of Risk Assessments in Sentencing*. Responsive Communities Initiative, Berkman Klein Center for Internet & Society, Harvard Law School.
- Kim, Hohan. 2020. "Entrapment, Culpability, and Legitimacy." *Law and Philosophy* 39: 67-91.
- Kleinberg, J., J. Ludwig, S. Mullainathan, and C. Sunstein. 2018. "Discrimination in the Age of Algorithms." *Journal of Legal Analysis* 10: 113-174.
- Kleinberg, J., S. Ludwig, C. Mullainathan, and C. Sunstein. 2020. "Algorithms as Discrimination Detectors." *Proceedings of the National Academy of Sciences*.
- Lippert-Rasmussen, K. 2011. "'We are all Different': Statistical Discrimination and the Right to be Treated as an Individual." *Journal of Ethics* 15: 47-59.
- Marušić, Berislav, and Stephen White. 2018. "How Can Beliefs Wrong? -- A Strawsonian Epistemology." *Philosophical Topics* (46): 97-114.
- Moss, Sarah. 2018. "Moral Encroachment." *Proceedings of the Aristotelian Society* 118 (2): 177-205.
- . 2018. *Probabilistic Knowledge*. Oxford: Oxford University Press.
- Murphy, Colleen. 2005. "Lon Fuller and the Moral Value of the Rule of Law." *Law and Philosophy* 24 (3): 239-262.
- Narayanan, A. 2019 MS. "How to Recognize AI Snake Oil." *Working paper, on file with author*.
- O'Neil, Cathy. 2016. *Weapons of Math Destruction*. New York: Crown Publishers.
- Rabuy, Bernadette, and Daniel Kopf. 2016. *Detaining the Poor: How money bail perpetuates an endless cycle of poverty and jail time*. Prison Policy Initiative.
- Richardson, Schultz, and Crawford. 2019. "Dirty Data, Bad Predictions: How Civil Rights Violations Impact Police Data, Predictive Policing Systems, and Justice." *NYU Law Review* 94: 192.
- Risinger, D. M. 2004. "Unsafe verdicts: The need for reformed standards for the trial and review of factual innocence claims." *Houston Law Review* 21: 1281.
- Schroeder, Mark. 2018. "When Beliefs Wrong." *Philosophical Topics* 46: 115-127.
- Shelby, Tommie. 2007. "Justice, Deviance, and the Dark Ghetto." *Philosophy & Public Affairs* 35 (2): 126-160.
- Taslitz, Andrew. 2010. "Police are People Too: Cognitive Obstacles to, and Opportunities for, Police Getting the Individualized Suspicion Judgment Right." *Ohio State Journal of Criminal Law* 8 (1): 7-78.
- Thomson, Judith Jarvis. 1986. "Liability and Individualized Evidence." *Law and Contemporary Problems* 49 (3): 199-219.

- Tillers, Peter. 2005. "If wishes were horses: Discursive comments on attempts to prevent individuals from being unfairly burdened by their reference classes." *Law, Probability, and Risk* 4: 33-49.
- Tribe, L. 1971. "Trial by Mathematics: Precision and Ritual in the Legal Process." *Harvard Law Review* 84 (6): 1329-1393.
- Underwood, Barbara. 1979. "Law and the Crystal Ball: Predicting Behavior with Statistical Inference and Individualized Judgment." *Yale Law Journal* 88: 1408.
- Voigt, Camp, Prabhakaran, Hamilton, Hetey, Griffiths, Jurgens, Jurafsky, and Eberhardt. 2017. "Racial Disparities in Police Language." *Proceedings of the National Academy of Sciences* 114 (25): 6521-6526.
- Vredenberg, Kate. working paper. "The Right to an Explanation." *on file with author*.
- Walén, Alec. 2011. "A Punitive Precondition for Preventative Detention: Lost Status as a Foundation for a Lost Immunity." *San Diego Law Review* 48: 1229.
- Walén, Alec. 2011. "A Unified Theory of Detention, with Application to Preventative Detention for Suspected Terrorists." *Maryland Law Review* 70: 871.
- Wasserman, D. 1992. "The Morality of Statistical Proof and the Risk of Mistaken Liability." *Cardozo Law Review* 13: 935.
- Watson, Gary. 2015. "A Moral Predicament in the Criminal Law." *Inquiry* 58 (2): 168-188.
- Yang, M, S.C. Wong, and J Coid. 2010. "The Efficacy of Violence prediction: a meta-analytic comparison of nine risk assessment tools." *Psychological Bulletin* 136: 740-767.