

Escaping the “Impossibility of Fairness”: From Formal to Substantive Algorithmic Fairness

Ben Green

bzgreen@umich.edu

Michigan Society of Fellows

Gerald R. Ford School of Public Policy

Abstract

In the face of compounding crises of social and economic inequality, many have turned to algorithmic decision-making to achieve greater fairness in society. As these efforts intensify, reasoning within the burgeoning field of “algorithmic fairness” increasingly shapes how fairness manifests in practice. This paper interrogates whether algorithmic fairness provides the appropriate conceptual and practical tools for enhancing social equality. I argue that the dominant, “formal” approach to algorithmic fairness is ill-equipped as a framework for pursuing equality, as its narrow frame of analysis generates restrictive approaches to reform. In light of these shortcomings, I propose an alternative: a “substantive” approach to algorithmic fairness that centers opposition to social hierarchies and provides a more expansive analysis of how to address inequality. This substantive approach enables more fruitful theorizing about the role of algorithms in combatting oppression. The distinction between formal and substantive algorithmic fairness is exemplified by each approach’s responses to the “impossibility of fairness” (an incompatibility between mathematical definitions of algorithmic fairness). While the formal approach requires us to accept the “impossibility of fairness” as a harsh limit on efforts to enhance equality, the substantive approach allows us to escape the “impossibility of fairness” by suggesting reforms that are not subject to this false dilemma and that are better equipped to ameliorate conditions of social oppression.

1 Introduction

The United States is home to overlapping crises of social and economic inequality. Black Americans possess far fewer economic and political resources than their white counterparts and suffer disproportionately from the injustices of policing and mass incarceration. Women are underrepresented in prominent jobs and leadership roles across society and face a significant pay gap relative to men. And as economic inequality rises, fewer people have access to the material resources necessary for a comfortable and dignified life. These and other forms of inequality intersect in compounding ways.

As debates about how to promote a more egalitarian society have become increasingly salient, one approach to achieving greater fairness that has gained traction is to make (or aid) socially consequential decisions using algorithms—in particular, machine learning systems that infer patterns from historical data to make predictions about the future. To scholars across a range of

fields, algorithms represent a novel approach to overcoming the cognitive limits and social biases of human decision-makers (Kleinberg et al., 2019; Miller, 2018; Sunstein, 2019). Proponents describe how algorithms could “disparate[ly] benefit” historically disadvantaged groups, which are typically judged unfavorably by stereotype-prone human decision-makers (Kleinberg et al., 2019: 160). Policymakers and advocates praise algorithms as being able to replace biased human decisions with “objective” data-driven ones (Arnold Ventures, 2019: 1; Harris and Paul, 2017: 10). This general optimism toward algorithms by prominent and influential proponents has led to algorithms being used to improve the fairness of decision-making in contexts such as criminal sentencing, hiring, and social services.

Yet alongside the increased use of algorithms to combat discriminatory decision-making is accumulating evidence that algorithms themselves can discriminate. A growing list of examples demonstrate algorithms making more favorable decisions for men than women and white individuals than Black individuals, for instance (Dastin, 2018; Obermeyer et al., 2019). Although scholars and policymakers have long considered how algorithms, software, and quantitative assessments can discriminate (Friedman and Nissenbaum, 1996; Hutchinson and Mitchell, 2019; Ochigame, 2020), concerns about algorithmic bias have gained particular traction as data-driven algorithms became more pervasive in government, business, and daily life (Barocas and Selbst, 2016; Noble, 2018; Sweeney, 2013). These concerns have prompted the emerging field of “algorithmic fairness.” A central component of algorithmic fairness is developing and comparing mathematical definitions of fairness (Barocas et al., 2019). Other lines of research include developing methods to optimize for various articulations of algorithmic fairness (Feldman et al., 2015; Hardt et al., 2016) and auditing existing algorithms to test for biases (Angwin et al., 2016; Obermeyer et al., 2019; Raji and Buolamwini, 2019).

As algorithms become common tools for achieving greater equality across social and legal domains, reasoning about algorithmic fairness increasingly informs broader reasoning about what fairness entails and how to achieve it. Because “[t]he way in which [a] problem is conceived decides what specific suggestions are entertained and which are dismissed” (Dewey, 1938: 108), the turn to algorithms for social interventions makes algorithmic methods influential in shaping how fairness is pursued. It is therefore essential to consider not just whether specific algorithms are “fair” or “biased,” but whether algorithmic fairness provides the appropriate conceptual and practical tools for enhancing social equality. Without this capacity, algorithmic fairness will be unable to provide a productive guide for efforts to achieve a more equal society.

This paper argues that the current approach to algorithmic fairness is ill-equipped as a framework for pursuing equality, and in its place proposes an alternative approach that can more rigorously inform the use of algorithms to help combat social hierarchies. I argue that the dominant approach to algorithmic fairness is fundamentally *formal*: it relies on a restrictive frame of analysis that considers only the explicit functioning of an algorithm when rendering decisions. This formal

algorithmic fairness is closely aligned with formal theories of equality that emphasize equal treatment for individuals based on their attributes or behavior at a particular decision point. Because of its narrow scope of analysis, formal algorithmic fairness overlooks a great deal, such as the social contexts represented by data, the impacts of policies enhanced by algorithms, and the full suite of possible social reforms. This leads to an impoverished theory of social change that involves applying algorithms to decisions where concerns about inequality manifest, treating the rest of the social context as static. Although this approach can sometimes yield benefits, it also leads to the application of algorithms in ways that fail to reduce—and may even exacerbate—oppression, despite appearing to be fair.

In place of the formal approach to algorithmic fairness, I propose a *substantive* approach to algorithmic fairness. This approach draws on theories of substantive equality from law and philosophy that expand the frame of analysis beyond specific decisions and emphasize the elimination of social hierarchies. Because of its broader scope of analysis, substantive algorithmic fairness provides a more productive guide regarding what roles algorithms can (and cannot) play in enhancing equality. The substantive approach suggests a tripartite strategy for incorporating algorithms into social reforms. First, diagnose the conditions of social hierarchy in the given social context and how the structure of decisions may exacerbate those conditions. Second, develop a theory of change regarding which reforms can address the forms and mechanisms of inequality that were identified. Third, consider whether and how algorithms can enhance the desired reforms.

Although the substantive approach provides arguments against certain forms of algorithmic interventions (such as pretrial risk assessments), its deeper diagnosis of inequality and how to remedy it also suggests new directions for algorithmic work. A central insight of this approach is that algorithmic injustice often results from accurate predictions that reproduce existing inequities rather than inaccurate predictions that misjudge already-disadvantaged individuals. Thus, rather than simply applying algorithms to make or inform a particular decision where inequality manifests, as formal algorithmic fairness suggests, we must strive to alleviate social hierarchies and lower the stakes of decisions that act on those hierarchies, while considering how algorithms might help achieve these goals.

I explore the contrast between the formal and substantive approaches to algorithmic fairness by analyzing a central insight in the field of algorithmic fairness: the “impossibility of fairness.” This mathematical result proves that, aside from rare edge cases, it is impossible for algorithms to simultaneously satisfy every mathematical definition of fair decision-making (Chouldechova, 2017; Kleinberg et al., 2016). In particular, if algorithms satisfy notions of fairness that emphasize treating people similarly based on their *probability* for a given behavior, they will violate notions of fairness that emphasize treating people similarly based on their *eventual* behavior (and vice versa). The “impossibility of fairness” appears to leave efforts to promote algorithmic fairness in

an impossible bind: no matter how algorithms make decisions, they are guaranteed to be unfair along some criterion.

The “impossibility of fairness” exemplifies the limits of the formal approach to algorithmic fairness and the benefits of the substantive approach. Within the formal approach’s narrow frame of analysis, the mathematical incompatibility between different definitions of fairness implies that fairness is “impossible” to achieve. This suggests that any efforts to achieve fairness are subject to a zero-sum tradeoff between distinct notions of fairness; the best that we can do is to navigate this tradeoff. In contrast, the substantive approach reveals the “impossibility of fairness” to be a misnomer. What is strictly “impossible” is the simultaneous achievement of two different mathematical notions of fair decision-making. Yet egalitarian goals of equality cannot be fully realized by considering particular decision points in isolation. We must instead consider the broader context of social hierarchies and how decisions exacerbate those hierarchies. By pursuing this expanded analysis, substantive algorithmic fairness suggests strategies for reform that are not subject to the “impossibility of fairness” and that are better equipped to ameliorate social oppression.

This article builds on prior critical scholarship on algorithmic fairness, which has emphasized two primary concerns. The first concern involves the myriad of issues associated with formalizing the broad and contested concept of “fairness” in terms of mathematical metrics. Work on algorithmic fairness emphasizes mathematical definitions of fairness, often treating these articulations of fairness as equivalent to the concept itself (Green and Hu, 2018; Selbst et al., 2019). This approach overlooks not just the contextual and philosophical underpinnings of fairness (Binns, 2018; Selbst et al., 2019), but also the limits of these formalizations to capture the realization of fairness in practice (Green and Chen, 2019). The second major critique focuses on the normative framings embodied in algorithmic fairness. Approaches to algorithmic fairness rely on narrow conceptions of discrimination and disadvantage, thus “mirroring some of [liberal] antidiscrimination discourse’s most problematic tendencies” (Hoffmann, 2019: 901). Efforts to achieve algorithmic fairness in practice can therefore subvert more substantial reforms by depoliticizing and legitimating the unjust systems in which algorithms are embedded (Green, 2020; Ochigame, 2020; Powles and Nissenbaum, 2018).

My analysis extends these critiques, articulating how these issues result from the formal approach to algorithmic fairness, how this formal approach limits algorithmic fairness’ ability to enhance equality, and how a substantive approach to algorithmic fairness can enable algorithms to play a more productive role in combatting social hierarchies. To the extent that work on algorithmic fairness is motivated by egalitarian goals of equality, the field must expand its frame of analysis to consider more than the outcomes of isolated decision-making processes. Efforts to instantiate fairness at a single decision point, while potentially beneficial, are fundamentally limited. Combatting social hierarchies requires a substantive approach that places decisions in the context

of social relations and institutional structures and that develops reforms (and algorithms) in light of this analysis.

2 A Seemingly Intractable Dilemma: The “Impossibility of Fairness”

2.1 The “Impossibility of Fairness”

In May 2016, investigative journalists at ProPublica exposed racial bias in an algorithm used to judge pretrial defendants in Broward County, Florida (Angwin et al., 2016). This “risk assessment” algorithm, known as COMPAS (short for Correctional Offender Management Profiling for Alternative Sanctions), was created by the company Northpointe (since rebranded as Equivant). Such risk assessments predict the likelihood that pretrial defendants will be arrested in the near future or will fail to appear for trial. The predicted probabilities are classified into scores (e.g., 1-10), categories (e.g., “high risk”), or recommendations (e.g., “detain”) that guide judicial decisions about which defendants to detain in jail until their trial. Risk assessments like COMPAS have proliferated across the United States over the last decade as tools for enhancing objectivity and reducing discrimination in the criminal justice system (Green, 2020).

ProPublica found that COMPAS was “biased against blacks” (Angwin et al., 2016). Among defendants who were not arrested in the two years after being evaluated, Black defendants were 1.9 times more likely than white defendants to be misclassified as “high risk” (i.e., subjected to false positive predictions). Conversely, among defendants who were arrested in the two years after being scored, white defendants were 1.7 times more likely than Black defendants to be misclassified as “low risk” (i.e., subjected to false negative predictions). Overall, 58.8% of Black defendants were labeled “high risk” while only 34.8% of white defendants were labeled “high risk.”

Although tech critics responded to ProPublica’s article with outrage about racist algorithms (Doctorow, 2016; O’Neil, 2016), others challenged ProPublica’s methods and conclusions. The debate centered on competing statistical measures for determining whether an algorithm is biased. Northpointe defended COMPAS, arguing that ProPublica failed to account for the different rates at which Black and white defendants were actually arrested (Dieterich et al., 2016). Northpointe claimed that COMPAS was unbiased because the probability of being arrested, given any particular risk score, was similar across race. Numerous researchers unaffiliated with Northpointe critiqued ProPublica’s claim that COMPAS is racially biased on similar grounds (Corbett-Davies et al., 2017; Flores et al., 2016; Gong, 2016).

This debate over whether COMPAS is biased illustrates a conflict between different mathematical definitions of fairness (Barocas et al., 2019). These definitions emphasize distinct normative claims regarding what fairness entails and, in turn, how to mathematically evaluate the fairness of algorithms. Two notions of algorithmic fairness are particularly salient. The first is “separation,” which is satisfied if all groups subject to an algorithm’s decisions experience the same false

positive rate and the same false negative rate. Separation expresses the idea that people who behave the same way should be evaluated similarly. For instance, under the separation definition of fairness, it would be unfair if Black non-recidivists are more likely to be labeled “high risk” than white non-recidivists. This is what ProPublica argued regarding COMPAS.

The second notion of algorithmic fairness is “sufficiency,” which is satisfied if, among those who receive a particular prediction or classification, all groups exhibit the outcome being predicted at the same rate. Sufficiency measures express the idea that individuals who have a similar likelihood to exhibit a particular behavior of interest should be treated similarly. For instance, under the sufficiency definition of fairness, it would be unfair if low-risk Black defendants are labeled “high risk” while low-risk white defendants are labeled “low risk.” Northpointe and others argued that the absence of this type of racial disparity meant that COMPAS is fair.

The central distinction between separation and sufficiency fairness measures is whether they condition on the *outcomes exhibited* by individuals or on the *predictions made* about individuals. Separation conditions on outcomes and evaluates fairness based on whether those who exhibit the same outcome receive the same prediction. Sufficiency conditions on predictions and evaluates fairness based on whether those who receive the same prediction actually exhibit the same outcome.

This debate about whether COMPAS is fair raised a fundamental question about algorithmic fairness: is it possible for predictions to simultaneously satisfy both separation and sufficiency measures of fairness? The “impossibility of fairness,” which was discovered in the wake of the COMPAS debate, proves that the answer is no (Barocas et al., 2019; Chouldechova, 2017; Kleinberg et al., 2016).¹ A risk assessment that satisfies sufficiency will necessarily violate separation, and vice versa. The only exceptions to this rule involve two unlikely scenarios: when the algorithm is able to perfectly predict every outcome or when the groups in question exhibit the outcome being predicted at the same “base rate” (Kleinberg et al., 2016: 5). This result reveals a serious dilemma for algorithmic fairness: how can algorithms enable fair decision-making if we are inevitably caught between two compelling, yet conflicting, notions of fairness? As one article summarizing the tension between separation and sufficiency notes, “the tradeoff between [these] two different kinds of fairness has real bite” and means that “total fairness cannot be achieved” (Berk et al., 2018: 32, 15).

2.2 Responses to the “Impossibility of Fairness”

Two common paths for responding to the “impossibility of fairness” have emerged. The first response is what I call the “formal equality response.” This response defends sufficiency as the proper instantiation of fair decision-making, justifying a lack of separation as the inevitable

¹ This tension between separation and sufficiency arises with predictions made by both humans and algorithms (Mayson, 2019).

byproduct of groups exhibiting the outcome in question at different base rates. This argument adheres to the logic of formal equality, emphasizing the need to treat people the same based on their (estimated) likelihood to exhibit the outcome in question, and aligns with much of U.S. anti-discrimination law. Most critiques of ProPublica followed the formal equality response (Corbett-Davies et al., 2017; Dieterich et al., 2016; Flores et al., 2016; Gong, 2016). Most notably, Northpointe emphasized that the violation of separation presented by ProPublica “does *not* show evidence of bias, but rather is a natural consequence of using unbiased scoring rules for groups that happen to have different distributions of scores” (Dieterich et al., 2016: 8).

The second response, which I call the “formalism response,” posits what appears to be a more sophisticated view of the dilemma. Rather than choosing a single fairness metric, the formalism response suggests using the explicit mathematical formalization required by algorithms to rigorously consider the tradeoffs between separation and sufficiency in any given context. In this view, the formalism of algorithms provides a reality check by revealing the difficult tradeoff between these notions of fairness that might otherwise remain murky and unarticulated (Barocas et al., 2019; Berk et al., 2018). Algorithms therefore provide “clarity” to help us identify and make unavoidable tradeoffs between competing goals (Kleinberg et al., 2019: 163; Sunstein, 2019: 504). Proponents of this view argue that algorithms can “be a positive force for social justice” because they “let us precisely *quantify tradeoffs* among society’s different goals” and “force us to make more explicit judgments about underlying principles” (Kleinberg et al., 2019: 120, 151, 119).

Both the formal equality and the formalism responses are intentionally limited in their scope and ambition. Because these responses derive from the formal approach to algorithmic fairness, they focus on the functioning of an algorithm at a specific decision point and thus accept the statistical conflict between separation and sufficiency as irreconcilable. Yet we should not give up so easily.

The “impossibility of fairness” is not the first seemingly intractable dilemma that has arisen from clashes between differing notions of equality. Such conflicts are commonplace in law and philosophy. It is therefore worth considering how approaches for managing such dilemmas, developed in these fields, can inform our responses to the “impossibility of fairness” and our approaches to algorithmic fairness more generally. This is the task of the next several sections. First, I consider how philosophers and legal scholars suggest responding to seemingly intractable moral dilemmas. I then apply these lessons to consider the limits of the formal approach to algorithmic fairness and to introduce an alternative, substantive approach. I use pretrial risk assessments in the United States as a representative case study throughout the paper, as these algorithms figure most centrally in debates about the “impossibility of fairness.”

3 Formal and Substantive Approaches to Equality Dilemmas

3.1 *Escaping Intractable Dilemmas*

What can we do when faced with an inherent conflict between different notions of equality? Three egalitarian approaches to escaping dilemmas that appear unresolvable are particularly relevant for informing our responses to the “impossibility of fairness”: philosopher Elizabeth Anderson’s theory of “democratic equality” (1999), legal scholar Martha Minow’s “social-relations approach” to managing social differences (1991), and legal scholar Joseph Fishkin’s theory of “opportunity pluralism” (2014).

Elizabeth Anderson, in her theory of “democratic equality,” proposes that egalitarianism should focus on equality of relationships—fostering communities in which people express mutual respect and obligation to one another—rather than equality of distributions (Anderson, 1999). She critiques the luck egalitarian focus on equality of fortune for its narrow emphasis on distributions. Among other flaws, this approach leads to a notable “dilemma”: not providing aid means blaming individuals for their misfortune, while providing special treatment on account of one’s inferiority means expressing contempt for the disadvantaged (1999: 334). Anderson escapes this dilemma by expanding the frame of analysis. Given that the goal of egalitarianism is “to end oppression, which by definition is socially imposed” (1999: 288), it is necessary to look beyond the distribution of particular goods (both tangible and intangible) and also consider the social relations that shape material distributions, cultural norms, and the structure of opportunities.

Consider, as Anderson does, someone who is stigmatized because of their physical appearance. While some egalitarians would propose subsidizing plastic surgery as a remedy for this individual, doing so upholds oppressive beauty norms even if it provides aid for the specific person in question. Democratic equality, on the other hand, “lets us see that the injustice lies not in the natural misfortune of the ugly but in the social fact that people shun others on account of their appearance” (Anderson, 1999: 336). Anderson thus suggests a more desirable remedy: altering social norms so that no one is prevented from participating in civil society due to their appearance. While such changes may be difficult to achieve, thus necessitating a more individualized remedy (at least in the immediate term), democratic equality “enables us to see that we have a choice between redistributing material resources and changing other aspects of society to meet the demands of equality” (Anderson, 1999: 336).

Martha Minow maintains a similar perspective in her description of a “social-relations approach” for how the law should address differences between people and social groups (1991). When attempting to deal with social differences, the law typically confronts what Minow calls the “dilemma of difference” (1991: 20). On the one hand, giving similar treatment to everyone regardless of their circumstances can “freeze in place the past consequences of differences” (1991: 21). On the other hand, giving special treatment to those deemed “different” risks further entrenching and stigmatizing that difference. This tension forces the law into “either/or” thinking

that “make[s] the difference dilemma seem intractable” (1991: 73-74). Minow argues, however, that the “[t]he dilemma becomes less paralyzing if we [...] look at the issues from another point of view” and question “the social arrangements that make [particular] traits seem to matter” (1991: 375, 389). Restructuring relationships and institutions to more equitably distribute the burdens and benefits of differences enables the law to “escape or transcend the dilemmas of difference” (1991: 50). For instance, rather than denying women maternity leave on the grounds of equal treatment with men or providing women with maternity leave as a form of special treatment, the law could restructure employment such that all employees obtain collectively financed parental leave and childcare support.

Joseph Fishkin’s theory of “opportunity pluralism” (2014), concerning equality of opportunity, complements Anderson and Minow. At their most basic level, conceptions of equal opportunity adhere to a “formal equal opportunity” that centers around a “fair contest”—judge people for an opportunity based on their performance or attributes at a particular moment in time (2014: 25). This approach perpetuates inequalities in different groups’ development opportunities and life chances. Although other theories of equal opportunity (such as Rawlsian equal opportunity and luck egalitarianism) attempt to account for existing inequalities, it is impossible to truly level the playing field in the manner these approaches prescribe. As a result, debates about equality of opportunity center on “zero-sum, high-stakes competitions” where “trench warfare is a certainty, and any successes will be incremental” (2014: 130, 257).

Rather than equalizing specific competitions, Fishkin (2014) escapes this fraught terrain by considering the broader structure of opportunities. From this perspective, the problem is not merely that groups face vastly different development opportunities; the problem is that we must adjudicate between these differences when making decisions that drastically impact people’s lives. Fishkin thus proposes “renovat[ing] the structure [of opportunities] itself, in ways large and small, to open up a broader range of paths,” such that any one decision point represents less of a high-stakes “bottleneck” that cuts people off from opportunities (2014: 23). In the case of college admissions, for instance, in addition to making it easier for disadvantaged individuals to obtain college degrees, we should also reduce the stakes of college admissions by creating more paths for people to lead comfortable and fulfilling lives without a college degree.

These three theories demonstrate how an approach to equality attentive to social relations and structures provides an escape from dilemmas that appear impossible to resolve when these dimensions are ignored. Each scholar presents a dilemma between treating everyone the same regardless of circumstances and providing a direct remedy to those in less fortunate circumstances. Attempting to navigate the intractable tradeoff between these options leads to dead ends that prevent meaningful advances in equality. Yet in all three cases, expanding the analysis to consider the structure of social relationships and decisions enables more effective reforms not restricted by these dilemmas.

3.2 *Formal and Substantive Equality*

What can we learn from these three theories? Broadly, they highlight the contrast between two distinct notions of equality: formal equality and substantive equality. “Formal equality” defines equality as equal treatment or equal process at a particular moment in time, with everyone judged according to the same standard (Fishkin, 2014; MacKinnon, 2011). It stipulates that similar people should be treated similarly and is often associated with neutrality. In the United States, disparate treatment law is grounded in notions of formal equality, ensuring that people are not mistreated on the basis of protected attributes such as race or gender.

By contrast, “substantive equality” defines equality more broadly, accounting for oppression and inequality. Although the precise definition of substantive equality is contested, its core principle is to account for disadvantage and exclusion instead of treating everyone the same (Fredman, 2016). Here, I follow MacKinnon’s notion of substantive equality as opposition to social hierarchies—“social relation[s] of rank ordering, typically on a group or categorical basis,” leading to both material and dignitary inequalities (2011: 11). In other words, “hierarchy identifies the substance of substantive equality” (MacKinnon, 2016: 744). In this formulation, substantive equality is aligned with relational egalitarianism—of which Elizabeth Anderson is a prominent theorist—which asserts that “people should relate to one another as equals or should enjoy the same fundamental status” (Arneson, 2013). Disparate impact law is grounded in notions of substantive equality (albeit partially (MacKinnon, 2011; MacKinnon, 2016)), ensuring that formally neutral rules do not disproportionately burden protected groups.

Due to their differing normative bases regarding what equality entails, formal and substantive equality require distinct frames of analysis. Because formal equality emphasizes equality at particular decision points, it relies on a narrow frame of analysis that focuses on specific decisions at specific points in time, in isolation. In contrast, because substantive equality accounts for existing inequalities, it requires a broader frame of analysis that considers social hierarchies and institutional structures. When confronted with instances of inequality, rather than focus only on the processes and tradeoffs at a single decision point, “[a] substantive equality approach [...] begins by asking, what is the substance of this particular inequality, and are these facts an instance of that substance?”, emphasizing that “it is the hierarchy itself that defines the core inequality problem” (MacKinnon, 2011: 11-12). Anderson, Minow, and Fishkin all provide distinctly substantive analyses by examining equality dilemmas in light of broader social relationships and institutional conditions, demonstrating the benefits of accounting for these factors. For without the ability to represent “normatively relevant political facts,” political theories will fail to identify the causes of injustice, leading to misplaced attempts at reform (Anderson, 2009: 132).

Formal and substantive notions of equality come into conflict through their different emphases and scopes of analysis. For instance, Crenshaw notes that anti-discrimination law contains a “definitional tension” between “conflicting and contradictory” visions: a “restrictive view” which

“treats equality as a process” (i.e., formal equality) and an “expansive view [which] stresses equality as a result,” aiming to eradicate racial subordination (i.e., substantive equality) (1988: 1336, 1348, 1342, 1341). While an approach of formal equality may be desirable in an already equitable society, it “would make no sense at all in a society in which identifiable groups had actually been treated differently historically and in which the effects of this difference in treatment continued into the present” (1988: 1345). In the context of academic admissions, for instance, evaluating all students according to the same standard perpetuates patterns of racial injustice because white students receive much greater educational opportunities and resources than non-white students (EdBuild, 2019).

These notions of equality have direct implications for algorithmic fairness. The dominant approach to fairness is fundamentally formal, particularly in the sense that the scope of analysis is limited to the functioning of algorithms at particular decision points. This restrictive frame leads to heavy reliance on mathematical definitions of fairness and narrow notions of how to remedy inequality, as critics have emphasized. Even when algorithmic fairness research does draw on substantive theories of equality (such as luck egalitarianism), it typically incorporates these notions into an isolated decision-making process without considering broader social relationships and structures (Binns, 2018; Heidari et al., 2019).

A substantive approach to algorithmic fairness remedies these flaws with a more expansive analysis. When faced with disparities in data, a substantive approach would ask: do these disparities reflect social conditions of hierarchy? Similarly, when faced with particular decision points, a substantive approach would ask: do these decisions (and the interventions that they facilitate) exacerbate or ameliorate social hierarchies? This substantive analysis allows for a deeper diagnosis of inequality in any given scenario, which in turn provides new terrain for responding to the “impossibility of fairness.” Following the paths of Anderson, Minow, and Fishkin, reforming social conditions and the structure of decisions provides new avenues for social change that are not subject to the “impossibility of fairness” and that more robustly challenge interlocking inequalities.

4 The Formal Approach: Navigating the “Impossibility of Fairness”

In this section, I return to the “impossibility of fairness,” interrogating the two responses to this dilemma that arise from the formal approach to algorithmic fairness.

4.1 The Formal Equality Response

The “formal equality response” chooses a side within the “impossibility of fairness.” It commits to sufficiency as the appropriate definition of fairness, aligning with notions of formal equality by emphasizing that fairness entails treating people the same based on their (estimated) likelihood for future crimes. Despite its broad appeal, the formal equality response overlooks both the social

context of why different people pose different crime risks and the morality of existing policy responses to high-risk individuals.

To see these limitations more clearly, it is necessary to distinguish between two distinct (but not mutually exclusive) sources of algorithmic injustice.

1. **Measurement bias:** One form of algorithmic injustice occurs when an algorithm is trained on data that presents a distorted view of the underlying empirical reality. A common instance of measurement bias involves data that reflects the decisions of biased humans. For instance, a risk assessment would be subject to measurement bias if its training data overestimated the recidivism rates of Black defendants due to disproportionate policing in Black neighborhoods. In other words, measurement bias introduces algorithmic unfairness through *inaccuracy*: if data presents a flawed representation of the underlying social reality, the algorithm will reproduce the flawed representation.
2. **Social inequity:** A second form of injustice occurs when an algorithm is trained on data that accurately depicts population-level disparities. For instance, a risk assessment would be subject to social inequity if its training data reflects accurately that Black defendants are more likely than white defendants to recidivate. In other words, social inequity introduces algorithmic unfairness through *accuracy*: if data accurately depicts a social reality that is the product of oppression, the algorithm will reproduce those unjust conditions.

Because of its emphasis on sufficiency, the formal equality response is concerned with alleviating measurement bias but not social inequity. When applied to risk assessments, the fundamental logic of the formal equality response is that people should be treated similarly based on their crime risk (Corbett-Davies et al., 2017; Dieterich et al., 2016; Flores et al., 2016; Gong, 2016). On this logic, the best way to advance racial justice in algorithmic decision-making is to increase prediction accuracy and thereby ensure that decisions are based on accurate judgments about each individual (Hellman, 2020; Kleinberg et al., 2019). Measurement bias thus needs to be remedied because it generates inaccuracy (typically in ways that magnify existing inequality), while social inequity does not because it enables accurate predictions. As critiques of ProPublica expressed, given accurate enough predictions, a lack of separation can be attributed to the fact that different groups commit crimes at different rates rather than to algorithmic bias.

Yet the problem of discrimination is not so neatly explained away by reference to the accuracy of predictions and group differences in base rates, because the racial disparities reflected in population-level outcomes are *themselves* the product of discrimination. In the case of risk assessments, Black and white defendants do not just “happen to have different distributions of scores,” as adherents of sufficiency assert (Dieterich et al., 2016: 8). Instead, the disparity in Black and white crime rates is the product of longstanding anti-Black oppression that includes redlining and segregation (Rothstein, 2017), racial criminalization (Muhammad, 2011), and severe underfunding of schools (EdBuild, 2019)—all of which increase crime (Krivo et al., 2009; Lochner

and Moretti, 2004; Rose and Clear, 1998). Taking Black crime rates at face value, rather as the product of oppression, extends a long legacy of justifying further racial inequities through associations of Blackness with dangerousness (Muhammad, 2011).

The racial disparity in crime risk is particularly salient in light of the policy that pretrial risk assessments facilitate: preventative detention (i.e., detaining criminal defendants before trial due to concerns about crime). When the United States Supreme Court determined that this policy was constitutional in 1987, Justice Thurgood Marshall deemed it to be “incompatible with the fundamental human rights protected by our Constitution” (U.S. Supreme Court, 1987). Preventative detention has faced continued scrutiny for undermining the rights of the accused and exacerbating mass incarceration (Koepke and Robinson, 2018). It imposes severe costs on defendants including the loss of freedom, an increased likelihood of conviction, and a reduction in future employment (Dobbie et al., 2018).

This combination of social inequity and carceral policy mean that even the complete elimination of measurement bias would not prevent risk assessments from reproducing oppression. Suppose that we eliminated all measurement bias and developed a risk assessment that predicts perfectly whether or not each defendant will commit a crime. Even this perfect risk assessment would disproportionately label Black defendants as “high risk” compared to white ones. For although measurement bias was removed from the algorithm, the social inequity remains. Past and present discrimination have created social conditions in which Black people in the U.S. *are* empirically at higher risk to commit crimes, above and beyond differences in arrest and enforcement patterns (Sampson et al., 2005). This disparity is the product of numerous injustices and inequities rather than any sort of disparity in inherent criminality, as has often been assumed historically (Muhammad, 2011). A perfectly accurate risk assessment would therefore identify that proportionately more Black defendants than white defendants will recidivate, leading to higher pretrial detention rates for Black defendants and in effect punishing Black communities for having been subjected to criminogenic circumstances.

In dismissing the significance of social inequity and instead emphasizing that people be treated equally based only on accurate predictions, the formal equality response embodies what Minow calls the “equal rights” approach (1991: 12) and what Fishkin calls the “formal-plus” approach (2014: 33). These approaches acknowledge that evaluations of merit and behavior can be biased and therefore seek more accurate methods for making such judgments. Although these enhanced tests can benefit some members of groups that systematically receive unfavorably mistaken judgments, they also exhibit serious conceptual limits. First, adhering to formally equal judgments in the face of social inequities “perpetuates a caste system” and “mask[s] an unjust, and permanently unequal, social order” (Fishkin, 2014: 35). Because risk is distributed along racial hierarchies, even a perfectly accurate risk assessment will lead to the detention of more Black defendants than white defendants. Second, striving for greater accuracy “preserves the either/or

construction of the problem,” reifying existing social categories and hierarchies as “natural and inevitable” (Minow, 1991: 214, 119). Risk assessments uphold the use of risk to determine whether defendants are released or detained, despite the racialized nature of “risk” as a category—and the more accurate a risk assessment’s predictions, the more legitimate detention based on risk appears.

In sum, the formal equality response relies on a simplistic view of where oppression comes from and how to remedy it. The central source of unfairness with respect to risk assessments is not that flawed data might lead an algorithm to make erroneous predictions of someone’s crime risk, but that racial stratification makes Black defendants higher risk than white ones and that the consequences of being deemed high risk include the loss of liberty. The formal equality response of treating defendants the same based on accurate estimates of risk may help Black defendants who are relatively low risk and who might otherwise be stereotyped as high risk, but does nothing for the many Black defendants who actually are high risk and who become incarcerated as a result. In fact, more accurate predictions may undercut the normative salience of why some groups are higher risk than others and how pretrial detention harms defendants.

4.2 The Formalism Response

The “formalism response” takes a more measured view than the “formal equality response.” It emphasizes that—because of their reliance on mathematical formalism—algorithmic approaches to enhancing fairness are valuable because they provide “clarity” and precision in determining the balance between different considerations (Kleinberg et al., 2019: 163; Sunstein, 2019: 504). Although the formalism response provides a way to navigate the “impossibility of fairness” without committing to sufficiency as the only appropriate notion of fairness, it suffers from a fundamental shortcoming of its own: because the formalism response considers only the precise functioning of the algorithm when making decisions, it relies on a highly constrained theory of social change.

The supposed clarity of mathematical formalism depends on an incredibly restrictive analysis focused on isolated decision points. Although this formalism enables precise reasoning about considerations that are applicable to the decision point and rendered legible in mathematical terms, it ignores all other considerations (Green and Viljoen, 2020). Notably, arguments for the formalism response assume that implementing an algorithm is the only possible (or, at least, relevant) alternative to the status quo (Berk et al., 2018; Kleinberg et al., 2019). Furthermore, research on fairness in risk assessments explicitly places structural disadvantage and existing racial disparities outside the scope of algorithms and the responsibility of their designers (Chouldechova, 2017; Corbett-Davies et al., 2017; Kleinberg et al., 2019), such that “total fairness” refers only to the outcomes of a specific decision (Berk et al., 2018: 15).

The narrow, formal frame thus locks us into the “impossibility of fairness” as a severe constraint on efforts to achieve greater equality. When the analysis of fairness is limited to an isolated

decision point, the statistical incompatibility between separation and sufficiency—a dilemma within a specific decision point—comes to represent a fundamental limit on our ability to achieve greater fairness writ large. Because other aspects of society are treated as static, an impossibility of achieving fairness in a specific decision-making process become conflated with an impossibility of achieving fairness more broadly. As a result, the supposed clarity of algorithmic formalism actually obscures other dimensions of analysis and potential reform.

This intractable tradeoff might be acceptable if either separation or sufficiency was appropriate for redressing social subordination. Instead, formal algorithmic fairness leaves us stuck making zero-sum choices between two flawed options. For although separation may seem appealing as an alternative to sufficiency that provides special treatment to Black defendants (e.g., in the form of higher risk thresholds for being labeled “high risk”), it fails to hold up as a solution on closer inspection. First, deviating from sufficiency and creating race-specific thresholds seems to obviate the point of using algorithmic risk predictions at all, as risk scores would have different meanings based on the defendant in question (Flores et al., 2016; Mayson, 2019). Such explicit differential treatment based on race would be illegal to implement in many instances (Corbett-Davies et al., 2017; Hellman, 2020). Second, although a lack of separation demonstrates that different groups face disparate burdens from mistaken judgments (Chouldechova, 2017; Hellman, 2020), separation cannot fully capture the extent and impacts of social hierarchies. Even the perfectly accurate risk assessment described above (which satisfies separation) would identify that proportionately more Black defendants than white defendants are “high risk,” thus reproducing the racial hierarchy that leads to racially disparate crime rates.

In sum, the formalism response leaves us navigating a zero-sum tradeoff between two flawed notions of fairness. If our interventions are indeed limited to tuning how an algorithm makes decisions, the formalism response is correct in praising the clarity that algorithms provide to at least precisely balance sufficiency and separation. Yet it is only because of the narrow, formal frame of analysis that such a balancing act is seen as the vanguard of reform. If we want to move beyond the harsh limits of the “impossibility of fairness,” we need a new approach to algorithmic fairness.

5 The Substantive Approach: Escaping the “Impossibility of Fairness”

Given the notable flaws of the formal approach responses that attempt to *navigate* the “impossibility of fairness,” we should turn our attention to pursuing more substantive responses that allow us to *escape* it. The egalitarian accounts of Anderson, Minow, and Fishkin all demonstrate productive responses to seemingly intractable social dilemmas. Anderson’s notion of democratic equality “lets us see how injustices may be better remedied by changing social norms and the structure of public goods than by redistributing resources” (1999: 336). Minow similarly demonstrates how recognizing social categories as relational (rather than intrinsic to individuals) and social arrangements as political and mutable (rather than neutral and static) “introduces new

angles of vision, new possibilities for change” (1991: 78). Finally, Fishkin’s treatise on “bottlenecks” suggests that altering the structure of opportunities enables us to escape many of the “extraordinarily contentious” and “zero-sum” debates about how to distribute opportunities (2014: 256). While at first glance these proposals may seem too broad or far removed from algorithms, they provide valuable guidance to orient social reforms and, in turn, algorithmic interventions that support those reforms.

Escaping the “impossibility of fairness” requires a substantive approach to algorithmic fairness that analyzes inequality and decision-making within the broader context of social relations of hierarchy. Doing so highlights the central—but undervalued—role of existing social inequities as the driving force behind algorithmic injustice. With respect to risk assessments, social hierarchies mean that disadvantaged groups are not just perceived as being higher-risk—society has shaped them to actually be higher-risk, on average. Therefore, the fundamental concern regarding algorithmic decision-making is not *mistaken* predictions that misclassify the disadvantaged, but *accurate* predictions that reproduce and exacerbate existing disadvantage. This diagnosis suggests a shift in how we respond to the “impossibility of fairness” and approach the task of improving social equality. Although we certainly should not strive for inaccurate predictions, we must also avoid decisions that, through their accuracy, simply recreate existing conditions of inequality (Mayson, 2019). Rather than following the formal approach of only enhancing the accuracy or fairness of particular decisions, we are better served searching for ways to a) ameliorate the underlying conditions of social inequity and b) reduce the scope and harms of decisions that hinge on attributes which are unequally distributed due to social inequity. Because these reforms consider a scope beyond individual decision points, they are not subject to the intractable dilemma of the “impossibility of fairness.”

The core problem at the heart of the “impossibility of fairness” is social inequity: underlying disparities between groups that are the product of oppression. We must therefore start with what I call the “relational response”: altering the social relations that form the hierarchy represented in the data. In the case of risk assessments, this involves restructuring the relationships that define “risk” and shape its unequal distribution across the population. The relational response thus requires, first and foremost, that we reduce the crime risk of Black populations by reducing criminogenic social conditions of disadvantage. For instance, public policies that extend access to education (Lochner and Moretti, 2004), welfare (Tuttle, 2019), and affordable housing (Diamond and McQuade, 2019) all reduce crime, and therefore could reduce the racial disparity in crime risk. Also essential, although admittedly more challenging, are efforts to combat the association of Blackness with criminality and the effects of this association. As scholars have shown, societal choices about what behaviors to punish are based on racial hierarchies (Harcourt, 2015; Muhammad, 2011), such that “risk [is] a proxy for race” (Harcourt, 2015: 238). The relational response thus requires not merely challenging stereotypes that link Blackness with crime, but decriminalizing behaviors that have been criminalized in the past to subjugate minorities.

The secondary problem at the heart of the “impossibility of fairness” is that consequential (and often punitive) decisions hinge on characteristics that are disparately distributed across social groups. We must therefore complement the relational response with what I call the “structural response”: reshaping the structure of decisions to avoid or lower the stakes of decisions that act on social hierarchies. In the context of risk assessments, this means eliminating (or reducing the scope and harms of) decisions that determine one’s freedom and opportunities based on their risk of future crime. When fewer people are subjected to decisions in which liberty and well-being hinge on having low levels of crime risk, existing racial disparities in the distribution of risk become less consequential. Most directly, such an approach could entail abolishing (or drastically reducing the scope of) pretrial detention on the basis of risk, such that fewer people would stand to be incarcerated, regardless of their risk level. Reforms targeted at decoupling pretrial detention from increased conviction, and conviction from reduced future employment, could diminish the downstream harms of pretrial detention and thereby further reduce the consequences of disparate risk levels. Another reform along these lines would be to shift from responding to risk with punishment to responding with social or material support. Because these structural response reforms target the nature and impacts of decisions themselves, this approach can yield benefits even absent the relational response reforms to reduce racial disparities.

These substantive responses suggest both limits of existing approaches to algorithmic fairness and new opportunities for algorithms to enhance social equality. On the one hand, a substantive analysis reveals clear limitations to algorithms like pretrial risk assessments. Even in the best-case scenario of perfectly accurate predictions, pretrial risk assessments reinforce racial disparities in risk that are the product of oppression and legitimize the policy of detaining defendants due to their risk levels. These algorithms are ill-suited to the task of facilitating decarceral criminal justice reform (Green, 2020). On the other hand, a substantive analysis reveals new and more fruitful pathways for algorithms to play a role in addressing social inequities. Following the relational response, algorithms could be used to improve the social well-being (and thereby reduce the crime risk) of disadvantaged groups by enhancing policies that increase access to education (Lakkaraju et al., 2015), welfare (DataSF, 2018), and affordable housing (Ye et al., 2019). Following the structural response, algorithms could be used to enhance supportive responses to risk (Mayson, 2019) and to enable a more systemic view of how the criminal justice system exacerbates racial inequalities (Crespo, 2015; Goel et al., 2016). These represent just a few examples that have been proposed or developed; many more applications along these lines could become possible with further research on substantive responses.

The substantive responses to the “impossibility of fairness” demonstrate how an expansive analysis of social conditions and institutions can lead to more rigorous theories of social change, and how, in turn, those theories of change can inform work on algorithms. Although the relational and structural proposals resist easy implementation, they present a compass with which to develop and evaluate interventions, thus helping us distinguish between reforms that uphold social

hierarchies and reforms that challenge such hierarchies. And because the substantive responses target social relations and the structure of decisions (rather than isolated decision-making procedures), they enable us to escape the harsh tradeoffs imposed by the “impossibility of fairness.” While there may be cases where we cannot (or do not want to) pursue the relational and structural reforms—and therefore must revert to the formalism response—that conclusion should be arrived at following a thorough substantive analysis rather than treated as the default.

In sum, substantive algorithmic fairness suggests a new framework for algorithmic fairness that aligns with recent calls to shift the field’s emphasis from “fairness” to “justice” (Bui and Noble, 2020; Green, 2018). This substantive approach follows a tripartite strategy. The first step is to consider the substance of the social conditions in question, looking for conditions of hierarchy and questioning how existing policies and institutions reinforce those conditions. The second step is to consider what types of social reforms will appropriately address the diagnosed inequities. This step can reveal opportunities for reform beyond simply adjusting the mechanisms of specific decision-making processes; instead, we can reform social relations and institutional structures. This then raises the third step of analyzing whether and how algorithms can be used to enhance or facilitate the desired reforms. Much of what these theories of change suggest will not involve algorithms, at least not directly, but algorithms can still play a role in combatting hierarchy and oppression. Algorithmic interventions therefore need to be pursued through an “agnostic approach” that emphasizes social impacts without prioritizing any necessary role for algorithms (Green and Viljoen, 2020: 28). An important area for future work will be to develop theories and methods for incorporating algorithms into broader efforts for social change, rather than using algorithms as centerpieces that act in isolation. Doing so will suggest new roles for algorithms, as well as new methodological questions for how to develop and evaluate algorithms. Achieving these new approaches requires actively combatting the exclusion of women and minorities, as well as the influence of private companies, in efforts to build algorithms and define algorithmic fairness (Bui and Noble, 2020; West, 2020).

6 Conclusion

Algorithms have become active sites of theorizing about questions of equality that have long been within the purview of normative theory, yet today appear anew under the guise of technical discussions about computational systems. If we fail to approach these questions with the appropriate conceptual and practical tools, we will fail to fully comprehend the nature of injustices and to identify effective paths for remediating those injustices. The current, formal approach to algorithmic fairness is poorly equipped to guide strategies for enhancing social equality. When the focus of analysis is directed only to the functioning of an algorithm when making consequential decisions, we are left with the intractable dilemma of the “impossibility of fairness” and proposals for reform that exacerbate existing social hierarchies. We can escape the supposed “impossibility” of fairness by taking a substantive approach to algorithmic fairness, expanding our gaze from decision-making procedures alone to social hierarchies and to the structure of decisions that

magnify the stakes of those hierarchies. Within the criminal justice system and beyond, this substantive approach suggests both alternative reforms worth pursuing and new roles for algorithms in combatting social hierarchies—thus aligning algorithmic interventions with the demands for social justice that motivate much of the work on algorithmic fairness.

Acknowledgements

I am grateful to Salomé Viljoen, Lily Hu, Will Holub-Moorman, Zach Wehrwein, and the Michigan Society of Fellows postdoctoral scholars for their detailed and insightful comments on early drafts. I also thank Andrew Schrock for his expert editing advice that greatly improved the article.

References

- Anderson E (2009) Toward a Non-Ideal, Relational Methodology for Political Philosophy: Comments on Schwartzman's *Challenging Liberalism*. *Hypatia* 24: 130-145.
- Anderson ES (1999) What is the Point of Equality? *Ethics* 109: 287-337.
- Angwin J, Larson J, Mattu S, et al. (2016) Machine Bias. *ProPublica*.
- Arneson R (2013) Egalitarianism. *The Stanford Encyclopedia of Philosophy*.
- Arnold Ventures (2019) Statement of Principles on Pretrial Justice and Use of Pretrial Risk Assessment.
- Barocas S, Hardt M and Narayanan A (2019) *Fairness and Machine Learning*: fairmlbook.org.
- Barocas S and Selbst AD (2016) Big Data's Disparate Impact. *California Law Review* 104: 671-732.
- Berk R, Heidari H, Jabbari S, et al. (2018) Fairness in Criminal Justice Risk Assessments: The State of the Art. *Sociological Methods & Research*: 1-42.
- Binns R (2018) Fairness in Machine Learning: Lessons from Political Philosophy. In: *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*. PMLR, 149--159.
- Bui ML and Noble SU (2020) We're Missing a Moral Framework of Justice in Artificial Intelligence: On the Limits, Failings, and Ethics of Fairness. In: Dubber MD, Pasquale F and Das S (eds) *The Oxford Handbook of Ethics of AI*. Oxford University Press.
- Chouldechova A (2017) Fair Prediction with Disparate Impact: A Study of Bias in Recidivism Prediction Instruments. *Big Data* 5: 153-163.
- Corbett-Davies S, Pierson E, Feller A, et al. (2017) Algorithmic Decision Making and the Cost of Fairness. In: *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 797-806.
- Crenshaw KW (1988) Race, Reform, and Retrenchment: Transformation and Legitimation in Antidiscrimination Law. *Harvard Law Review* 101: 1331-1387.
- Crespo AM (2015) Systemic Facts: Toward Institutional Awareness in Criminal Courts. *Harvard Law Review* 129: 2049-2117.
- Dastin J (2018) Amazon scraps secret AI recruiting tool that showed bias against women. *Reuters*.
- DataSF (2018) Keeping Moms and Babies in Nutrition Program.

- Dewey J (1938) *Logic: The Theory of Inquiry*: Henry Holt and Company.
- Diamond R and McQuade T (2019) Who Wants Affordable Housing in Their Backyard? An Equilibrium Analysis of Low-Income Property Development. *Journal of Political Economy* 127: 1063-1117.
- Dieterich W, Mendoza C and Brennan T (2016) COMPAS Risk Scales: Demonstrating Accuracy Equity and Predictive Parity. *Northpoint Inc. Research Department*.
- Dobbie W, Goldin J and Yang CS (2018) The Effects of Pretrial Detention on Conviction, Future Crime, and Employment: Evidence from Randomly Assigned Judges. *American Economic Review* 108: 201-240.
- Doctorow C (2016) Algorithmic risk-assessment: hiding racism behind "empirical" black boxes. *Boing Boing*.
- EdBuild (2019) \$23 Billion.
- Feldman M, Friedler SA, Moeller J, et al. (2015) Certifying and Removing Disparate Impact. In: *21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 259-268.
- Fishkin J (2014) *Bottlenecks: A New Theory of Equal Opportunity*: Oxford University Press.
- Flores AW, Bechtel K and Lowenkamp CT (2016) 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.”. *Federal Probation* 80: 38-46.
- Fredman S (2016) Substantive equality revisited. *International Journal of Constitutional Law* 14: 712-738.
- Friedman B and Nissenbaum H (1996) Bias in Computer Systems. *ACM Transactions on Information Systems* 14: 330–347.
- Goel S, Rao JM and Shroff R (2016) Precinct or prejudice? Understanding racial disparities in New York City’s stop-and-frisk policy. *The Annals of Applied Statistics* 10: 365-394.
- Gong A (2016) Ethics for powerful algorithms (1 of 4). *Medium*.
- Green B (2018) Putting the J(ustice) in FAT. *Berkman Klein Center Collection - Medium*.
- Green B (2020) The False Promise of Risk Assessments: Epistemic Reform and the Limits of Fairness. In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. ACM, 594–606.
- Green B and Chen Y (2019) Disparate Interactions: An Algorithm-in-the-Loop Analysis of Fairness in Risk Assessments. In: *Proceedings of the Conference on Fairness, Accountability, and Transparency*. ACM, 90–99.
- Green B and Hu L (2018) The Myth in the Methodology: Towards a Recontextualization of Fairness in Machine Learning. In: *Machine Learning: The Debates workshop at the 35th International Conference on Machine Learning*.
- Green B and Viljoen S (2020) Algorithmic Realism: Expanding the Boundaries of Algorithmic Thought. In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. ACM, 19–31.

- Harcourt BE (2015) Risk as a Proxy for Race: The Dangers of Risk Assessment. *Federal Sentencing Reporter* 27: 237-243.
- Hardt M, Price E and Srebro N (2016) Equality of Opportunity in Supervised Learning. In: *30th Conference on Neural Information Processing Systems (NIPS 2016)*. 3315-3323.
- Harris K and Paul R (2017) Pretrial Integrity and Safety Act of 2017. *115th Congress*.
- Heidari H, Loi M, Gummadi KP, et al. (2019) A Moral Framework for Understanding Fair ML through Economic Models of Equality of Opportunity. In: *Proceedings of the Conference on Fairness, Accountability, and Transparency*. ACM, 181–190.
- Hellman D (2020) Measuring Algorithmic Fairness. *Virginia Law Review* 106: 811-866.
- Hoffmann AL (2019) Where fairness fails: data, algorithms, and the limits of antidiscrimination discourse. *Information, Communication & Society* 22: 900-915.
- Hutchinson B and Mitchell M (2019) 50 Years of Test (Un)fairness: Lessons for Machine Learning. In: *Proceedings of the Conference on Fairness, Accountability, and Transparency*. ACM, 49–58.
- Kleinberg J, Ludwig J, Mullainathan S, et al. (2019) Discrimination in the Age of Algorithms. *Journal of Legal Analysis* 10: 113-174.
- Kleinberg J, Mullainathan S and Raghavan M (2016) Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807*.
- Koepke JL and Robinson DG (2018) Danger Ahead: Risk Assessment and the Future of Bail Reform. *Washington Law Review* 93: 1725-1807.
- Krivo LJ, Peterson RD and Kuhl DC (2009) Segregation, Racial Structure, and Neighborhood Violent Crime. *American Journal of Sociology* 114: 1765-1802.
- Lakkaraju H, Aguiar E, Shan C, et al. (2015) A Machine Learning Framework to Identify Students at Risk of Adverse Academic Outcomes. In: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1909–1918.
- Lochner L and Moretti E (2004) The Effect of Education on Crime: Evidence from Prison Inmates, Arrests, and Self-Reports. *American Economic Review* 94: 155-189.
- MacKinnon CA (2011) Substantive Equality: A Perspective. *Minnesota Law Review* 96: 1.
- MacKinnon CA (2016) Substantive equality revisited: A reply to Sandra Fredman. *International Journal of Constitutional Law* 14: 739-746.
- Mayson SG (2019) Bias In, Bias Out. *Yale Law Journal* 128: 2218-2300.
- Miller AP (2018) Want Less-Biased Decisions? Use Algorithms. *Harvard Business Review*.
- Minow M (1991) *Making All the Difference: Inclusion, Exclusion, and American Law*: Cornell University Press.
- Muhammad KG (2011) *The Condemnation of Blackness: Race, Crime, and the Making of Modern Urban America*: Harvard University Press.
- Noble SU (2018) *Algorithms of Oppression: How Search Engines Reinforce Racism*: NYU Press.
- O’Neil C (2016) ProPublica report: recidivism risk models are racist. *mathbabe*.
- Obermeyer Z, Powers B, Vogeli C, et al. (2019) Dissecting racial bias in an algorithm used to manage the health of populations. *Science* 366: 447-453.

- Ochigame R (2020) The Long History of Algorithmic Fairness. *Phenomenal World*.
- Powles J and Nissenbaum H (2018) The Seductive Diversion of ‘Solving’ Bias in Artificial Intelligence. *OneZero*.
- Raji ID and Buolamwini J (2019) Actionable Auditing: Investigating the Impact of Publicly Naming Biased Performance Results of Commercial AI Products. In: *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, 429-435.
- Rose DR and Clear TR (1998) Incarceration, Social Capital, and Crime: Implications for Social Disorganization Theory. *Criminology* 36: 441-480.
- Rothstein R (2017) *The Color of Law: A Forgotten History of How Our Government Segregated America*: Liveright Publishing Corporation.
- Sampson RJ, Morenoff JD and Raudenbush S (2005) Social Anatomy of Racial and Ethnic Disparities in Violence. *American Journal of Public Health* 95: 224-232.
- Selbst AD, Boyd D, Friedler SA, et al. (2019) Fairness and Abstraction in Sociotechnical Systems. In: *Proceedings of the Conference on Fairness, Accountability, and Transparency*. ACM, 59-68.
- Sunstein CR (2019) Algorithms, Correcting Biases. *Social Research* 86: 499-511.
- Sweeney L (2013) Discrimination in Online Ad Delivery. *Queue* 11: 10-29.
- Tuttle C (2019) Snapping Back: Food Stamp Bans and Criminal Recidivism. *American Economic Journal: Economic Policy* 11: 301-327.
- U.S. Supreme Court (1987) United States v. Salerno. *481 U.S.* 739.
- West SM (2020) Redistribution and Recognition. *Catalyst: Feminism, Theory, Technoscience* 6.
- Ye T, Johnson R, Fu S, et al. (2019) Using machine learning to help vulnerable tenants in New York City. In: *Proceedings of the 2nd ACM SIGCAS Conference on Computing and Sustainable Societies*. ACM, 248–258.