


Algorithmic reparation

Jenny L. Davis¹ , Apryl Williams^{2,3}  and Michael W. Yang⁴

Big Data & Society
July-December: 1–12
© The Author(s) 2021
Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/20539517211044808
journals.sagepub.com/home/bds


Abstract

Machine learning algorithms pervade contemporary society. They are integral to social institutions, inform processes of governance, and animate the mundane technologies of daily life. Consistently, the outcomes of machine learning reflect, reproduce, and amplify structural inequalities. The field of fair machine learning has emerged in response, developing mathematical techniques that increase fairness based on anti-classification, classification parity, and calibration standards. In practice, these computational correctives invariably fall short, operating from an *algorithmic idealism* that does not, and cannot, address systemic, Intersectional stratifications. Taking present fair machine learning methods as our point of departure, we suggest instead the notion and practice of *algorithmic reparation*. Rooted in theories of Intersectionality, reparative algorithms name, unmask, and undo allocative and representational harms as they materialize in sociotechnical form. We propose algorithmic reparation as a foundation for building, evaluating, adjusting, and when necessary, omitting and eradicating machine learning systems.

Keywords

Machine learning, algorithmic bias, Intersectionality, artificial intelligence, critical race theory, fair machine learning

Introduction

In socially stratified societies, power concentrates but its mechanisms are diffuse. Power flows through governing bodies, social institutions, and micro-interactions, all of which entangle with technologies of the time. By default, technologies reflect and reinforce existing social orders, expressing and materializing hierarchical relations. However, technologies can also be tools of liberation. They can expose, undo, and reshape status quos. This latter project necessitates concerted and targeted efforts, underpinned by socially informed perspectives. In service of such efforts, we present *algorithmic reparation* as a concept and a scaffold for Intersectional¹ approaches to machine learning (ML) systems, displacing fairness in favor of redress. Beyond improving code, a reparative approach uses computational tools for social intervention, while critically assessing when and where computation does not belong.

Algorithmic reparation is a transdisciplinary, sociotechnical proposal that converges theories of Intersectionality with acts of reparation, together applied to ML, with the goal of recognizing and rectifying structural inequality. Both Intersectionality and reparation have legal historical foundations, and each address systemic discrimination. Both have also now expanded beyond their legal origins via intellectual and activist movements. We continue

these expansions, fusing Intersectionality and reparation into a cogent framework for critical algorithmic reform.

Algorithmic reform requires both social and technical expertise. Transdisciplinary collaboration is thus central to this proposal. Social theorists and computer scientists are equally vital for the design, production, and evaluation of equitable algorithmic systems, best achieved through tandem work. This does not mean perfunctory partnerships in which technicians work on one thing and theoreticians another, but meaningful collaboration and cross-training (and cross-training *through* collaboration)² such that reforms emerge from the pools of multiple knowledge.

Our argument proceeds as follows: first, we review the problem of algorithmic inequality in ML – what it is, why

¹School of Sociology, The Australian National University, Canberra, Australia

²Department of Communication and Media, University of Michigan, Ann Arbor, MI, USA

³University of Notre Dame Institute for Advanced Study & Technology Ethics Center, Notre Dame, IN, USA

⁴School of Computer Science, The Australian National University, Canberra, Australia

Corresponding author:

Jenny L. Davis, School of Sociology, The Australian National University, Canberra, Australia.

Email: Jennifer.davis@anu.edu.au

it persists and how technologists have attempted to address the issue. Next, we summarize key tenets of Intersectionality, link it to ML, and delineate how its pairing with reparation produces a critical orienting framework. With this foundation, we dig into the central techniques that drive the field of fair machine learning (FML), analyzing how and why these techniques are ineffective at combatting algorithmic inequality, and thus making the case for an alternative, reparative approach. Finally, we discuss methods for, and barriers to, implementing algorithmic reparation, addressing opportunities and constraints for a reparative algorithmic praxis.

Algorithmic inequality in ML

An algorithm is simply a set of rules for completing a task. In computation, these are encoded mathematical directives which traditionally have been written manually by computer programmers. ML uses a special type of algorithm developed via automated statistical inference procedures over large datasets (Barocas et al., 2017b; Kearns and Roth, 2019). ML is utilized by major institutions to guide criminal sentencing, welfare distributions, access to loans, hiring processes, and other resource allocations that shape opportunity structures for individuals and groups. ML also pervades everyday practices through search engines, dating applications, social media platforms, and entertainment streaming services. ML thus informs governance, shapes organizations, and weaves through the mundanities of daily life.

The rationale for ML is pleasantly benevolent – to make institutional decisions fairer and to make tasks more convenient. However, the implementation of these systems consistently results in data-driven outcomes that reflect and augment patterns of inequality (Amoore, 2020; Benjamin, 2019; Costanza-Chock, 2020; Crawford, 2021; Crawford et al., 2019; D'Ignazio and Klein, 2020; Noble, 2018; O'Neil, 2016). These patterns have been documented by journalists, academics, and activists over the past decade, exemplified by high-profile cases of automation gone awry, such as Google's racist image labels (Simonite, 2018; Kayser-Bril, 2020), pricing algorithms that overcharge Asian communities for college test prep services (Angwin et al., 2015), and facial recognition tools that result in wrongful arrests due to poor fidelity with dark skin combined with racist patterns of over-policing (Hill, 2020). These harms are both allocative and representational, creating material divisions and reinforcing cultural stereotypes that devalue marginalized individuals and groups (Barocas et al., 2017a).

Why does ML reproduce inequality?

The fundamental reason that ML algorithms continue to reproduce inequality is because these technical systems are intrinsically and fundamentally social (Ames, 2018;

Bucher, 2018; Kitchin, 2017; Seaver, 2017). Put simply, algorithms are animated by data, data come from people, people make up society, and society is unequal. Algorithms thus arc towards existing patterns of power and privilege, marginalization and disadvantage (Benjamin, 2016, 2019; Broussard, 2018; Browne, 2015; Costanza-Chock, 2020; Davis, 2020; D'Ignazio and Klein, 2020).

Barocas et al. (2017b) summarize the ML process as a pipeline that proceeds in four steps: capture and quantify what is (measure)→model generalizations from the training data (learn)→apply the model to novel inputs (action)→collect feedback and refine. Through the course of this pipeline, there are several specific, overlapping ways algorithmic inequalities materialize. They can be a product of unjust goals rooted in racist, sexist, heteronormative, ableist, nationalist, and/or colonialist priorities; they can derive from biased, non-representational data; they can use biased proxies (e.g. arrest rates as an indicator of actual crime rates); and they can take real population differences that have been created through structural oppression and treat these differences as unproblematic and essential (e.g. health insurance pricing that penalizes Black men and rewards White women based on differential rates of chronic illness) (Caplan et al., 2018; Hoffmann, 2019).

FML and algorithmic idealism

The problem of algorithmic inequality is not lost on computer scientists and engineers. Indeed, a vibrant field of FML has emerged with the shared goal of rectifying biases in ML systems (e.g. Barocas et al., 2017b; Chouldechova and Roth, 2020; Corbett-Davies and Goel, 2018; Kearns and Roth, 2019; Suresh and Gutttag, 2019). A recent review categorizes technical FML solutions into three categories, which map onto distinct definitions of fairness: anti-classification, classification parity, and calibration (Corbett-Davies and Goel, 2018). We define and discuss each of these in a subsequent section. For now, the relevant point is that each of these solutions proposes a computational path towards fair algorithmic outcomes. However, despite laudable aims, the proposed solutions consistently fall short.

FML approaches fall short because they stem from what we refer to as *algorithmic idealism*, enacting computation that assumes a meritocratic society and seeks to neutralize demographic disparities. Such an approach will always be inadequate in a context that is fundamentally unjust (Fazelpour and Lipton, 2020; Green and Viljoen, 2020). Algorithmic idealism begins with a base belief in equal opportunity, defining the problem of stratification as one caused by fallible human biases on the one hand, and imperfect statistical procedures, on the other. This perspective derives from illusory cultural narratives that misalign with the world that is – a world in which discrimination is

entrenched, elemental and compounding at the intersections of multiple marginalizations. From their current theoretical packaging, FML proposals emerge disinterested and objective; they seek optimal precision to apportion risks and rewards evenly across neatly bounded identity-based groups. Such proposals are consistently eluded by the fairness they mean to achieve.

We take FML's idealism as our point of departure, proposing instead *algorithmic reparation*, which re-conceives society through a critical Intersectional lens. This approach strives not for social equality, which treats everyone the same, but for social equity, which provides resources based on differential need, thus accounting for axes of historical (dis)advantage (Cook and Hegtvædt, 1983; Deutsch, 1975; Rawls, 1971). This means doing away with fairness and instead, coursing resources to those who have been systematically denied. This approach pairs the logic of Intersectionality with the praxis of reparation.

Algorithmic reparation

Intersectionality as a lens on ML

Intersectionality is not a singular theory, but an approach and a prism with a set of orienting assertions, goals and tools. It undergirds critical theories across subfields – critical race theory, critical feminist studies, queer theory – all of which share a fundamental focus on systemic power relations that privilege and penalize centralize and silence (Cho et al., 2013; Collins, 2019; Crenshaw, 1990; Hooks, 2000; Rahman, 2010). An Intersectional orientation is premised on the notion that identities are multiple and interrelated, shaped by and filtered through, societal structures and institutions. These structures and institutions concentrate and compound opportunities and constraints in ways that reflect and reinforce essentialized hierarchical arrangements. However, these hierarchical arrangements are not predetermined, and practitioners of Intersectionality task themselves with revealing and undoing, systems of injustice (Chepp and Collins, 2013; Collins, 2002; Collins and Bilge, 2020).

Intersectionality has taken on various meanings and been deployed towards varied ends while sustaining a core set of tenets (Cho et al., 2013; Collins, 2019; Ferree, 2018; McCall, 2005). The main tenets of Intersectionality are that inequalities are systemic and entangled, meaning that identities cannot be understood apart from their interrelation with each other and from their imbrication with socio-structural systems; 'objectivity' is never neutral, meaning positionality matters and marginal subjects provide a necessary but undervalued lens; that inequalities manifest through legal, personal, and professional (dis)advantage; and that hierarchies of power and privilege hide behind essentialisms, rendering their mechanisms imperceptible by default. These tenets combine with imperatives to

expose and negate essentialisms; empower the marginalized; and to name, highlight, and challenge agents and structures of domination (Carastathis, 2016; Collins and Bilge, 2020; Ferree, 2018).

Although Intersectionality has become embedded in academic texts and activist movements, it originates in the legal sector. Intersectionality arose in response to legal codes that erased and ignored co-occurring identity axes (e.g. Black women), working to account for discriminatory policies and practices that affect doubly marginalized legal subjects. With these legal foundations, proponents of Intersectionality emphasize the approach as an active political project (Cho et al., 2013; Collins and Bilge, 2020). Intersectionality is not just something to think with, but something to *do*. It is an intellectual method, but also, and in the first instance, a tool for empowering people and fostering social justice (Collins and Bilge, 2020). Thus, beyond identifying cases of systemic disadvantage, an Intersectional project also works to surge resources to those who are marginalized and deprived. This imperative to treat Intersectionality as a grounded, practical, material endeavor, can be served through the application of Intersectionality to ML evaluation and design.

As an approach to ML, our deployment of Intersectionality joins with and builds on a growing body of work attending to socio-historical power relations within computational systems. These include proposals for critical race methodologies for algorithmic fairness (Hanna et al., 2020), critical race theories applied to human-computer interaction (Ogbonnaya-Ogburu et al., 2020), decolonial AI (Mohamed et al., 2020), decolonial computer science (Birhane and Guest), computing for social change (Abebe et al., 2020), and affirmative action in algorithmic policing and criminal sentencing (Humerick, 2019). Inspired by, and combining elements from each of these projects, algorithmic reparation has a fundamental foundation in praxis, an emphasis on the multiplex of intersecting identities, and an explicit position of compensatory resource redistributions accomplished proactively through a reparative approach.

A reparative approach

Bringing Intersectionality to bear on ML, and bringing ML to bear on Intersectionality, grounds Intersectional politics in material conditions that interplay with contemporary lived experience through computational forms of governance and mundane technical engagements. That is, the *doing* (and undoing) that drives Intersectionality converges directly with issues of algorithmic inequality (Benjamin, 2019; Costanza-Chock, 2020; Mann and Matzner, 2019). We suggest animating Intersectional politics through practices of reparation.

'Reparation' is a historically grounded mechanism by which offending parties symbolically and materially mend wrongdoings enacted against individuals and groups (Torpey, 2006). Reparations have been assigned in the

context of war (Lu, 2017; Young, 2010), in acknowledgement of and apology for acts of colonialism (Gunstone, 2016; Lenzerini, 2008), and they remain a point of mobilization for Black civil rights activists in the United States, demanding material recompense for the multigenerational damages of slavery and segregation (Bittker, 2018 [1972]; Coates, 2014; Henry, 2009). Reparative acts are not just backward-looking, but also proactive, aiming to address the way historical wrongdoings affect current and future opportunity structures by channeling resources to make up for and overcome existing deficits.

Although traditionally applied in a legislative, often geopolitical context, we use ‘reparation’ in a broader sense, arguing for structural redress through algorithmic reform. This is more than the conceptual loosening of a legal term. Legal and political systems hinge reparation on identifiable culprits and victims along with demonstrable links between the wrongdoing of one party and the consequences of wrongful actions upon the aggrieved. However, this is rarely how structural, Intersectional oppressions operate. What makes Intersectional oppressions so pervasive and pernicious is their diffusion through institutional infrastructures, policies of governance, language, culture, individual attitudes and interpersonal dynamics. The systematic, multifaceted, often subtle nature of Intersectional inequality is at odds with linear relations of harm and blame. Algorithmic reparation thus incorporates redress into the assemblage of technologies that interweave macro institutions and micro-interactions, embedding an equitable agenda into the material systems that govern daily life³.

Our call for reparative algorithms is motivated by a broader mandate for equity and social justice, but it is also motivated by the specific conditions of automation that leave no neutral option (Broussard, 2018; Bucher, 2018; Mann and Matzner, 2019; Noble, 2018). In general, the distribution of resources can either reinforce inequalities, make them worse, or make them better. However, ML systems are intrinsically self-perpetuating in ways that ossify and intensify the outcomes they engender. This is because algorithms render decisions seemingly objective and divorced from human discretion; because they are opaque and inscrutable; and because their outcomes often have no technical means of undoing, even if circumstances call for correction (Bucher, 2018; Eubanks, 2018; Gillespie, 2014, 2018; Pasquale, 2015; Vaidhyanathan, 2018). Our proposal for algorithmic reparation assumes a moral duty to ameliorate, rather than aggravate, structural and historical stratifications as they manifest in computational code. This proposal sits in direct opposition to the prevailing logic of FML, which seeks to de-bias algorithms and make them fairer. In contrast, a reparative approach assumes and leverages bias to make algorithms more equitable and just.

A critical read on FML: from fair to reparative

The field of FML is dedicated to making algorithms fairer for the people whom ML systems affect. In a review of the field, Corbett-Davies and Goel (2018) catalogue FML strategies, distinguishing between three definitions of fairness that underpin various computational solutions: *anti-classification*, *classification parity*, and *calibration*.⁴ As lamented by the authors, these efforts have been largely unsuccessful, reconstituting the unjust social conditions they were designed to alleviate (Corbett-Davies and Goel, 2018: 2).

FML’s troubles, we argue, stem from the field’s foundation in *algorithmic idealism* – a meritocratic misconception of the world and a political ambivalence that this fallacy permits.

In this section, we describe existing FML solutions and the definitions of fairness to which they ascribe, highlight empirical instances in which these solutions proved lacking, and reimagine for each instance an alternative starting point derived from an Intersectional reparative approach. In doing so, we advance the case for algorithmic reparation in juxtaposition to the idealism embedded in aspirations towards ‘fair.’

Anti-classification

Anti-classification stipulates that algorithmic estimates do not consider protected class attributes such as race, class, gender, or (dis)ability. This includes direct consideration of these characteristics as well as proxies for them. Corbett-Davies and Goel (2018) equate this to principles of equal protection under the law (Karst, 1977) and ‘taste-based’ discrimination in economics (Becker, 2010 [1957]), by which advantages and disadvantages cannot be assigned based on demographic preference. Algorithmically, anti-classification systems strive to encode indifference to the identities of individuals who will be subject to automated outcomes.

Anti-classification principles underlie automated employment programs that aim to circumvent managerial biases in candidate selection, avoiding the historical race–class–gender–age–nationality (dis)advantages that have historically shaped which candidates make it past initial screenings (Lahey and Oxley, 2018; Oreopoulos, 2011; Quillian et al., 2017). In practice, these algorithmic systems reproduce social hierarchies pervasive to the populations from which they select. Technology conglomerate Amazon’s use of anti-classification algorithms exemplifies this point.

In 2014, Amazon developed a recruitment tool to aid in its own hiring processes. The tool used ML to sort applicants based on optimal fit for each position, removing social identity characteristics from consideration (Dastin, 2018).

The trifold purpose was to increase efficiency, select the best candidates, and avoid implicit biases, especially against women, as this group has been (and remains) under-represented in the technology sector (Beede et al., 2011; Harrison, 2019). However, by 2015, it became evident that the automated system was not operating as planned. Consistently, the recruitment algorithms assigned higher scores to men and lower scores to women. The reason for this is that the system was trained on the company's previous 10 years of employment data, which reflected a male-dominated sector. That is, Amazon's workforce, like the broader technology workforce, was populated disproportionately by men. Consequently, using existing data, the hiring system learned that men were the preferred candidates. This self-perpetuating cycle was so pronounced that any indicator of feminine gender identity in an application lowered the applicant's score. A degree from a women's college, participation in women-focused organizations, and feminized language patterns all reduced the evaluative outcome. Although Amazon attempted to adjust for these issues, the system continued to find proxies for gender and reward men at the expense of women. Amazon eventually retired the program (Dastin, 2018).

From an Intersectional perspective, anti-classification systems are intrinsically faulty. These systems are premised on the erasure of difference, a flattening of demographic traits. Such an approach ideologically sidesteps the empirical reality of systemic inequality, but it cannot statistically or mathematically address it. The data that feeds these systems and the people who are subject to them, operate from hierarchically differentiated positions. These distinctions are, and will continue to be, captured and reproduced through computation.

In contrast, a reparative approach would highlight, name and encode hierarchical distinctions as they manifest across social identity categories. From this foundation, Amazon's algorithms would not invisibilize gender but would instead define gender as a primary variable on which to optimize. This could mean weighting women, trans, and non-binary applicants in ways that mathematically bolster their candidacy, and potentially deflating scores that map onto stereotypical indicators of White cisgender masculinity, thus elevating women, trans, and non-binary folks in accordance with, and in rectification of, the social conditions that have gendered (and raced) the high-tech workforce. Moreover, it would not treat 'woman' as a homogenous (binary) category, but would label and correct for intersections of age, race, ability and other relevant variables that shape gendered experiences and opportunity structures.

This reparative system would literally value the contributions underrepresented applicants bring to the company while normalizing Intersectional gender diversity in tech, such that high-level positions and the pathways to them, are recast as plausible and expected across gender groups. The technical solution (women, trans, and non-binary

individuals get a statistical boost) would thus have direct effects on the company's work environment (more women, trans, and non-binary employees are hired at Amazon) and broader social effects on Intersectional gendered social relations (women, trans, and non-binary folks are normalized in the technology sector and the pathways to technology careers more seamless for these individuals to pursue). If these ends remain untenable, a reparative approach would indicate that ML ought not to be used in hiring decisions⁵.

Classification parity

Classification parity is defined in terms of equal errors in classification across social identity groups. This aims to achieve parity in the error rates of predictive performance measures. Corbett-Davies and Goel (2018) identify several measures of classification error: false-positive rates, false-negative rates, precision, recall, the proportion of decisions that are positive, and the area under the ROC curve (AUC) (see Berk et al., 2018; Skeem and Lowenkamp, 2016). They focus in particular on false positives and the proportion of decisions that are positive, as these are the error metrics that FML researchers have given the most attention (Corbett-Davies and Goel, 2018). We also focus on those metrics here, along with false negatives, as these are relevant to high-profile cases of algorithmic inequality.

False positives and false negatives are errors in predicting how likely it is that something will (or will not) happen. Proportion of positive decisions, also known as 'demographic parity' (Feldman et al., 2015), means that a given outcome distributes equally across social identity groups. These measures – false positives, false negatives, and demographic parity – have been central to debates about (and critiques of) ML in criminal sentencing, the most notable case of which is the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) recidivism risk assessment tool.

COMPAS is a widely used and commercially available tool designed to predict the likelihood that a criminal defendant will reoffend. In 2016, a ProPublica report analyzed pre-sentencing data from Broward County, Florida, a large jurisdiction using the COMPAS system. The report found that Black defendants were systematically assigned higher risk scores than White defendants, and that risk was overpredicted for Black defendants and underpredicted for White defendants (i.e. Black defendants recidivated at a lower rate than what the algorithm predicted, and White defendants recidivated at a higher rate than what the algorithm predicted) (Angwin et al., 2016a, 2016b). The overall disparity in Black–White risk assessments represents an error of demographic parity, while the over- and under prediction of Black–White recidivism

represents errors of false positives (Black defendants) and false negatives (White defendants).

Classification parity is rooted in the assumption that if errors distribute evenly, then decisions will be fair, and that if data are accurate and representative enough, fair distributions can be achieved. Intersectionality indicates that these assumptions are misguided. They are misguided because the *effects* of risk are not equivalent across groups and because the data that feed into the ML system are infused with social policies and practices that shape statistical inputs and outputs. In terms of sentencing decisions, a criminal record and time in prison are undesirable for anyone. However, the negative effects of conviction and incarceration compound for individual Black defendants, flow on to their families and communities, and reinforce racial disparities in wealth, health, family stability, and mental well-being (Pettit and Western, 2004; Travis et al., 2014; Western and Pettit, 2010; Western and Sirois, 2019). Moreover, the likelihood of contact with police and of conviction is significantly higher for poor Black men than any other group (Alexander, 2010). Indeed, criminal justice data are self-perpetuating, such that those groups defined as criminally ‘risky’ are in fact, at disproportionate risk of ensnarement by the criminal justice system (Brayne, 2017; Brayne et al., 2015; Christin et al., 2015; Ferguson, 2017; Richardson et al., 2019).

A reparative approach would supplant the goal of ‘parity’ with, instead, systemic redress, beginning with the social facts of disproportionate risk between racial groups and the history of race–class dynamics that inform training data. From this, reparative decision aids would work to actively protect poor communities of color, especially poor Black men, over and above other subpopulations. This means the production and deployment of algorithms that keep Black men out of prison and keep police out of Black communities, defending against the criminalization of Blackness and rectifying racialized prison pipelines.

Calibration

Calibration specifies that ‘outcomes should be independent of protected attributes conditional on risk scores’ (Corbett-Davies and Goel, 2018: 6). Calibration can be thought of as a more nuanced take on anti-classification. The calibration approach is such that identity characteristics should only be considered by an algorithmic equation if those characteristics have demonstrable, empirical effects on the outcome under consideration. That is, the system calibrates to differential risk levels between groups and assigns scores according to those base-level differences.

To illustrate calibration, we remain with the COMPAS example. We do so because Northpointe, the company behind COMPAS, has responded to critics by claiming that in fact, their algorithms are fair because they satisfy

calibration. What they mean is that a Black defendant classified as high risk by COMPAS is equally likely to recidivate as a White defendant classified as high risk. In an open letter to ProPublica, the company states:

ProPublica focused on classification statistics that did not take into account the *different base rates of recidivism for blacks and whites*. Their use of these statistics resulted in false assertions in their article that were repeated subsequently in interviews and in articles in the national media (Dietrich et al., 2016: 1) (emphasis added).

Defending itself, Northpointe justifies its product based on calibration standards. Their defense is inadequate on both technical and social grounds.

On a technical level, although errors calibrate for base differences between groups, the *kinds* of errors are inconsistent. As detailed by ProPublica, Black defendants remain subject to disproportionate false positives and White defendants are rewarded with disproportionate false negatives (Angwin et al., 2016a). Black people are mis-assessed with overly strong risk scores and White people are mis-assessed with overly lenient risk scores. Concretely, this means that more Black people end up in jail and more White people remain free.

There are also non-technical reasons to be dissatisfied with Northpointe’s response and in turn, dissatisfied with calibration as an algorithmic standard. In particular, the data used by Northpointe to train their algorithms reflect racist policing tendencies in the United States that over-indict Black men, creating (not just reflecting) different base rates between raced, classed, and gendered groups (Brayne, 2017; Brayne et al., 2015; Ferguson, 2017; Richardson et al., 2019). Moreover, the carceral system not only responds to criminality, but through a constellation of mechanisms, also begets further violations (Alexander, 2010). Thus, Northpointe’s reliance on calibration as a technical justification affirms and entrenches a system in which existing injustices act as the basis for their own amplified reproduction.

Like anti-classification, calibration seeks to remove identity from the decision equation (though in qualified form). Like classification parity, calibration works to achieve equivalence between groups (by adjusting for differential base risk). As detailed in the subsections above, both of these objectives are ineffective for reducing inequality and may intensify inequitable social arrangements. In contrast, algorithmic reparation rejects the notion of identity erasure, even on the grounds of empirically distinct risk rates. It instead takes stock of social disparities as they map along axes of identity, exposing the underlying causes of differential risk and undoing the stratification that those differences both represent and produce. In practice, this could mean higher recidivism risk thresholds for Black defendants, lower thresholds for bail and parole,

and weighted statistical adjustments that account for over-policing in poor communities of color. Evaluated through a reparative Intersectional lens, any algorithm that did not address these base inequities would be deemed inadequate.

Methods and barriers

The technical means of algorithmic reparation are already computationally viable, but its social effects can only take hold through meaningful implementation. It is thus to implementation that we now turn. Rather than reinvent the wheel, we select two recently proposed methods of algorithmic praxis that serve as possible tools of application for the reparative strategies discussed herein: archivist data curation and distributed AI power. Both of these methods are founded in transdisciplinarity and require mutual collaborations between academic and non-academic actors. We also identify and discuss three challenges to implementing algorithmic reparation, including social, legal, and institutional barriers. Together, these methods and barriers ground algorithmic reparation within a context of both possibility and constraint.

Methods of implementation

Archivist curation is one promising approach to implementing algorithmic reparation. This draws on the professional expertise of archival practice, honed by librarians and museum curators, applying these skills to ML data (Donovan, 2020; Jo and Gebru, 2020). Unjust algorithmic outputs are inextricable from problems with source data. These problems can be a function of representation in datasets and/or social factors that crystalize in data form. Managing these data issues can be prohibitively complex. However, professionals trained in collection and curation have skill sets that are transferrable to the ML sector, with Jo and Gebru (2020) noting *consent, inclusivity, power, transparency, ethics, and privacy* as data-relevant issues that have been well addressed in library sciences.

Drawing on their extant skill sets, curation professionals are capable of managing, collecting, arranging, and auditing data in ways that not only avoid re-entrenched inequalities, but optimize for marginal elevation, enacting targeted precision unachievable by those who are not professionally trained in curatorial methods. This includes the capacity to account for complex identity configurations in which advantages and disadvantages are in simultaneous operation, and the insight to determine which pieces of data are relevant to collect and, more importantly, what data ought not be collected. Such skills and practices are well suited to the problems discussed above, such as hiring and criminal sentencing, in which the complexity of the data and its entanglement with a multitude of confounding and compounding variables have proven intractable for data practitioners alone.

Distributed AI power is a second potential method. This method is premised on undoing standard power asymmetries between those who make, and those who are affected by, ML systems. The approach argues for tools that are legible to, and co-created with, impacted communities, especially those communities with histories of vulnerability prior to, and re-entrenched with, automation (Kalluri, 2020). Distributed AI power tactics rely on reciprocal engagement between developers and community stakeholders, with reverse pedagogies by which community stakeholders serve as experts in their lived experiences (Mohamed et al., 2020).

This method is exemplified by academic-activist collaborative projects, undertaken by groups such as the Algorithmic Justice League, Data for Black Lives, and the Carceral Tech Resistance Network, among others. Each of these organizations leverages community knowledge to challenge and partner with commercial, governance, and regulatory bodies to enact technical, social, and policy changes. The Algorithmic Justice League, for example, has performed audits of race and gender in facial recognition technologies, leading several companies to revamp their programming in ways that improve the classification accuracy for dark-skinned women in image search tasks (Raji and Buolamwini, 2019). Data for Black Lives, which coordinates thousands of engineers, mathematicians and activists, is training former inmates in data science so that this directly affected population can actively participate in the reform of the criminal justice system (Heaven, 2020). In turn, the Carceral Tech Resistance Network (2020) trains in and with communities, mobilizing towards the abolition of carceral tech and reparations for these systems' racialized damages. All three organizations have joined with others to activate against the use of facial recognition in policing, demonstrating the fundamental incongruity between these tools and racial justice, eventuating a cascade of corporate and legislative moratoria (Flynn, 2020; Heilweil, 2020; Lazar et al., 2020). These projects begin with, are led by, and develop through, affected communities, with a record that demonstrates the capacity to enact reparative approaches to ML evaluation and design.

Barriers

Grounding algorithmic reparation means identifying both opportunities and challenges. The methods just discussed represent encouraging prospects, but there are empirical reasons that ML keeps reproducing inequality, and these realities are robust and obdurate. Enacting algorithmic reform requires unvarnished realism about the conditions under which any sociotechnical intervention will go into effect. For algorithmic reparation, implementation will face interrelated social, legal, and institutional barriers. Although addressing each barrier is beyond the scope of

the present work, we lay them out to set clear terms for the path ahead.

Socially, reparation relies on a base logic that diverges from normative conceptions of fairness, opting instead for uneven resource allocations targeted at the margins. As evidenced by the backlash against affirmative action policies and resistance to critical race curricula (Ray and Gibbons, 2021; Vought, 2020), an intentional reallocation of resources will, undoubtedly, come up against significant friction. Rectificatory tactics will be difficult to accept for those who ascribe to an image of society that is functionally meritocratic, and this baseline assumption is indeed, deep-seated.

There will also be legal and institutional challenges. Reparation calls for centralized knowledge about and action based upon protected class attributes. This is difficult under legal conditions that prohibit the collection of such data and/or its consideration in consequential decisions like employment, lending, school admissions, and criminal sentencing (Lieberwitz, 2008; Long and Batemen, 2020; Skeem and Lowenkamp, 2020). Similar prohibitions written into institutional policies will create blockades against algorithmic reparation within organizational settings.

There are also real challenges to the kinds of interdisciplinary and socially engaged collaborations necessary for reparative algorithmic projects. Power and compensation disparities persist between computer scientists and social scientists, and between academic and non-academic organizations (Carrigan and Bardini, 2021; Hackett and Rhoten, 2011; Stavrianakis, 2015; Viseu, 2015), along with epistemological schisms that are difficult to reconcile (Bauer, 1990; Richter and Paretti, 2009). These impediments to meaningful inter/trans/non-disciplinary collaboration are exacerbated by academic incentive structures that reward traditional intra-disciplinary outputs over and above hybrid and expansively defined research products (Woelert and Millar, 2013), despite widespread statements about the value of disciplinary blending and community-engaged science (Hackett and Rhoten, 2011; Viseu, 2015). Contending with these institutional challenges means considering not only who will do the work of algorithmic reparation, but also how it can be done across sectors, with the support of leadership, mechanisms of accountability, democratic oversight, and equitable returns for practitioners' labor.

Conclusions

Summary

Technologies reflect and create the societies from which they stem and in which they proliferate. By default, technologies will embody the values of the powerful and reconstitute the stratified hierarchies those values represent

(Benjamin, 2016, 2019; Broussard, 2018; Browne, 2015; Costanza-Chock, 2020; Davis, 2020). These patterns of reflection, reconstitution, and in turn, amplification of structural inequality have borne out in spectacular fashion with the integration of ML systems into personal and institutional life.

The field of FML has emerged in response, with computer scientists and engineers proposing myriad technical fixes to the injustices of automation. Yet, algorithmic inequalities persist. In their efforts to hide, distribute evenly between, and calibrate social identity traits, FML practitioners operate with a goal of fairness and equality when instead, equity and reparation are required. We make this case in the body of the text above, suggesting a move away from fairness, replaced by an anti-oppressive, Intersectional approach. We intend for this approach to guide algorithmic design and to act as an evaluative standard by which existing algorithmic systems are judged, adjusted, and where necessary, omitted or dismantled. Our proposal is thus geared towards building better systems and holding existing ones to account.

We highlight two possible methods of implementation – professionalized archival data curation and distributed AI power. Both methods are consonant with the base assumptions and objectives of algorithmic reparation and they both show promise as practical means for algorithmic reform. We also take stock of social, legal, and institutional barriers to implementation, providing a realistic perspective on the work ahead.

Next steps

Continuing this focus on the work ahead, we conclude by considering next steps in the ongoing project towards social and technical restructuring. Here, we emphasize the need for context-specific attention, more and multiple tools, and multipronged approaches that converge technical, social, and institutional efforts.

Instruments of social change – technical or otherwise – never operate in a vacuum. In the final substantive section of this paper, we selected two newly introduced mechanisms by which algorithmic reparation might be implemented. Testing these in diverse empirical settings will reveal how they function, where they fall short, and what kinds of infrastructural conditions will be required for these methods to take meaningful effect.

It will also be vital to explore and create a cache of methods and tools, addressing specific needs, specific conditions, and creating interoperability between social and technical systems. The acute need for a constellation of methods and tools becomes clear when we consider the varied and engrained structural reasons why inequalities continue to manifest in algorithmic form. Algorithmic reparation will, necessarily, run against the grain of multiple

status quos, requiring numerous iterations, agile applications, and persistent adjustments for this uphill endeavor.

In service of creating a robust toolbox, this paper's third author (Yang) is currently leading a project to devise technical instruments that audit and optimize for inequality reduction in decision systems. This is a computational mechanism that centers impact estimations that most reduce inequality in automated decision outputs. These auditing tools are intended specifically for institutional decision aids, such as those used in hiring processes, loan allocations, and admission decisions, calibrated to the particular inequalities of the communities affected. Projects such as this, which are currently in development, portend a new and critically informed landscape of sociotechnical relations.

We also note that 'next steps' cannot be technical alone. Any algorithmic solution to social problems is necessarily partial and incomplete, requiring complementary social, legal, and institutional evolutions. Concretely, this means rethinking discrimination policies that erase and thus ignore identity attributes; reworking institutional incentive structures and power arrangements that silo academic disciplines from each other and from the public sector; introducing regulatory implements that capture and censure discriminatory algorithmic outputs; and forming organizational bodies dedicated to auditing technical systems and assuring their allocative and representational ends.

In practice, the problems of algorithmic systems are the problems of social systems, and meaningful solutions will be technical *and* social in nature. These solutions will not come easy, nor fast, nor with absolute finitude. Just as undoing racism, sexism, classism, and colonialism are continuous, evolving, non-linear projects, so too is the journey to unmake the inequalities of algorithms and code. Algorithmic reparation is not a final or encompassing answer, but a critical, equitable, Intersectional foundation.

Acknowledgements

The authors would like to thank members of the Humanising Machine Intelligence project at the Australian National University and affiliates with the Berkman Klein Center for Internet & Society at Harvard University. We are especially grateful to Professor Toni Erskine, Dr Claire Benn, and Dr Sarah Logan for their comments and ideas during early stages of this paper's development.


Declaration of conflicting interests


The authors declared no potential conflicts of interest with respect to the research, authorship and/or publication of this article.

Funding

The authors received no financial support for the research, authorship and/or publication of this article.

ORCID iDs

Jenny L. Davis  <https://orcid.org/0000-0003-0952-5842>

Apryl Williams  <https://orcid.org/0000-0003-4896-9366>

Notes

1. Throughout we capitalize 'Intersectional' when referencing the theoretical paradigm, as is convention. We use a lower case 'i' in all other circumstances.
2. We embody this call in the present work, co-authored by two sociologists and a computer scientist, with combined backgrounds in critical race theory, critical technology studies, communication theory, and ML.
3. Some theorists contest the use of reparation in a structural sense due to its uneasy fit with the specified relations of harm, and the adjudication of the specific wrongdoings that currently define reparative outcomes in legal settings (Young, 2010; Lu, 2017). We do not disagree with this point but depart from it, challenging the specified nature of reparation as a tool of redress when in practice, harms are often structural and diffuse, operating outside the scope of legal-political institutions.
4. In addition to the fairness models reviewed here, causality-based notions of fairness are of burgeoning interest in FML. These models have not yet been implemented in mainstream ways and are thus not included in Corbett-Davies and Goel's (2018) review or in our paper. However, recent critiques of causal fairness in ML show similar vulnerabilities to existing approaches (see Hu and Kohler-Hausmann, 2020).
5. See Glazebrook and Sundaram (2020) 'Why we don't use AI for hiring decisions'. Available at <https://www.beapplied.com/post/why-we-dont-use-ai-for-hiring-decisions>.

References

- Abebe R, Barocas S, Kleinberg J, et al. (2020) Roles for computing in social change. In: Proceedings of the 2020 conference on fairness, accountability, and transparency, January 27–30, 2020, Barcelona, Spain, pp. 252–260.
- Alexander M (2010) *The New Jim Crow: Mass Incarceration in the Age of Colorblindness*. New York: The New Press.
- Ames MG (2018) Deconstructing the algorithmic sublime. *Big Data & Society* 5(1): 1–4.
- Amoore L (2020) *Cloud Ethics*. Durham, NC: Duke University Press.
- Angwin J, Larson J, Mattu S, et al. (2016a) How we analyzed the compas recidivism algorithm. *ProPublica*. Available at: <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>
- Angwin J, Larson J, Mattu S, et al. (2016b) Machine bias: There's a software used across the country to predict future criminals. And it's biased against Blacks. *ProPublica*. Available at: <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
- Angwin J, Mattu S and Larson J (2015) The tiger mom tax: Asians are nearly twice as likely to get a higher price from Princeton review. *ProPublica*. Available at: <https://www.propublica.org/article/asians-nearly-twice-as-likely-to-get-higher-price-from-princeton-review>

- Barocas S, Crawford K, Shapiro A, et al. (2017a) The problem with bias: Allocative versus representational harms in machine learning. In: 9th annual conference of the special interest group for computing, information and society October 29, 2017. Philadelphia, PA.
- Barocas S, Hardt M and Narayanan A (2017b) Fairness in machine learning: Limitations and opportunities. *NIPS Tutorial 1*. <https://fairmlbook.org/pdf/fairmlbook.pdf>.
- Bauer HH (1990) Barriers against interdisciplinarity: Implications for studies of science, technology, and society (STS). *Science, Technology, & Human Values* 15(1): 105–119.
- Becker GS (2010 [1957]) *The Economics of Discrimination*. Chicago, IL: University of Chicago Press.
- Beede DN, Julian TA, Langdon D, et al. (2011) Women in STEM: A gender gap to innovation. *Economics and Statistics Administration Issue Brief* 4(11).
- Benjamin R (2016) Catching our breath: Critical race STS and the carceral imagination. *Engaging Science, Technology, and Society* 2: 145–156.
- Benjamin R (2019) *Race After Technology: Abolitionist Tools for the New Jim Code*. Cambridge, UK: Polity.
- Berk R, Heidari H, Jabbari S, et al. (2018) Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research*. 50(1): 3–44.
- Birhane A and Guest O (2020) Towards decolonising computational sciences. *arXiv preprint arXiv:2009.14258*
- Bittker B (2018 [1972]) *The Case for Black Reparations*. Boston, MA: Beacon Press.
- Brayne S (2017) Big data surveillance: The case of policing. *American Sociological Review* 82(5): 977–1008.
- Brayne S, Rosenblat A and Boyd D (2015) Predictive policing. *Data & Civil Rights: A New Era Of Policing And Justice*. Available at: https://datacivilrights.org/pubs/2015-1027/Predictive_Policing.pdf
- Broussard M (2018) *Artificial Unintelligence: How Computers Misunderstand the World*. Cambridge, MA: MIT Press.
- Browne S (2015) *Dark Matters: On the Surveillance of Blackness*. Durham, NC: Duke University Press.
- Bucher T (2018) *If... Then: Algorithmic Power and Politics*. NY: Oxford University Press.
- Caplan R, Donovan J, Hanson L, et al. (2018) *Algorithmic Accountability: A Primer*. NY: Data & Society. Available at: <https://datasociety.net/library/algorithmic-accountability-a-primer/>
- Carastathis A (2016) *Intersectionality: Origins, Contestations, Horizons*. Lincoln, NE: University of Nebraska Press.
- Carceral Tech Resistance Network (2020, March 30). Available at: <http://carceral.tech/practice>.
- Carrigan C and Bardini M (2021) Majorism: Neoliberalism in student culture. *Anthropology & Education Quarterly*. In press. 52(1): 42–62.
- Chepp V, Collins PH (2013) Intersectionality. In: Celis K, Kantola J, Waylen G, et al. (eds) *The Oxford Handbook of Gender and Politics*. New York: Oxford University Press, pp. 57–87.
- Cho S, Crenshaw KW and McCall L (2013) Toward a field of intersectionality studies: Theory, applications, and praxis. *Signs: Journal of Women and Culture in Society* 38(4): 785–810.
- Chouldechova A and Roth A (2020) A snapshot of the frontiers of fairness in machine learning. *Communications of the ACM* 63(5): 82–89.
- Christin A, Rosenblat A and Boyd D (2015) *Courts and Predictive Algorithms*. NY: Data & Society. Available at: <https://datasociety.net/library/data-civil-rights-courts-and-predictive-algorithms/>
- Coates T-N (2014) The case for reparations. *The Atlantic* 313. Available at: <https://www.theatlantic.com/magazine/archive/2014/06/the-case-for-reparations/361631/>
- Collins PH (2002) *Black Feminist Thought: Knowledge, Consciousness, and the Politics of Empowerment*. London, UK: Routledge.
- Collins PH (2019) *Intersectionality as Critical Social Theory*. Durham, NC: Duke University Press.
- Collins PH and Bilge S (2020) *Intersectionality*. Medford, MA: Polity Press.
- Cook KS and Hegtvedt KA (1983) Distributive justice, equity, and equality. *Annual Review of Sociology* 9(1): 217–241.
- Corbett-Davies S and Goel S (2018) The measure and mismeasure of fairness: A critical review of fair machine learning. *arXiv preprint arXiv:1808.00023*. Available at: <https://arxiv.org/abs/1808.00023>
- Costanza-Chock S (2020) *Design Justice: Community-led Practices to Build the Worlds we Need*. Cambridge, MA: MIT Press.
- Crawford K (2021) *Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence*. New Haven, CT: Yale University Press.
- Crawford K, Dobbe R, Dryer T, et al. (2019) *AI NOW 2019 Report*. NY: AI NOW Institute. Available at: https://ainowinstitute.org/AI_Now_2019_Report.pdf
- Crenshaw K (1990) Mapping the margins: Intersectionality, identity politics, and violence against women of color. *Stanford Law Review* 43(6): 1241–1299.
- Dastin J (2018) Amazon Scraps secret AI recruiting tool that showed bias against women. *Reuters* 9. Available at: <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>
- Davis JL (2020) *How Artifacts Afford: The Power and Politics of Everyday Things*. Cambridge, MA: MIT Press.
- Deutsch M (1975) Equity, equality, and need: What determines which value will be used as the basis of distributive justice? *Journal of Social Issues* 31(3): 137–149.
- Dieterich W, Mendoza C and Brennan T (2016) COMPAS risk scales: Demonstrating accuracy equity and predictive parity. Northpointe Inc. Available at: http://go.volarisgroup.com/rs/430-MBX-989/images/ProPublica_Commentary_Final_070616.pdf.
- D'Ignazio C and Klein LF (2020) *Data Feminism*. Cambridge, MA: MIT Press.
- Donovan J (2020) You purged racists from your website? Great, now get to work. *Wired*. Available at: <https://www.wired.com/story/you-purged-racists-from-your-website-great-now-get-to-work/>
- Eubanks V (2018) *Automating Inequality: How High-Tech Tools Profile, Police, and Punish the Poor*. New York: St Martin's Press.
- Fazelpour S and Lipton ZC (2020) Algorithmic fairness from a non-ideal perspective. In: Proceedings of the AAAI/ACM conference on AI, ethics, and society, February 7–9 New York, pp. 57–63. Association for Computing Machinery.

- Feldman M, Friedler SA, Moeller J, et al. (2015) Certifying and removing disparate impact. In: Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining, August 10-13th Sydney, Australia, pp. 259–268.
- Ferguson AG (2017) Illuminating black data policing. *Ohio State Journal of Criminal Law* 15: 503–525.
- Ferree MM (2018) Intersectionality as theory and practice. *Contemporary Sociology: A Journal of Reviews* 47(2): 127–132.
- Flynn S (2020) 13 Cities where police are banned from using facial recognition tech. *Innovation and Tech Today*. Available at: <https://innotechtoday.com/13-cities-where-police-are-banned-from-using-facial-recognition-tech/>
- Gillespie T (2014) The relevance of algorithms. In: Gillespie T, Boczkowski PJ and Foot KA (eds) *Media Technologies: Essays on Communication, Materiality, and Society*. Cambridge, MA: MIT Press, pp. 167–193.
- Gillespie T (2018) *Custodians of the Internet: Platforms, Content Moderation, and the Hidden Decisions That Shape Social Media*. New Haven, CT: Yale University Press.
- Glazebrook K and Sundaram K (2020) Why we don't use AI in hiring decision. *Applied*. Available at: <https://www.beapplied.com/post/why-we-dont-use-ai-for-hiring-decisions>
- Green B and Viljoen A (2020) Algorithmic realism: expanding the boundaries of algorithmic thought. In: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, pp. 19–31. January 27-30, Barcelona, Spain.
- Gunstone A (2016) Reconciliation, reparations and rights. In: Short CLaD (ed) *Handbook of Indigenous Peoples' Rights*. London, UK: Routledge pp. 301–312.
- Hackett EJ and Rhoten DR (2011) Engaged, embedded, enjoined: Science and technology studies in the National Science Foundation. *Science and Engineering Ethics* 17(4): 823–838.
- Hanna A, Denton E, Smart A, et al. (2020) Towards a critical race methodology in algorithmic fairness. In: Proceedings of the 2020 conference on fairness, accountability, and transparency, pp. 501–512. January 27-30th. Barcelona, Spain.
- Harrison S (2019) Five years of tech diversity reports—and little progress. *Wired*. Available at: <https://www.wired.com/story/five-years-tech-diversity-reports-little-progress/>
- Heaven WD (2020) Predictive policing algorithms are racist they need to be dismantled. *MIT Technology Review*. Available at: <https://www.technologyreview.com/2020/07/17/1005396/predictive-policing-algorithms-racist-dismantled-machine-learning-bias-criminal-justice/>
- Heilweil R (2020) Big tech companies back away from selling facial recognition to police. That's progress. *Vox*. Available at: <https://www.vox.com/recode/2020/6/10/21287194/amazon-microsoft-ibm-facial-recognition-moratorium-police>
- Henry CP (2009) *Long Overdue: The Politics of Racial Reparations*. New York: NYU Press.
- Hill K (2020) Another arrest, and jail time, due to a bad facial recognition match. The New York Times. Available at: <https://www.nytimes.com/2020/12/29/technology/facial-recognition-misidentify-jail.html>
- Hoffmann AL (2019) Where fairness fails: Data, algorithms, and the limits of antidiscrimination discourse. *Information, Communication & Society* 22(7): 900–915.
- Hooks B (2000) *Feminist Theory: From Margin to Center*. London, UK: Pluto Press.
- Hu L and Kohler-Hausmann I (2020) What's sex got to do with machine learning. *arXiv preprint arXiv:2006.01770*.
- Humerick JD (2019) Reprogramming fairness: Affirmative action in algorithmic criminal sentencing. *HRLR Online* 4(2): 213–244.
- Jo ES and Gebru T (2020) Lessons from archives: Strategies for collecting sociocultural data in machine learning. In: Proceedings of the 2020 conference on fairness, accountability, and transparency, pp. 306–316. January 27-30, Barcelona, Spain.
- Kalluri P (2020) Don't ask if AI is good or fair, ask how it shifts power. *Nature World View*. 583(169) Available at: <https://media.nature.com/original/magazine-assets/d41586-020-02003-2/d41586-020-02003-2.pdf>
- Karst KL (1977) Foreword: Equal citizenship under the fourteenth amendment. *Harvard Law Review* 91: 1–301.
- Kayser-Bril N (2020) Google apologizes after its vision AI produced racist results. *Algorithm Watch*. Available at: <https://algorithmwatch.org/en/story/google-vision-racism/>
- Kearns M and Roth A (2019) *The Ethical Algorithm: The Science of Socially Aware Algorithm Design*. Oxford, UK: Oxford University Press.
- Kitchin R (2017) Thinking critically about and researching algorithms. *Information, Communication & Society* 20(1): 14–29.
- Lahey JN and Oxley DR (2018) Discrimination at the intersection of age, race, and gender: Evidence from a lab-in-the-field experiment. Working Paper 25357 National Bureau of Economic Research Cambridge, MA December 2018 <https://www.nber.org/papers/w25357>
- Lazar S, Benn C and Günther M (2020) Large scale facial recognition is incompatible with a free society. *The Conversation*. Available at: <https://theconversation.com/large-scale-facial-recognition-is-incompatible-with-a-free-society-126282>
- Lenzerini F (2008) *Reparations for Indigenous Peoples: International and Comparative Perspectives*. Oxford, UK: Oxford University Press.
- Lieberwitz R (2008) Employment discrimination law in the United States: on the road to equality? New developments in Employment Discrimination Law, Bulletin of Comparative Labour Relations. In: New developments in employment discrimination law conference, Tokyo, Japan. Available at: https://www.jil.go.jp/event/ro_forum/resume/080220/USA_.pdf
- Long MC and Batemen NA (2020) Long-run changes in underrepresentation after affirmative action bans in public universities. *Educational Evaluation and Policy Analysis* 42(2): 188–207.
- Lu C (2017) *Justice and Reconciliation in World Politics*. Cambridge, UK: Cambridge University Press.
- McCall L (2005) The complexity of intersectionality. *Signs: Journal of Women in Culture and Society* 30(3): 1771–1800.
- Mann M and Matzner T (2019) Challenging algorithmic profiling: The limits of data protection and anti-discrimination in responding to emergent discrimination. *Big Data & Society* 6(2): 1–11.
- Mohamed S, Png M-T and Isaac W (2020) Decolonial AI: Decolonial theory as sociotechnical foresight in artificial intelligence. *Philosophy & Technology* 33(4): 659–684.

- Noble SU (2018) *Algorithms of Oppression: How Search Engines Reinforce Racism*. NY: NYU Press.
- Ogbonnaya-Ogburu IF, Smith AD, To A, et al. (2020) Critical race theory for HCI. In: Proceedings of the 2020 CHI conference on human factors in computing systems, pp. 1–16.
- O’Neil C (2016) *Weapons of Math Destruction: How big Data Increases Inequality and Threatens Democracy*. NY: Broadway Books.
- Oreopoulos P (2011) Why do skilled immigrants struggle in the labor market? A field experiment with thirteen thousand resumes. *American Economic Journal: Economic Policy* 3(4): 148–171.
- Pasquale F (2015) *The Black box Society: The Secret Algorithms That Control Money and Information*. Cambridge, MA: Harvard University Press.
- Pettit B and Western B (2004) Mass imprisonment and the life course: Race and class inequality in US incarceration. *American Sociological Review* 69(2): 151–169.
- Quillian L, Pager D, Hexel O, et al. (2017) Meta-analysis of field experiments shows no change in racial discrimination in hiring over time. *Proceedings of the National Academy of Sciences* 114(41): 10870–10875.
- Rahman M (2010) Queer as intersectionality: Theorizing gay Muslim identities. *Sociology* 44(5): 944–961.
- 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, pp. 429–435.
- Rawls J (1971) *A Theory of Justice*. Cambridge, MA: Harvard university press.
- Ray R and Gibbons A (2021) Why are states banning critical race theory? *Brookings Institute*. Available at: <https://www.brookings.edu/blog/fixgov/2021/07/02/why-are-states-banning-critical-race-theory/>
- Richardson R, Schultz JM and Crawford K (2019) Dirty data, bad predictions: How civil rights violations impact police data, predictive policing systems, and justice. *New York University Law Review Online* 94: 15–55.
- Richter DM and Paretto MC (2009) Identifying barriers to and outcomes of interdisciplinarity in the engineering classroom. *European Journal of Engineering Education* 34(1): 29–45.
- Seaver N (2017) Algorithms as culture: Some tactics for the ethnography of algorithmic systems. *Big Data & Society* 4(2): 1–12.
- Simonite T (2018) When it comes to gorillas, Google photos remains blind. *Wired*. Available at: <https://www.wired.com/story/when-it-comes-to-gorillas-google-photos-remains-blind/>
- Skeem JL and Lowenkamp CT (2016) Risk, race, and recidivism: Predictive bias and disparate impact. *Criminology: An Interdisciplinary Journal* 54(4): 680–712.
- Skeem JL and Lowenkamp CT (2020). Using algorithms to address trade-offs inherent in predictive recidivism. *Social Sciences & the Law*, 38(3): 259–278.
- Stavrianakis A (2015) From anthropologist to actant (and back to anthropology): Position, impasse, and observation in socio-technical collaboration. *Cultural Anthropology* 30(1): 169–189.
- Suresh H and Gutttag JV (2019) A framework for understanding unintended consequences of machine learning. *arXiv preprint arXiv:1901.10002*. Available at: <https://arxiv.org/abs/1901.10002>
- Torpey JC (2006) *Making Whole What has Been Smashed: On Reparations Politics*. Cambridge, MA: Harvard University Press.
- Travis J, Western B and Redburn FS (2014) *The Growth of Incarceration in the United States: Exploring Causes and Consequences*. Washington, DC: National Research Council. Available at: <https://www.nap.edu/catalog/18613/the-growth-of-incarceration-in-the-united-states-exploring-causes>
- Vaidhyanathan S (2018) *Antisocial Media: How Facebook Disconnects us and Undermines Democracy*. New York: Oxford University Press.
- Viseu A (2015) Integration of social science into research is crucial. *Nature News* 525(7569): 291.
- Vought R (2020) Memorandum for the heads of executive departments and agencies: Training in the federal government. *United States Office of Management and Budget*. Available at: <https://www.whitehouse.gov/wp-content/uploads/2020/09/M-20-34.pdf>
- Western B and Pettit B (2010) Incarceration & social inequality. *Daedalus* 139(3): 8–19.
- Western B and Sirois C (2019) Racialized re-entry: Labor market inequality after incarceration. *Social Forces* 97(4): 1517–1542.
- Woelert P and Millar V (2013) The ‘paradox of interdisciplinarity’ in Australian research governance. *Higher Education* 66(6): 755–767.
- Young IM (2010) *Responsibility for Justice*. Oxford, UK: Oxford University Press.