

# Fairness & Bias

(in peer review & other sociotechnical systems)

Nihar Shah & Zachary Lipton  
Carnegie Mellon University

Email: [zlipton@cmu.edu](mailto:zlipton@cmu.edu), [nihars@cs.cmu.edu](mailto:nihars@cs.cmu.edu)

Twitter: [@zacharylipton](https://twitter.com/zacharylipton)

# Are these vectors fair?

-8.1, 4.1, 9.6, -3.8, -2.5

-0.8, 3.4, -7.0, 8.8, -0.8

-0.8, 3.4, -7.0, 8.8, -0.8

7.6, 1.3, 1.1, -3.2, 8.5

# Are these vectors fair?

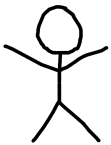
-8.1, 4.1, 9.6, -3.8, -2.5



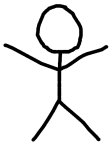
-0.8, 3.4, -7.0, 8.8, -0.8



-0.8, 3.4, -7.0, 8.8, -0.8



7.6, 1.3, 1.1, -3.2, 8.5



# Part 1 Goals

- Put the “fairness” back in “fair machine learning”
- Put the “ethics” back in AI Ethics
- Provide conceptual clarity to avoid category errors
- Give broad background via philosophy, economics, and the law
- Examine injustices due to ML & current approaches in fair ML literature
- Critically examine proposed mitigation strategies
- Re-focus attention on context required to determine just actions

In Part 2, we will focus on peer review considering how CS/stats might offer tools to diagnose and remediate various notions of unfairness.



# Part 1: Fairness and Machine Learning

- **What is fairness?**
  - Philosophical perspectives
  - Economic perspectives
  - Legal Perspectives
- Fair Machine Learning
  - Motivating problems
  - Statistical fairness metrics
  - Mitigation Strategies
- Limitations and Dangers
  - Potential harms of misguided interventions
  - Revisiting philosophical foundations
  - A non-ideal perspective
- Causal approaches to fair ML
  - Causal parities
  - Mitigation Strategies
  - Limitations
- Can Interpretability help?
  - Techniques
  - Limitations
  - Contestability

What is  $\{fairness \rightarrow justice\}$ ?

# Institutes of Justinian

*The most plausible candidate for a core definition comes from the Institutes of Justinian, a codification of Roman Law from the sixth century AD, where justice is defined as ‘the constant and perpetual will to render to each his due’.*

*— Stanford Encyclopedia of Philosophy*

## **Key points:**

- Concerns treatment of individuals
- Arises to resolve conflicts when interests clash
- Justice concerns one's due (an obligation, in contrast to charity)
- Invokes impartiality—two cases relevantly alike should be treated similarly
- Centers on an agent “whose will alters circumstances of its objects”

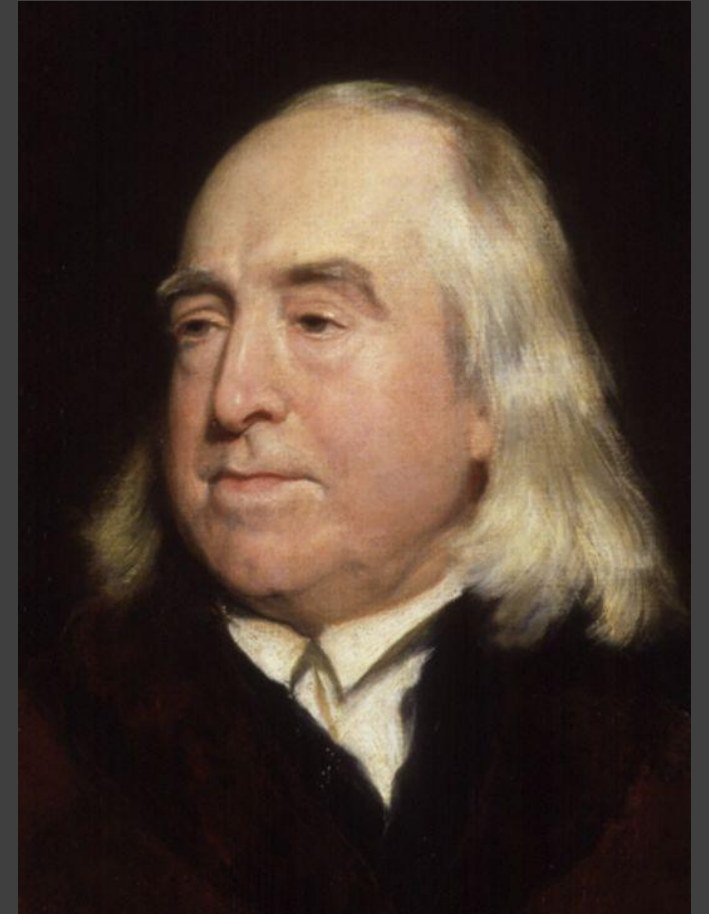
# Aristotle — Nichomachean ethics

- Earliest ethical treatises in philosophy
- Concerns development of *virtuous character*
- **Book V addresses “Justice & Fairness”**
- Justice “in relation to another person”
- Distinguishes b/w justice as
  1. **law-abiding vs. fair**
  2. **on the whole vs. particular**
- Distributive (“equal shares for equals”) vs. corrective (“subtract unjust gain”)



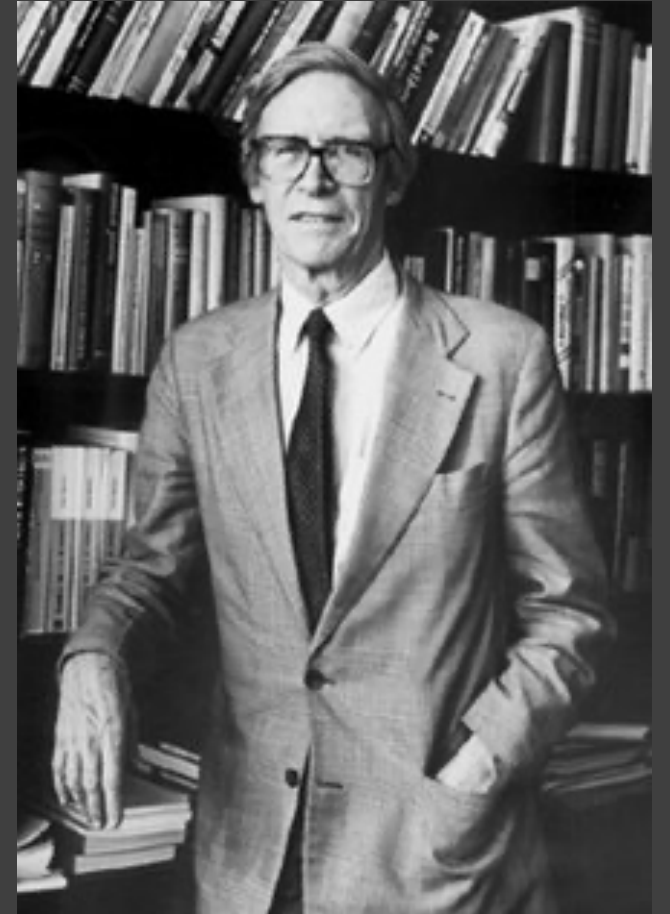
# Bentham

- Founder of modern utilitarianism
- Fundamental axiom:  
*"it is the greatest happiness of the greatest number that is the measure of right and wrong."*
- "Advocated individual and economic freedoms, the separation of church and state, freedom of expression, equal rights for women, the right to divorce, and (in an unpublished essay) the decriminalising of homosexual acts."



# Rawls

- Revived normative political philosophy with *A Theory of Justice* in 1971 (—Will Kymlicka)
- *Theory of “justice as fairness” recommends equal basic rights, equality of opportunity, and promoting the interests of the least advantaged members of society.*
- Pioneered thought experiment of “original position” (veil of ignorance)
- Origin of **ideal/non-ideal** distinction



Four distinctions

# Conservative vs Ideal

- Should justice be views “conservative of existing norms and practices” or “demanding reform of these norms and practices”
- Conservative:
  - Respect people’s rights under existing laws, rules & expectations
- Ideal:
  - “Reason to change laws, practices and conventions quite radically, thereby creating new entitlements and expectations”
- The ideal specifies a notion of equality, dismisses claims of justice that do not arise from / accord with the principle.



# Corrective vs. Distributive

- Distributive: justice is a principle for allocating good to individuals
  - **Multilateral**, assumes a **distributing agent**
- Corrective: “remedial principle that applies when one person interferes with another’s legitimate holdings”
  - **Bilateral**, concerns relationship between wrong-doer and the wronged
- Idea: theft of a rich person’s property ought to be remediated via corrective justice, but is not demanded by distributive justice
- Philosophers and lawyers disagree about standard of responsibility to mandate corrective justice

# Procedural vs. Substantive

- Distinction between the virtue of the method by which benefits and burdens are allocated vs. the final allocation itself.
- Coin tosses may yield equal allocations but be procedurally unjust.
- Some (e.g. Nozick) suggest final distribution is irrelevant, only “sequence of prior events that created it” matters
- Some suggest justice of a procedure is determined by its outcomes

# Comparative vs. Non-comparative

- When does determining justice require looking at what others can claim?
- Comparative harms:
  - E.g., denied a job that was offered to a less qualified candidate
- Non-comparative harms:
  - Rights to free speech, religion, etc.
  - Whether or not these rights are denied others, they are still one's right.
- May face trade-offs between comparative/non-comparative harms.
- Focusing myopically on one category can blind us to the other.
- Denying **everyone** a good may have a comparative (but not NC) harm

*(Fair ML literature typically focuses on comparative justice)*

# E.g.: Universal Declaration of Human Rights

- A3—Everyone has the right to life, liberty and security of person. (NC)
- A4—No one shall be held in slavery or servitude; slavery and the slave trade shall be prohibited in all their forms. (NC)
- A5—No one shall be subjected to torture or to cruel, inhuman or degrading treatment or punishment. (NC)
- A7—All are equal before the law and are entitled without any discrimination to equal protection of the law. All are entitled to equal protection against any discrimination in violation of this Declaration and against any incitement to such discrimination. (C)
- A9—No one shall be subjected to arbitrary arrest, detention or exile. (NC)
- A16—Men and women of full age, without any limitation due to race, nationality or religion, have the right to marry and to found a family. They are entitled to equal rights as to marriage, during marriage and at its dissolution. (NC & C)

# The Scope of Justice

- To who or what does justice apply?
- When & among whom do principles of justice take effect?
- “Who can make claims of justice?”
- “Who might have the corresponding obligation to meet them?”
- If comparative principles are being applied, who should be counted as part of the comparison group?
- Which principles are universal vs contextual?

# Ideal and Non-Ideal Theorizing about Justice

- Key distinction in Rawls and subsequent theorizing on justice/fairness:
- The ideal approach:
  - Imagine a perfectly just world.
  - Try to minimize discrepancy between our world and the ideal.
  - Has been used to argue against affirmative action—*ideal world is color-blind*
- The non-ideal approach:
  - *[Non-ideal theorists] ... seek a **causal explanation** of the problem to determine what can and ought to be done about it, and who should be charged with correcting it. This requires an evaluation of the mechanisms causing the problem, as well as responsibilities of different agents to alter these mechanisms.*

— “The imperative of integration” Elizabeth Anderson 2019

Economic perspectives

# Becker—“The Economics of Discrimination”

- Considers workers belonging to two groups (say whites & blacks)
- Introduces “taste-based discrimination” a model of outright prejudice
- Employer acts as though there is a cost associated w. hiring blacks
- However, profit function  $\pi_i$  regards two groups as perfect substitutes
- Each employer’s utility function assigns “disutility”  $d_i$  per black worker

$$V_i = \pi_i - n_i^b \cdot d_i$$

- Market equilibrium results in
  1. Induces a sorting of workers → firms hire only blacks or only whites
  2. Different wages for white and black workers
  3. Wages determined by the marginal discriminator



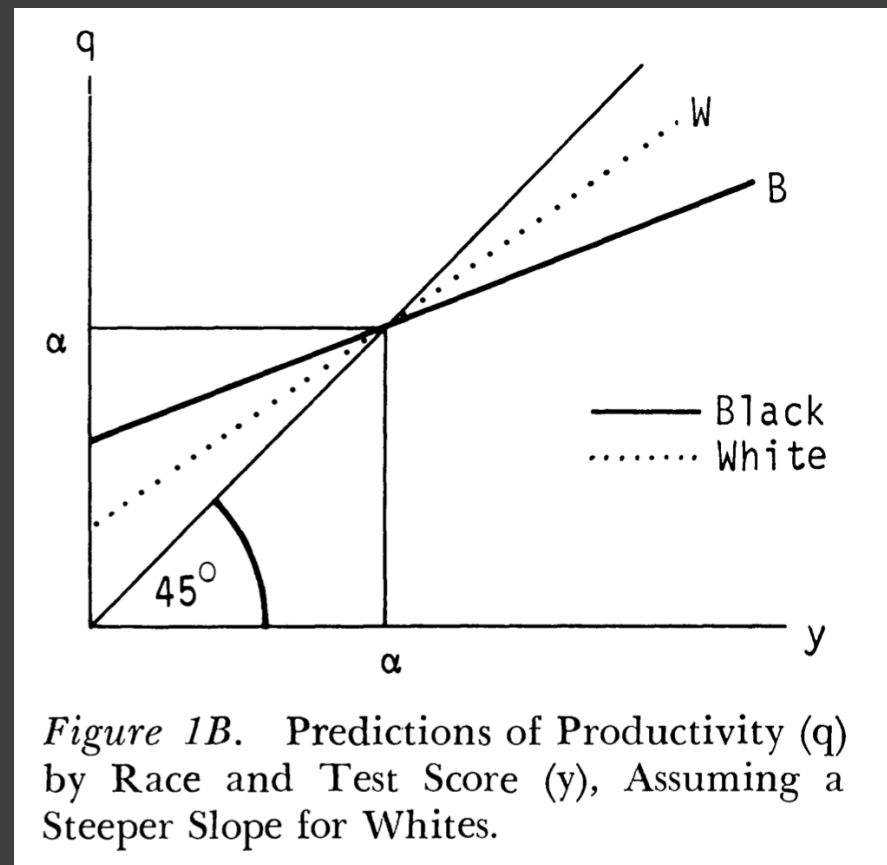
# Arrow's Rebuttal of Becker (1973)

- Argues that taste-based discrimination will fail because discriminating employers will be driven from market by inefficiency
- Discusses situations with actual productivity differences among groups due to discrimination in other spheres of life (e.g., education)
- Argues that consequence of forcing identical wages may be that employers stop employing from minority group
- Suggests imperfect information as alternative cause of disparities

*"I believe these results are only the barest fragment of what could be found with better and more detailed systems in which there is an interaction between reality and perceptions of it"*

# “The statistical theory of racism and sexism”

- Introduced by Phelps (1972)
- Models how disparities arise absent disutility, and w. identically dist. skills.
- Requires only signal more difficult to obtain for minority workers.
- Simplified by Aigner & Cain (1977)
  - Worker quality  $q$  normally distributed, group-conditioned noise levels  $u$ .
  - Observed test results  $y = q + u, u \sim N(0, \sigma_g^2)$



# Modeling dynamics of affirmative-action

- Several papers and a book (the anatomy of racial inequality) by Glenn Loury investigate discrimination in hiring, notable for richer considerations of the interplay of policies and the behavior of agents.
- Coates & Loury (1993) look at long term effects of affirmative action.
- Consider interplay of interventions, investment in education.
- One key insight: even when groups are equal ex ante, equilibrium outcomes following some interventions can appear to confirm negative stereotypes.

Legal thinking

# What distinguishes legal thinking?

- Some argue legal profession's "special skill" is combined abilities to:
  - (i) address facts and evidence
  - (ii) understand the "full context" of a particular event, dispute or decision
- More specific forms of legal reasoning:
  - making decisions according to rules, determining which sources are authoritative, respecting precedent, notions of burdens of proof, resolving questions of decision-making jurisdiction
- "Law is inevitably and especially subject to the unforeseeable complexity of the human condition"

# Generality

*“Although disputes, in court and out, involve particular people with particular problems engaged in particular controversies, the law tends to treat the particulars it confronts as members of larger categories. Rather than attempting to reach the best result for each controversy in a wholly particularistic and contextual way, law’s goal is often to make sure that the outcome for all or at least most of the particulars in a given category is the right one.”*

- “It is better saith the Law to suffer a mischief (that is particular to one) than an inconvenience that may prejudice many.” — Lord Coke
- Legal thinking attempts to extract from particular cases general rules
- Pragmatic exercise, accepts that result is approximate, not always just.
- Legal decisions in appellate court require reasoning not just about a particular case, but hypothetical scenarios.

# Formalistic vs. Interpretative Arguments

- Rules can be over- or under-inclusive w.r.t. background justification.
- Toy examples:
  - Speed limit is always 55, not varying w. road conditions.
  - Law prohibiting naturalized citizens from becoming president (b/c loyalty)
- Formalistic arguments prefer letter of law to intent.
- Injustices allowed in view that legislature (not court) sets the rules
- Interpretative approach more flexible, but open to arbitrariness
- Who decides what speed is *reasonable* to drive?

# Precedent

- “Law characteristically faces backward.”
- “The legal system [is] particularly concerned with precedent—with doing the same thing that has been done before just because it has been done before.”
- Precedent flows
  1. Vertically — (lower courts obey higher)
  2. Horizontally — courts default to respecting past decisions

*Related to philosophical notions:  
Impartiality, anti-arbitrariness, conservative justice*



# Authority-based reasoning

- “taking the source of a directive rather than the reasons behind it as a justification for following it”
- Limit interpretation to make application of law impartial, consistent.
- What a rule *says* matters

# The fringes of a law

- Rule often has *core* purpose, defs. can get fuzzy w. wide application
- Example: what is prohibited by a “no vehicles in park” rule?
  - Cars and motorcycles surely. But what about Scooters? Wheelchairs? Roller skates? Skateboards?
- Fringe cases require interpretation, establishment of precedent.
- Application dominated by “easy cases”, adjudication by “hard cases”

# Anti-discrimination law



President Lyndon B. Johnson shakes hands with Martin Luther King after signing the Civil Rights Act of 1964

# Disparate treatment

- Addresses intentional discrimination
- Includes decisions explicitly based on a *protected characteristic*
- Also intentional discrimination via proxy variables



# Disparate impact

- Facially neutral practices that might nevertheless have an “*unjustified adverse impact on members of a protected class*”
- Complicated doctrine w 3 tests
  1. Plaintiff must demonstrate **statistical disparity** (e.g. 4/5 rule)
  2. Defendant must show that decisions are justified by ‘**business necessity**’
  3. Plaintiff must show defendant can achieve goal w ‘**alternative practice**’

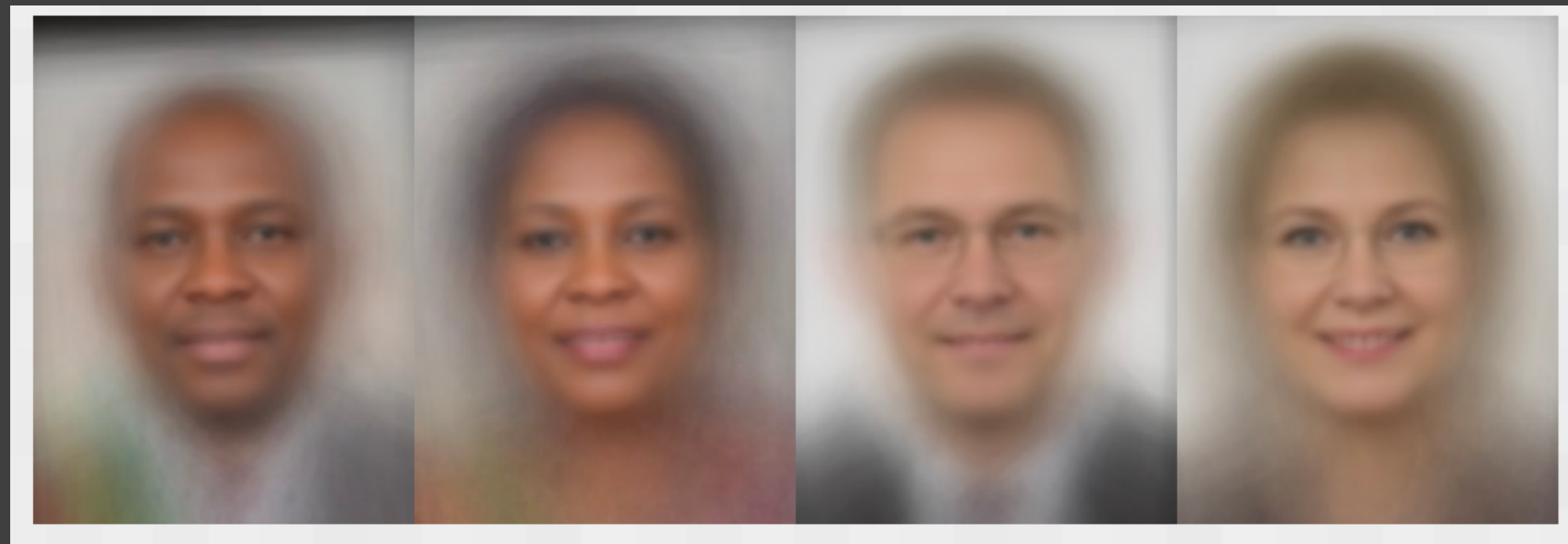
# Part 1: Fairness and Machine Learning

- What is fairness?
  - Philosophical perspectives
  - Economic perspectives
  - Legal Perspectives
- **Fair Machine Learning**
  - Motivating problems
  - Statistical fairness metrics
  - Mitigation Strategies
- Limitations and Dangers
  - Potential harms of misguided interventions
  - Revisiting philosophical foundations
  - A non-ideal perspective
- Causal approaches to fair ML
  - Causal parities
  - Mitigation Strategies
  - Limitations
- Can Interpretability help?
  - Techniques
  - Limitations
  - Contestability

# ProPublica — Machine Bias, 2016



# Gender Shades—2018





# Bias in word embeddings, 2016

---

## Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

---

Tolga Bolukbasi<sup>1</sup>, Kai-Wei Chang<sup>2</sup>, James Zou<sup>2</sup>, Venkatesh Saligrama<sup>1,2</sup>, Adam Kalai<sup>2</sup>

<sup>1</sup>Boston University, 8 Saint Mary's Street, Boston, MA

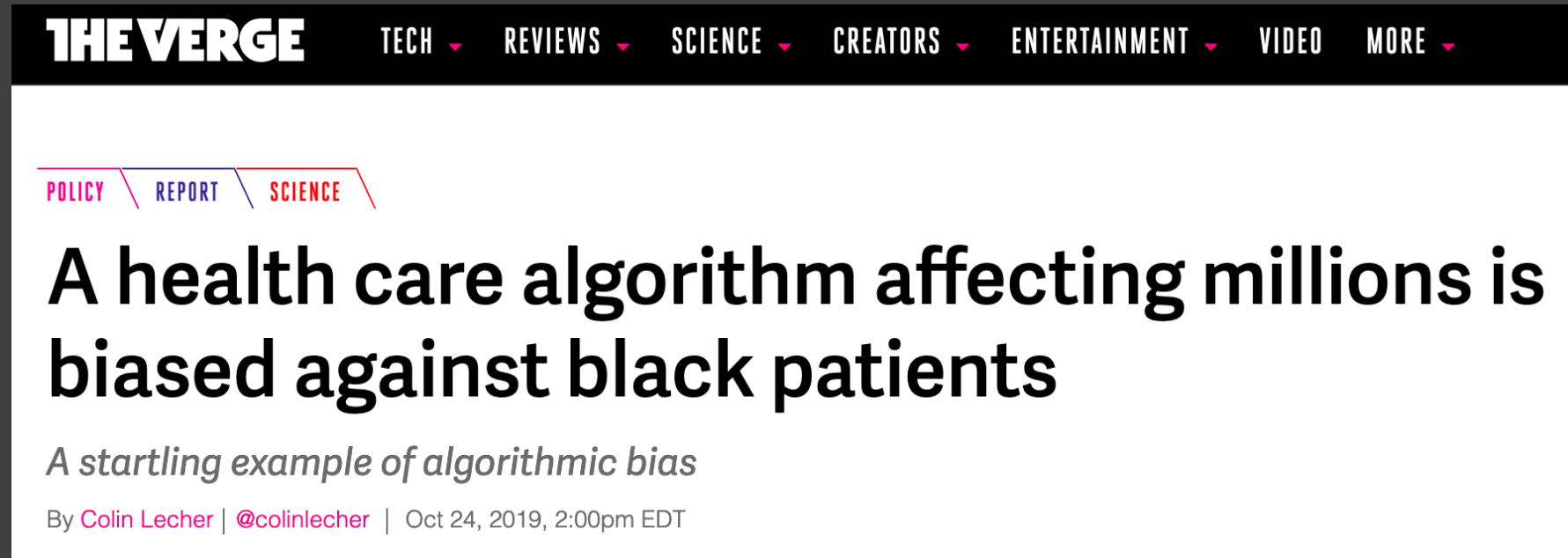
<sup>2</sup>Microsoft Research New England, 1 Memorial Drive, Cambridge, MA

tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

### Abstract

The blind application of machine learning runs the risk of amplifying biases present in data. Such a danger is facing us with *word embedding*, a popular framework to represent text data as vectors which has been used in many machine learning and natural language processing tasks. We show that even word embeddings trained on Google News articles exhibit female/male gender stereotypes to a disturbing extent. This raises concerns because their widespread use, as we describe, often tends to amplify these biases. Geometrically, gender bias is first shown to be captured by a direction in the word embedding. Second, gender neutral words are shown to be linearly separable from gender definition words in the word embedding. Using

# Biased allocation of healthcare (2019)




“The authors estimated that this racial bias reduces the number of Black patients identified for extra care by more than half. Bias occurs because the algorithm uses health costs as a proxy for health needs. Less money is spent on Black patients who have the same level of need, and the algorithm thus falsely concludes that Black patients are healthier than equally sick White patients.”

# A Pernicious Pattern

1. Take a problem ill-described as statistical prediction.
2. Fashion a surrogate prediction problem anyway.
3. Define metrics of success, e.g. accuracy, assuming prediction as task.
4. Trouble arises due to insufficiency of problem description.
5. Work to “solve” the problem while working entirely within the paradigm whose insufficiencies are themselves the root cause.
6. Mislead the public by purporting to have addressed the problem, often by redefining the objective.

# Some examples:



## 3 BIAS-FREE ALGORITHMS

We know that a diverse workforce is critical for a company's success. However, not all algorithms are created equal and they are not inherently objective or fair. If an algorithm is trained on a biased training set, it will simply codify human biases, and often worsen the bias that exists. pymetrics has developed an algorithmic auditing technique that uses a reference set of tens of thousands of potential biases, and we do this to ensure that we produce a bias-free algorithm. We believe that a 4/5ths rule. We believe that a certain demographic population signal.

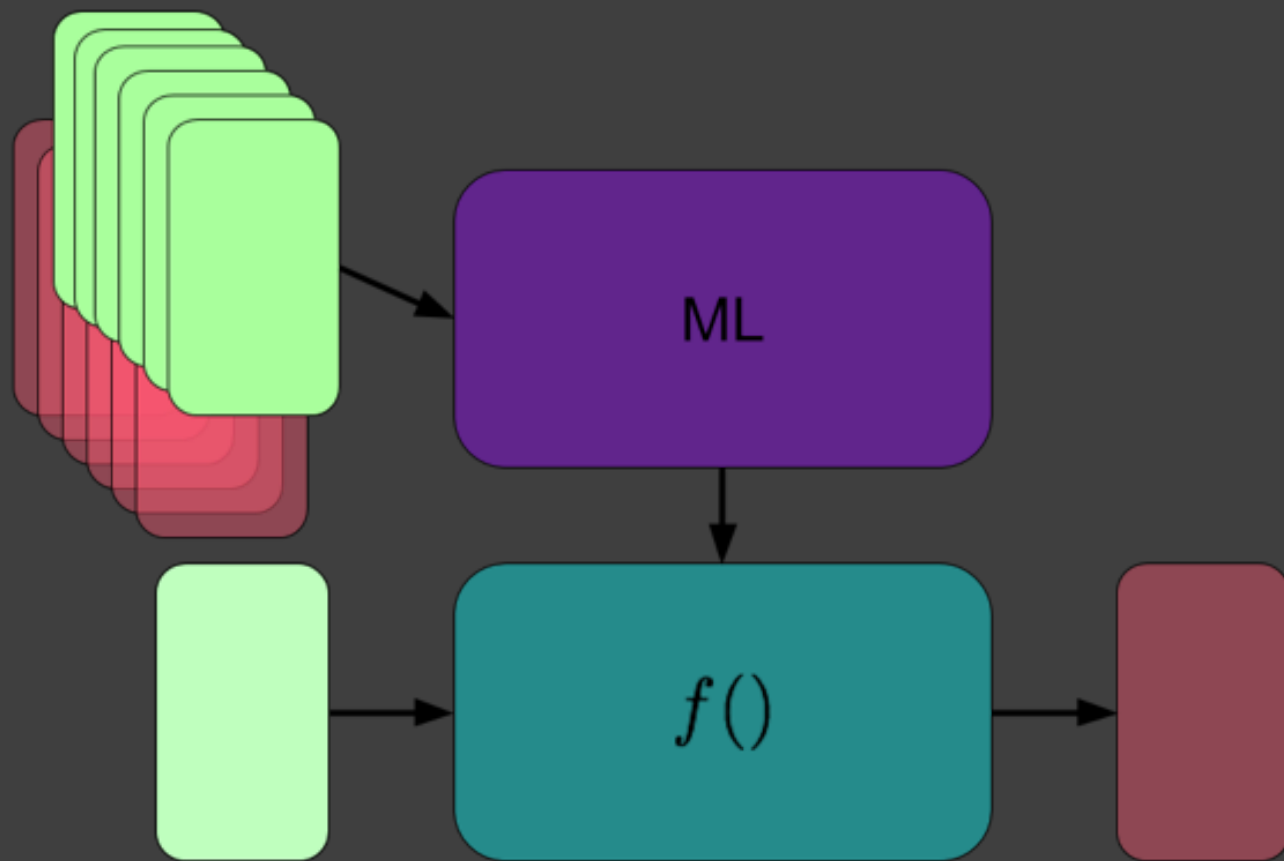
[Request Whitepaper](#)

The field is a bit more sophisticated than this. There are many excellent papers on bias, eg using powerful tools for causal reasoning on arxiv. My colleagues are making good progress and not giving up.

3:46 PM - 30 May 2019

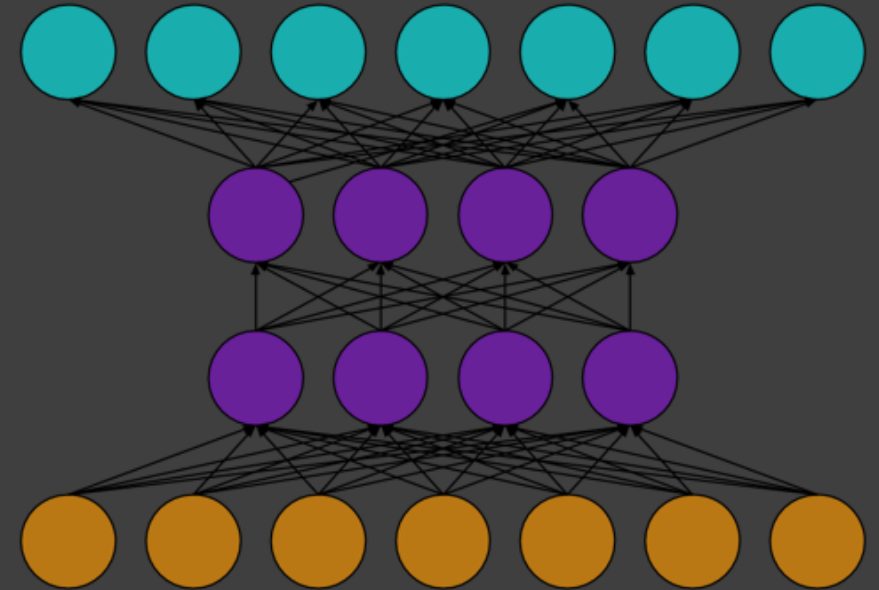
# Fair Machine Learning (Statistical)

# Supervised Learning

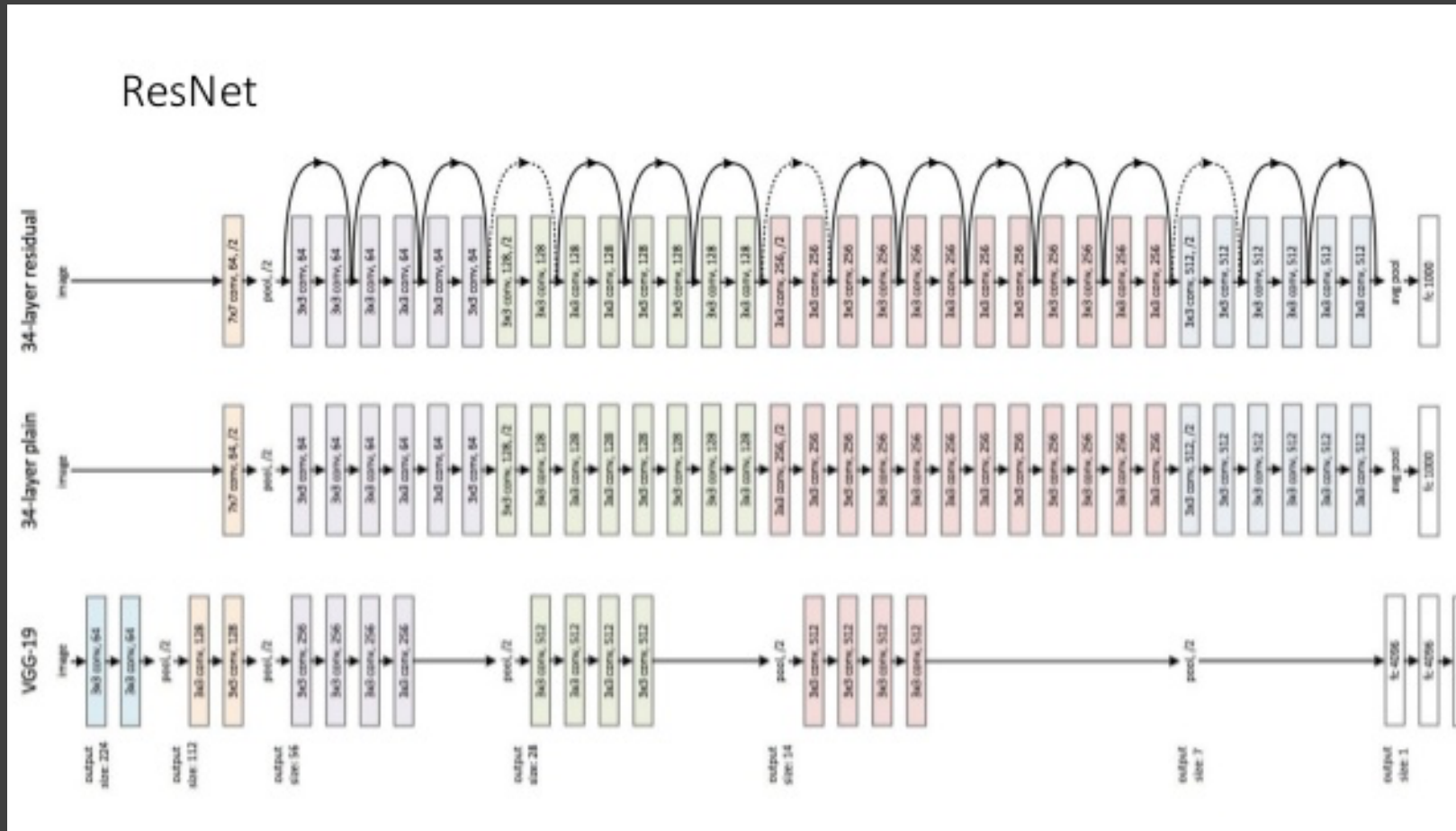


# Neural Networks

- Composed of artificial **neurons**
- Connected by weighted edges (*like synapses*)
- Each takes **input**, emits **output**
- Can approximate complex functions



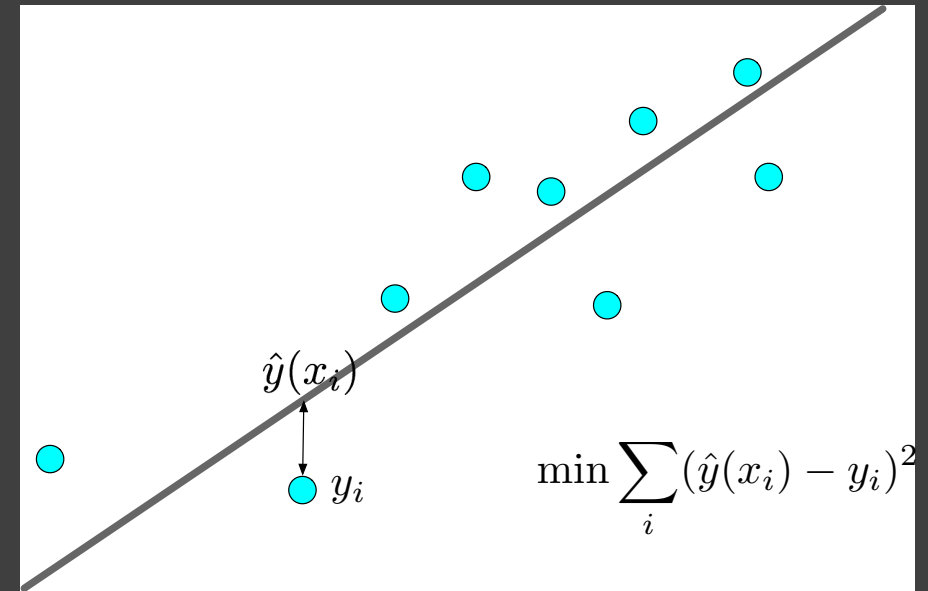
# Deep Learning (2005+)





# Curve-fitting

- [How a pioneer of artificial intelligence became one of its sharpest critics](#)
- “ML is stuck on... learning associations”
- ...we did not expect...so many problems could be solved by pure curve fitting”
- Learns associations, not causal relations
- Sometimes, that’s enough



# Wrong at the Root

- In real world supervised learning is applied in scenarios **that are not really prediction problems**
- Applications of ML riddled with mismatches between real-world desiderata and the formalism provided by supervised learning
- Growing literatures on *fair algorithms* and *interpretability* purport to *take on these problems*.
- Due to flawed formulations, **may give only appearance of solutions.**

# The foundations of *algorithmic bias*

Even if we truly were addressing a prediction problem, things go wrong:

- Some groups under-represented, benefits of automation unequal.
- The training labels themselves may be noisy or *biased*.
- Models often optimized for wrong task altogether (choice of surrogate task may have disparate effects).
- Task may be *easier* for one group.

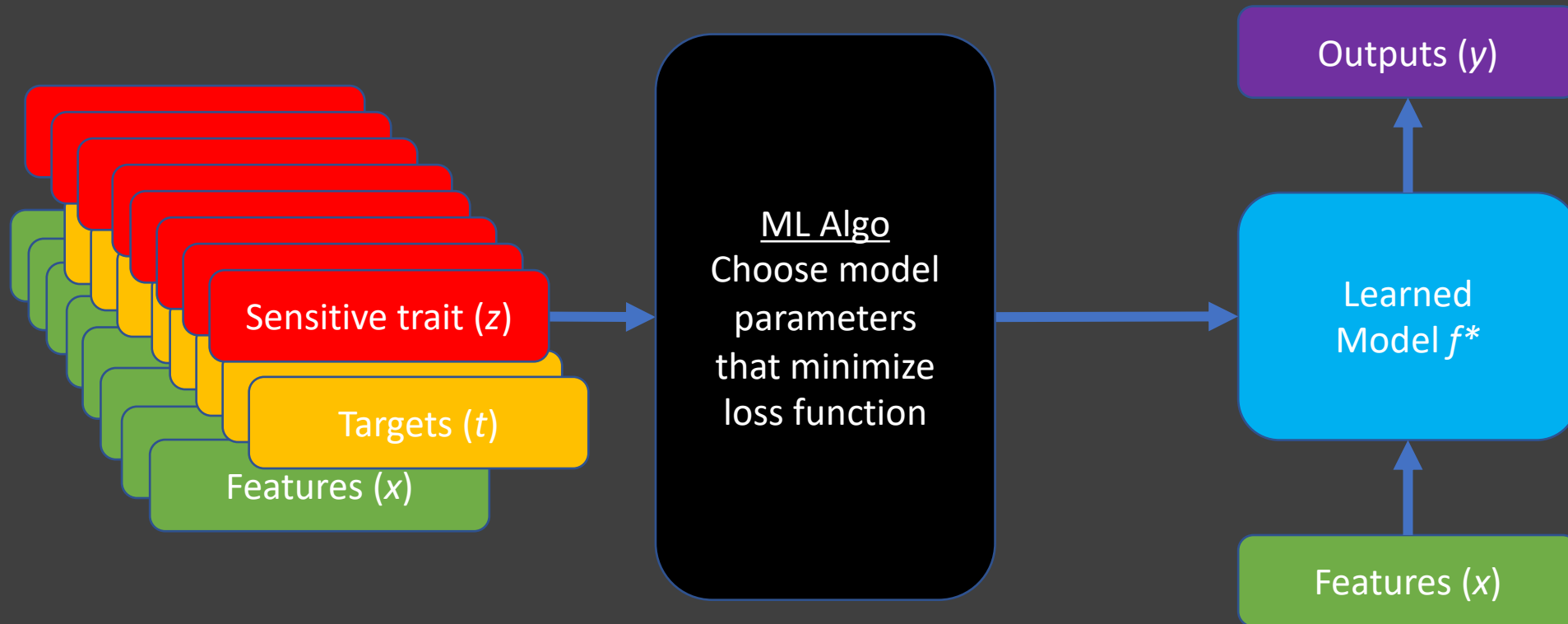
## Complications

- All of our features are correlated.
- And many subject to measurement error.

# Formal setup

- $\mathbf{x} \in R^d$  — Features characterizing an individual
- $z \in \{0, 1\}$  — Protected group membership
- $t \in \{0, 1\}$  — “Ground truth” target
- $f(\mathbf{x})$  — Prediction, an estimate of  $p(y|\mathbf{x})$
- $c(\mathbf{x}, z)$  — Decision, often by scoring/thresholding with  $f(\mathbf{x})$

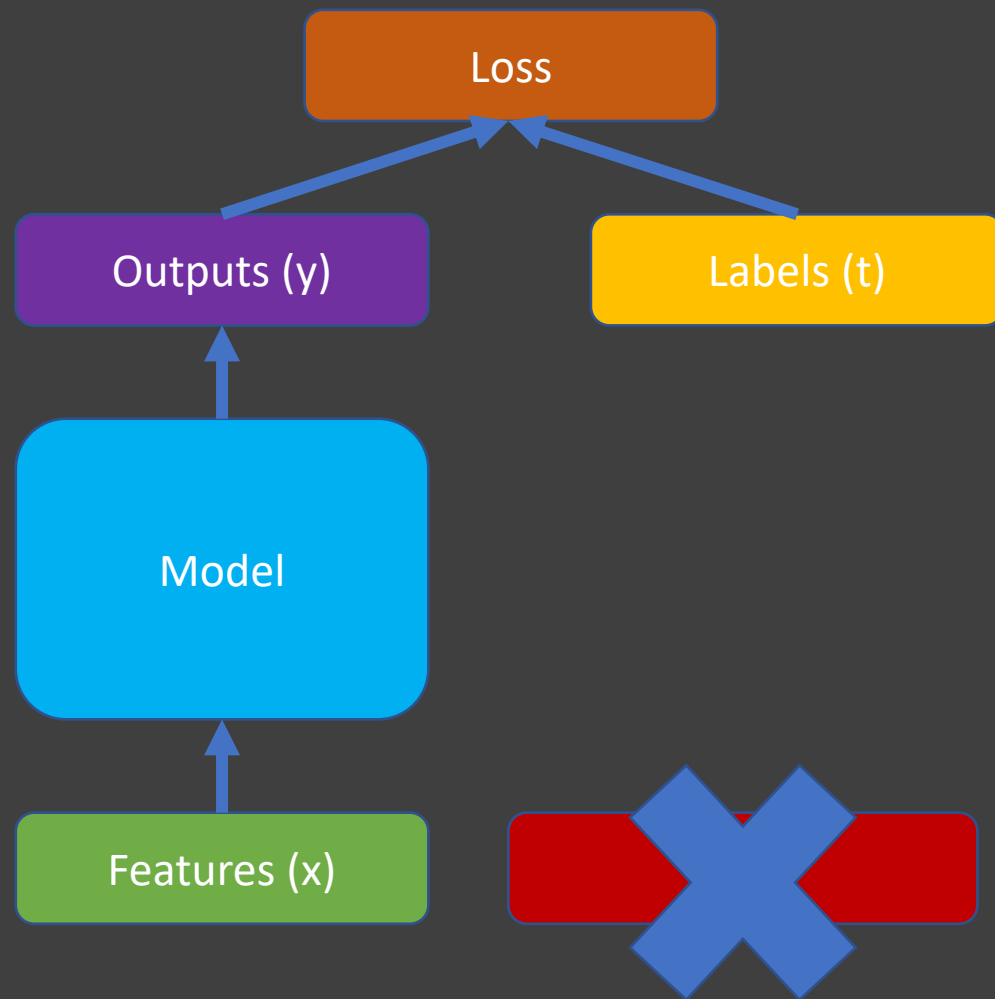
# Fair supervised learning



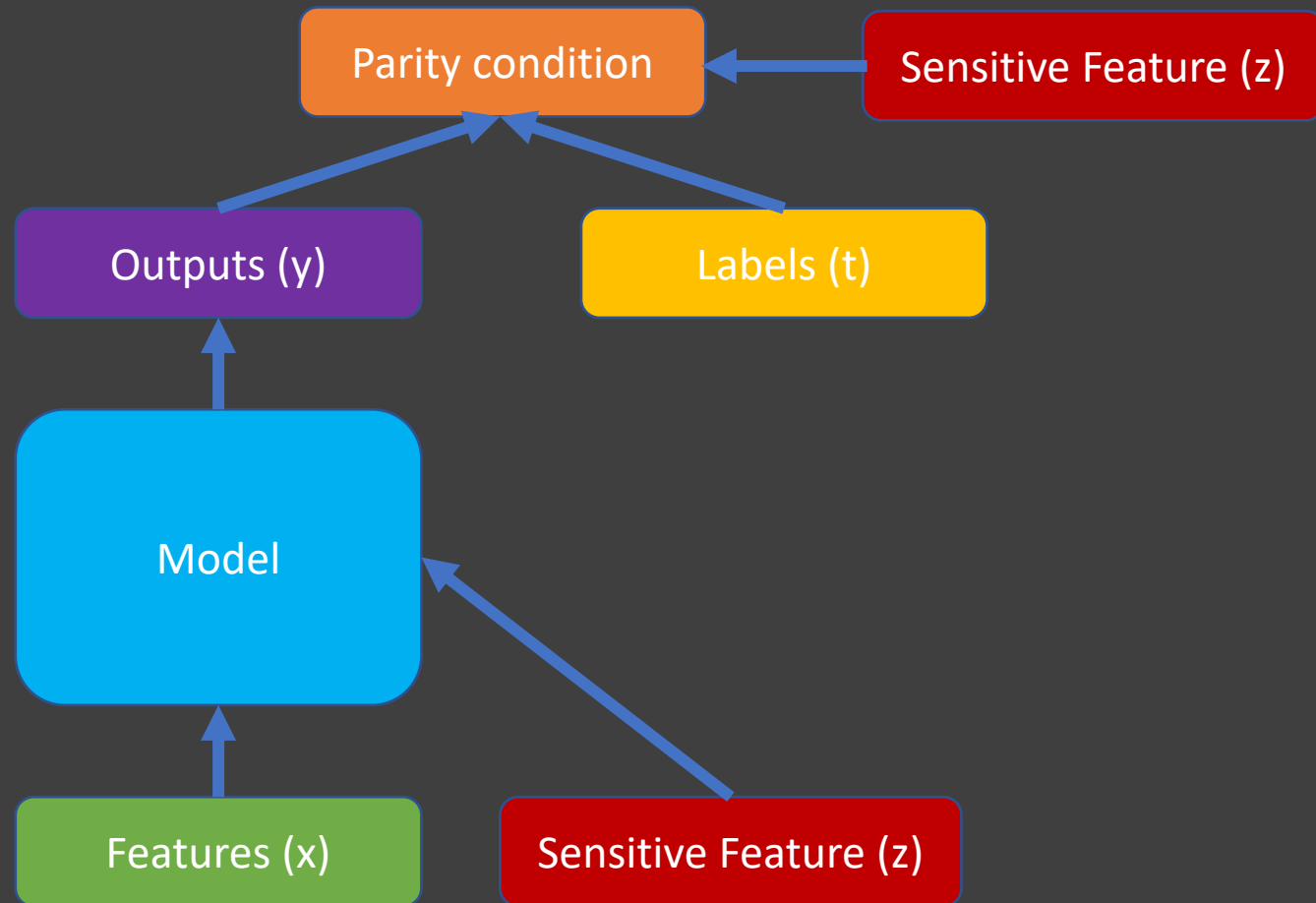
# Make groups equal but how?

- Impact parity
  - Outcome independent of group status  $y \perp z$
- Treatment parity
  - The output  $y$  depends only on  $x$ , not on  $z$
- Representational parity
  - Map  $x$  to  $r(x)$  such that  $r(x) \perp Z$
  - Entails impact parity
- Calibration:
  - Independence of truth and demographic for predicted value —  $(T \perp Z \mid Y)$
- Equalized Odds / “Opportunity” parity
  - Equal false negative and/or false positive rates

# Treatment parity / blindness



# Impact Parity / equal outcomes





# Impossibility Theorems

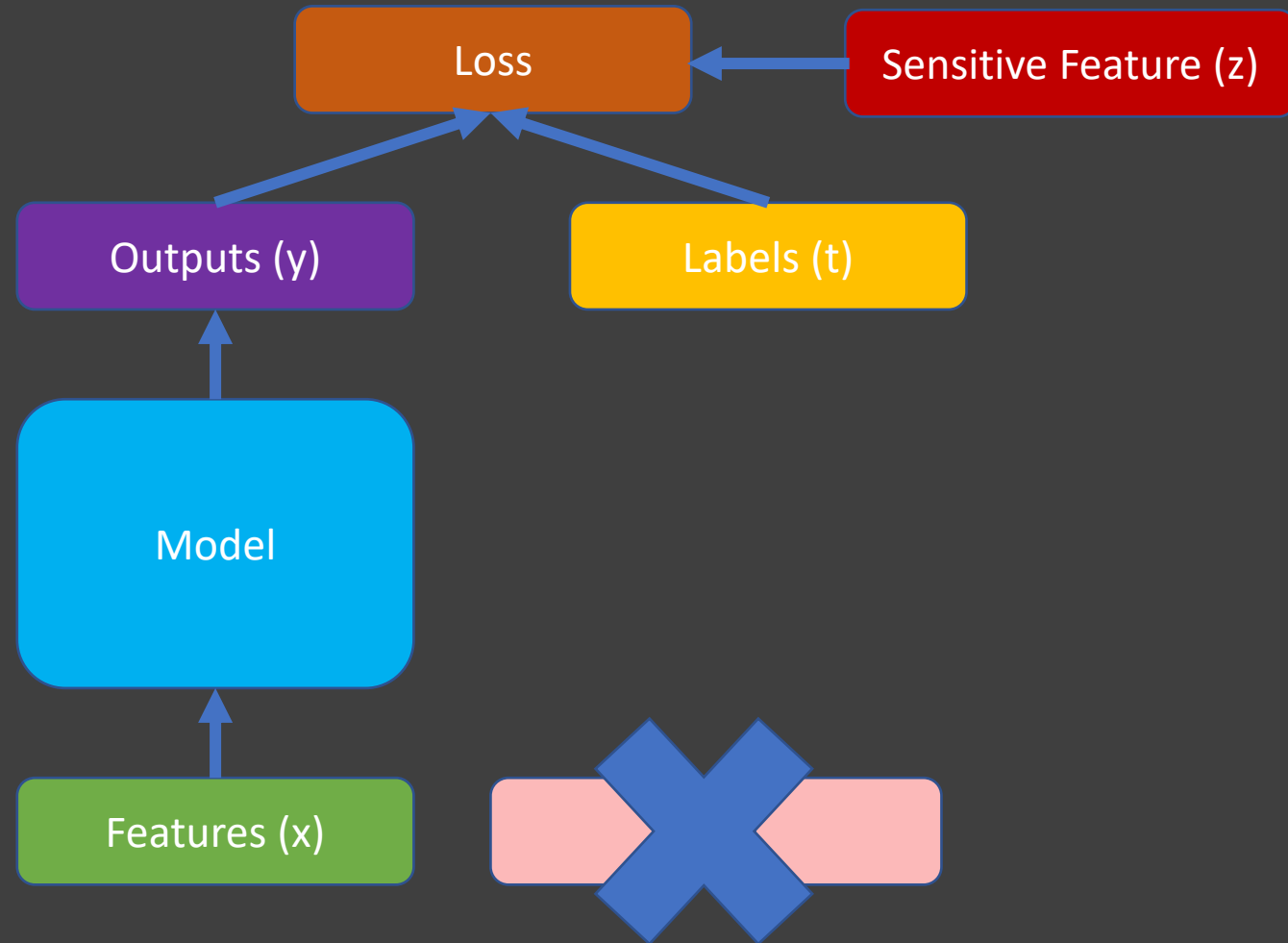
- The following 3 conditions cannot (in general) hold simultaneously:
  - Demographic parity ( $Y \perp Z$ )
  - Separation ( $Y \perp Z \mid T$ )
  - Calibration ( $T \perp Z \mid Y$ )
- Characterized by
  - [Chouldechova \(2016\)](#)
  - [Kleinberg, Mullainathan, Raghavan \(2016\)](#).
- Trade-offs among parities unavoidable.

# Limitations & Dangers

# Technical vs legal (vs ethics) terminology



# Disparate Learning Processes



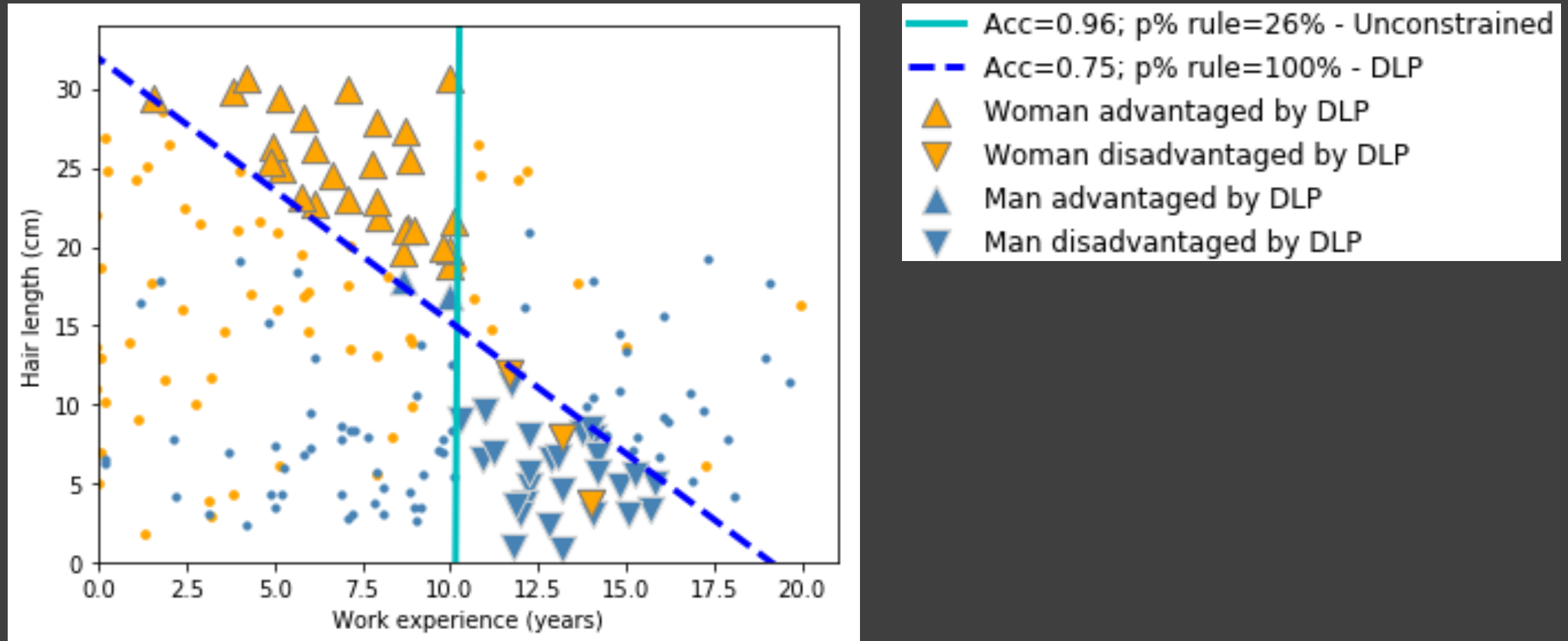
# Problems

- If all groups are the same in every way, easy
- Otherwise various parities are mutually irreconcilable
- Statistical parities don't capture legal /philosophical notions
- Do not address whether decisions are justified
- Lacks **even the ingredients** required to determine just action:
  - How did the data come to be / did disparities arise?
  - What are the impacts of decisions?
  - What are responsibilities of the decision-maker?

# Findings

1. For reconciling impact disparity and treatment disparity, **treatment disparity is optimal** (theoretical)
2. When  $x$  fully encodes  $z$ , for sufficiently powerful model, **DLP indistinguishable from treatment disparity** (theoretical)
3. When  $x$  partially encodes  $z$ , DLP results in side effects (empirical)
  - A. Re-orders within-group based on otherwise irrelevant characteristics
  - B. Produces potentially bizarre incentive to conform to stereotype

# Toy example

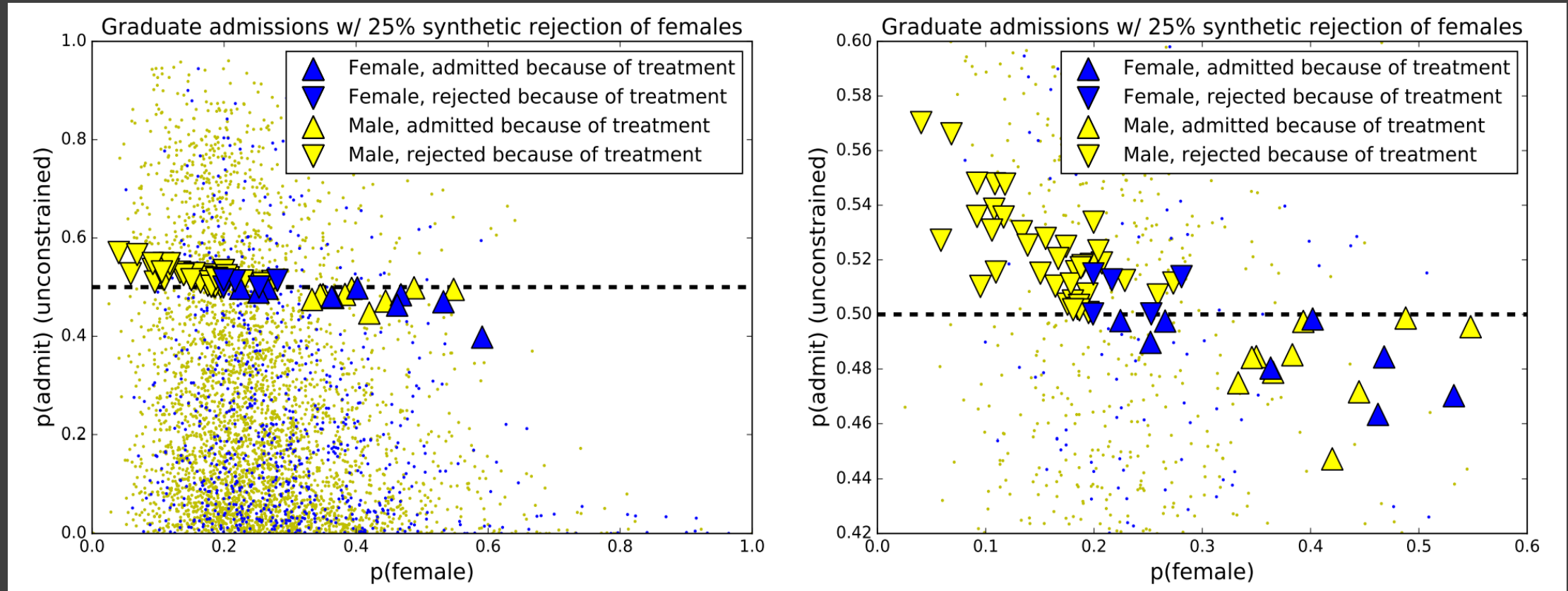


# Case study: Gender bias in CS admissions

- **Dataset:** sample of ~9,000 students considered for admission to the MS program of a large US university over an 11-year period
- **Labels:** admissions decisions provided by a faculty admissions committee
- **Attributes:** **Gender** the **protected attribute**. Country of origin, interest area, and GRE, etc. are used as features
- **Synthetic discrimination:** applied to mimic biased training data: of all women who were admitted, we flip 25% of their labels to 0



# Effects of DLP in CS admissions



# Algorithmic fairness from a non-ideal perspective

(work w. Sina Fazelpour)

<https://arxiv.org/abs/2001.09773>

# Ideal and Non-Ideal Theorizing about Justice

- The ideal approach:
  - Imagine a perfectly just world.
  - Try to minimize discrepancy between our world and the ideal.
  - Has been used to argue against affirmative action—*ideal world is color-blind*
- The non-ideal approach:
  - *[Non-ideal theorists] ... seek a **causal explanation** of the problem to determine what can and ought to be done about it, and who should be charged with correcting it. This requires an evaluation of the mechanisms causing the problem, as well as responsibilities of different agents to alter these mechanisms.*
    - Elizabeth Anderson 2019 The imperative of integration

# Solutions or Solutionism?

- From the perspective of stakeholders caught in the tension between (i) the potential profit to be gained from deploying machine learning in socially-consequential domains, and (ii) the increased scrutiny of a public concerned with algorithmic harms, these metrics offer an alluring solution: continue to deploy machine learning systems per the status quo, but use some chosen parity metric to claim a certificate of fairness, seemingly inoculating the actor against claims that they have not taken the moral concerns seriously, and weaponizing the half-baked tools produced by academics in the early stages of formalizing fairness as a shield against criticism.

# A new perspective on impossibility theorems?

- One potential contribution of ML impossibility theorems to philosophy is that they make evident an often-overlooked shortcoming with the ideal approach. These impossibility results make clear that in general, if we start from a non-ideal world, no set of actions (by a single agent) can instantaneously achieve the ideal world in every respect. Moreover, matching the ideal in a particular respect, may only be possible at the expense of widening gaps in others. Thus this naive form of an ideal approach appears to be fundamentally underspecified. If matching the ideal in various respects simultaneously is impossible, then we require, in addition to an ideal, a basis for deciding which among competing discrepancies to focus on. In this manner, the impossibility results in fair ML provide a novel lens to approach the philosophical debate about the extent to which normative theorizing on matters of justice can proceed in isolation from empirical sociohistorical facts.

# Or... an old perspective on impossibility theorems

Many other problems of applied equity follow a similar pattern. What seems simple at first turns out to be riddled with puzzles and contradictions. Inevitably, we must turn to logical analysis to sort them out. The study of equity turns out, therefore, to have close ties with the axiomatic method in mathematics. From simple and intuitively plausible propositions about the meaning of equity, one draws general and sometimes surprising conclusions about the form that an equitable rule must take.

The axiomatic method has two weaknesses however. The first is that, while each axiom seems reasonable by itself, when piled on top of one another they almost inevitably lead to “impossibility” theorems. This confirms the skeptic’s predisposition to believe that the problem had no solution anyway. The proper conclusion, however, is that not all desirable conditions can be satisfied simultaneously. Some choice must be made. A second difficulty with the axiomatic method is that it can easily become disengaged from the problem that it was intended to solve. The invention of axioms and conditions is a fascinating business. The danger is that the exercise can take on a life of its own and lead to results that are mathematically elegant, but that have little or no relation to the realities of the underlying situation. To guard against this tendency I have tried to mix formal definitions and theorems with informal arguments and examples

weaknesses  
of axiomatic  
method  
interpreting  
“impossibility”  
theorems

When should ML be off the table altogether?

# Original Works and Collaborators

- *Does Mitigating ML's Impact Disparity Require Treatment Disparity?*  
<https://arxiv.org/abs/1711.07076> (NeurIPS 2018)
- *Algorithmic Fairness from a Non-Ideal Perspective*  
<http://zacklipton.com/media/papers/fairness-non-ideal-fazelpour-lipton-2020.pdf> (AIES 2020)
- The Mythos of Model Interpretability  
(<https://arxiv.org/abs/1606.03490>) CACM 2018 (& ICML WHI workshop 2016)





# Thanks!

- **Co-authors:**

Alex Chouldechova (CMU), Julian McAuley (UCSD),  
Sina Fazelpour (UBC/CMU)

- **Papers:**

- *Does Mitigating ML's Impact Disparity Require Treatment Disparity?*  
<https://arxiv.org/abs/1711.07076> (NeurIPS 2018)

- *Algorithmic Fairness from a Non-Ideal Perspective*  
<http://zacklipton.com/media/papers/fairness-non-ideal-fazelpour-lipton-2020.pdf>  
(AIES 2020)

- The Mythos of Model Interpretability  
(<https://arxiv.org/abs/1606.03490>) CACM 2018 (& ICML WHI workshop 2016)

- Stay in touch! [zlipton@cmu.edu](mailto:zlipton@cmu.edu)