



# Fairness

---

March 12, 2020

Material from Berkeley's CS 294: Fairness in Machine Learning (<https://fairmlclass.github.io>)

and

N[eur]IPS2017 tutorial (<https://vimeo.com/248490141>)

by Solon Barocas (Cornell)  
and Moritz Hardt (Berkeley)



**Massachusetts  
Institute of  
Technology**

# Bias in Optum's Algorithm to Predict Healthcare Utilization

## Racial bias in a medical algorithm favors white patients over sicker black patients



“... black patients who were ranked by the algorithm as equally as in need of extra care as white patients were much sicker: They collectively suffered from 48,772 additional chronic diseases.

By **Carolyn Y. Johnson**

Oct. 24, 2019 at 2:00 p.m. EDT

Scientists discovered racial bias in a widely used medical algorithm that predicts which patients will have complex health needs. (iStock)

# NASEM Committee on Science, Technology, and Law

March, 2018

---

- Blockchain and Distributed Trust
- **Artificial Intelligence and Decision-Making**
  - **Hank Greely, Stanford**
  - **Cherise Fanno Burdee, Pretrial Justice Institute**
  - **Matthew Lungren, Stanford**
  - **Peter Szolovits, MIT**
  - **Suresh Venkatasubramanian, U. Utah**
- Privacy and Informed Consent in an Era of Big Data
- Science Curriculum for Law School
- Emerging Issues in Science, Technology, and Law
- Using Litigation to Target Scientists
- Communicating Advances in the Life Sciences to a Skeptical Public

Co-Chairs:

David Baltimore, Caltech  
David S. Tatel, U.S. Court  
of Appeals for the District  
of Columbia Circuit

# Algorithms and Justice

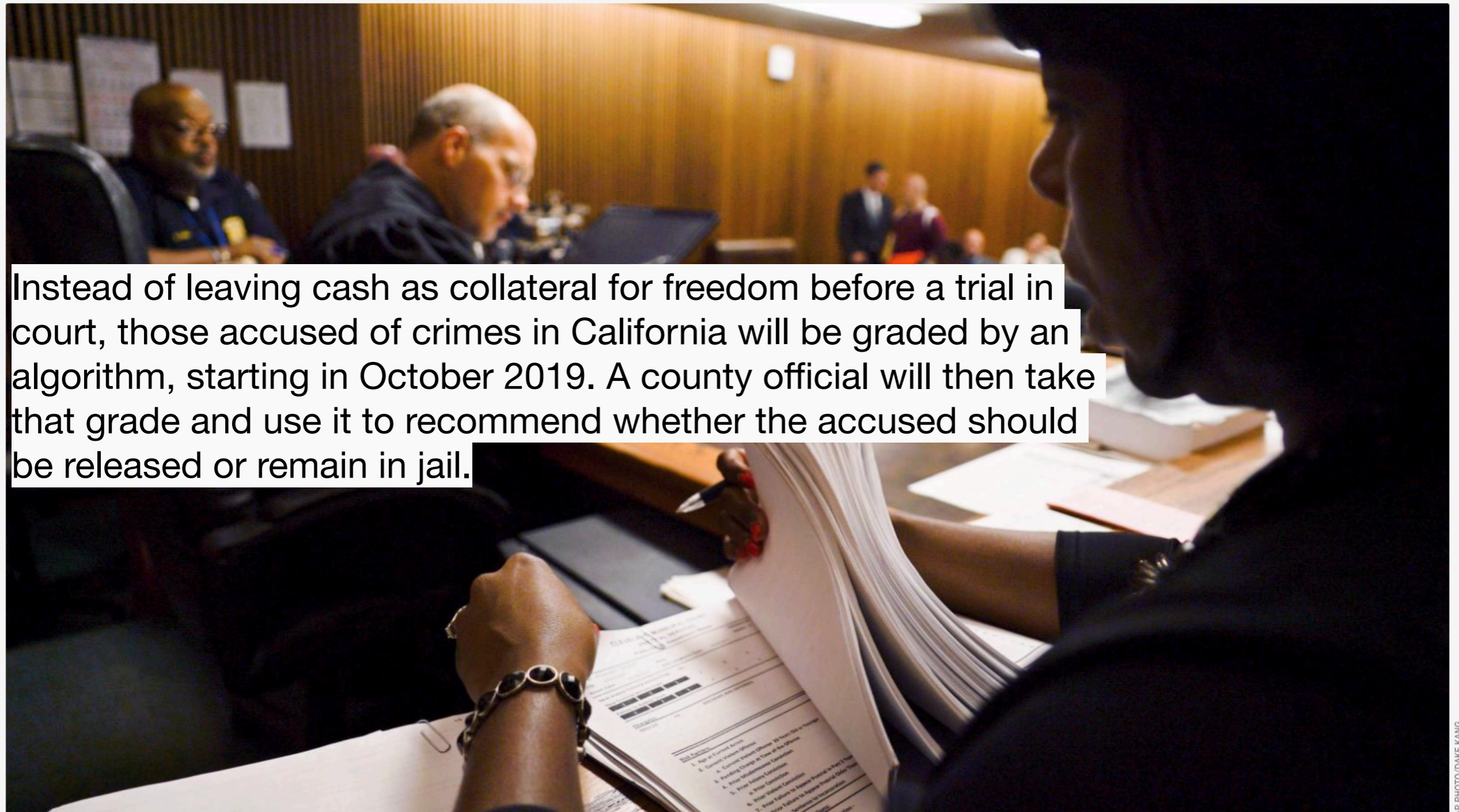
---

- Government use of decision automation for
  - determining eligibility for services
  - evaluating where to deploy health inspectors and law enforcement personnel
  - defining boundaries around voting districts
- In the law
  - “To the extent they inject clarity and precision into bail, parole, and sentencing decisions, algorithmic technologies may minimize harms that are the products of human judgment.”
  - “Conversely, the use of technology to determine whose liberty is deprived and on what terms raises significant concerns about transparency and interpretability.”



FRYING PAN, FIRE, ETC

# California just replaced cash bail with algorithms

By [Dave Gershgorn](#) · September 4, 2018

Instead of leaving cash as collateral for freedom before a trial in court, those accused of crimes in California will be graded by an algorithm, starting in October 2019. A county official will then take that grade and use it to recommend whether the accused should be released or remain in jail.

A probation officer sorts through automated risk scores in Cleveland.

# Critique of Bail Algorithms

---

- “... the machine learning systems used to calculate these risk scores throughout the criminal justice system, have been shown to hold severe racial biases, scoring people of color more likely to commit future crimes.”
- “Furthermore, since private companies have been typically contracted to offer these services, the formulas derived by machine learning algorithms to calculate these scores are generally withheld as intellectually property that would tip competitors to the company’s technology.”
- “... you have data collection that’s flawed with a lot of the same biases as the criminal justice system.”

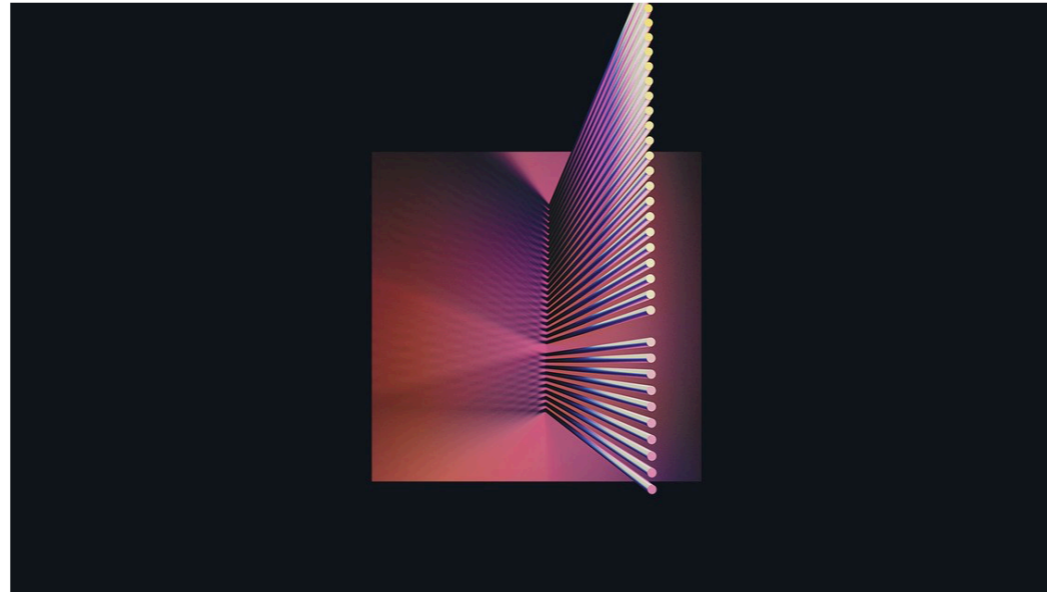
# COURTS ARE USING AI TO SENTENCE CRIMINALS. THAT MUST STOP NOW



- In the case of *Wisconsin v. Loomis*, defendant Eric Loomis was found guilty for his role in a drive-by shooting.
- During intake, Loomis answered a series of questions that were then entered into Compas, a risk-assessment tool developed by a privately held company and used by the Wisconsin Department of Corrections.
- The trial judge gave Loomis a long sentence partially because of the "high risk" score the defendant received from this black box risk-assessment tool.
- Loomis challenged his sentence, because he was not allowed to assess the algorithm.
- Last summer, the state supreme court ruled against Loomis, reasoning that **knowledge of the algorithm's output was a sufficient level of transparency.**



# HOW ALGORITHMS COULD HELP KEEP PEOPLE OUT OF JAIL



MARIO HUGO

- “... trying to use data to keep low-level offenders out of jail, figure out who needs psychiatric help, and even set bail and parole. In the same way that law enforcement uses data to deploy resources—so-called predictive policing—cities are using techniques borrowed from public health and machine learning to figure out what to do with people after they get arrested”



# Can an Algorithm Hire Better Than a Human?

---



**Claire Cain Miller** @clairecm JUNE 25, 2015

---

Hiring and recruiting might seem like some of the least likely jobs to be automated. The whole process seems to need human skills that computers lack, like making conversation and reading social cues.

But people have biases and predilections. They make hiring decisions, often unconsciously, based on similarities that have nothing to do with the job requirements — like whether an applicant has a friend in common, went to the same school or likes the same sports.



11/20/2015

## Why Machines Discriminate—and How to Fix Them

🕒 27:50 minutes



- Some believers in big data have claimed that, in big data sets, “the numbers speak for themselves.” Or in other words, the more data available to them, the closer machines can get to achieving objectivity in their decision-making. But data researcher Kate Crawford says that’s not always the case, because big data sets can perpetuate the same biases present in our culture, teaching machines to discriminate when scanning resumes or approving loans, for example.
- And when algorithms do discriminate, computer scientist Suresh Venkatasubramanian says he tends to hear expressions of disbelief, such as, “Algorithms are just code—they only do what you tell them.” But the decisions that machine-learning algorithms spit out are a lot more complicated and opaque than people think, he says, which makes tracking down an offending line of code a near impossibility.

# What is Fairness?

---

- *your ideas...*

# COMPAS

## Correctional Offender Management Profiling for Alternative Sanctions

---

- Used on >1M offenders since 2000
- Predicts risk of subject committing a misdemeanor or felony within 2 years from 137 features, not including race
- But, FP for blacks was 44.9%, whites 23.5%, FN was 47.7% for whites, 28% for blacks in a Broward County study.
- But, COMPAS satisfies
  - predictive parity (likelihood of recidivism among high-risk offenders regardless of race)
  - very similar AUC for predicting recidivism among whites and blacks
  - provides a well-calibrated score regardless of race
- But recidivism rate among blacks in Broward is 51% vs. 39% for whites
- Human judgment on 7 features achieves nearly identical results
- LR or SVM/RBF on same 7 features also does; *so much for proprietary superiority*

J. Angwin, J. Larson, S. Mattu, L. Kirchner, "Machine bias: There's software used across the country to predict future criminals. And it's biased against blacks," ProPublica, 23 May 2016; [www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing](http://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing).

Dressel, J., & Farid, H. (2018). The accuracy, fairness, and limits of predicting recidivism. *Science Advances*, 4(1), eaao5580. <http://doi.org/10.1126/sciadv.aao5580>



## Bias, Technically

---


### Genetic Misdiagnoses and the Potential for Health Disparities

Arjun K. Manrai, Ph.D., Birgit H. Funke, Ph.D., Heidi L. Rehm, Ph.D., Morten S. Olesen, Ph.D., Bradley A. Maron, M.D., Peter Szolovits, Ph.D., David M. Margulies, M.D., Joseph Loscalzo, M.D., Ph.D., and Isaac S. Kohane, M.D., Ph.D.

- {Selection, Sampling, Reporting} bias
- Case of Hypertrophic Cardiomyopathy
  - ... risk stratification for hypertrophic cardiomyopathy has been enhanced by targeted genetic testing
  - Multiple patients, all of whom were of African or unspecified ancestry, received positive reports, with variants misclassified as pathogenic on the basis of the understanding at the time of testing.
  - Subsequently, all reported variants were re-categorized as benign.
  - The mutations that were most common in the general population were significantly more common among black Americans than among white Americans ( $P < 0.001$ ).
  - Simulations showed that the inclusion of even small numbers of black Americans in control cohorts probably would have prevented these misclassifications.

Article | [OPEN](#) | Published: 15 April 2019

# Genetic risk factors identified in populations of European descent do not improve the prediction of osteoporotic fracture and bone mineral density in Chinese populations

Yu-Mei Li , Cheng Peng, Ji-Gang Zhang, Wei Zhu, Chao Xu, Yong Lin, Xiao-Ying Fu, Qing Tian, Lei Zhang, Yang Xiang, Victor Sheng & Hong-Wen Deng 

*Scientific Reports* **9**, Article number: 6086 (2019) | [Download Citation](#) ↓

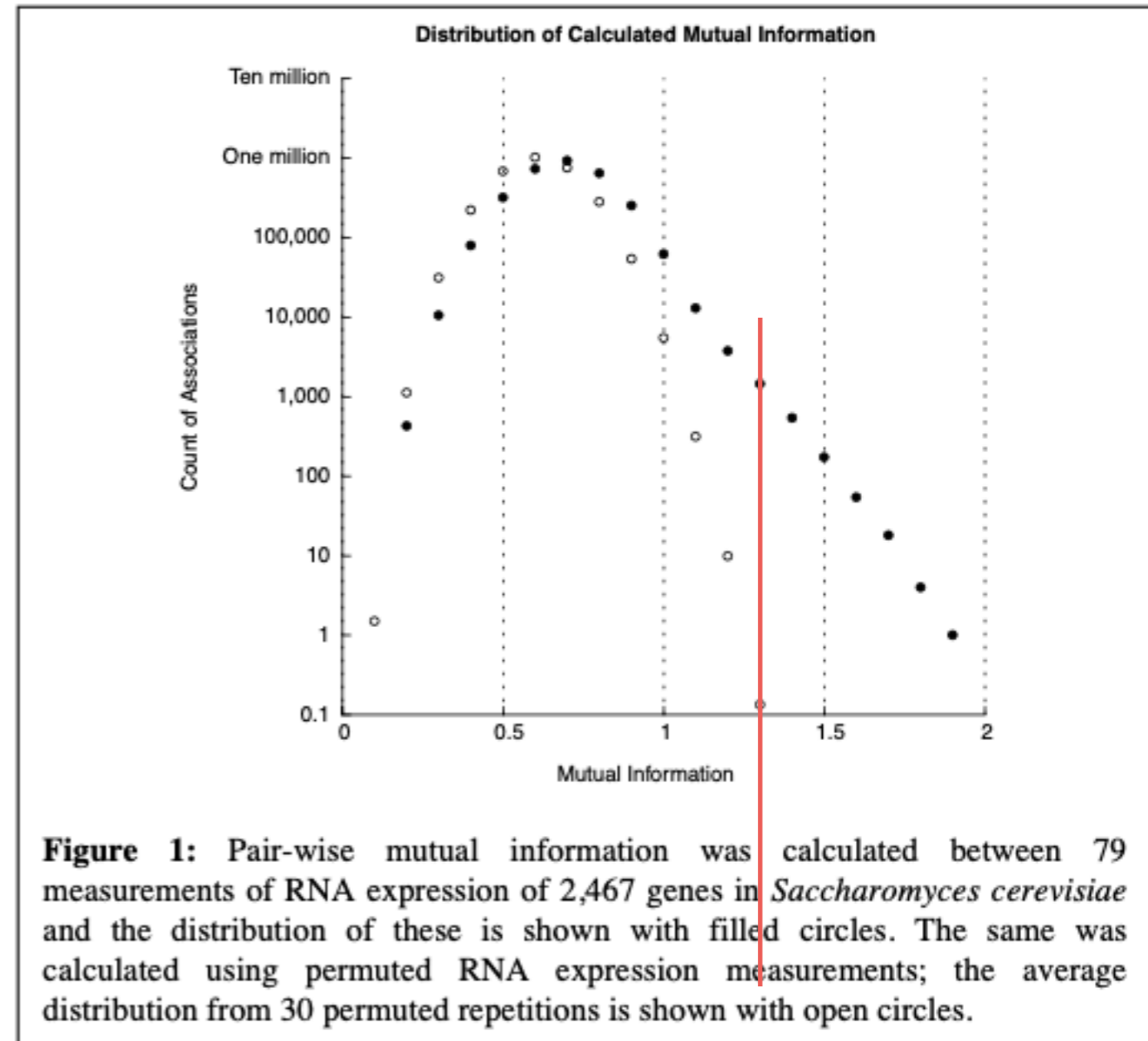
# Bias, Technically

---

- {Selection, Sampling, Reporting} bias
- Bias of an Estimator
  - Generally, we have bias, variance, and noise
  - $O$  = optimal possible model over all possible learners (model family)
  - $L$  = best model learnable by this learner
  - $A$  = actual model learned
  - Bias =  $O - L$  (limitation of learning method or target model)
  - Variance =  $L - A$  (error due to sampling of training cases)
    - Estimate significance by comparing against learning from randomly permuted data
- Inductive Bias — assumptions made by the learning algorithm about regularities that allow prediction on unseen cases

# Aside, on Permutation Testing of Significance

- If null hypothesis is that the input data is no better at predicting output than random chance
  - Permute the input data (to “disconnect” it from the outputs)
  - See if the results are the same
- E.g., Butte, A. J., & Kohane, I. S. (2000). Mutual information relevance networks: functional genomic clustering using pairwise entropy measurements. *Pacific Symposium on Biocomputing*, 418–429.





# Isn't Discrimination the Very Point of Machine Learning?

---

- *Unjustified* basis for differentiation
- Practical irrelevance
- Moral irrelevance
  
- Fairness focuses on ethical concerns
  
- Discrimination is
  - domain specific — how it influences people's life chances
  - feature specific — socially salient qualities that have served as the basis for unjustified and systematically adverse treatment in the past

# Regulated Domains

---

- Credit (Equal Credit Opportunity Act)
- Education (Civil Rights Act of 1964; Education Amendments of 1972)
- Employment (Civil Rights Act of 1964)
- Housing (Fair Housing Act)
- ‘Public Accommodation’ (Civil Rights Act of 1964)
- Marriage (Defense of Marriage Act, 1996, struck down by Supreme Court in 2013; also 1967 landmark civil rights case of Loving v. Virginia)
- Extends to marketing and advertising; not limited to final decision
- This list sets aside complex web of laws that regulates the government

# Legally recognized 'protected classes'

---

- Race (Civil Rights Act of 1964)
- Color (Civil Rights Act of 1964)
- Sex (Equal Pay Act of 1963; Civil Rights Act of 1964)
- Religion (Civil Rights Act of 1964)
- National origin (Civil Rights Act of 1964)
- Citizenship (Immigration Reform and Control Act)
- Age (Age Discrimination in Employment Act of 1967)
- Pregnancy (Pregnancy Discrimination Act)
- Familial status (Civil Rights Act of 1968)
- Disability status (Rehabilitation Act of 1973; Americans with Disabilities Act of 1990)
- Veteran status (Vietnam Era Veterans' Readjustment Assistance Act of 1974; Uniformed Services Employment and Reemployment Rights Act); Genetic information (Genetic Information Nondiscrimination Act)
- Sexual orientation (Mass SJC 2004, SCOTUS 2015)

# Two Doctrines of Discrimination Law

---

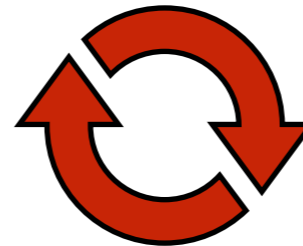
- Disparate Treatment
  - Formal — considering class membership
    - E.g., country club exclusion based on race or religion,
  - Intentional — without explicit reference to class, but with same effect
    - E.g., red-lining (mortgage availability based on geographic location)
- Disparate Impact
  - Unjustified, Avoidable
  - How to demonstrate: “4/5 rule” (20% difference establishes it)
  - How to defend: business necessity, job-related
  - Alternative practice: can we achieve the same goal but with less disparity?



# Goals of (Anti-)Discrimination Law

---

- Disparate Treatment
  - Procedural fairness
  - Equality of opportunity
- Disparate Impact
  - Distributive justice
  - Minimize inequality of outcome
- Non-discrimination:
  - ensuring that decision-making treats similar people similarly on the basis of relevant features, given their current degree of similarity
- Equality of opportunity:
  - organizing society in such a way that people of equal talents and ambition can achieve equal outcomes over the course of their lives
- Equality of outcome:
  - treat seemingly dissimilar people similarly, on the belief that their current dissimilarity is the result of past injustice



Conflict

E.g., affirmative action

# Discrimination Persists in Many Areas

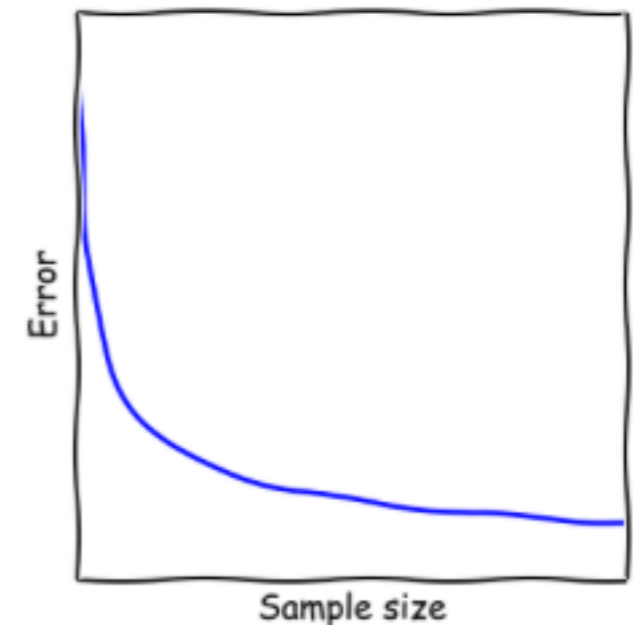
---

- Criminal justice — “Predictive Policing”
  - Police records measure “some complex interaction between criminality, policing strategy, and community-policing relations”
  - Future observations of crime confirm predictions
  - Fewer opportunities to observe crime that contradicts predictions
  - Initial bias may compound over time
- Housing
- Employment
- Health care
- ...

# Ongoing Problems

---

- Limited features
  - Features may be less informative or less reliably collected for certain parts of the population
  - A feature set that supports accurate predictions for the majority group may not for a minority group
  - Different models with the same reported accuracy can have a very different distribution of error across population
- Sample size disparity
- Leakage
  - With rich data, protected class membership will be unavoidably encoded across other features
  - No self-evident way to determine when a relevant attribute is too correlated with proscribed features



# Bias in Data

- E.g., Pasta consumption vs. BMI
- Three populations: slugs, normal, athletes

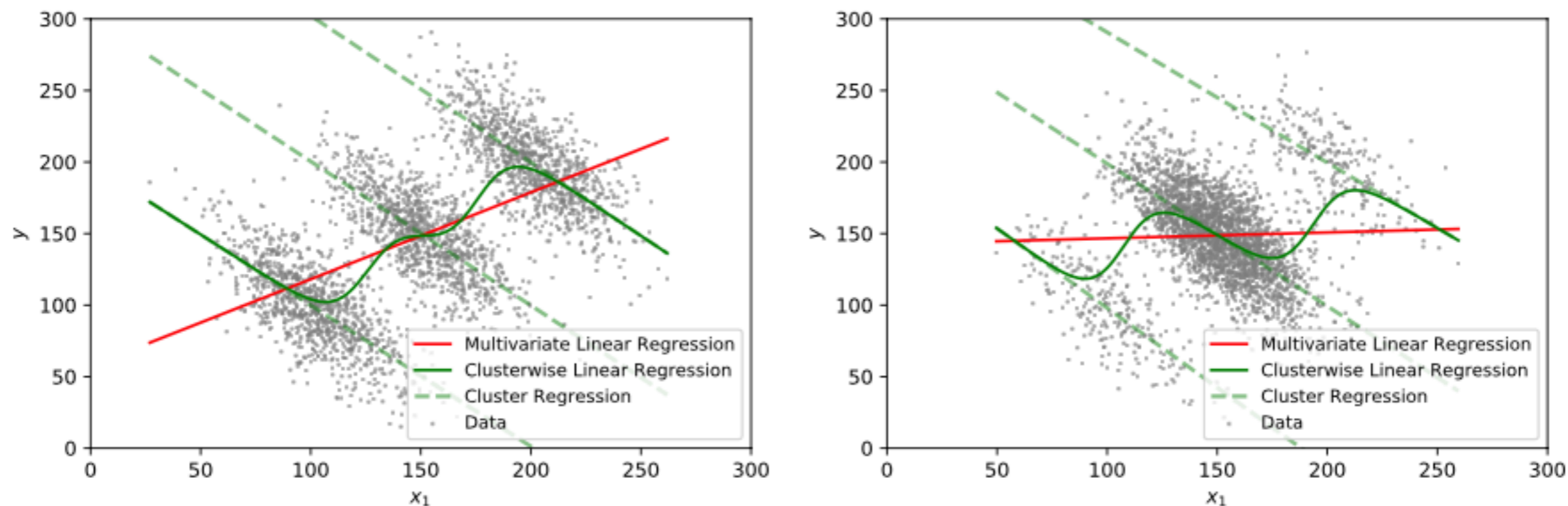


Fig. 1. Illustration of biases in data. Red line shows the regression (MLR) for the entire population, while dashed green lines are regressions for each subgroup, and the solid green line is the unbiased regression. (a) When all subgroups are of equal size, then MLR shows a positive relationship between the outcome and the independent variable. (b) Regression shows almost no relationship in less balanced data. The relationships between variables within each subgroup, however, remain the same. (Credit: Nazanin Alipourfard)

# Many Forms of Bias

---

- Historical
- Representation
- Measurement
- Evaluation
- Aggregation
- Population
- Simpson's Paradox
  - a trend, association, or characteristic observed in underlying subgroups may be quite different from association or characteristic observed when these subgroups are aggregated.
- Longitudinal Data Fallacy
- Sampling
- Behavioral
- Content Production
- Linking
- Temporal
- Popularity
- Algorithmic
- User Interaction/Presentation/Ranking
- Social
- Emergent
- Self-Selection
- Omitted Variable
- Cause-Effect
- Observer
- Funding

# Formalizing Fairness Discussion

---

- Hardt's example: advertising for a software engineer, question of gender bias

- Notation:

$$\mathbb{P}_a \{E\} = \mathbb{P}\{E \mid A=a\}$$

<b><math>X</math></b>	features of an individual (browsing history)
<b><math>A</math></b>	sensitive attribute (gender)
<b><math>R = r(\mathbf{X}, \mathbf{A})</math></b> <b><math>C = c(\mathbf{X}, \mathbf{A})</math></b>	score/predictor (show ad) [classify by thresholding score]
<b><math>Y</math></b>	hire software engineer



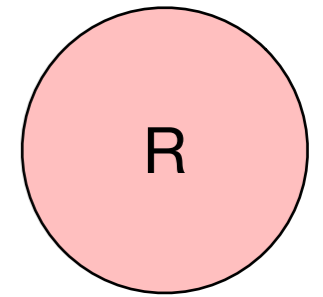
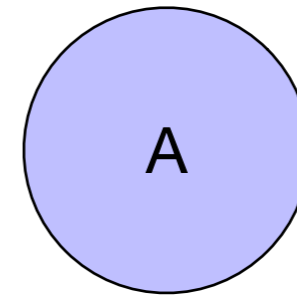
# Proposed Criteria of Fairness

---

- **Independence** of scoring function from sensitive attributes
  - $R \perp A$
- **Separation** of score and sensitive attribute given outcome
  - $R \perp A \mid Y$
- **Sufficiency**
  - $Y \perp A \mid R$

# Independence

$R \perp A$



- 
- Also called demographic parity, statistical parity, group fairness, disparate impact
  - $P\{R = 1 \mid A = a\} = P\{R = 1 \mid A = b\}$  for all groups A
  - thus, unfair if
    - $|P\{R = 1 \mid A = a\} - P\{R = 1 \mid A = b\}| > \epsilon$
    - $\left| \frac{P\{R = 1 \mid A = a\}}{P\{R = 1 \mid A = b\}} - 1 \right| \geq \epsilon$
    - $\epsilon = 0.2$  relates to 4/5 rule

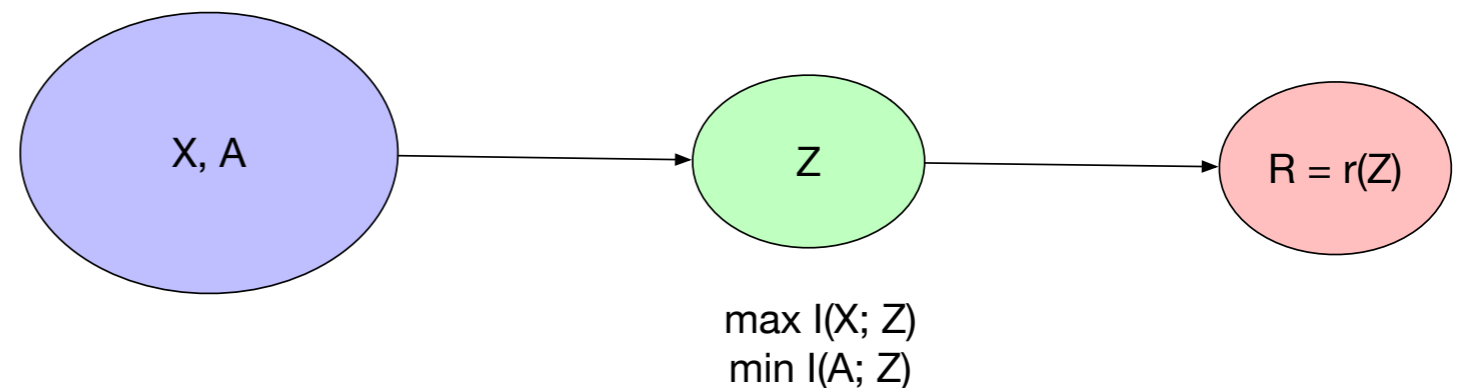
# Problems with Independence

---

- Only requires equal rates of decisions (hiring, liver transplants, etc.)
  - But, what if hiring is based on a good score in group a, but random in b, though with same probability?
  - Outcomes will (most likely) be better for group a, establishing problems for the future!
  - Could be caused by malice, or by better information about group a.
- What if A is a perfect predictor of Y?
  - ... or at least is strongly correlated?
  - How much are you willing to decrease the effectiveness of the predictor to achieve fairness?

# Potential Fixes to Achieve Independence

- Pre-processing:
  - Adjust the feature space to be uncorrelated with the sensitive attribute
    - Domain-specific
  - Representation learning



Zemel, R. S., Wu, Y., Swersky, K., Pitassi, T., & Dwork, C. (2013). Learning Fair Representations. ICML.

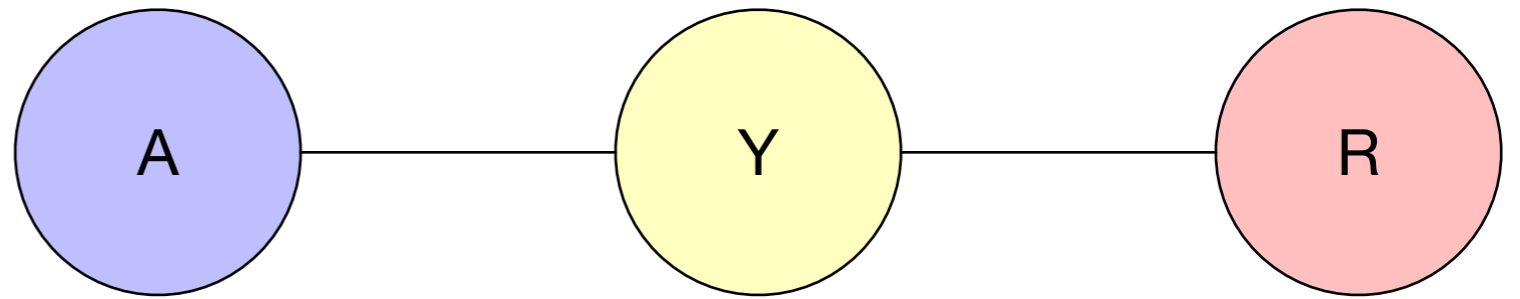
- Impose independence constraints at training time (for a given data set)  
E.g., include dependence in the loss function, differential sampling, ...

Calders, T., Kamiran, F., & Pechenizkiy, M. (2010). Building Classifiers with Independency Constraints (pp. 13–18). Presented at the 2009 IEEE International Conference on Data Mining Workshops (ICDMW), IEEE. <http://doi.org/10.1109/ICDMW.2009.83>

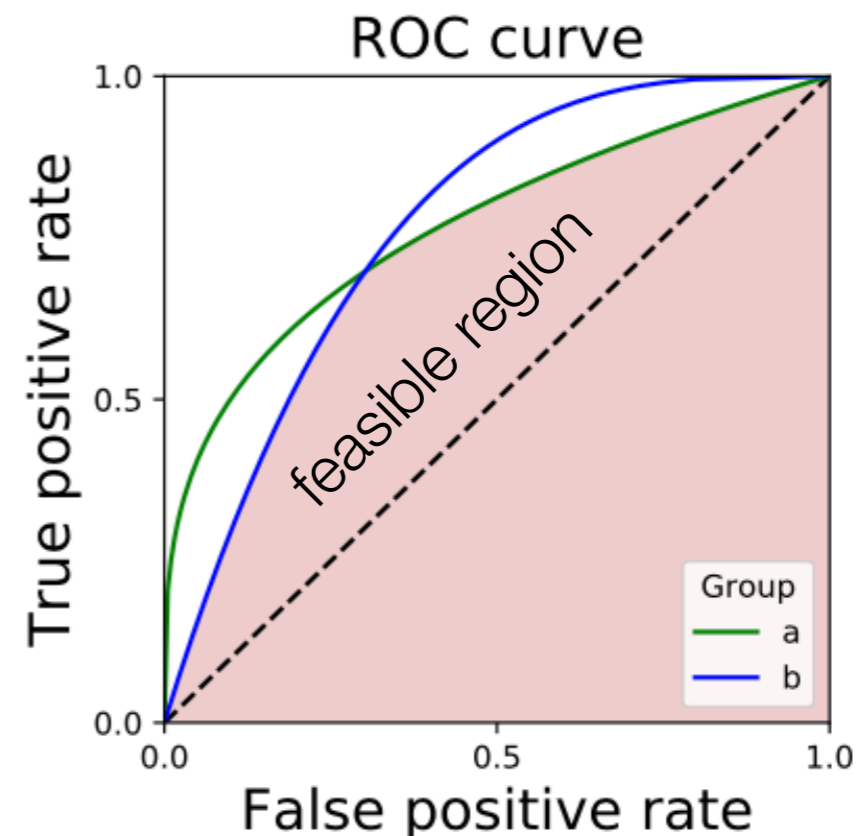
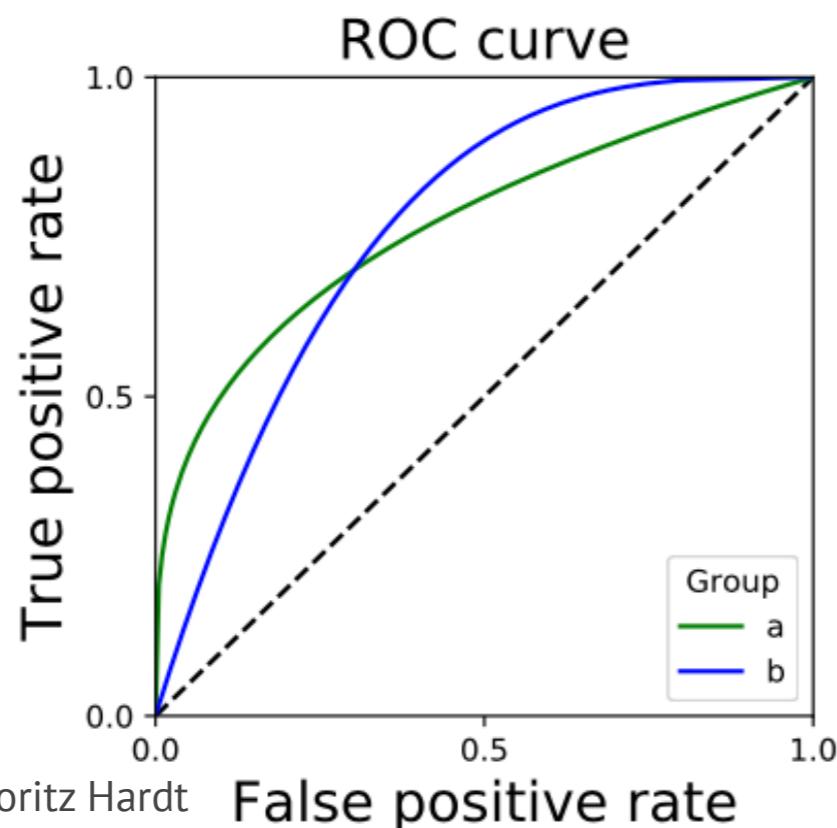
- Post-processing
  - Create a new classifier  $F$ ,  $\hat{Y} = F(R, A)$
  - minimize cost of misclassification, perhaps more strongly for protected  $A$

# Separation

$R \perp A | Y$



- Recognizes that A may be correlated with the target variable
  - E.g., different success rates in a drug trial for different ethnic populations
- $P\{R = 1 | Y = 1, A = a\} = P\{R = 1 | Y = 1, A = b\}$   
 $P\{R = 1 | Y = 0, A = a\} = P\{R = 1 | Y = 0, A = b\}$ 
  - i.e., true and false positive rates for both classes must be the same
- Can choose any true positive/false positive tradeoff in the feasible region, depending on relative costs



# Advantages of Separation over Independence

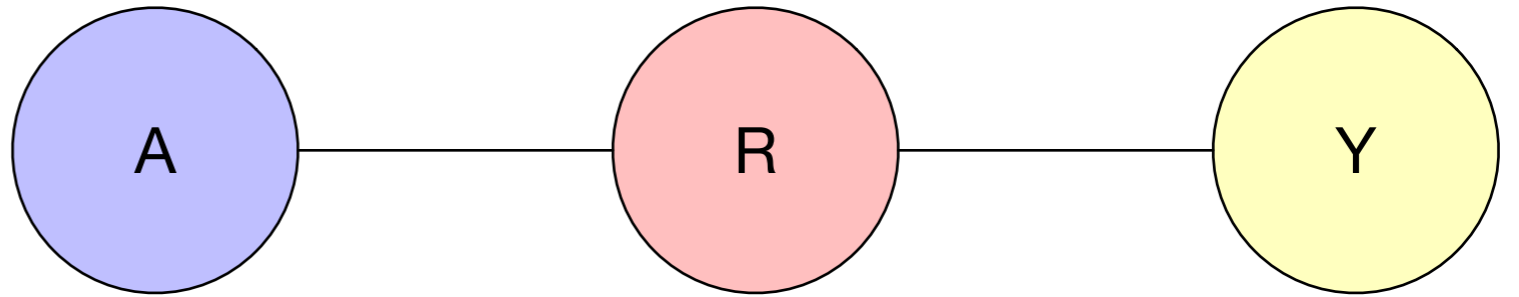
---

- Allows correlation between  $R$  and  $Y$  (even perfect predictor)
- Incentive to reduce errors uniformly in all groups



# Sufficiency

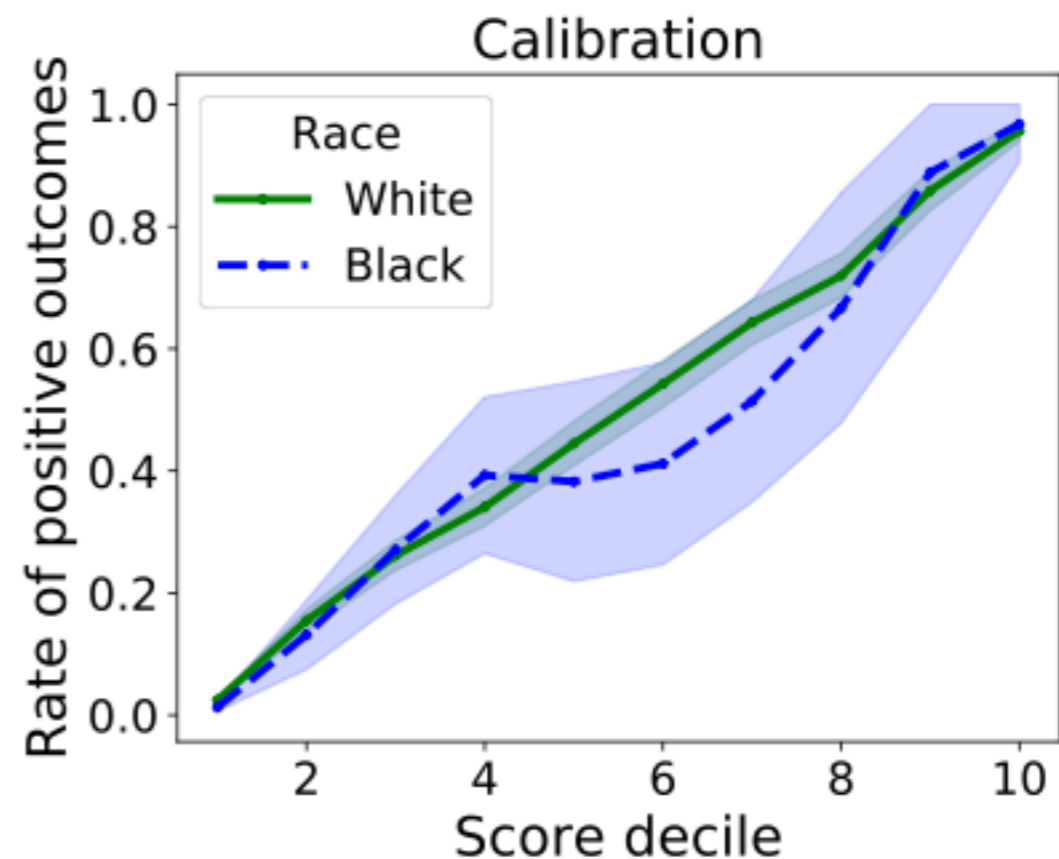
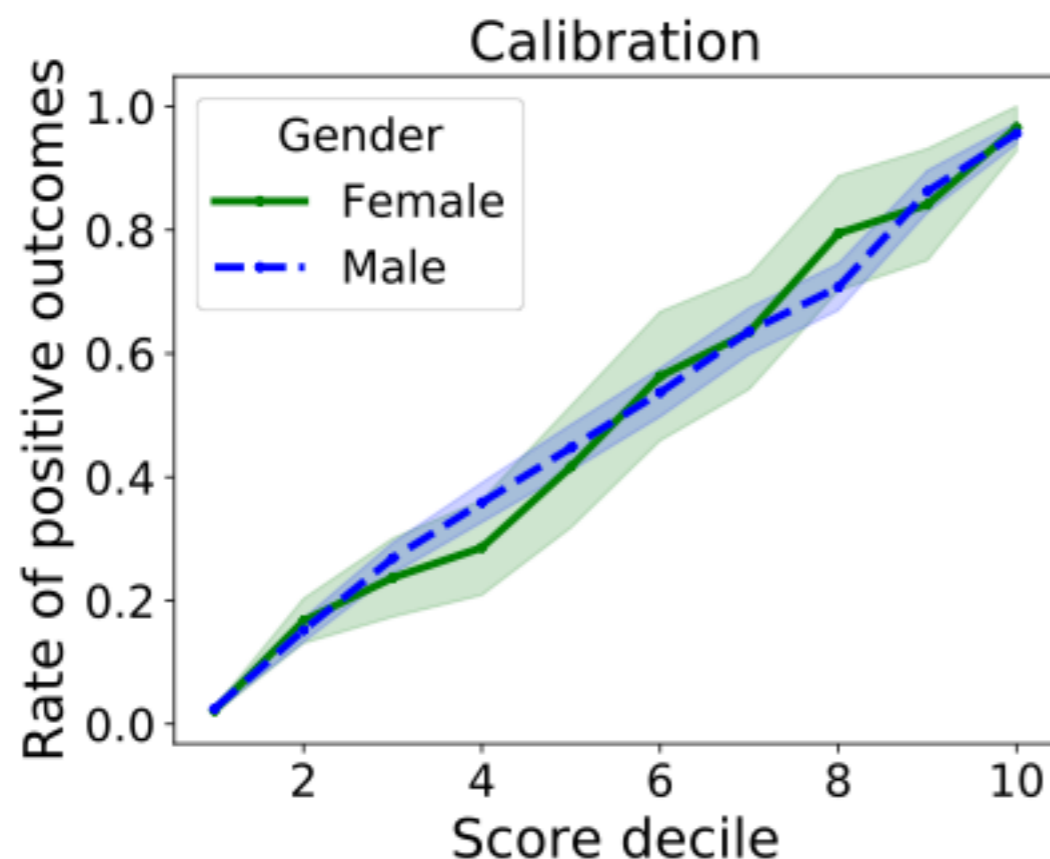
$Y \perp A \mid R$



- $P\{Y = 1 \mid R = r, A = a\} = P\{Y = 1 \mid R = r, A = b\}$
- Requires parity of positive and negative predictive values across groups
- R is *calibrated* if  $P\{Y = 1 \mid R = r, A = a\} = r$ 
  - I.e., if the scoring function is a probability of outcome, or
  - “the set of all instances assigned a score value  $r$  has an  $r$  fraction of positive instances among them”
- Can recalibrate a scoring function R by fitting a sigmoid
  - $S = \frac{1}{1 + e^{aR+b}}$
  - and optimizing log loss  $-\mathbb{E}[Y \log S + (1 - Y) \log(1 - S)]$
- Calibration by group implies sufficiency

# Calibration Can be Good Without Even Trying

- E.g., UCI census data set, predicting income  $> \$50,000/\text{year}$  for those over 16yo with some income
- Features (14): age, type of work, weight of sample, education, marital status, occupation, military service, race, sex, capital gain/loss, hours per week of work, native country, ...

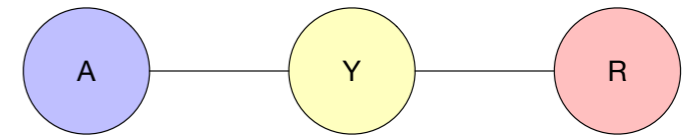


# Bad News!

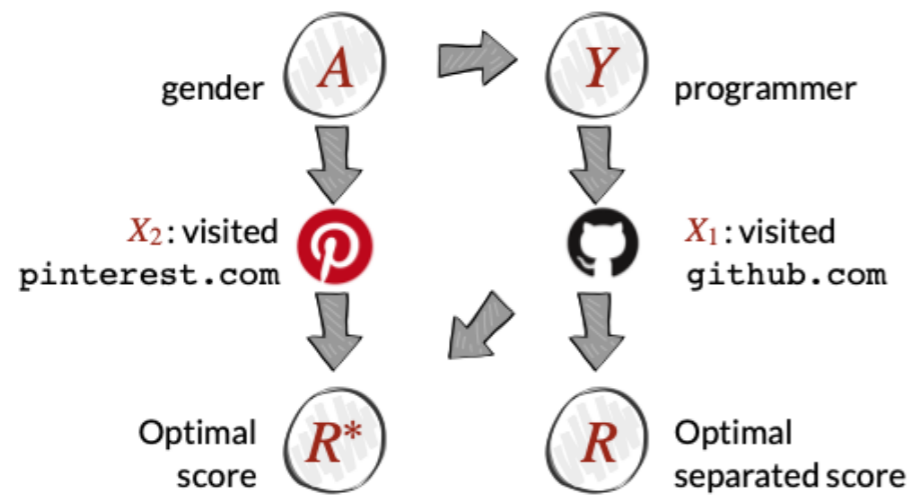
---

- It is not possible to jointly achieve any pair of these conditions
  - Independence *xor* Separation
  - Independence *xor* Sufficiency
  - Separation *xor* Sufficiency
- Nice illustration at
  - <https://research.google.com/bigpicture/attacking-discrimination-in-ml/>

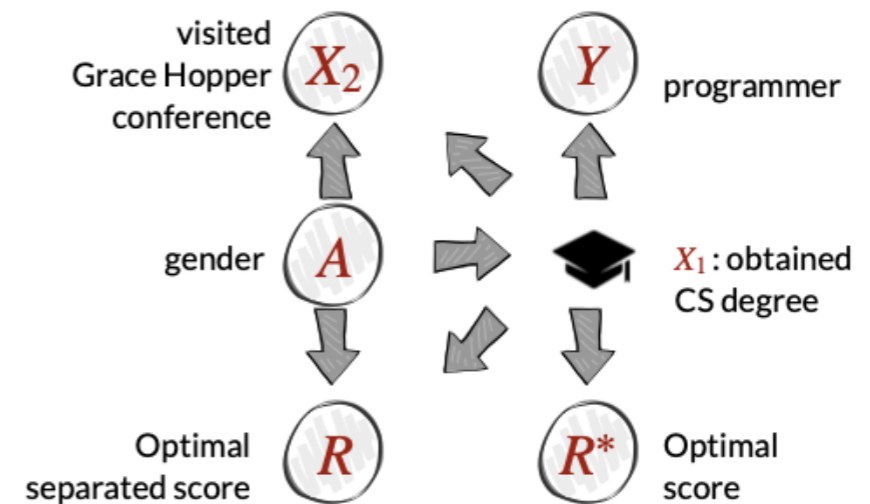
# Different Scenarios Can Lead to Same Observed Distributions



## Scenario I



## Scenario II



- The distributions of  $A$ ,  $R$ ,  $Y$ ,  $X_1$  and  $X_2$  can be identical in the two scenarios
- In Scenario II, gender is used directly to adjust separated score



ORIGINAL RESEARCH  
FEB 2019

## Can AI Help Reduce Disparities in General Medical and Mental Health Care?

Irene Y. Chen, Peter Szolovits, PhD, and Marzyeh Ghassemi, PhD

# Examined Error Rates in Two Data Sets

---

- Data: de-identified unstructured notes
  - MIMIC-III, predict ICU mortality
  - Psych inpatient data, predict 30-day psych readmission
- Is there bias, based on race, gender, insurance type (as proxy for socio-economic status)?
- Topic modeling on notes: 50 topics



# Interpreting Notes by Topic Modeling

Topic Name <sup>a</sup>	Characteristic Words
Cancer	Mass, cancer, metastatic
Heart flow	Afib, atrial, Coumadin <sup>®</sup> , fibrillation
Kidney	Renal, dialysis, ESRD, line
Orthopedic	Liver, cirrhosis, hepatic, ascites
Pulmonary	COPD, home, BiPAP, chronic
Substance abuse	EtOH, abuse, CIWA, withdrawal
Abbreviations: afib, atrial fibrillation; BiPAP, bilevel positive airway pressure; CIWA, Clinical Institute Withdrawal Assessment; COPD, chronic obstructive pulmonary disease; ESRD, end-stage renal disease; EtOH, ethanol.	
<sup>a</sup> Topic name was inferred based on algorithmically found top words.	
Anxiety	Anxiety, depression, disorder
Bipolar disorder	Bipolar, lithium, manic, episode
Chronic pain	Pain, chronic, mg
Depression	Depression, suicidal, depressive
Psychosis	Psychotic, psychosis, paranoia
Substance abuse	Use, substance, abuse, cocaine

# Psychiatry Results

---

- Race:
  - White patients had higher topic enrichment values for the **anxiety** and **chronic pain** topics
  - Black, Hispanic, and Asian patients had higher topic enrichment values for the **psychosis** topic
- Gender:
  - Male patients had higher topic enrichment values for **substance abuse** (0.024 v 0.015)
  - Female patients had higher topic enrichment values for **general depression** (0.021 v 0.019) and **treatment resistant depression** (0.025 v 0.015)
- Insurance:
  - private insurance patients have higher topic enrichment values than public insurance patients for **anxiety** (0.029 v 0.0156) and **general depression** (0.026 v 0.017)
  - public insurance patients have higher topic enrichment values for **substance abuse** (0.022 v 0.016)

# ICU Results

---

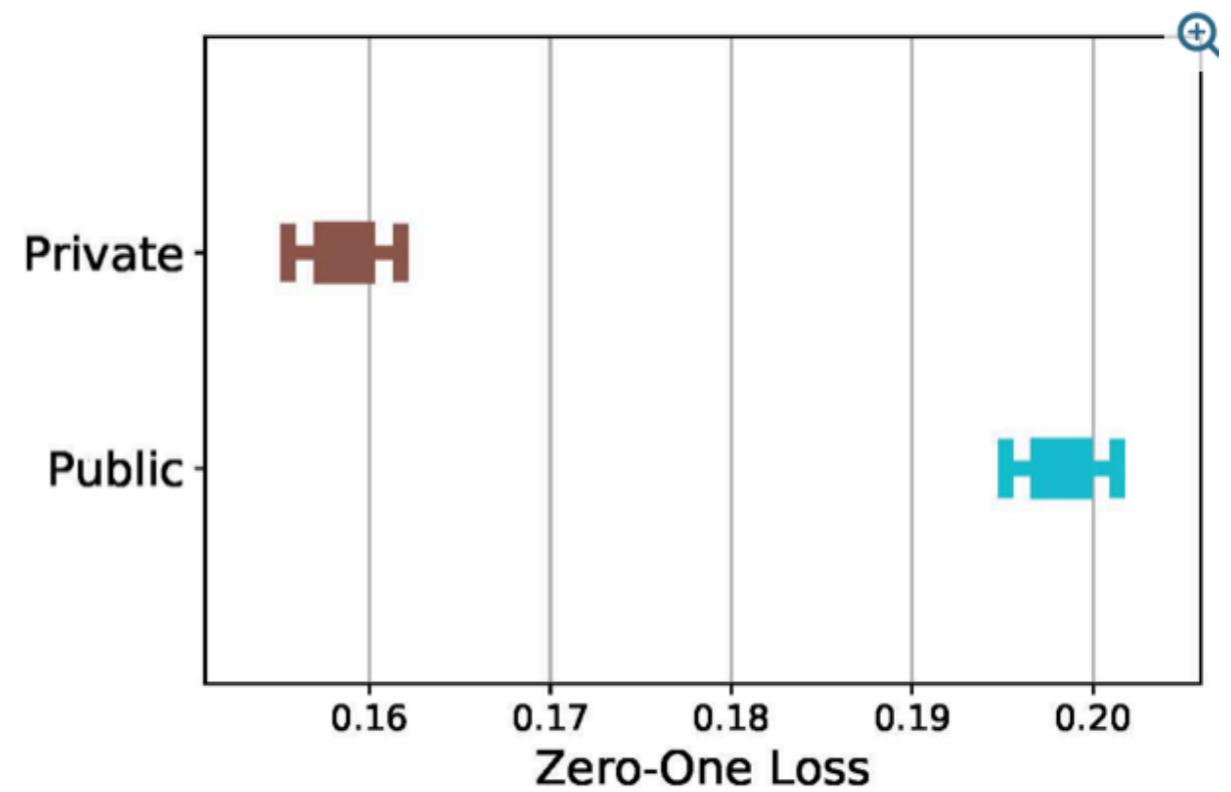
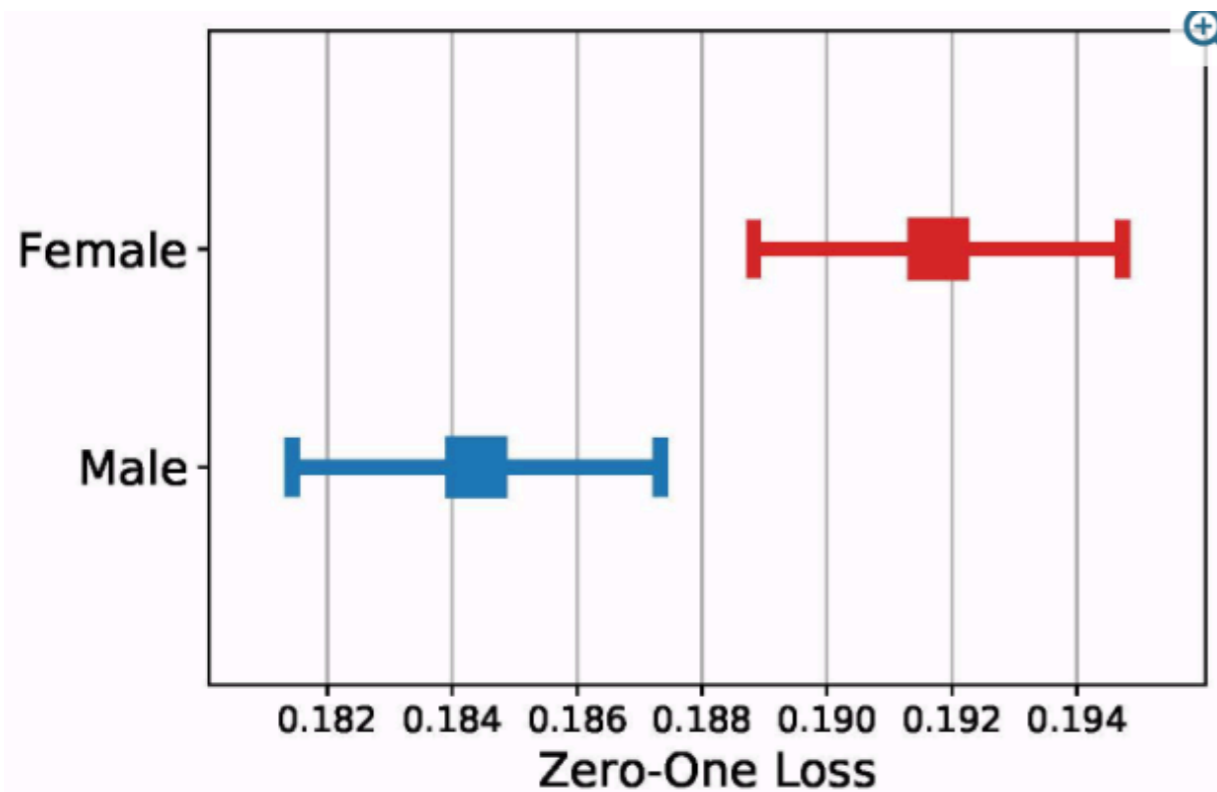
- Gender:
  - male patients have higher topic enrichment values for **substance use** (0.027 v 0.011)
  - female patients have higher topic enrichment values for **pulmonary disease** (0.026 v 0.016), potentially reflecting known underdiagnosis of chronic obstructive pulmonary disease in women
- Race:
  - Asian patients have the highest topic enrichment values for **cancer** (0.036), followed by white patients (0.021), other patients (0.016), and black and Hispanic patients (0.015)
  - Black patients have the highest topic enrichment values for **kidney problems** (0.061), followed by Hispanic patients (0.027), Asian patients (0.022), white patients (0.015), and other patients (0.014)
  - Hispanic patients have the highest topic enrichment values for **liver concerns** (0.034), followed by other patients (0.024), Asian patients (0.023), white patients (0.019), and black patients (0.014)
  - White patients have the highest topic enrichment values for **atrial fibrillation** (0.022), followed by other patients (0.017), Asian patients (0.015), black patients (0.013), and Hispanic patients (0.011)

# ICU Results, continued

---

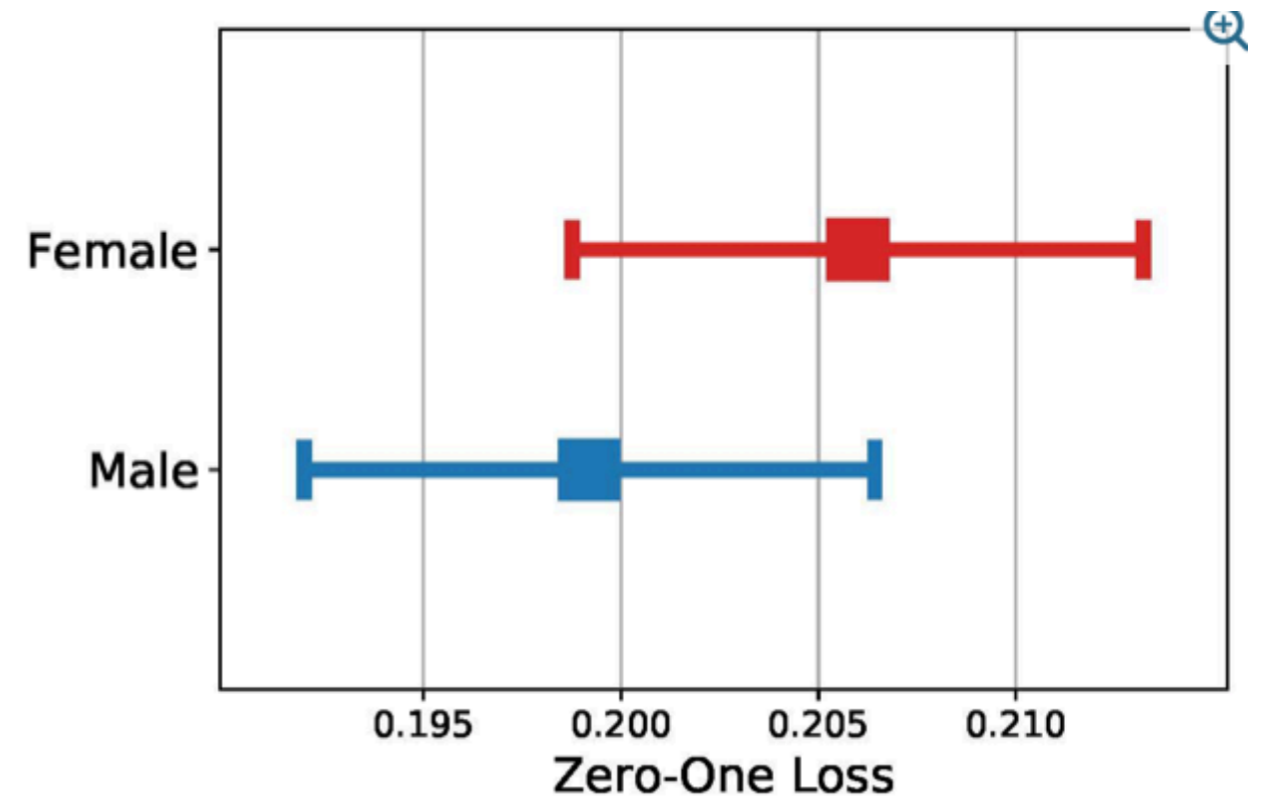
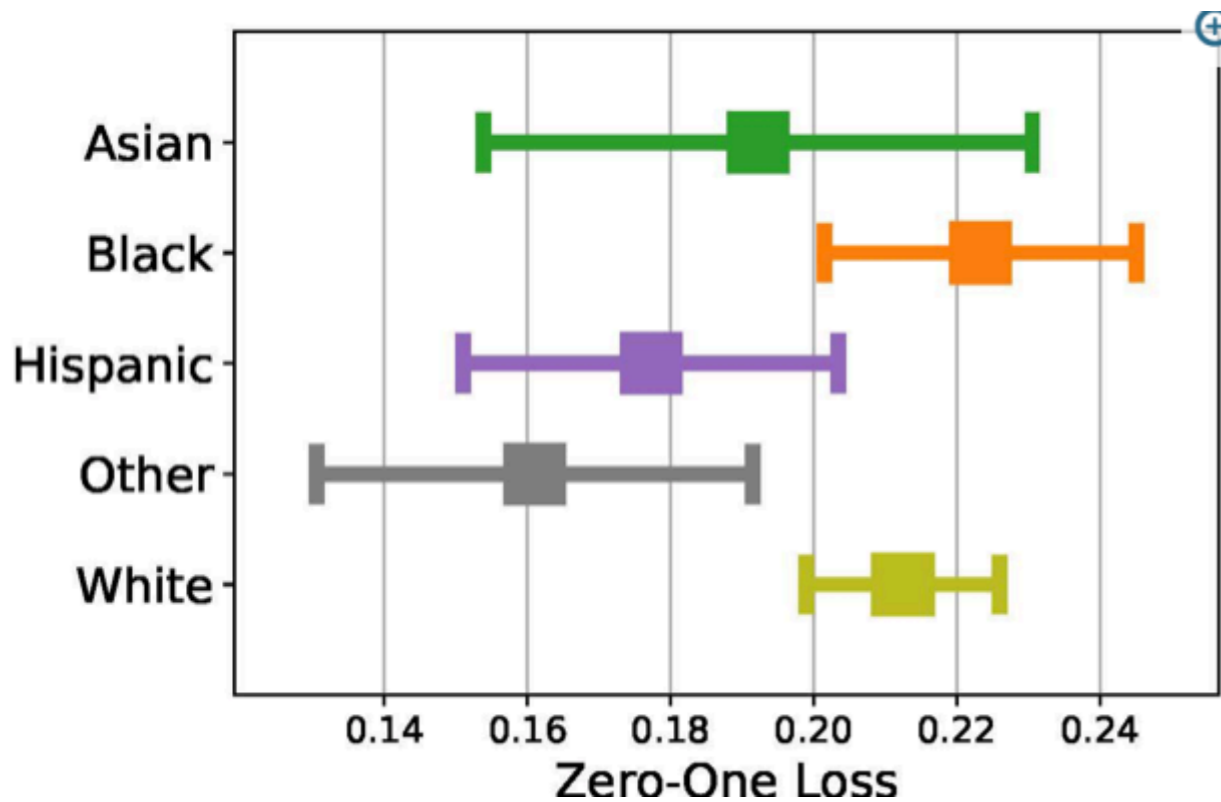
- Insurance:
  - Those with public insurance often have **multiple chronic conditions** that require regular care
  - Public insurance patients have higher topic enrichment values for **atrial fibrillation** (0.024 v 0.013), **pacemakers** (0.023 v 0.014), and **dialysis** (0.023 v 0.013)
  - private insurance patients have higher topic enrichment values for **fractures** (0.035 v 0.012), **lymphoma** (0.030 v 0.015), and **aneurysms** (0.028 v 0.016)
- These results are consistent with known disparities from literature

# Prediction Errors in ICU (violation of Separation)

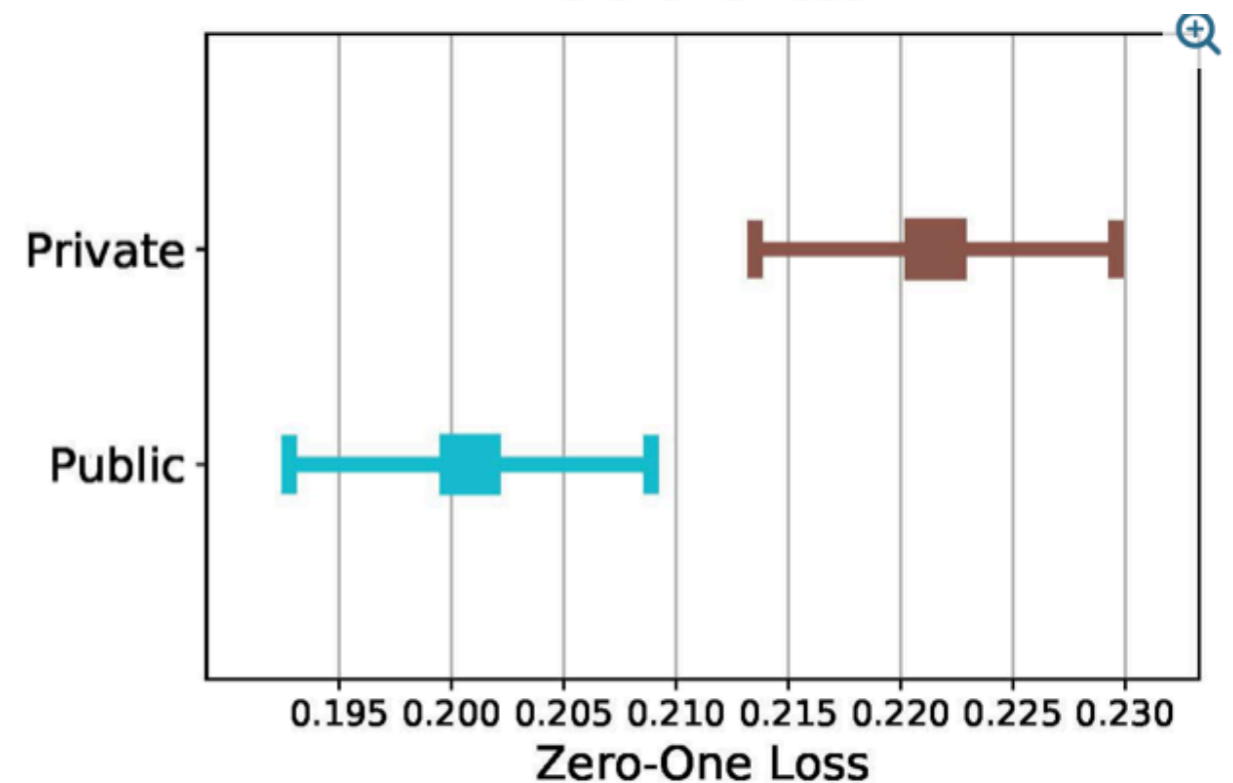


- 95% confidence intervals for zero-one loss differences across gender and insurance type

# Prediction Errors in Psychiatry (violation of Separation)



- 95% confidence intervals for zero-one loss differences across race, gender and insurance type





---

## Modeling Mistrust in End-of-Life Care

---

Willie Boag<sup>1</sup> Harini Suresh<sup>1</sup> Leo Anthony Celi<sup>1</sup> Peter Szolovits<sup>1</sup> Marzyeh Ghassemi<sup>1 2 3 4</sup>

Based on Boag, W. (2018, June). Quantifying Racial Disparities in End-of-Life Care. Master's Thesis, MIT EECS. Cambridge, MA.

- Replicate in MIMIC Racial Disparities expectation from previous studies
- Model Mistrust Algorithmically
- Compare Racial and Mistrust Disparities

# Racial Disparities in End-of-Life Care

## African American patients receive longer durations of aggressive treatment during end-of-life care

Figure 3-1: **Mechanical Ventilation:** CDF of ventilation duration by race, where dotted lines represent the median duration treatment for a population. In multiple datasets, the median black patient receives statistically significant longer ventilation durations than the median white patient.

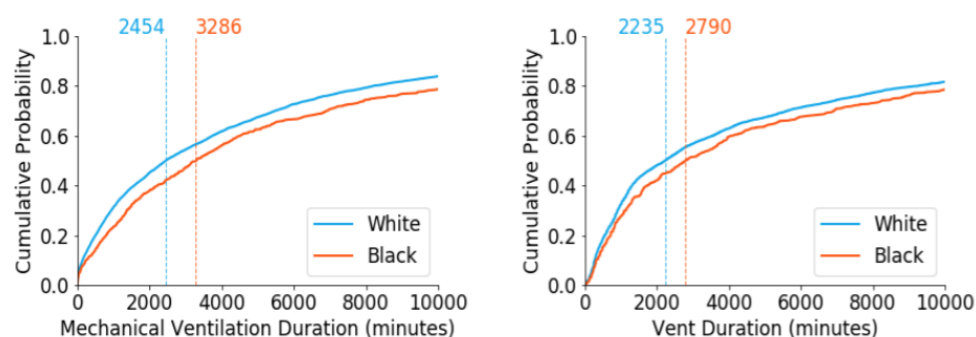
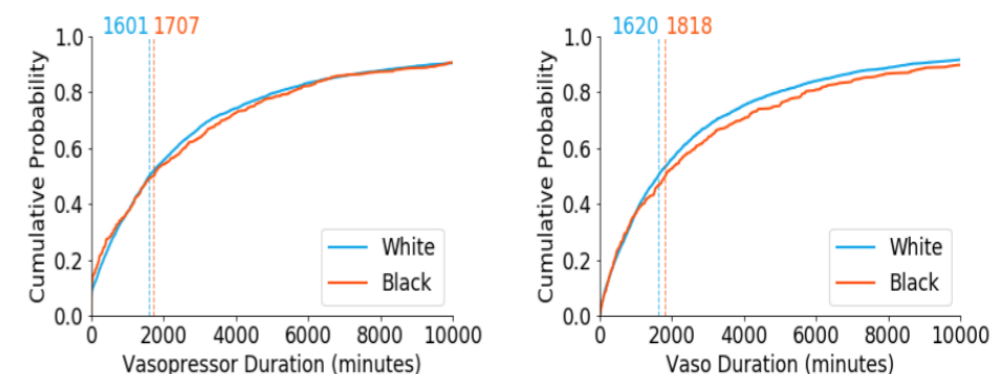


Figure 3-2: **Vasopressors:** In both datasets, the median black patient receives a longer duration of vasopressors than the median white patient. This trend is not statistically significant in either dataset..



## Could this be the result of mistrust?

(e.g. If your doctor recommends hospice, do you accept their advice?)

# Clues of Mistrust

---

## Noncompliance in Clinical Notes

# Social: Pt **refused to sign ICU consent** and expressed wishes to be DNR/DNI, seemingly **very frustrated** and **mistrusting of healthcare system** in relation to [REDACTED]. Also, w/ hx of **poor medication compliance and follow-up**

## Autopsy Rates

Table 4.3: Autopsy rates by race in MIMIC III.

population	consent	decline	% consent
Asian	2	23	8.0%
White	161	505	24.2%
Other	56	102	32.9%
Black	32	51	38.6%
Hispanic	9	11	45.0%
ALL	260	692	27.3%

Problem: Not every patient has an “obvious” label.



Can we use the obvious examples as labels and train a model to interpolate every patient’s “mistrust” score onto the scale?

# Chart Events Give Clues About Patient State Relevant to True

Table 4.1: Coded interpersonal feature types from chartevents.

1:1 sitter present?	baseline pain level (0 to 10)	received bath?	bedside observer
behavioral intervent	currently experiencing pain	disease state	consults
education barrier	education learner	education method	feamily meeting?
education readiness	harm by partner?	education topic	judgement
follows commands?	family communication method	gcs - verbal response	informed?
hair washed?	goal richmond-ras scale	headache?	health care proxy?
pain management	non-violent restraints?	orientation	pain (0 to 10)
pain assess method	understand & agree with plan?	pain level acceptable?	reason for restraint
restraint device	richmond-ras scale (-5 to +4)	rsbi deferred	riker-sas scale
safety measures	violent restraints ordered?	security	security guard
side rails	status and comfort	sitter	skin care?
spiritual support	behavior during application	support systems	stress
verbal response	teaching directed toward	wrist restraints?	social work consult?

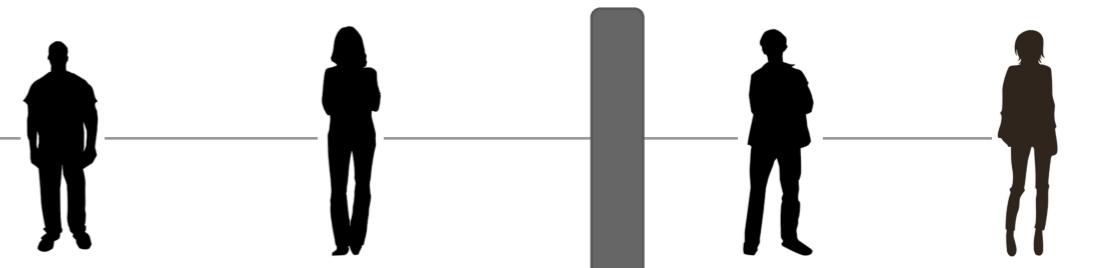
Structured data in the EHR documenting interpersonal variables, including:

- Is the patient's comfort being taken seriously?
- Is the patient being treated as a threat?
- Is the patient's pain being managed?
- Are there good communication between staff and the family?

# Modeling Mistrust

# Social: Pt **refused to sign ICU consent** and expressed wishes to be DNR/DNI, seemingly **very frustrated** and **mistrusting of healthcare system** in relation to [REDACTED]. Also, w/ hx of **poor medication compliance and follow-up**

620 binary indicators of trust
indication of family meetings
patient education engagement
patient needed to be restrained
pain is being monitored and treated
healthcare literacy
has a healthcare proxy
has a support system (such as family, social workers, and religion)
agitation scales (Riker-SAS and Richmond-RAS)



$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

“trustful” “mistrustful”

L1-reg Logistic Regression

Mistrust Scores:

0.30

0.80

0.10

0.72

Labeled Examples

Unlabeled Examples

# Inspecting the Mistrust Metrics

---

**Mistrustful patients:** Agitated & in pain  
**Trustful patients:** No pain & calm

**Mistrustful patients:** Restrained  
**Trustful patients:** No pain & healthcare literacy

Feature	Weight
state: alert	-1.0156
riker-sas scale: agitated	0.7013
pain: none	-0.5427
richmond-ras scale: 0 alert and calm	-0.3598
education readiness: no	0.2540
pain level: 7-mod to severe	0.2168

(a) Noncompliance-derived Mistrust

Feature	Weight
pain present: no	-0.2689
spokesperson is healthcare proxy	-0.2271
family communication: talked to m.d.	-0.1184
reapplied restraints	0.1153
restraint type: soft limb	0.0980
orientation: oriented 3x	0.0363

(b) Autopsy-derived Mistrust

**Treatment Disparities are much larger across trust cohorts than race.**

	Race-based Disparity	Trust-based Disparity
Mechanical ventilation	<p>(a) <i>MIMIC Mechanical Ventilation</i>  <b>White:</b> 4810 patients  <b>Black:</b> 510 patients  <math>p=0.005</math></p>	<p>(a) <b>Mechanical Ventilation</b>  <b>High Trust:</b> 4810 patients  <b>Low Trust:</b> 510 patients  <math>p &lt; 0.001</math></p>
Vasopressors	<p>(a) <i>MIMIC Vasopressors</i>  <b>White:</b> 4458 patients  <b>Black:</b> 453 patients  <math>p=0.122</math></p>	<p>(b) <b>Vasopressors</b>  <b>High Trust:</b> 4456 patients  <b>Low Trust:</b> 453 patients  <math>p=0.001</math></p>



# Mistrust is Not Just a Proxy for Severity

---

Table 4: Pairwise Pearson correlation coefficients between scores.

	OASIS	SAPS II	Noncompliance	Autopsy	Sentiment
OASIS	1.0	0.679	0.050	-0.012	0.075
SAPS II	0.679	1.0	0.013	-0.013	0.086
Noncompliance	0.050	0.013	1.0	0.262	0.058
Autopsy	-0.012	-0.013	0.262	1.0	0.044
Sentiment	0.075	0.086	0.058	0.044	1.0

# Population Mistrust

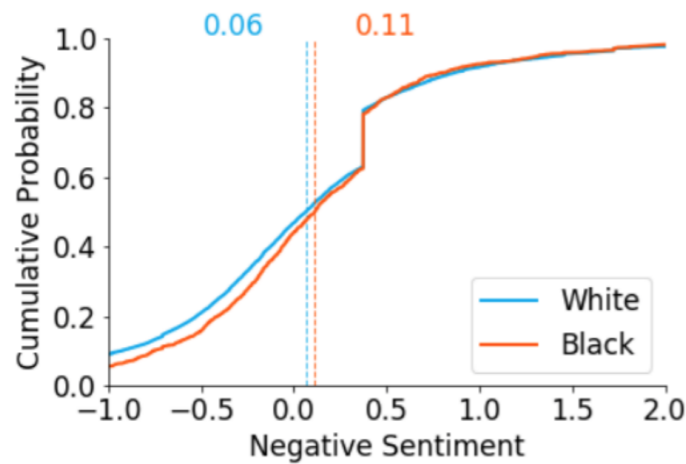


Figure 4-4: Racial disparity in (negative) sentiment.  
**White:** 9669 patients  
**Black:** 1173 patients  
 $p=0.007$

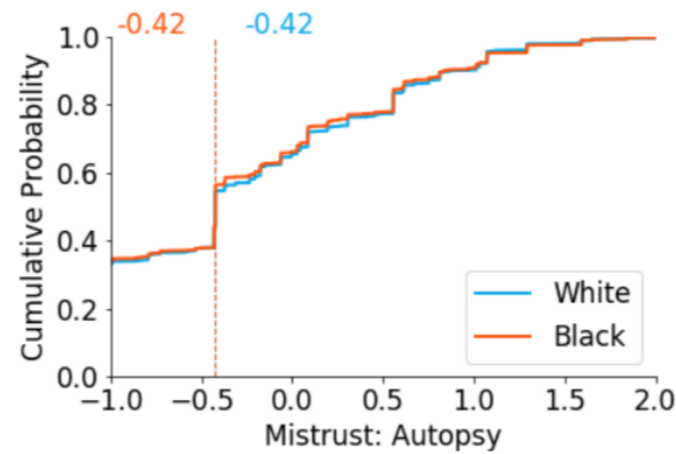


Figure 4-3: Racial disparity in autopsy-derived mistrust metric.  
**White:** 9923 patients  
**Black:** 1202 patients  
 $p=0.126$

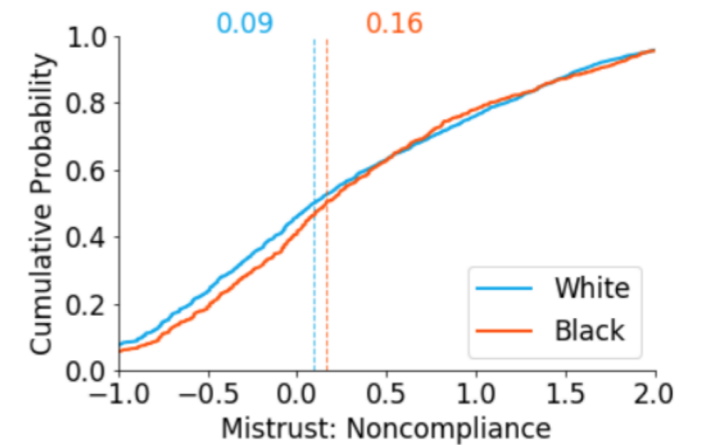


Figure 4-2: Racial disparity in noncompliance-derived mistrust metric.  
**White:** 9923 patients  
**Black:** 1202 patients  
 $p < 0.001$

For 2/3 metrics, the median black patient has a statistically significantly higher mistrust score than the median white patient.

# Much Work and Education to be Done

---

- Conferences and Workshops
  - Fairness, Accountability, and Transparency in Machine Learning (FAT/ML) Workshop
  - ACM Conference on Fairness, Accountability, and Transparency (ACM FAT\*)
  - Machine Learning for Healthcare Conference (MLHC)
  - ACM CHI Conference on Human Factors in Computing Systems (CHI)
- Popular Press
- Classes
  - Berkeley CS 294: Fairness in Machine Learning
  - U. Penn CIS 399 The Science of Data Ethics



# How to Achieve Algorithmic Fairness

- We've seen a few:
  - Defer only for cases where DM is “fair”
  - De-bias data; GANs
  - Calibrate outputs

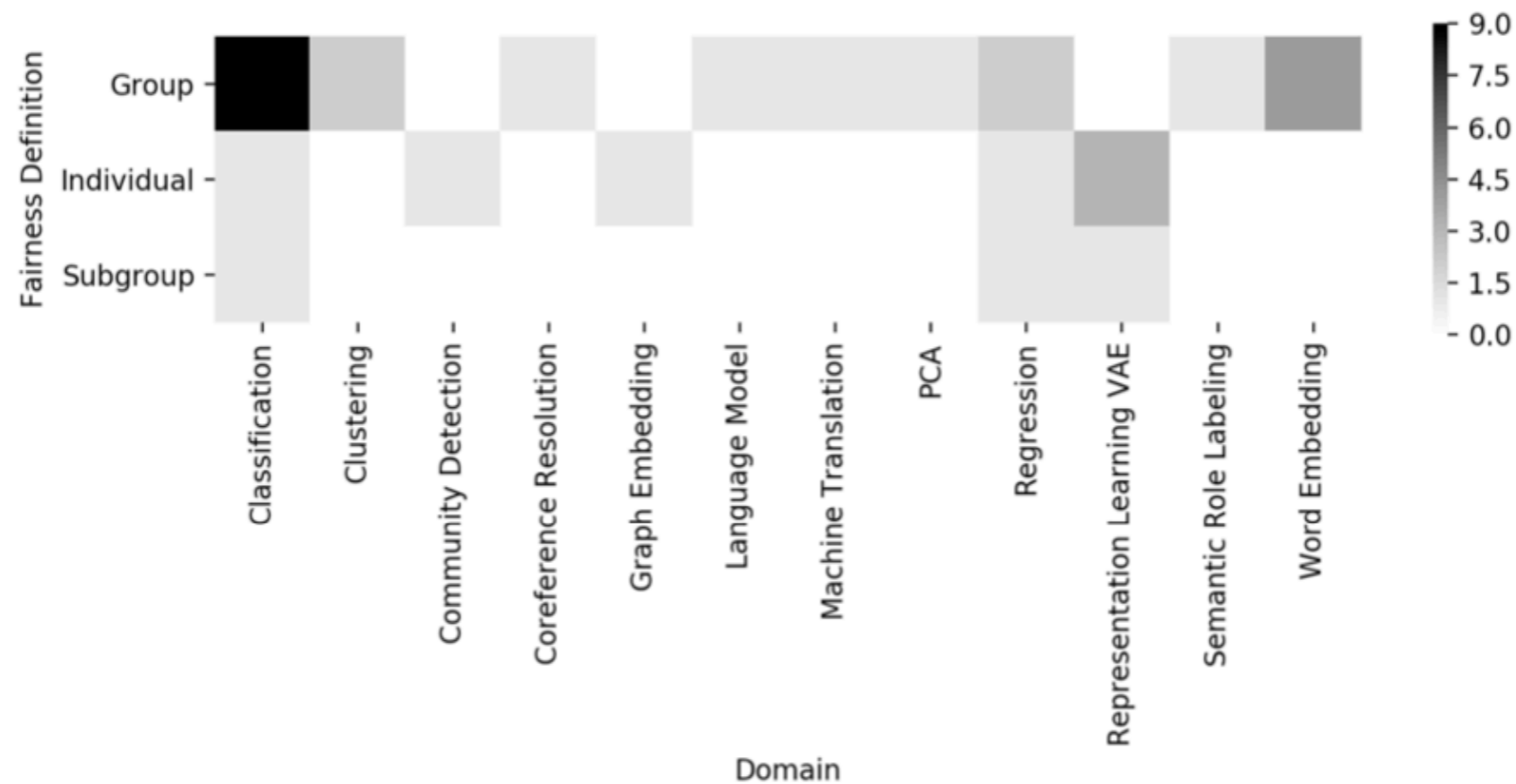


Fig. 7. Heatmap depicting distribution of previous work in fairness, grouped by domain and fairness definition.