# 10-301/601: Introduction to Machine Learning Lecture 28: Algorithmic Bias

Henry Chai

8/8/22

# Front Matter

- Announcements
  - HW9 released 8/3, due 8/9 (tomorrow) at 1 PM
    - Only one grace day allowed on HW9
  - Exam 3 on 8/12, this Friday!
    - Exam review recitation on 8/10, this Wednesday
    - Today's lecture is out-of-scope for Exam 3

- Recommended Supplementary Material
  - Solon Barocas and Mortiz Hardt's 2017 NeurIPS tutorial on Fairness in ML: https://vimeo.com/248490141

# Are Face-Detection Cameras Racist?

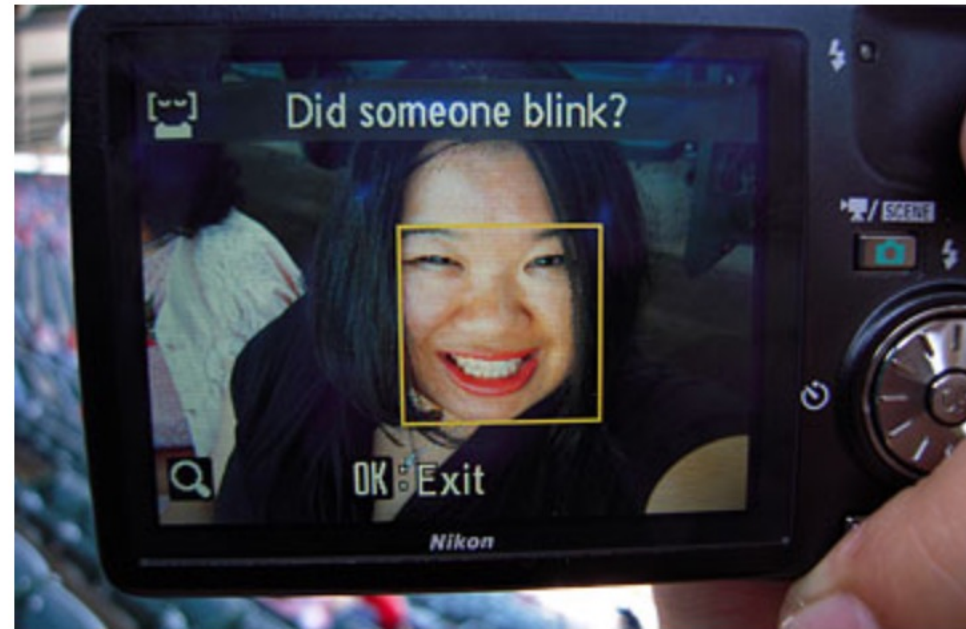By Adam Rose | Friday, Jan. 22, 2010

**Tweet**     **in Share**     **Read Later**

When Joz Wang and her brother bought their mom a Nikon Coolpix S630 digital camera for Mother's Day last year, they discovered what seemed to be a malfunction. Every time they took a portrait of each other smiling, a message flashed across the screen asking, "Did someone blink?" No one had. "I thought the camera was broken!" Wang, 33, recalls. But when her brother posed with his eyes open so wide that he looked "bug-eyed," the messages stopped.

Wang, a Taiwanese-American strategy consultant who goes by the Web handle "jozjozjoz," thought it was funny that the camera had difficulties figuring out when her family had their eyes open.



Joz Wang

## IS THE IPHONE X RACIST? APPLE REFUNDS DEVICE THAT CAN'T TELL CHINESE PEOPLE APART, WOMAN CLAIMS

BY **CHRISTINA ZHAO** ON 12/18/17 AT 12:24 PM EST

"A Chinese woman [surname Yan] was offered <u>two</u> refunds from Apple for her new iPhone X… [it] was unable to tell her and her other Chinese colleague apart."

"Thinking that a faulty camera was to blame, the store operator gave [Yan] a refund, which she used to purchase another iPhone X. But the new phone turned out to have the same problem, prompting the store worker to offer her another refund … <u>It is unclear whether she purchased a third phone</u>"

Source: https://www.newsweek.com/iphone-x-racist-apple-refunds-device-cant-tell-chinese-people-apart-woman-751263

"As facial recognition systems become more common, Amazon has emerged as a frontrunner in the field, courting customers around the US, including police departments and Immigration and Customs Enforcement (ICE)."

# Gender and racial bias found in Amazon's facial recognition technology (again)

*Research shows that Amazon's tech has a harder time identifying gender in darker-skinned and female faces*

By James Vincent | Jan 25, 2019, 9:45am EST

# Word embeddings and analogies

- https://lamyiowce.github.io/word2viz/

# Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
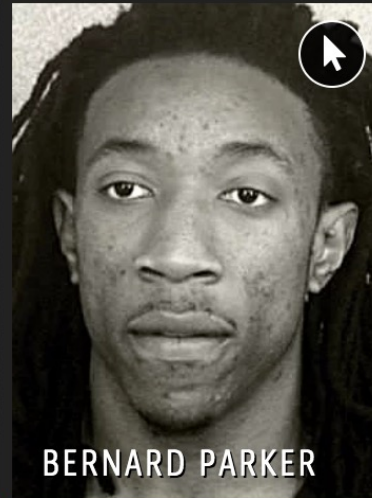
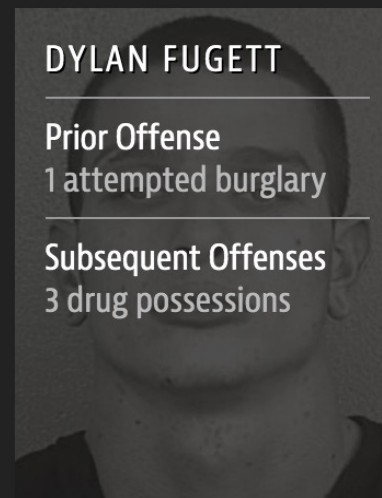May 23, 2016

## Two Drug Possession Arrests

**DYLAN FUGETT**

**BERNARD PARKER**

LOW RISK 3

HIGH RISK 10

*Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.*

## Two Drug Possession Arrests

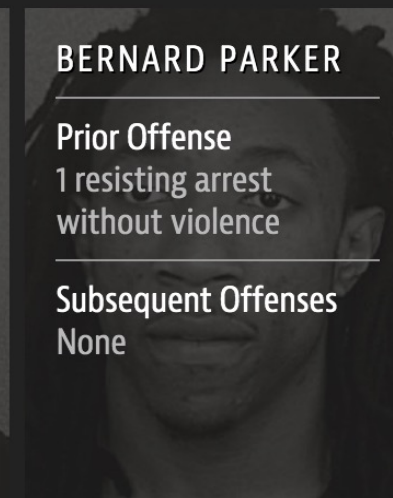**DYLAN FUGETT**

Prior Offense
1 attempted burglary

Subsequent Offenses
3 drug possessions

**BERNARD PARKER**

Prior Offense
1 resisting arrest without violence

Subsequent Offenses
None

LOW RISK 3

HIGH RISK 10

*Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.*

# Different Types of Errors

| | True label | Predicted label |
|---|---|---|
| True positive (TP) | $+1$ | $+1$ |
| False positive (FP) | $-1$ | $+1$ |
| True negative (TN) | $-1$ | $-1$ |
| False negative (FN) | $+1$ | $-1$ |

# How We Analyzed the COMPAS Recidivism Algorithm

*by Jeff Larson, Surya Mattu, Lauren Kirchner and Julia Angwin*

May 23, 2016

| All Defendants | Low | High |
| --- | --- | --- |
| Survived | 2681 | 1282 |
| Recidivated | 1216 | 2035 |

FP rate: 32.35
FN rate: 37.40

| Black Defendants | Low | High |
| --- | --- | --- |
| Survived | 990 | 805 |
| Recidivated | 532 | 1369 |

FP rate: 44.85
FN rate: 27.99

| White Defendants | Low | High |
| --- | --- | --- |
| Survived | 1139 | 349 |
| Recidivated | 461 | 505 |

FP rate: 23.45
FN rate: 47.72

This is one possible definition of unfairness.

We'll explore a few others and see how they relate to one another.

# Running Example

- Suppose you're an admissions officer for some program at CMU, deciding which applicants to admit

- $X$ are the non-protected features of an applicant (e.g., standardized test scores, GPA, etc...)

- $A$ is a protected feature (e.g., gender), usually categorical, i.e., $A \in \{a_1, \dots, a_C\}$

- $h(X, A) \in \{+1, -1\}$ is your model's prediction, usually corresponding to some decision or action (e.g., $+1 =$ admit to CMU)

- $Y \in \{+1, -1\}$ is the true, underlying target variable, usually some latent or hidden state (e.g., $+1 =$ this applicant would be "successful" at CMU)

## Attempt 1: Fairness through Unawareness

- Idea: build a model that only uses the non-protected features, $X$

- Achieves some notion of "individual fairness"
  - "Similar" individuals will receive "similar" predictions
  - Two individuals who are identical except for their protected feature $A$ would receive the same predictions

- Problem: the non-protected features $X$ might be affected by/dependent on $A$
  - In general, $X$ and $A$ are not independent

# Healthcare risk algorithm had 'significant racial bias'

It reportedly underestimated health needs for black patients.

Jon Fingas, @jonfingas
10.26.19 in Medicine

"While it [the algorithm] didn't directly consider ethnicity, its emphasis on medical costs as bellwethers for health led to the code routinely underestimating the needs of black patients. A sicker black person would receive the same risk score as a healthier white person simply because of how much they could spend."

# Three Definitions of Fairness

- **Independence**:

- **Separation**:

- **Sufficiency**:

# Three Definitions of Fairness

- **Independence**: $h(X, A) \perp A$

- **Separation**:

- **Sufficiency**:

# Independence

- Probability of being accepted is the same for all genders

$$P(h(X, A) = +1 | A = a_i) = P(h(X, A) = +1 | A = a_j) \, \forall \, a_i, a_j$$

or more generally,

$$P(h(X, A) = +1 | A = a_i) \approx P(h(X, A) = +1 | A = a_j) \, \forall \, a_i, a_j$$

$$\frac{P(h(X, A) = +1 | A = a_i)}{P(h(X, A) = +1 | A = a_j)} \geq 1 - \epsilon \, \forall \, a_i, a_j \text{ for some } \epsilon$$

# Achieving Fairness

- Pre-processing data

- Additional constraints during training

- Post-processing predictions

# Achieving Independence

- Massaging the dataset: strategically flip labels so that $Y \perp A$ in the training data

| $X$ | $A$ | $Y$ | Score | $Y'$ |
|-----|-----|-----|-------|------|
| | +1 | +1 | 0.98 | +1 |
| | +1 | +1 | 0.89 | +1 |
| | +1 | +1 | 0.61 | −1 |
| | +1 | −1 | 0.30 | −1 |
| ... | −1 | +1 | 0.96 | +1 |
| | −1 | −1 | 0.42 | +1 |
| | −1 | −1 | 0.31 | −1 |
| | −1 | −1 | 0.02 | −1 |

# Achieving Independence

- Reweighting the dataset: weight the training data points so that under the implied distribution, $Y \perp A$

| $X$ | $A$ | $Y$ | Score | $\Omega$ |
|---|---|---|---|---|
| | +1 | +1 | 0.98 | 1/12 |
| | +1 | +1 | 0.89 | 1/12 |
| | +1 | +1 | 0.61 | 1/12 |
| ... | +1 | −1 | 0.30 | 1/4 |
| | −1 | +1 | 0.96 | 1/4 |
| | −1 | −1 | 0.42 | 1/12 |
| | −1 | −1 | 0.31 | 1/12 |
| | −1 | −1 | 0.02 | 1/12 |

# Independence

- Probability of being accepted is the same for all genders

$$P(h(X, A) = +1 | A = a_i) = P(h(X, A) = +1 | A = a_j) \, \forall \, a_i, a_j$$

or more generally,

$$P(h(X, A) = +1 | A = a_i) \approx P(h(X, A) = +1 | A = a_j) \, \forall \, a_i, a_j$$

$$\frac{P(h(X, A) = +1 | A = a_i)}{P(h(X, A) = +1 | A = a_j)} \geq 1 - \epsilon \, \forall \, a_i, a_j \text{ for some } \epsilon$$

- Problem: permits laziness, i.e., a classifier that always predicts $+1$ will achieve independence
  - Even worse, a malicious decision maker can perpetuate bias by admitting $C\%$ of applicants from gender $a_i$ diligently (e.g., according to a model) and admitting $C\%$ of applicants from all other genders at random

# Three Definitions of Fairness

- **Independence**: $h(X, A) \perp A$
  - Probability of being accepted is the same for all genders

  - Permits laziness/is susceptible to adversarial decisions

- **Separation**:

- **Sufficiency**:

# Three Definitions of Fairness

- **Independence**: $h(X, A) \perp A$
  - Probability of being accepted is the same for all genders

  - Permits laziness/is susceptible to adversarial decisions

- **Separation**: $h(X, A) \perp A \mid Y$
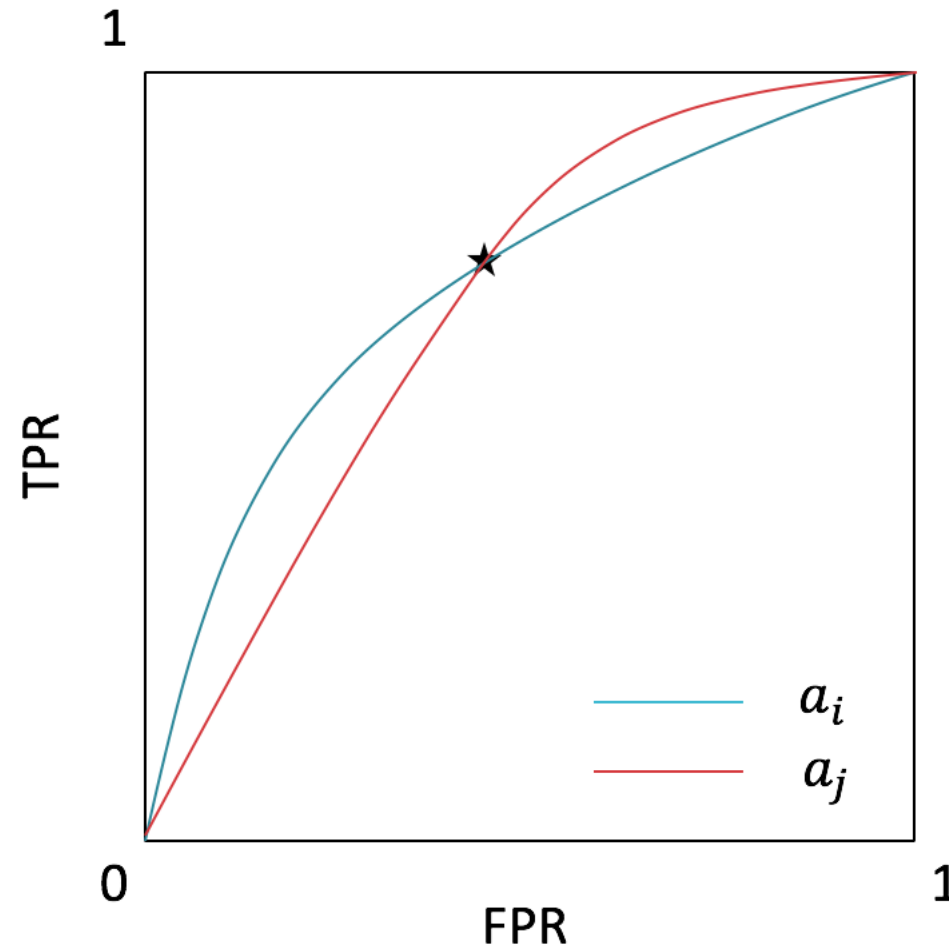
- **Sufficiency**:

## Separation

- Predictions and protected features can be correlated to the extent justified by the (latent) target variable

$$P(h(X, A) = +1 | Y = +1, A = a_i)$$
$$= P(h(X, A) = +1 | Y = +1, A = a_j) \ \&$$
$$P(h(X, A) = +1 | Y = -1, A = a_i)$$
$$= P(h(X, A) = +1 | Y = -1, A = a_j) \ \forall \ a_i, a_j$$

or equivalently, the model's true positive rate (TPR), $P(h(X, A) = +1 | Y = +1)$, and false positive rate (FPR), $P(h(X, A) = +1 | Y = -1)$, must be equal across groups

- Natural relaxations care about only one of these two

# Achieving Separation



- ROC curve plots TPR against FPR at different prediction thresholds, $\tau$:

$$h(X, A) = \mathbb{1}(\text{SCORE} \geq \tau)$$

- Can achieve separation by using different thresholds for different groups, corresponding to where their ROC curves intersect

# Separation

- Predictions and protected features can be correlated to the extent justified by the ~~(latent) target variable~~ training data

$$P(h(X,A) = +1 | Y = +1, A = a_i)$$
$$= P(h(X,A) = +1 | Y = +1, A = a_j) \, \&$$
$$P(h(X,A) = +1 | Y = -1, A = a_i)$$
$$= P(h(X,A) = +1 | Y = -1, A = a_j) \, \forall \, a_i, a_j$$

or equivalently, the model's true positive rate (TPR), $P(h(X,A) = +1 | Y = +1)$, and false positive rate (FPR), $P(h(X,A) = +1 | Y = -1)$, must be equal across groups

  - Natural relaxations care about only one of these two

- Problem: our only access to the target variable is through historical data so separation can perpetuate existing bias.

# Three Definitions of Fairness

- **Independence**: $h(X, A) \perp A$
  - Probability of being accepted is the same for all genders

  - Permits laziness/is susceptible to adversarial decisions

- **Separation**: $h(X, A) \perp A \mid Y$
  - All "good" applicants are accepted with the same probability, regardless of gender
  - Perpetuates existing biases in the training data

- **Sufficiency**:

# Three Definitions of Fairness

- **Independence**: $h(X, A) \perp A$
  - Probability of being accepted is the same for all genders

  - Permits laziness/is susceptible to adversarial decisions

- **Separation**: $h(X, A) \perp A \mid Y$
  - All "good" applicants are accepted with the same probability, regardless of gender
  - Perpetuates existing biases in the training data

- **Sufficiency**: $Y \perp A \mid h(X, A)$

# Sufficiency

- Knowing the prediction is *sufficient* for decorrelating the (latent) target variable and the protected feature

$$P(Y = +1 | h(X, A) = +1, A = a_i)$$
$$= P(Y = +1 | h(X, A) = +1, A = a_j) \ \&$$
$$P(Y = +1 | h(X, A) = -1, A = a_i)$$
$$= P(Y = +1 | h(X, A) = -1, A = a_j) \ \forall \ a_i, a_j$$

If a model uses some score to make predictions, then that score is *calibrated* if

$$P(Y = +1 | \text{SCORE}) = \text{SCORE}$$

# Sufficiency

- Knowing the prediction is *sufficient* for decorrelating the (latent) target variable and the protected feature

$$P(Y = +1 | h(X, A) = +1, A = a_i)$$
$$= P(Y = +1 | h(X, A) = +1, A = a_j) \&$$
$$P(Y = +1 | h(X, A) = -1, A = a_i)$$
$$= P(Y = +1 | h(X, A) = -1, A = a_j) \, \forall \, a_i, a_j$$

If a model uses some score to make predictions, then that score is *calibrated across groups* if

$$P(Y = +1 | \text{SCORE}, A = a_i) = \text{SCORE} \, \forall \, a_i$$

A model being calibrated across groups implies sufficiency

- In general, most off-the-shelf ML models can achieve sufficiency without intervention

# Three Definitions of Fairness

- **Independence**: $h(X, A) \perp A$
  - Probability of being accepted is the same for all genders

  - Permits laziness/is susceptible to adversarial decisions

- **Separation**: $h(X, A) \perp A \mid Y$
  - All "good"/"bad" applicants are accepted with the same probability, regardless of gender
  - Perpetuates existing biases in the training data

- **Sufficiency**: $Y \perp A \mid h(X, A)$
  - For the purposes of predicting $Y$, the information contained in $h(X, A)$ is "sufficient", $A$ becomes irrelevant

# Many Definitions of Fairness (Barocas et al., 2019)

| Name | Closest relative | Note |
|---|---|---|
| Statistical parity | Independence | Equivalent |
| Group fairness | Independence | Equivalent |
| Demographic parity | Independence | Equivalent |
| Conditional statistical parity | Independence | Relaxation |
| Darlington criterion (4) | Independence | Equivalent |
| Equal opportunity | Separation | Relaxation |
| Equalized odds | Separation | Equivalent |
| Conditional procedure accuracy | Separation | Equivalent |
| Avoiding disparate mistreatment | Separation | Equivalent |
| Balance for the negative class | Separation | Relaxation |
| Balance for the positive class | Separation | Relaxation |
| Predictive equality | Separation | Relaxation |
| Equalized correlations | Separation | Relaxation |
| Darlington criterion (3) | Separation | Relaxation |
| Cleary model | Sufficiency | Equivalent |
| Conditional use accuracy | Sufficiency | Equivalent |
| Predictive parity | Sufficiency | Relaxation |
| Calibration within groups | Sufficiency | Equivalent |
| Darlington criterion (1), (2) | Sufficiency | Relaxation |

Source: https://fairmlbook.org/pdf/fairmlbook.pdf
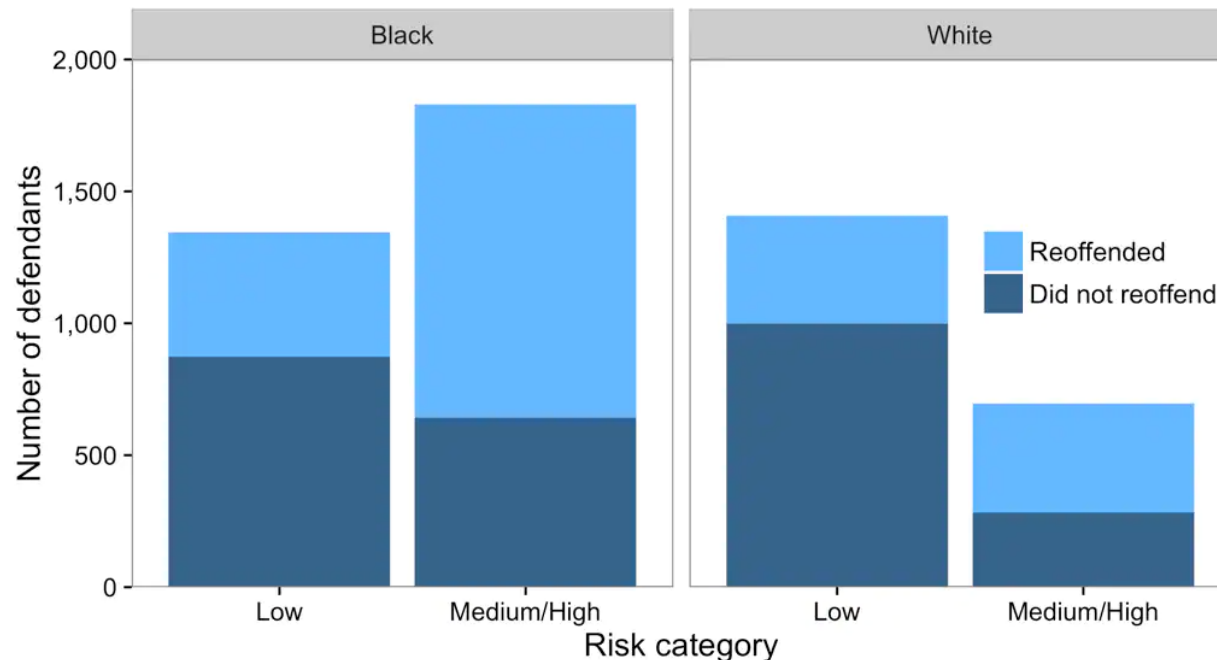
# Three Definitions of Fairness

- **Independence**: $h(X, A) \perp A$
  - Probability of being accepted is the same for all genders

  - Permits laziness/is susceptible to adversarial decisions

- **Separation**: $h(X, A) \perp A \mid Y$
  - All "good"/"bad" applicants are accepted with the same probability, regardless of gender
  - Perpetuates existing biases in the training data

- **Sufficiency**: $Y \perp A \mid h(X, A)$
  - For the purposes of predicting $Y$, the information contained in $h(X, A)$ is "sufficient", $A$ becomes irrelevant

# Three Incompatible Definitions of Fairness

- **Independence**: $h(X, A) \perp A$
  - Probability of being accepted is the same for all genders

  - Permits laziness/is susceptible to adversarial decisions

- **Separation**: $h(X, A) \perp A \mid Y$
  - All "good"/"bad" applicants are accepted with the same probability, regardless of gender
  - Perpetuates existing biases in the training data

- **Sufficiency**: $Y \perp A \mid h(X, A)$
  - For the purposes of predicting $Y$, the information contained in $h(X, A)$ is "sufficient", $A$ becomes irrelevant

Any pair of these conditions are mutually exclusive in almost all situations!

# A computer program used for bail and sentencing decisions was labeled biased against blacks. It's actually not that clear.

By Sam Corbett-Davies, Emma Pierson, Avi Feller and Sharad Goel
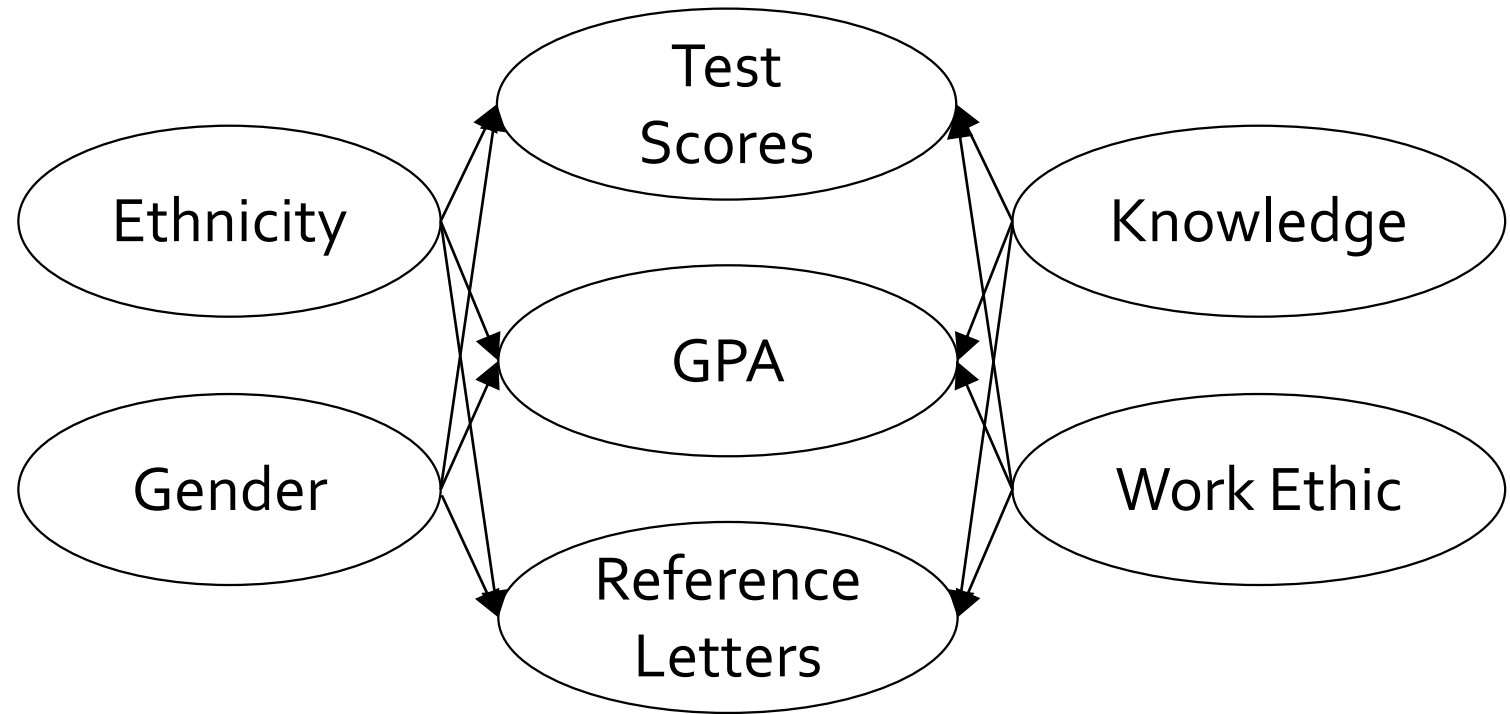
October 17, 2016



- Within each risk category, the proportion of defendants who reoffend is approximately the same regardless of race; this is Northpointe's definition of fairness.

- The overall recidivism rate for black defendants is higher than for white defendants (52 percent vs. 39 percent).

- Black defendants are more likely to be classified as medium or high risk (58 percent vs. 33 percent). While Northpointe's algorithm does not use race directly, many attributes that predict reoffending nonetheless vary by race. For example, black defendants are more likely to have prior arrests, and since prior arrests predict reoffending, the algorithm flags more black defendants as high risk even though it does not use race in the classification.

- Black defendants who don't reoffend are predicted to be riskier than white defendants who don't reoffend; this is ProPublica's criticism of the algorithm.

The key — but often overlooked — point is that the last two disparities in the list above are mathematically guaranteed given the first two observations.

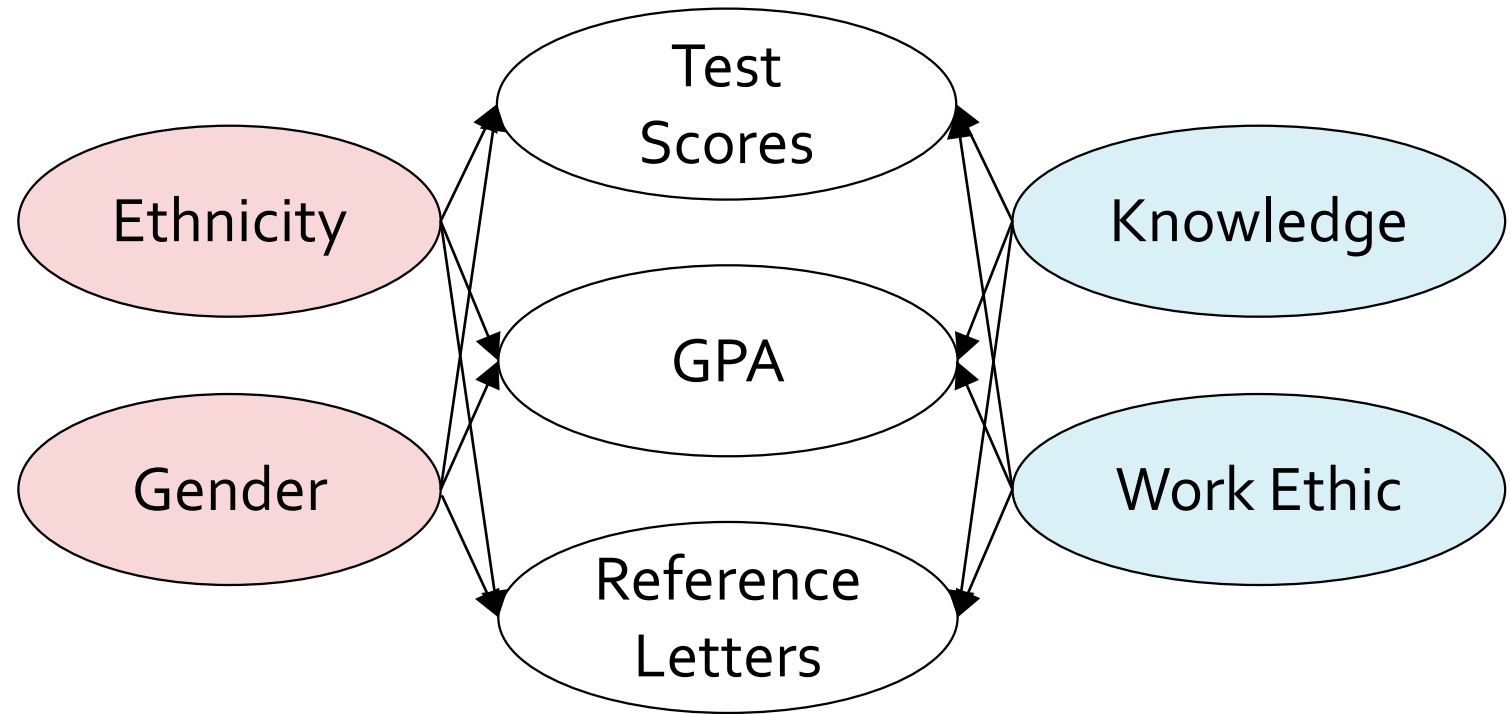# Yet another Definition of Fairness (Kusner et al., 2017)

- (Causal) Bayesian networks to the rescue!



- Counterfactual fairness: an applicant's probability of acceptance should not change if we were to change their gender

# Yet another Definition of Fairness (Kusner et al., 2017)

- (Causal) Bayesian networks to the rescue!



- Counterfactual fairness: any predictor that only relies on non-descendent of $A$ will be counterfactually fair

- Problem: how on earth do we specify this (causal) DAG?

# Key Takeaways

- High-profile cases of algorithmic bias are increasingly common as machine learning is applied more broadly in a variety of contexts

- Various definitions of fairness
  - Independence: $h(X, A) \perp A$
  - Separation: $h(X, A) \perp A \,|\, Y$
  - Sufficiency: $Y \perp A \,|\, h(X, A)$
    - In all but the simplest of cases, any two of these three are mutually exclusive
  - Counterfactual fairness via (causal) Bayesian networks