**Software Fairness**

April 11, 2023

---

TOM CRUISE

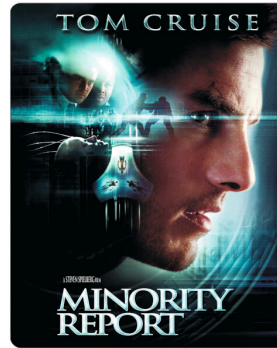MINORITY REPORT

---

Resilient cities  Cities

**Predicting crime, LAPD-style**

Cutting edge data-driven analysis directs Los Angeles patrol officers to likely future crime scenes – but critics worry that decision-making by machine will bring 'tyranny of the algorithm'

● Join our live Q&A with Homicide Watch this Friday

▲ PredPol co-developer P. Jeffrey Brantingham at the Unified Command Post in Los Angeles. 'This is not Minority Report,' he said. Photograph: Damian Dovarganes/AP

https://www.theguardian.com/cities/2014/jun/25/predicting-crime-lapd-los-angeles-police-data-analysis-algorithm-minority-report

---

ACLU                                   GET UPDATES / DONATE

**The Government Is Blacklisting People Based on Predictions of Future Crimes**

By Hina Shamsi, Director, ACLU National Security Project

**Modern software influences critical decisions**

https://www.aclu.org/blog/national-security/discriminatory-profiling/government-blacklisting-people-based-predictions

---

THE WALL STREET JOURNAL.

Home  World  U.S.  Politics  Economy  Business  Tech  Markets  Opinion  Arts  Life  Real Estate

TECH

**On Orbitz, Mac Users Steered to Pricier Hotels**

On Orbitz, Mac Users See Costlier Hotel Options

Orbitz has found that Apple users spend as much as 30% more a night on hotels, so the online travel site is starting to show them different, and sometimes costlier, options than Windows visitors see. Dana Mattioli has details on The News Hub. Photo: Bloomberg.

By Dana Mattioli

---

Forbes    LOG IN                          forbes.com

The Algorithm That Beats Your Bank Manager

HAAS NEWS > NEWS CATEGORIES > RESEARCH NEWS
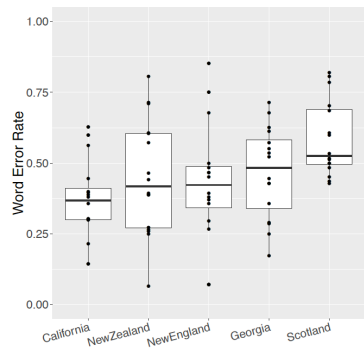
**Minority homebuyers face widespread statistical lending discrimination, study finds**

By Laura Counts  |  NOVEMBER 13, 2018

Face-to-face meetings between mortgage officers and homebuyers have been rapidly replaced by online applications and algorithms, but lending discrimination hasn't gone away.

A new University of California, Berkeley study has found that both online and face-to-face lenders charge higher interest rates to African American and Latino borrowers, earning 11 to 17 percent higher profits on such loans. All told, those homebuyers pay up to half a billion dollars more in interest every year than white borrowers with comparable credit scores do, researchers found.

The findings raise legal questions about the rise of statistical discrimination in the fintech era, and point to potentially widespread violations of U.S. fair lending laws, the researchers say. While lending discrimination has historically been caused by human prejudice, pricing disparities are increasingly the result of algorithms that use machine learning to target applicants who might shop around less for higher-priced loans.

"The mode of lending discrimination has shifted from human bias to algorithmic bias," said study co-author Adair Morse, a finance professor at UC Berkeley's Haas School of Business. "Even if the people writing the

Wisconsin Supreme Court allows state to continue using computer program to assist in sentencing

Machine Bias

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

Software can make bad decisions. Software can discriminate!

Google Translate

BING TRANSLATOR

Microsoft Bing

YouTube

YouTube automatic captions

Oh Jessica I am this stove I play the heroine me I am

## YouTube automatic captions

Word Error Rate

California  NewZealand  NewEngland  Georgia  Scotland
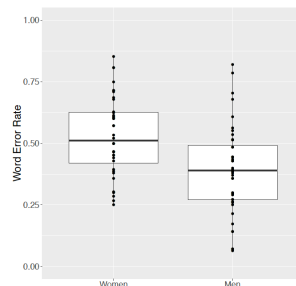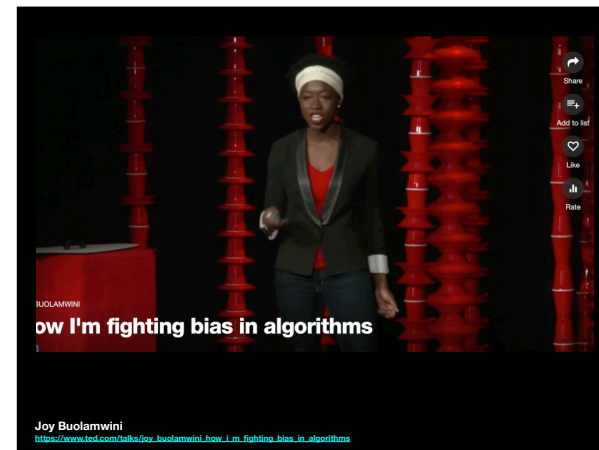
Rachael Tatman, "Gender and Dialect Bias in YouTube's Automatic Captions" in 2017 Workshop on Ethics in Natural Language Processing

---

## YouTube automatic captions

Word Error Rate

Women  Men

Rachael Tatman, "Gender and Dialect Bias in YouTube's Automatic Captions" in 2017 Workshop on Ethics in Natural Language Processing

---

ow I'm fighting bias in algorithms

Joy Buolamwini
https://www.ted.com/talks/joy_buolamwini_how_i_m_fighting_bias_in_algorithms

---

## how people want to use vision software

---

## how people want to use vision software

Hire-Vue
SEE A HIREVUE ASSESSMENTS DEMO

WANT TO SEE HIREVUE IN ACTION?
SEE A DEMO

The predictive power of a traditional assessment in a video interview.

- Predictive assessments tied to higher quality hires
- Assess and interview in one step and hire faster
- Better candidate experience
- Validated technology

JOIN 700+ CUSTOMERS THAT USE HIREVUE

vodafone   TJX   TIFFANY & CO.
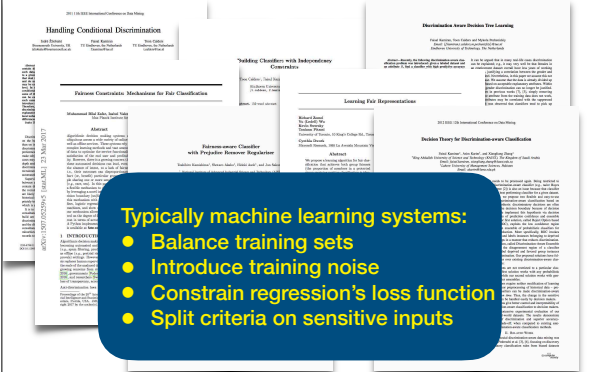
GEICO   DR PEPPER SNAPPLE   intel

---

## today's goals

Define software discrimination.

Operationalize measuring discrimination through causal software testing.

Provide provable fairness guarantees.

**Design software to be fair**

Typically machine learning systems:
- Balance training sets
- Introduce training noise
- Constrain regression's loss function
- Split criteria on sensitive inputs

---

**Design alone is not enough**

---

**possible causes**

biased data

poor design

implementation bugs

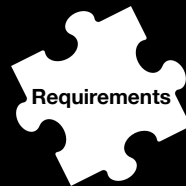unintended interactions and mismatched components

---

**Fairness is just like quality and security**

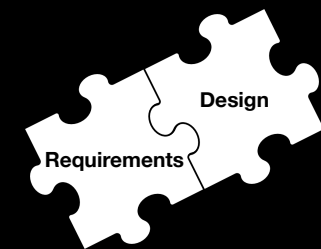Fairness must be part of the software engineering lifecycle

---

**Call to Action!**

Fairness must be part of the software engineering lifecycle

Requirements

We need methods for specifying **fairness requirements**

---

**Call to Action!**

Design

Requirements

We need fairness **design principles**

Call to Action!

We need automated fairness testing

Requirements

Testing



Call to Action!

We need fairness property verification

Requirements

Testing

Verification



Call to Action!

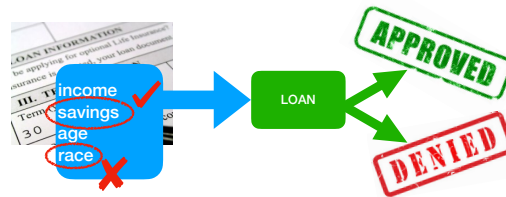Fairness must be part of the software engineering lifecycle

Design

Requirements

Testing

Verification



Let's talk about what it means for systems to discriminate



LOAN program

income
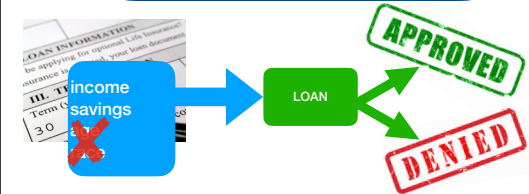savings
age
race

LOAN

APPROVED

DENIED

This talk is not about policy.



Fairness: Disparate Treatment

Hide the data

income
savings

LOAN

APPROVED

DENIED

Zafar et al. Fairness constraints: Mechanisms for fair classification. AISTATS 2017.

# Fairness: Disparate Treatment

**Hide the data**

Ads by Google

Latanya Sweeney, Arrested?
1) Enter Name and State. 2) Access Full Background Checks Instantly.
www.instantcheckmate.com/

Ineffective because of data correlation.
[Latanya Sweeney. Discrimination in online ad delivery. CACM 2013]

---



**BUSINESS INSIDER** — TECH NEWS

Amazon just showed us that 'unbiased' algorithms can be inadvertently racist

---

**BUSINESS INSIDER** — TECH · FINANCE · POLITICS · STRATEGY · LIFE · ALL — PRIME · INTELLIGENCE

Amazon built an AI tool to hire people but had to shut it down because it was discriminating against women

Isobel Asher Hamilton 4h

- Amazon tried building an artificial-intelligence tool to help with recruiting, but it showed a bias against women, Reuters reports.
- Engineers reportedly found the AI was unfavorable toward female candidates because it had combed through male-dominated résumés to accrue its data.
- Amazon reportedly abandoned the project at the beginning of 2017.

Amazon wo... with hiring...

**disparate treatment: still not fair**

https://www.businessinsider.com/amazon-built-ai-to-hire-people-discriminated-against-women-2018-10

---

# Fairness: Demographic Parity

**Compare subpopulation proportions**

APPROVED

35%   20%
65%   80%

**often called group discrimination**

Fails to identify discrimination against individuals.

Dwork et al. Fairness through awareness. ITCS 2012.
Calders and Verwer. Three naïve Bayes approaches for discrimination-free classification. DMKD 2010.

---

# How group discrimination can fail

## Europe        ## Asia

APPROVED

DENIED

approve loans to all **green** deny loans to all **purple** applicants

approve loans to all **purple** deny loans to all **green** applicants

European and Asian discriminations cancel each other out, and the group discrimination measure can be 0.

---

# Fairness: Disparate Impact

**Prohibits using a facially neutral practice that has an unjustified adverse impact on members of a protected class.**

**80% rule: Employer's hiring rates for protected groups may not differ by more than 80%.**

Zafar et al. Fairness constraints: Mechanisms for fair classification. AISTATS 2017.

# Fairness: Delayed Impact

Making seemingly fair decisions can (but shouldn't), in the long term, produce unfair consequences

Liu et al., Delayed impact of fair machine learning. ICML 2018

---

# Fairness: Predictive Equality
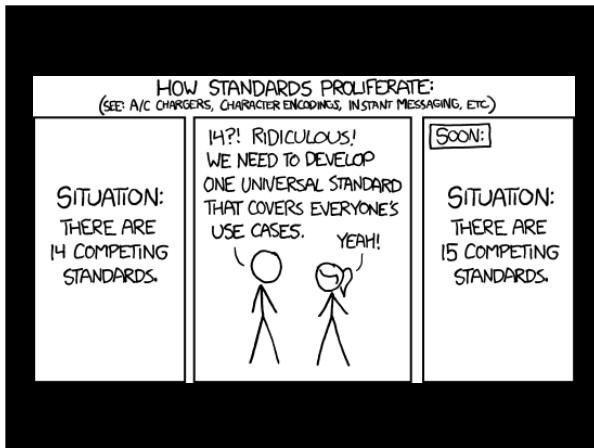
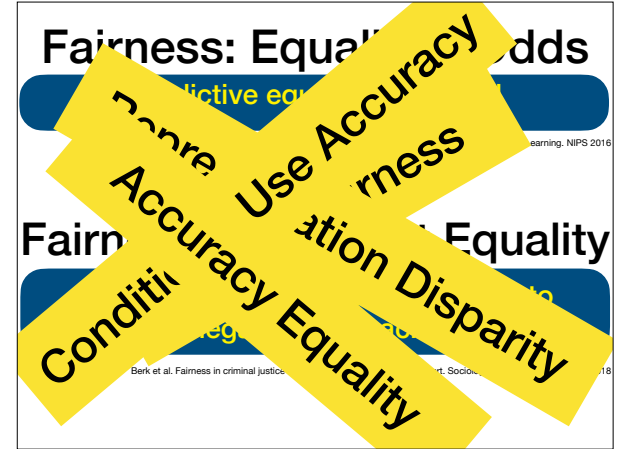False positive rates should not differ

Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. FATML 2016
Corbett-Davies. Algorithmic decision making and the cost of fairness. KDD 2017

# Fairness: Equal Opportunity

False negative rates should not differ

Hardt et al. Equality of Opportunity in Supervised Learning. NIPS 2016
Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments FATML 2016

---

# Fairness: Equalized Odds

Predictive equality

learning. NIPS 2016

# Fairness: ... ation Equality

Berk et al. Fairness in criminal justice ... ort. Sociol ... 18

Use Accuracy
Representation Fairness
Accuracy Disparity
Conditional Accuracy Equality

---



HOW STANDARDS PROLIFERATE:
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC)

SITUATION:
THERE ARE
14 COMPETING
STANDARDS.

14?! RIDICULOUS! WE NEED TO DEVELOP ONE UNIVERSAL STANDARD THAT COVERS EVERYONE'S USE CASES. YEAH!

SOON:

SITUATION:
THERE ARE
15 COMPETING
STANDARDS.

---

# Fairness: Correlation

correlation(race, APPROVED ) = 0.8

mutual information(race, APPROVED ) = 0.6

Correlation does not measure causation

Atlidakis et al. FairTest: Discovering unwarranted associations in data-driven applications. EuroS&P 2017

---

# What is fairness?

Sensitive inputs should not affect software behavior.

We want to measure causality!

Judea Pearl. Causal inference in statistics: An overview. Statistics Surveys 2009

# causal testing

Sensitive inputs should not affect software behavior.

hypothesis testing:

LOAN ?

Galhotra, Brun, and Meliou, Fairness Testing: Testing Software for Discrimination. ESEC/FSE 2017



# causal testing

No need for an oracle!



# causal testing



# Themis
automated test-suite generator

How much does my software discriminate with respect to …?

Does my software discriminate more than 10% of the time, and against

Themis generates a test suite or can use a manually written one

http://fairness.cs.umass.edu

Angell, Johnson, Brun, and Meliou, Themis: Automatically Testing Software for Discrimination. ESEC/FSE 2018 Demo



# discrimination measures

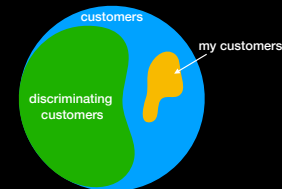causal discrimination

$$\text{LOAN}(\quad) \overset{?}{=} \text{LOAN}(\quad)$$

group discrimination

APPROVED
15%
DENIED



# apparent discrimination

customers

discriminating customers

my customers

customers

poor green

my customers

rich purple

Software may discriminate, but not for a given set of customers

Fair software may appear to discriminate
(e.g., Amazon same-day delivery)

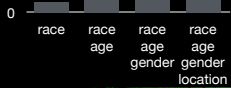★ Apparent discrimination can be group or causal, measured on a given test suite or operational profile.

# How does Themis work?

adaptive, confidence-driven sampling

input schema →
confidence →
error bound →

$$error = z^* \sqrt{\frac{p(1-p)}{r}}$$

Themis

sound pruning

1
0

race | race age | race age gender | race age gender location

---

# Evaluation

Eight open-source decision systems trained on two public data sets

| | |
|---|---|
| discrimination-aware logistic regression | [88] |
| discrimination-aware decision tree | [40] |
| discrimination-aware naive Bayes | [18] |
| discrimination-aware decision tree | [91] |
| naive Bayes | scikit-learn |
| decision tree | |
| logistic regression | |
| SVM | |

- Census income dataset:
  financial data
  45K people
  income > $50K?

- Statlog German credit dataset:
  credit data
  1K people
  "good" or "bad" credit?

---

# findings

**Group discrimination is not enough.**

More than 11% of the individuals had the output flipped just by altering the individual's gender.

Decision tree trained not to group discriminate against gender causal discriminated against gender: 0.11.

---

# findings

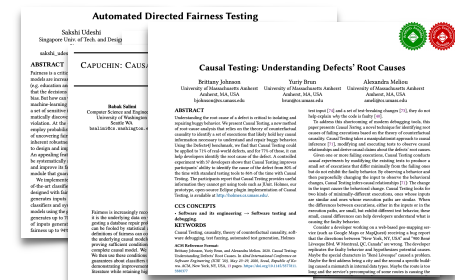**Trying to avoid group discrimination**

Training a decision tree not to discriminate against gender made it discriminate against race 38.4% of the time.

---

# findings

**Pruning is highly effective.**

- The more a system discriminates, the more efficient Themis is.

- On average, pruning reduced test suites by 148x for causal and 2,849x for group discrimination. Best improvement was 13,000x.

---

# Debugging