

Data science ethics

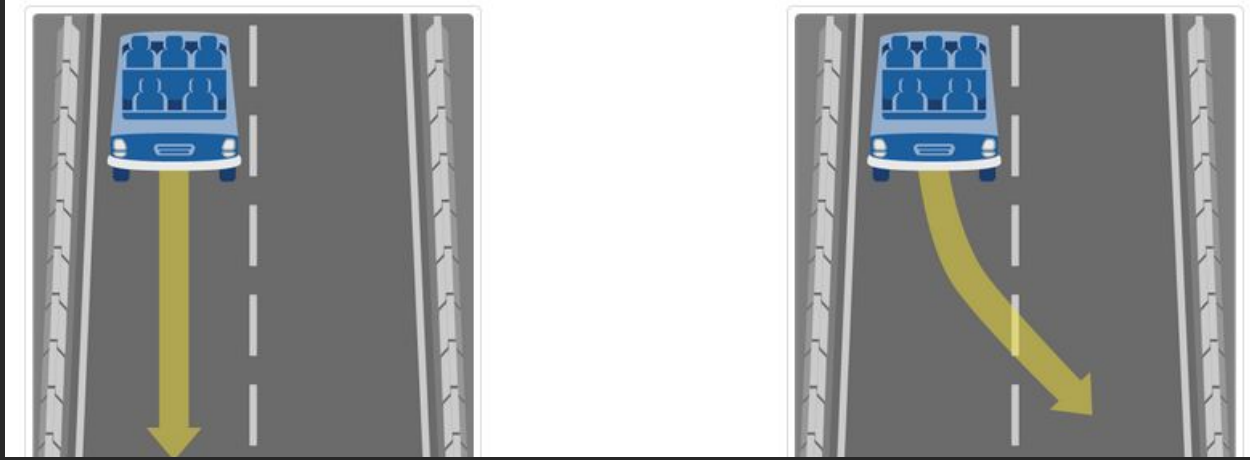
Data science ethics

Data science ethics

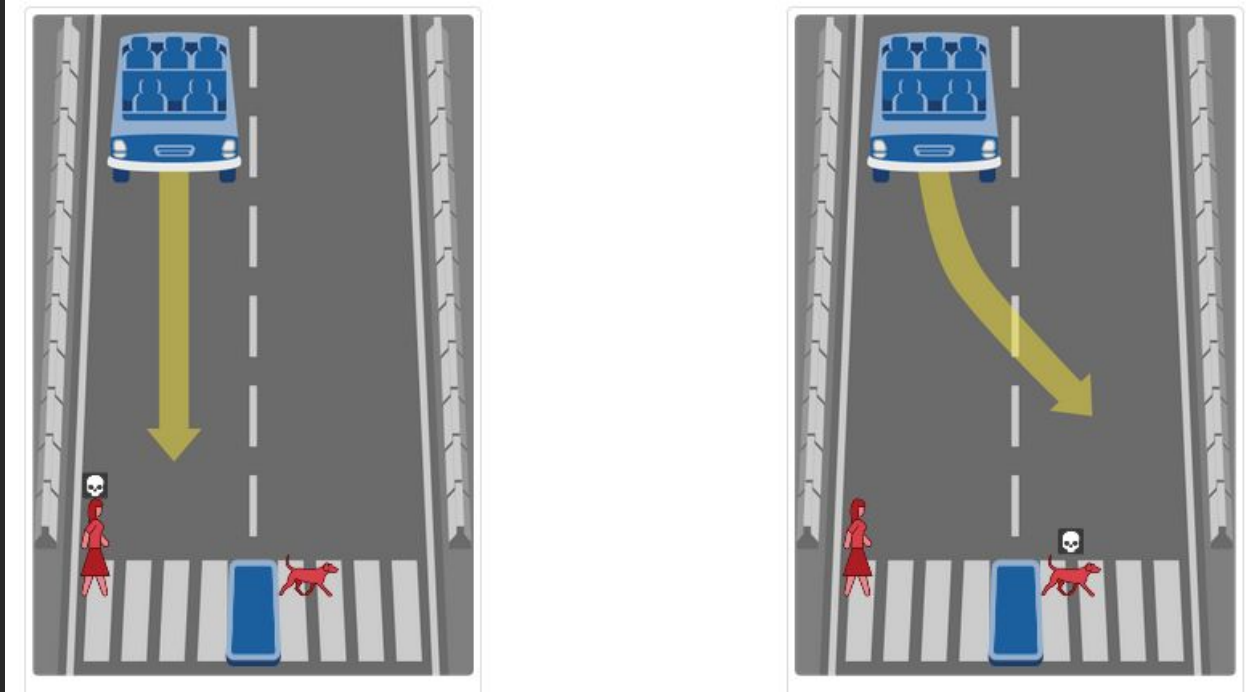
Topics

- Intro what data science ethics
- Bias & fairness

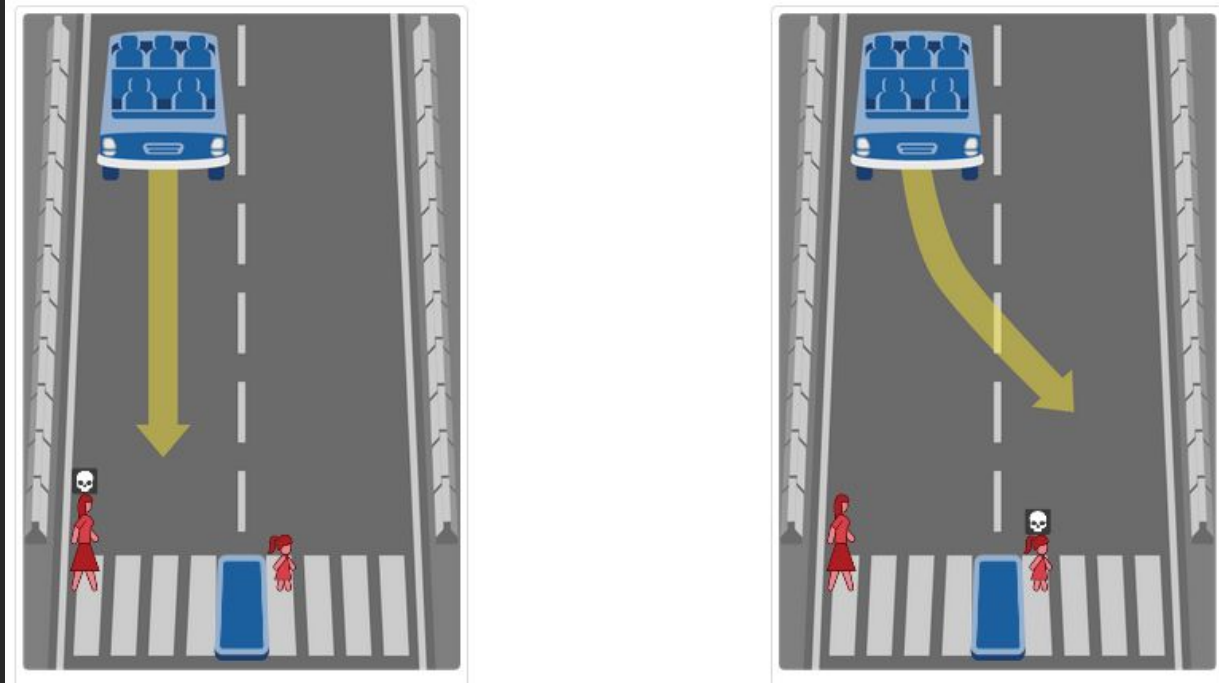
Ethical dilemma



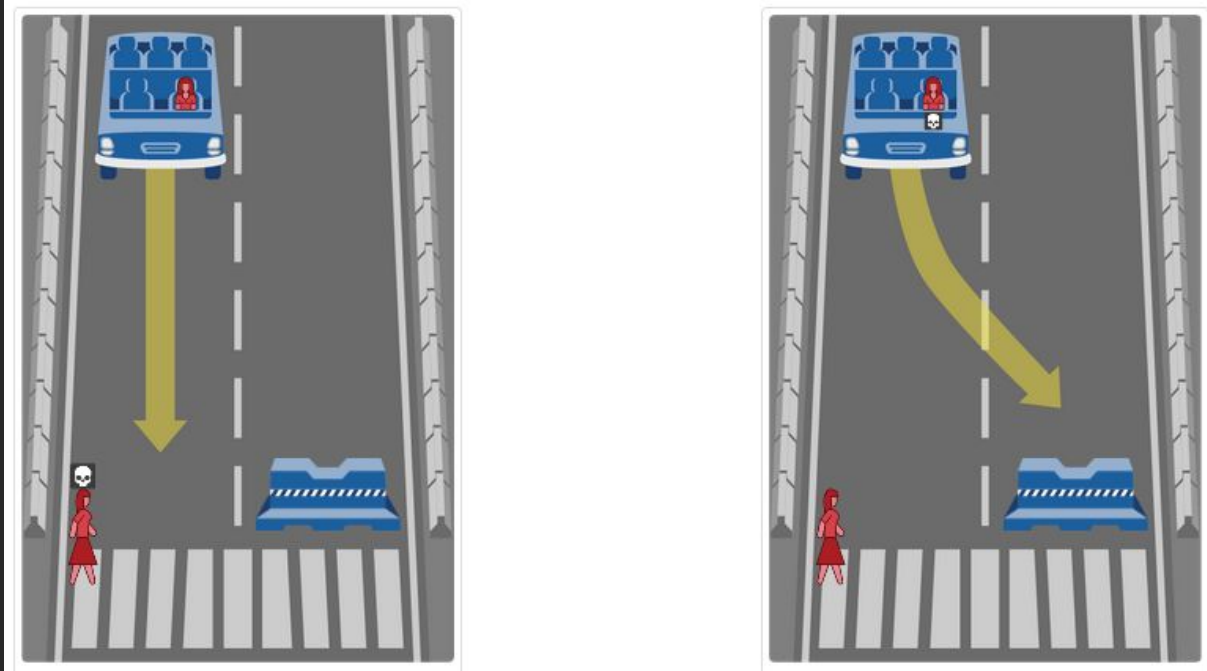
Ethical dilemmas



Ethical dilemmas



Ethical dilemmas



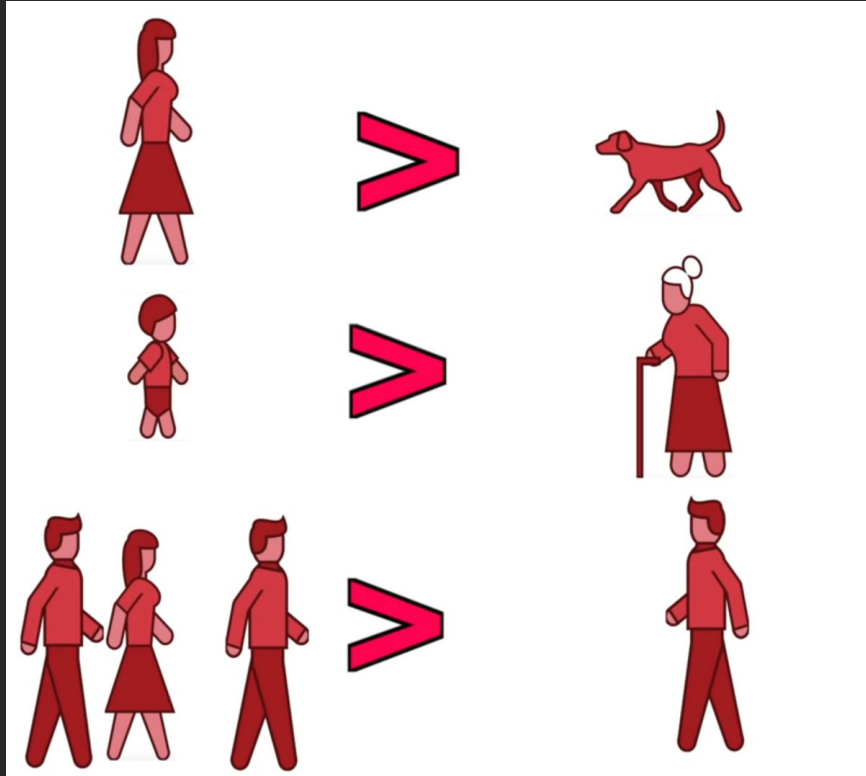
Self-driving cars



MIT self-driving car survey

- Online survey
- Millions of participants
- 233 countries

MIT self-driving car survey



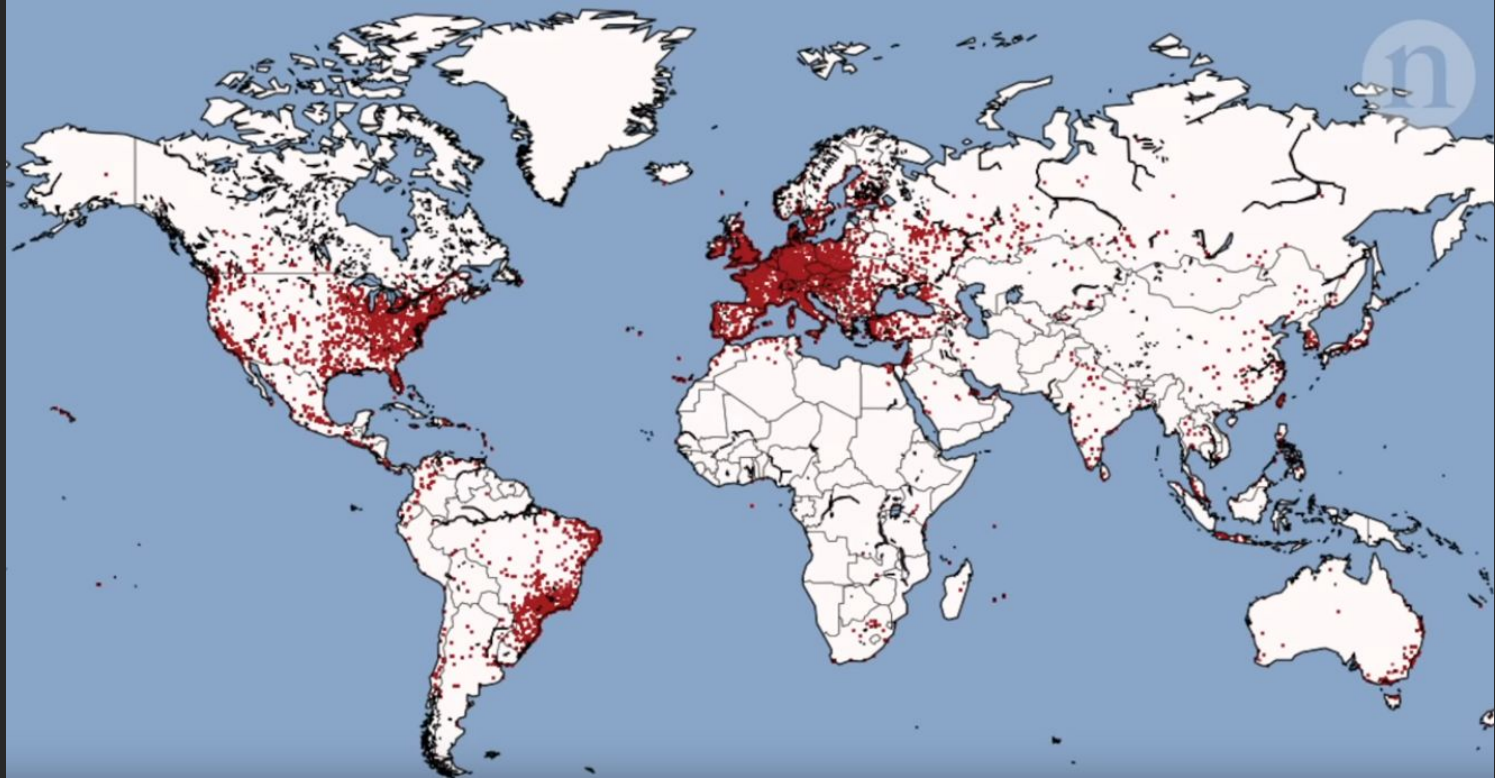
MIT self-driving car survey

- Many countries in east Asia, put a similar value on the elderly and younger lives.
- Many french speaking countries put a higher value on female lives.
- Countries with high income disparity generally valued the lives of “business people” higher than others.

Bias with the car survey data?

- Online survey
- Millions of participants
- 233 countries

Bias with the data?



Bias and fairness

Human bias

Sina Fazelpour
Carnegie Mellon University

Published online 11 April 2011 | Nature | doi:10.1038/news.2011.227

News

Hungry judges dispense rough justice

When they need a break, decision-makers gravitate towards the easy option.

Zoë Corbyn

Journal of Economic Perspectives—Volume 12, Number 2—Spring 1998—Pages 41–62

Evidence on Discrimination in Mortgage Lending

Helen F. Ladd

Science faculty's subtle gender biases favor male students

Corinne A. Moss-Racusin^{a,b}, John F. Dovidio^b, Victoria L. Brescoll^c, Mark J. Graham^{a,d}, and Jo Handelsman^{a,1}

^aDepartment of Molecular, Cellular and Developmental Biology, ^bDepartment of Psychology, ^cSchool of Management, and ^dDepartment of Psychiatry, Yale University, New Haven, CT 06520

Edited by Shirley Tilghman, Princeton University, Princeton, NJ, and approved August 21, 2012 (received for review July 2, 2012)



Health | Food | Fitness | Wellness | Parenting | Live Longer

Live TV | U.S. Edition | Search | Menu

Drowsy driving is a factor in almost 10% of crashes, study finds

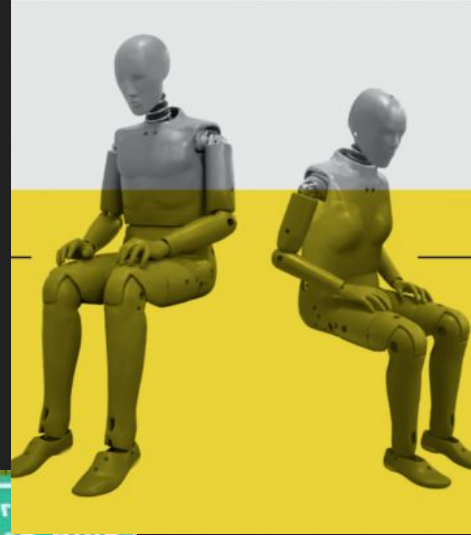
By Erin Gabriel, CNN

Updated 2:55 AM ET, Thu February 8, 2018



Can technologies
have bias?

Technological biases

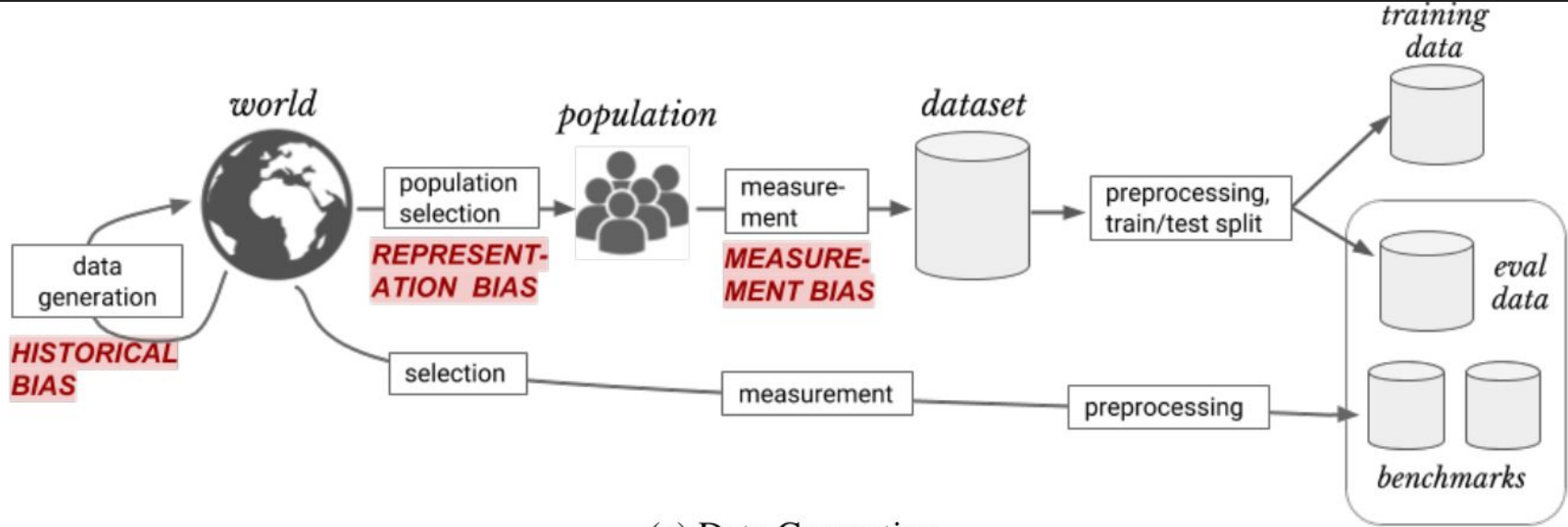


Humans have biases, so why does algorithmic bias matter?

- Cheap
- Scalable
- Automated
- Self-reinforcing
- Seemingly objective
- Often lacking appeals processes
- Does not just predict but also cause the future

Types of machine learning biases

Data & measurement biases



(a) Data Generation

Historical bias

Historical bias arises when there is a misalignment between world as it is and the values or objectives to be encoded and propagated in a model. It is a normative concern with the state of the world, and exists even given perfect sampling and feature selection.

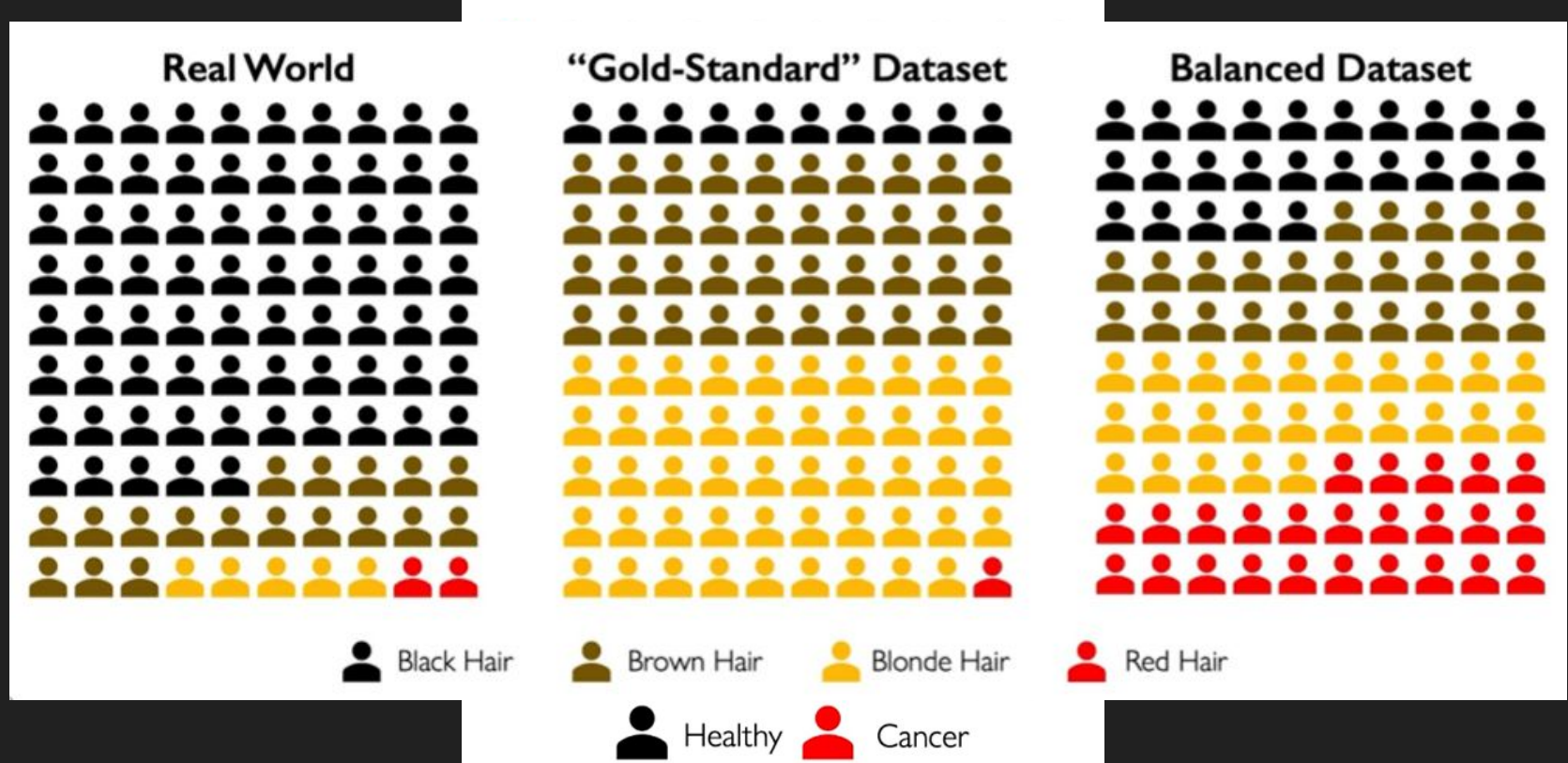
Example: image search In 2018, 5% of Fortune 500 CEOs were women (Zarya, 2018). Should image search results for “CEO” reflect that number? Ultimately, a variety of stakeholders, including affected members of society, should evaluate the particular harms that this result could cause and make a judgment. This decision may be at odds with the available data even if that data is a perfect reflection of the world. Indeed, Google has recently changed their Image Search results for “CEO” to display a higher proportion of women.

Representation bias

Representation bias arises while defining and sampling a development population. It occurs when the development population under-represents, and subsequently fails to generalize well, for some part of the use population.

1. **The sampling methods only reach a portion of the population.** For example, datasets collected through smartphone apps can under-represent lower-income or older groups, who are less likely to own smartphones. Similarly, medical data for a particular condition may be available only for the population of patients who were considered serious enough to bring in for further screening.
2. **The population of interest has changed or is distinct from the population used during model training.** Data that is representative of Boston, for example, may not be representative if used to analyze the population of Indianapolis. Similarly, data representative of Boston 30 years ago will likely not reflect today's population.

Representation bias

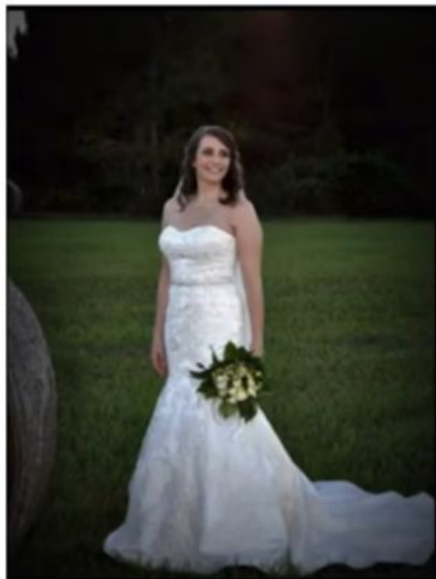


Representation bias

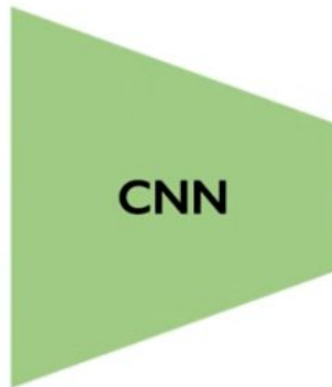


Representation bias

Bias in Image Classification



Ground Truth: Bride



CNN

CNN for image
classification.



Predicted Classes

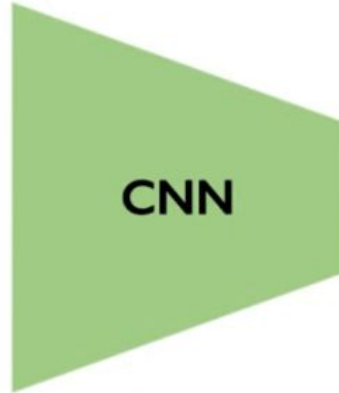
Bride
Dress
Ceremony
Woman
Wedding

Representation bias

Bias in Image Classification



Ground Truth: Bride



CNN

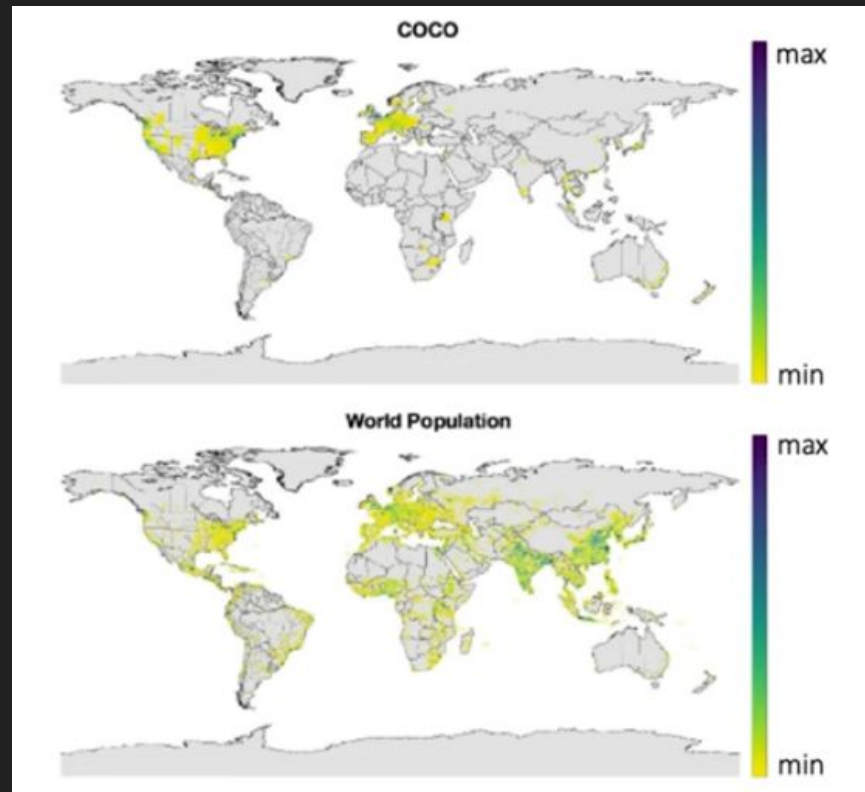
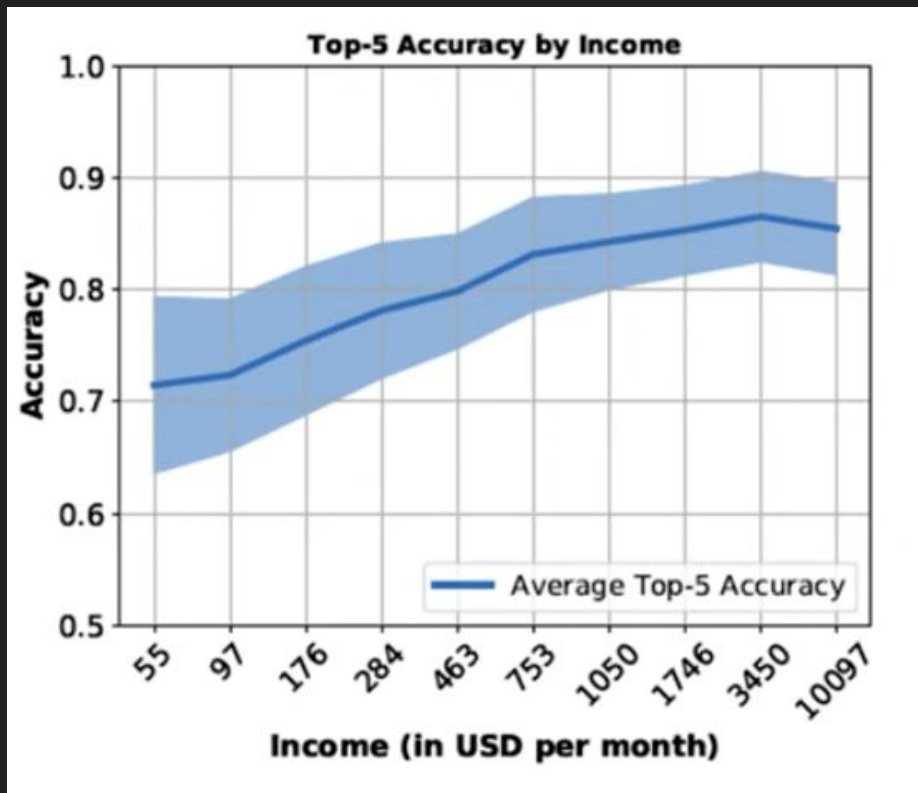
CNN for image
classification.



Predicted Classes

Clothing
Event
Costume
Red
Performance art

Representation bias



Measurement bias

Measurement Bias arises when choosing and measuring features and labels to use; these are often proxies for the desired quantities. The chosen set of features and labels may leave out important factors or introduce group- or input-dependent noise that leads to differential performance.

3. **The defined classification task is an oversimplification.** In order to build a supervised ML model, some label to predict must be chosen. Reducing a decision to a single attribute can create a biased proxy label because it only captures a particular aspect of what we really want to measure. Consider the prediction problem of deciding whether a student will be successful (e.g., in a college admissions context). Fully capturing the outcome of ‘successful student’ in terms of a single measurable attribute is impossible because of its complexity. In cases such as these, algorithm designers resort to some available label such as ‘GPA’ (Kleinberg et al. 2018), which ignores different indicators of success achieved by parts of the population.

1. **The measurement process varies across groups.** For example, if a group of factory workers is more stringently or frequently monitored, more errors will be observed in that group. This can also lead to a feedback loop wherein the group is subject to further monitoring because of the apparent higher rate of mistakes (Barocas and Selbst 2016).
2. **The quality of data varies across groups.** Structural discrimination can lead to systematically higher error rates in a certain group. For example, women are more likely to be misdiagnosed or not diagnosed for conditions where self-reported pain is a symptom (Calderone 1990). In this case, “*diagnosed* with condition X” is a biased proxy for “has condition X.”

Measurement bias (proxy metrics)

How to determine:

- A job candidate's fit?
- A student's potential?
- A patient's healthcare needs?
- A post's engagement?

Metrics as proxies

- Prior stroke
- Cardiovascular disease
- Accidental injury
- Benign breast lump
- Colonoscopy
- Sinusitis

Does Machine Learning Automate Moral Hazard and Error?[†]

By SENDHIL MULLAINATHAN AND ZIAD OBERMEYER*

Gaming/Short-term

- Cancelled scheduled operations to draft extra staff to ER
- Required patients to wait outside the ER, e.g. in ambulances
- Put stretchers in hallways and classified them as "beds"
- Hospital and patients reported different wait times.

WHAT'S MEASURED IS WHAT MATTERS: TARGETS AND GAMING IN
THE ENGLISH PUBLIC HEALTH CARE SYSTEM

GWYN BEVAN, CHRISTOPHER HOOD

Manipulation

TECHNOLOGY

How Facebook's Chaotic Push Into Video Cost Hundreds of Journalists Their Jobs

As media companies t
platform, they fired w

ALEXIS C. MADRIGAL AND ROBI

Als Are Designed to Maximize Watch Time

At YouTube, we used a complex AI to pursue a simple goal: maximize watch time. Google explains this focus in the following statement:

If viewers are watching more YouTube, it signals to us that they're happier with the content they've found. It means that creators are attracting more engaged audiences. It also opens up more opportunities to generate revenue for our partners.

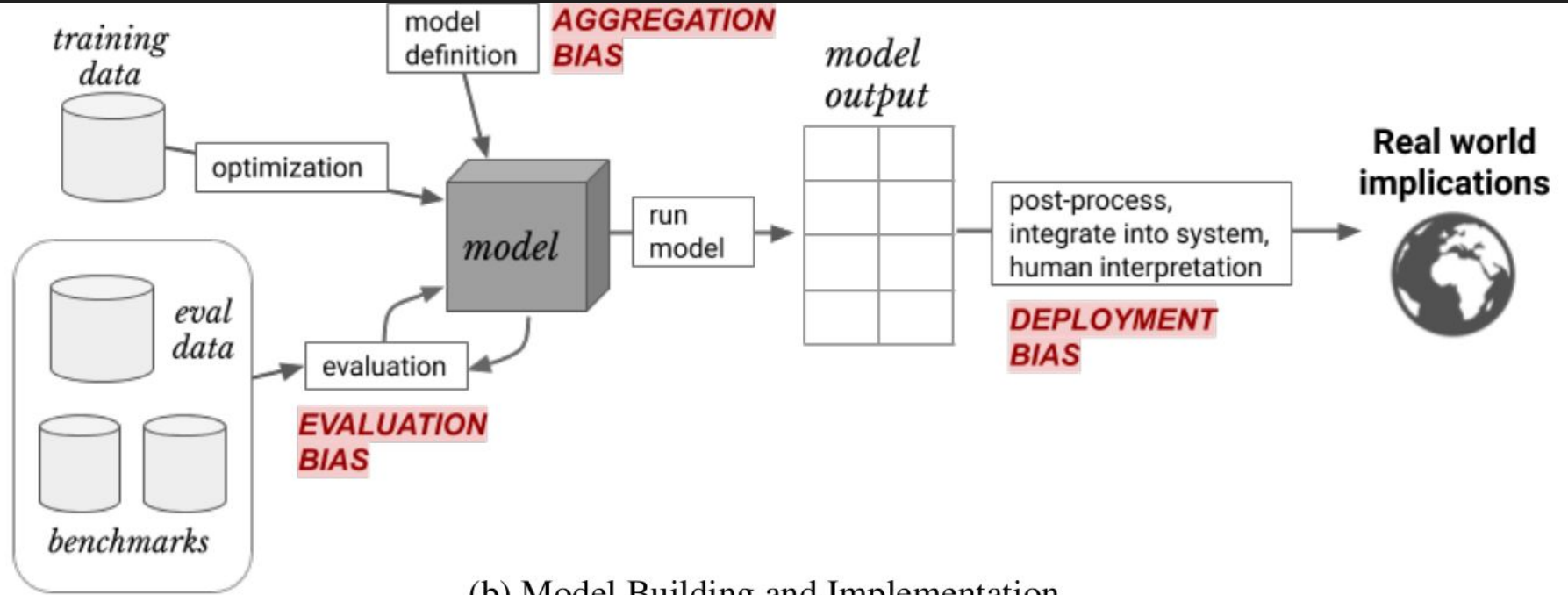
Goodhart's/Campbell's law

When a measure becomes a target
it ceases to be a good measure

When metrics are useful

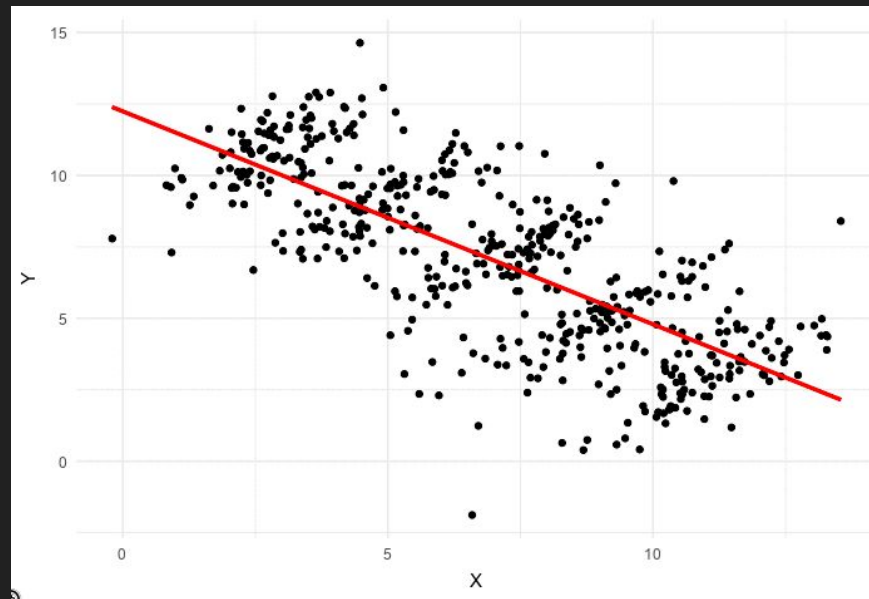
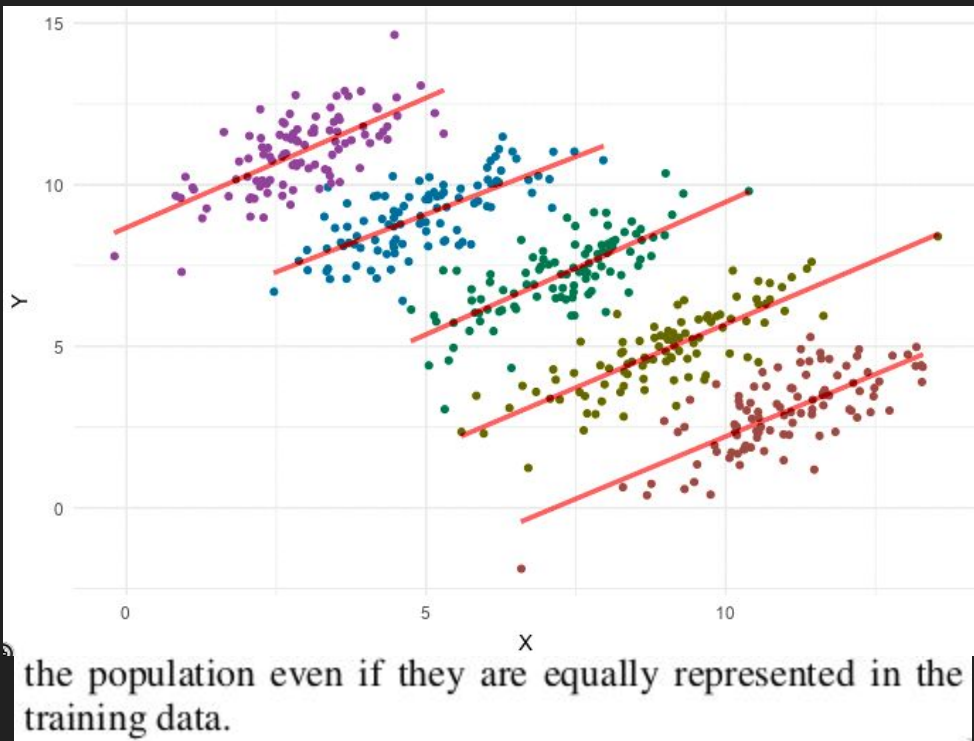
- Use multiple metrics instead of a single score
- Listening to qualitative first-person experiences
- Involve domain experts
- Consider biases and what can go wrong beforehand
- Transparent algorithms
- Assess and review your metrics regularly
- Legislative incentives for companies to adhere

Model & interpretation biases



(b) Model Building and Implementation

Aggregation bias

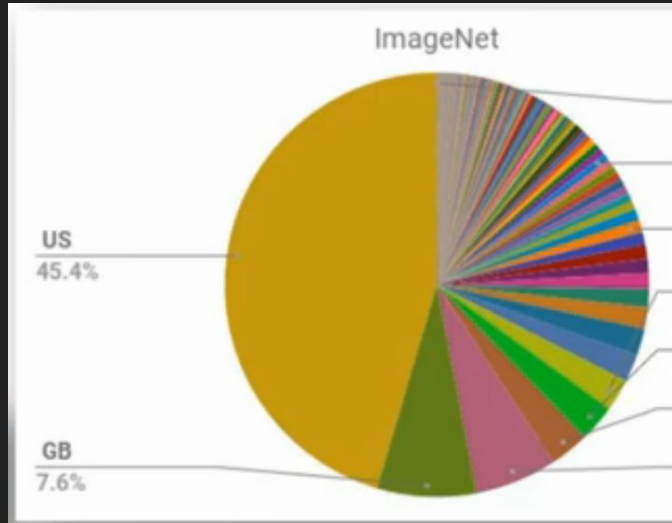


Evaluation bias

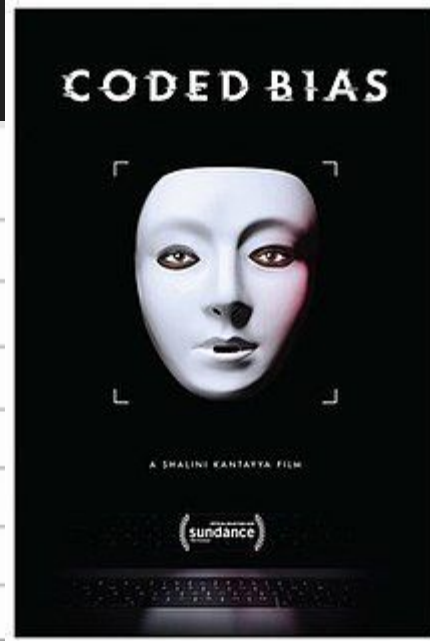
Evaluation bias occurs during model iteration and evaluation. It can arise when the testing or external benchmark populations do not equally represent the various parts of the use population. Evaluation bias can also arise from the use of performance metrics that are not appropriate for the way in which the model will be used.

		True condition	
	Total population	Condition positive	Condition negative
Predicted condition	Predicted condition positive	True positive	False positive, Type I error
	Predicted condition negative	False negative, Type II error	True negative

Evaluation bias



2/3 of ImageNet images from the West (Shankar et al, 2017)



Example: underperformance of commercial facial recognition algorithms Buolamwini and Gebru (2018) point out the drastically worse performance of commercially-used facial analysis algorithms (performing tasks such as gender-detection) on dark-skinned females. Looking at facial analysis benchmark datasets, it becomes evident why such algorithms were considered ap- use – just 7.4% and 4.4% of the images in datasets such as Adience and IJB-A are of dark-

of IJB-A images are skinned women



Deployment bias

Deployment Bias occurs after model deployment, when a system is used or interpreted in inappropriate

Example: Risk Assessment Tools in Practice

mic risk assessment tools are models intended

person's likelihood of c
tice, however, these tool
such as to help determin
(2018) describes the ha
ment tools for actuarial
creased incarceration on
tics. Stevenson (2018) t
in-depth study of the dep
Kentucky, demonstrates
tion created highly unre
efits and consequences a

Measure potential, not pedigree.

Don't judge a job seeker by their resume alone.

We use behavioral assessments to evaluate job seekers. Rather than focusing on backward-looking resumes or self-reported questionnaires, we collect objective behavioral data that measures a job seeker's true potential.

[See our assessments](#) ↗

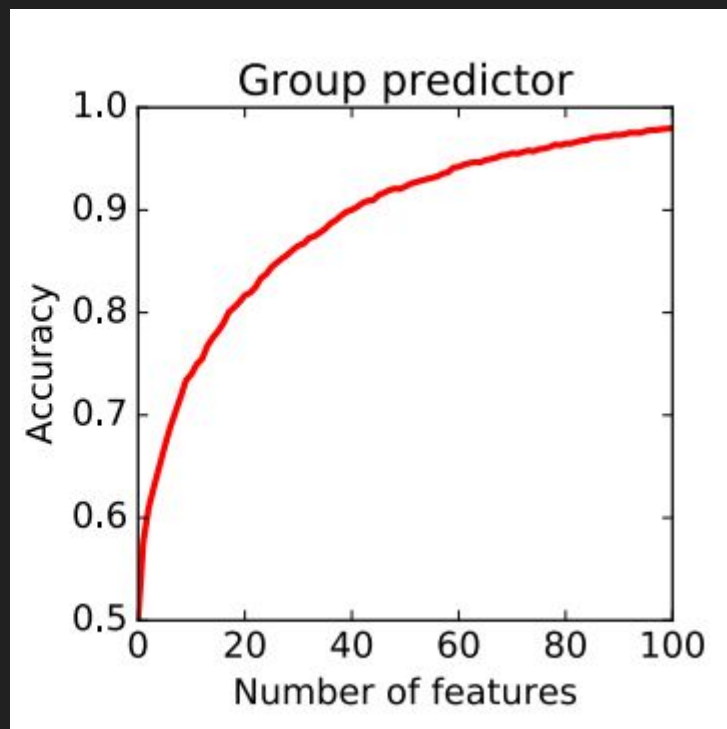
Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission

Johannes Gehrke
Microsoft
johannes@microsoft.com

Noémie Elhadad
Columbia University
noemie.elhadad@columbia.edu

Quantifying fairness

No fairness through unawareness



What is fair?

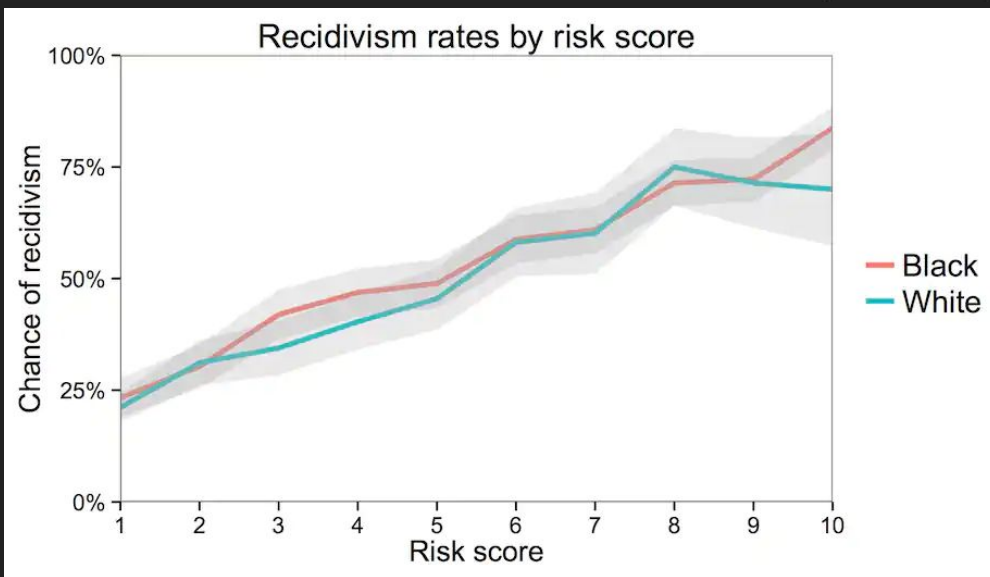
Propublica

COMPAS

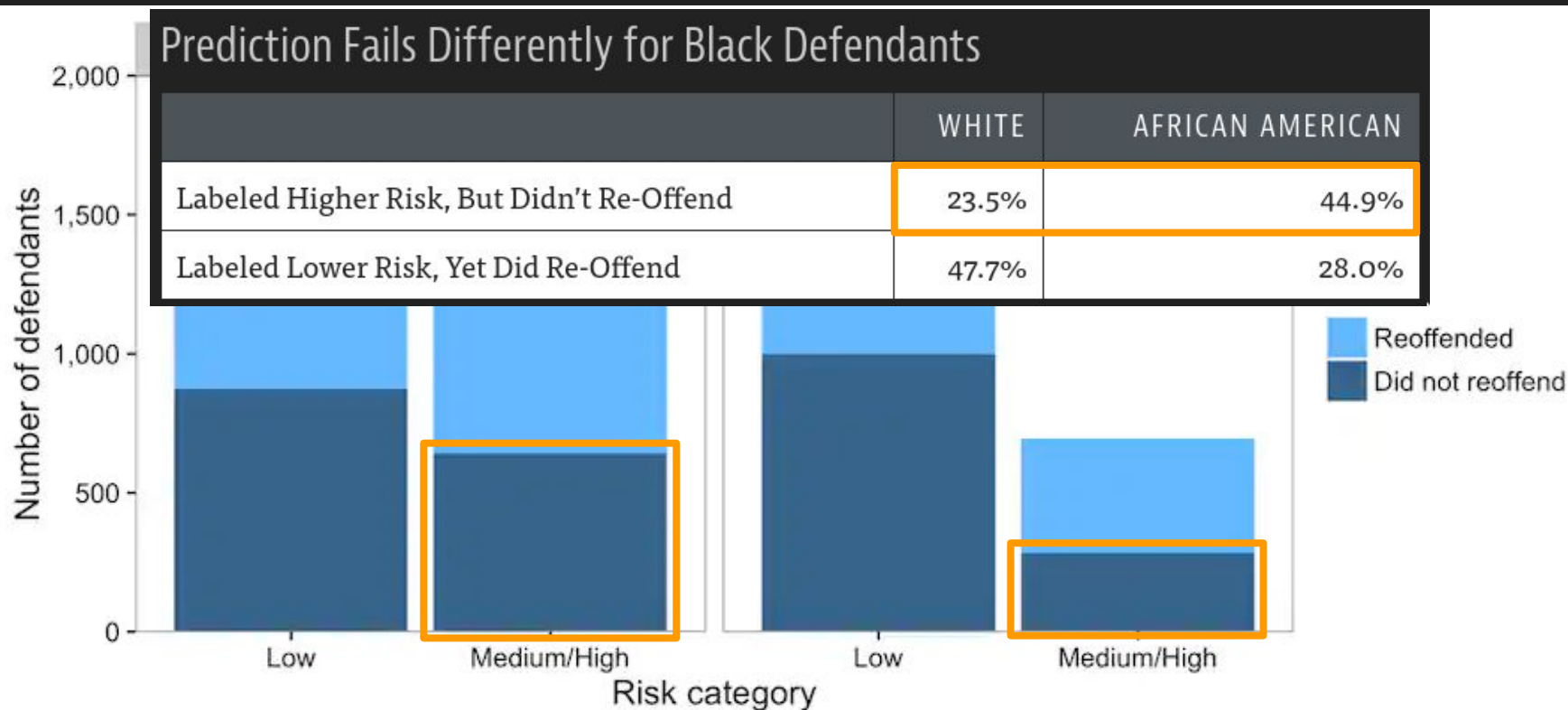
Northpointe

Prediction Fails Differently for Black Defendants

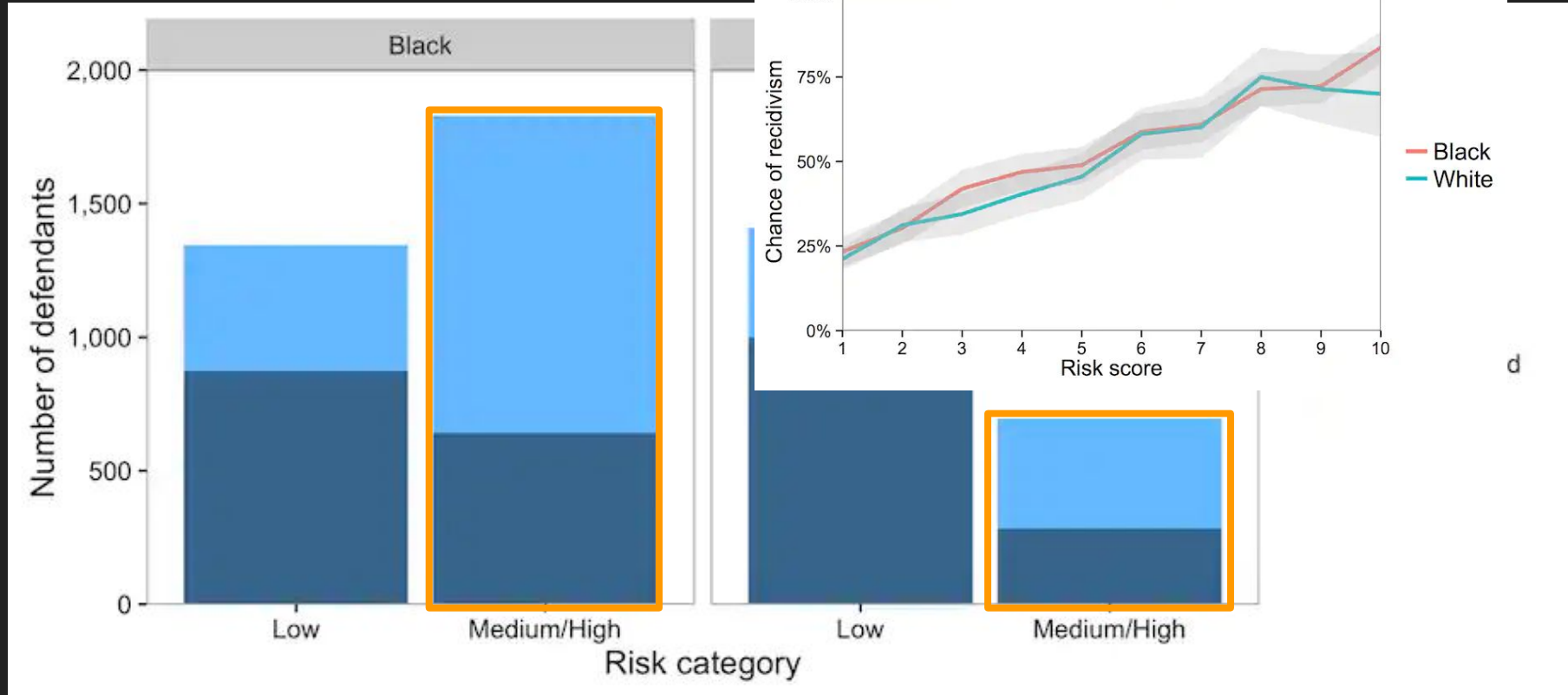
	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%



Propublica's fairness metric



Northpoint's fairness metric



Which fairness
metric is correct?

Confusion matrix & fairness metrics

		True condition			
Total population		Condition positive	Condition negative	Prevalence = $\frac{\Sigma \text{Condition positive}}{\Sigma \text{Total population}}$	Accuracy (ACC) = $\frac{\Sigma \text{ True positive} + \Sigma \text{ True negative}}{\Sigma \text{ Total population}}$
Predicted condition	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Predicted condition negative}}$
			True positive rate (TPR), Recall, Sensitivity, probability of detection, $\text{Power} = \frac{\Sigma \text{ True positive}}{\Sigma \text{ Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\Sigma \text{ False positive}}{\Sigma \text{ Condition negative}}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$
		False negative rate (FNR), Miss rate = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) $= \frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$	

Gaming fairness metrics

Detain everyone



Arrest more low
risk individuals
in orange group!

Detention rate	False pos. rate
38%	25%
61% 42%	42% 26%

Find the underlying problem

Failure to appear in court

One approach: Predict failure to appear, jail if risk is high.

Alternative: Recognize that people fail to appear in court due to lack of child care and transportation, work schedules, or too many court appointments. Implement steps to mitigate these issues.

Alternative is part of the Harris County Lawsuit settlement: *"require Harris County to provide free child care at courthouses, develop a two-way communication system between courts and defendants, give cell phones to poor defendants and pay for public transit or ride share services for defendants without access to transportation to court."* (Source: [Houston Chronicle, April 2019](#))



Take-home messages

- Rather than trying to understand **IF** your model is fair, try to understand **HOW** it is unfair.
- Look at multiple fairness metrics to diagnose potential issues among diverse stakeholders.
- Be careful when optimizing on fairness criteria.
- Use domain expertise to try to understand causal relationships underlying the observed results.

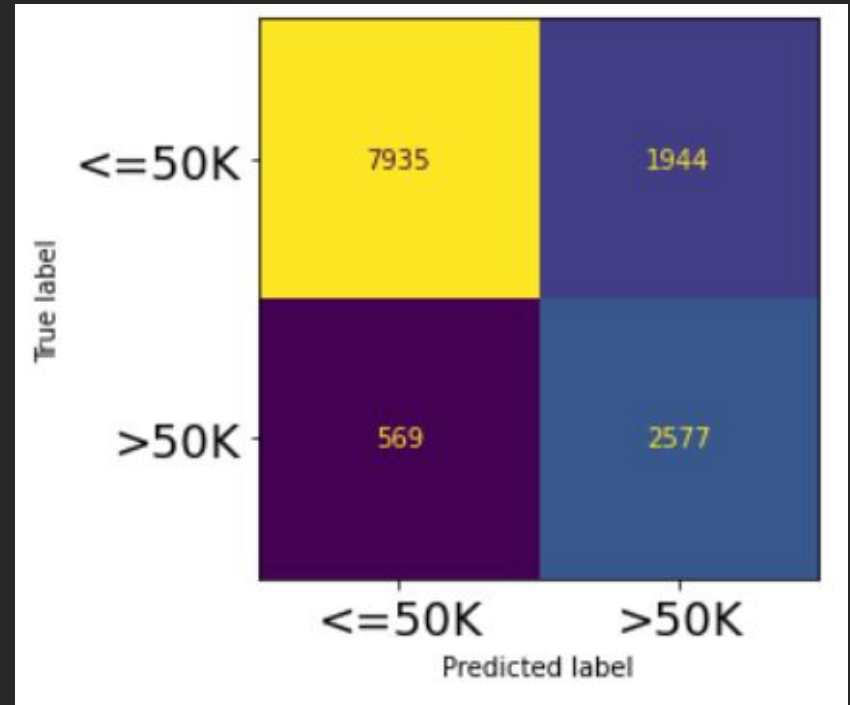
Break out room reflection

We are building a model in order to determine who to approve for a loan based on their predicted income. We are using US Census data that includes information about education, marital status, sex, race, etc. We want to predict who makes $>50k$ / year and ensure that our model doesn't have a bias towards either men or women in the prediction.

On the next few slides you will find the model accuracy and confusion matrix, both overall and for women and men separately. Do you think this model is fairly treating the two groups? Why/why not?

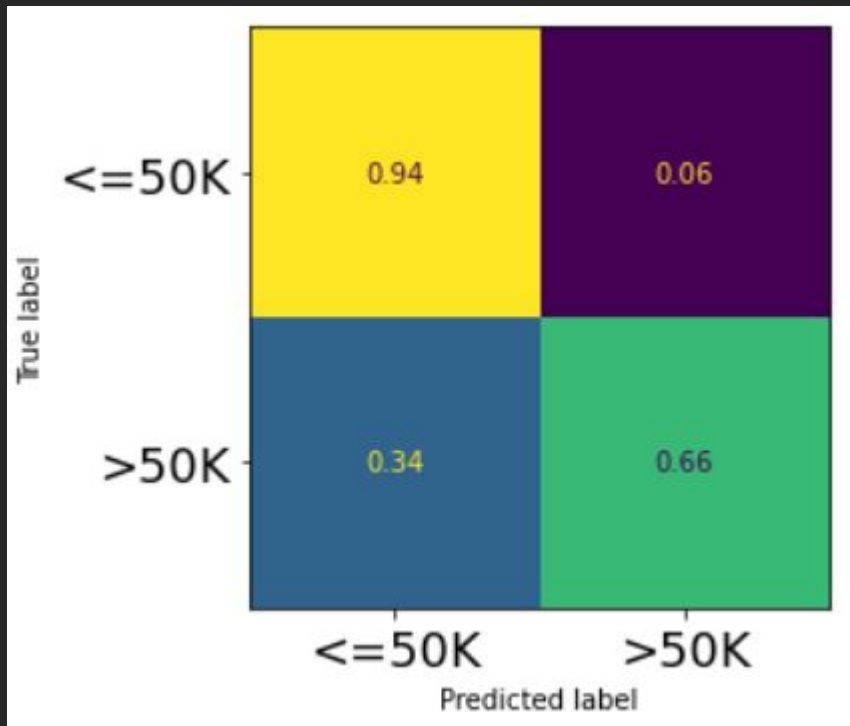
Accuracy & confusion matrix

- Overall accuracy ~85%
 - Men ~75%
 - Women ~90%



Confusion matrix by sex

Women



Men

