

AI in Medicine I

Trustworthy AI: Fairness & Bias

Prof. Christian Wachinger

Lab for AI in Medical Imaging (www.ai-med.de)

Institute of Radiology, Klinikum rechts der Isar, TUM

School of Computation, Information and Technology, TUM

Student projects in my group

The screenshot shows a website for "Artificial Intelligence in Medical Imaging". The header includes navigation links for Team, Publications, Research, Teaching (which is highlighted in red), Contact, and News. Below the header, a dark bar displays "Student Projects" and a search bar with a "Go" button. The main content area is divided into sections: "Open Projects" (with links to Master Thesis: Alzheimer's Disease Prediction, Master Thesis: Diabetes Prediction, and Master Thesis: Medical Image Segmentation), "Running Projects" (with a link to Master Thesis: Analysis of brain structure in premature-born adults), and "Finished Projects" (with links to Master Thesis: Explaining Deep Survival Analysis Models for Heterogenous Data and Master Thesis: Automatic Feature Interaction Learning for Alzheimer's Disease Diagnosis using Factorization Models).

Artificial Intelligence in Medical Imaging

Team Publications Research **Teaching** Contact News

Student Projects

Search Go

Open Projects:

[Master Thesis: Alzheimer's Disease Prediction](#)

[Master Thesis: Diabetes Prediction](#)

[Master Thesis: Medical Image Segmentation](#)

Running Projects:

[Master Thesis: Analysis of brain structure in premature-born adults](#)

Finished Projects:

[Master Thesis: Explaining Deep Survival Analysis Models for Heterogenous Data](#)

[Master Thesis: Automatic Feature Interaction Learning for Alzheimer's Disease Diagnosis using Factorization Models](#)

ai-med.de/research/student-projects/

Agenda

1. Trustworthy AI
2. Examples of unfair AI
3. Bias & Fairness
4. Fairness criteria
5. Metrics & Mitigation

Learning Outcomes

- **Understand importance** of trustworthy AI
- Awareness of **sources of bias** in AI
- **Quantify fairness** of prediction outcome
- Know strategies to **overcome bias**

Trustworthy

AI

Why should we care about Trustworthy AI?



👍 Maximize the benefits of AI systems

👎 Minimize the risks

AI use needs trust of the society

Ethics guidelines for trustworthy AI

European Commission



Framework for Trustworthy AI

Trustworthy AI

Lawful AI

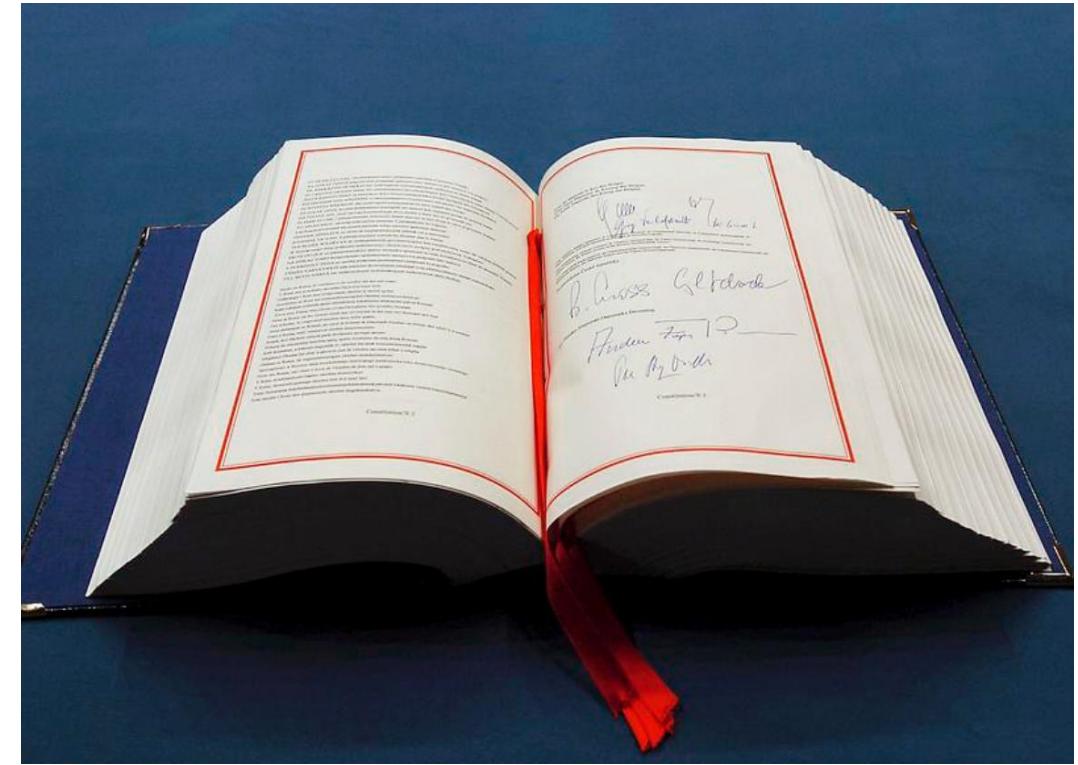
Ethical AI

Robust AI

EU Charter of Fundamental Rights

... brings together the most important personal freedoms and rights enjoyed by citizens of the EU

- Dignity
- Freedoms
- Equality
- Solidarity
- Citizen's rights
- Justice



Ethics guidelines for trustworthy AI

European Commission



Framework for Trustworthy AI

Trustworthy AI

Lawful AI

Ethical AI

Robust AI

Foundations of Trustworthy AI

Adhere to ethical principles based on fundamental rights



4 Ethical Principles

- Respect for human autonomy
- Prevention of harm
- Fairness
- Explicability

Ethics guidelines for trustworthy AI

European Commission



- Respect for human autonomy
- Prevention of harm
- Fairness
- Explicability

Processes need to be transparent, and decisions explainable

Humans must be able to keep self-determination.
AI systems should not deceive, or manipulate humans.

AI systems should not harm humans.
Protection of human dignity as well as mental and physical integrity.

Fairness

Ensuring equal and just distribution of both benefits and costs.

Ensuring that individuals and groups are free from unfair bias.

Unfair bias could lead to

- Marginalization of vulnerable groups
- Amplification of prejudice and discrimination.

Stated principles remain **abstract ethical prescriptions**. AI practitioners can hence not be expected to find the right solution based on the principles above, yet they should **approach ethical dilemmas via reasoned, evidence-based reflection** rather than intuition.

Discussion among students

Do you know UNFAIR or BIASED AI systems?

Consider general scenarios and medical ones.

5 minutes

Examples
of unfair AI

AI systems in the courtroom

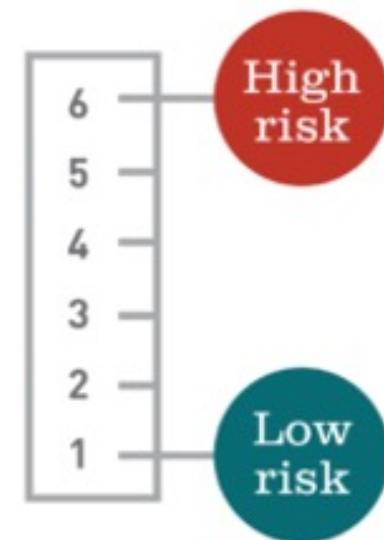
Released to await trial

Detained in jail to await trial



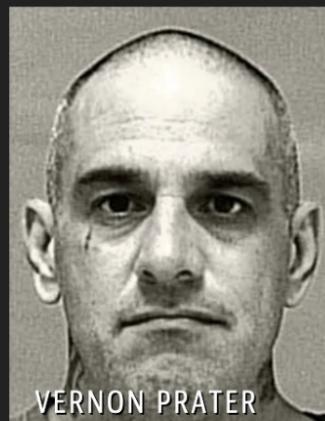
AI support: Predict risk for

- Future crimes
- Failure to appear in court



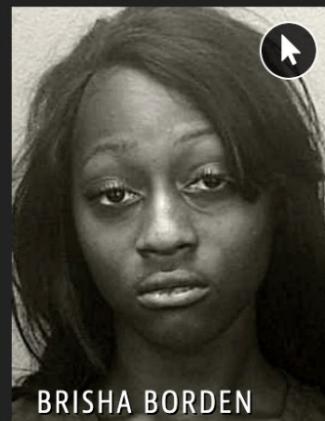
Controversy over the COMPAS Score (Propublica)

Two Petty Theft Arrests



VERNON PRATER

LOW RISK

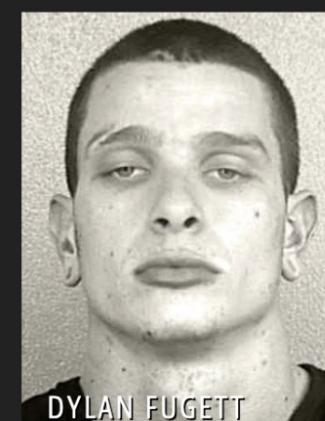


BRISHA BORDEN

HIGH RISK 8

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

Two Drug Possession Arrests



DYLAN FUGETT

LOW RISK 3

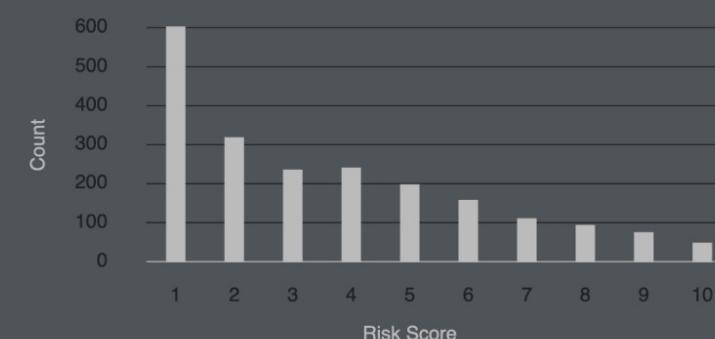


BERNARD PARKER

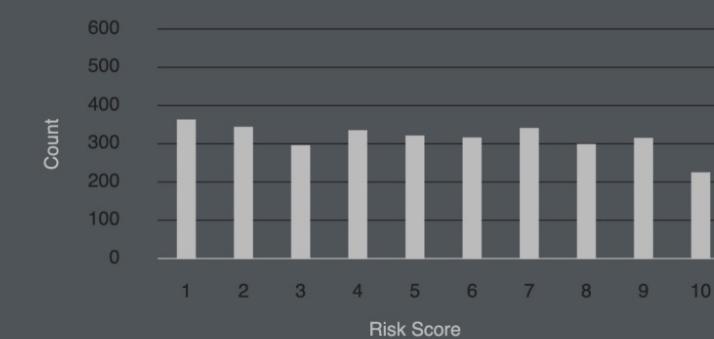
HIGH RISK 10

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

White Defendants' Risk Scores



Black Defendants' Risk Scores



AI systems in the courtroom

“... 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.”

The Public Safety Assessment (PSA)

Following a person's arrest, a judge must decide whether that person should:



A judge considers many factors in making this decision. One tool that judges may use to help make this decision is the PSA.



The PSA produces a score that represents the likelihood that a defendant who is released before trial will commit a new crime or will fail to appear for a future court appearance.

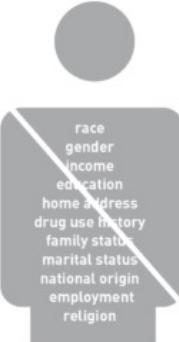


The PSA also flags the small number of defendants who pose an elevated risk of committing a crime of violence if released before trial.

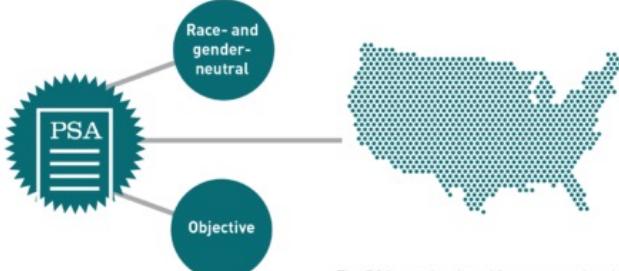
The PSA score is calculated based on nine factors.

Current violent offense	Pending charge at the time of the offense	Prior misdemeanor conviction
Prior felony conviction	Prior violent conviction	Prior failure to appear pretrial in past 2 years
Prior failure to appear pretrial older than 2 years	Prior sentence to incarceration	Age at current arrest

The PSA does NOT look at any of the following factors:



The PSA provides information that is race- and gender-neutral. It helps guide pretrial decision making in an effort to increase safety, reduce taxpayer costs, and enhance fairness and efficiency in the system.



The PSA score is not the only information that a judge considers, and the final decision will always be made by a judge.



The PSA was developed from research using data from across the United States.



BECOME
A MEMBER

11010
00111

01010
111001
010101
00101
01001

Illustration: The Intercept

A BAIL REFORM TOOL INTENDED TO CURB MASS INCARCERATION HAS ONLY REPLICATED BIASES IN THE CRIMINAL JUSTICE SYSTEM

Kentucky was an early adopter of risk assessments in an effort to release more people without bail. But the algorithms are reproducing systemic inequities.



Bryce Covert

July 12 2020, 2:00 p.m.



Recent stories

Major projects

Coverage areas

News414

Audio stories

Store

Be your own watchdog



BEYOND BAIL | JUSTICE & SAFETY

Debate rages over whether pretrial risk assessments are racially biased

Some researchers say there is not enough evidence to know whether the tools will reduce or exacerbate discrimination that already exists in the justice system

The Public Safety Assessment (PSA)

Following a person's arrest, a judge must decide whether that person should:



A judge considers many factors in making this decision. One tool that judges may use to help make this decision is the PSA.



The PSA produces a score that represents the likelihood that a defendant who is released before trial will commit a new crime or will fail to appear for a future court appearance.

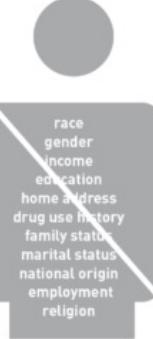


The PSA also flags the small number of defendants who pose an elevated risk of committing a crime of violence if released before trial.

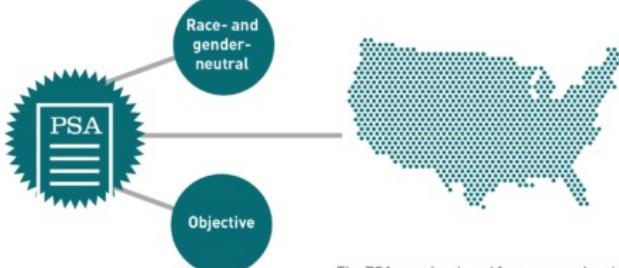
The PSA score is calculated based on nine factors.

Current violent offense	Pending charge at the time of the offense	Prior misdemeanor conviction
Prior felony conviction	Prior violent conviction	Prior failure to appear pretrial in past 2 years
Prior failure to appear pretrial older than 2 years	Prior sentence to incarceration	Age at current arrest

The PSA does NOT look at any of the following factors:



The PSA provides information that is race- and gender-neutral. It helps guide pretrial decision making in an effort to increase safety, reduce taxpayer costs, and enhance fairness and efficiency in the system.



The PSA score is not the only information that a judge considers, and the final decision will always be made by a judge.



The PSA was developed from research using data from across the United States.

"Prediction looks to the past to make guesses about the future. In a racially stratified world, any method of prediction will project the inequalities of the past into the future."

https://digitalcommons.law.uga.edu/fac_artchop/1293/

AI in recruiting

- Amazon hiring tool to mechanize search for top talent
- Found to be biased against women, penalizing candidates whose resumes included the word “women’s”.



Gender Shades Project

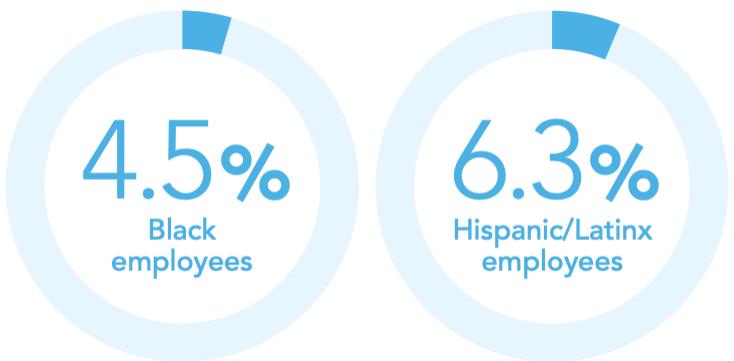


Lack of diversity in
the training data

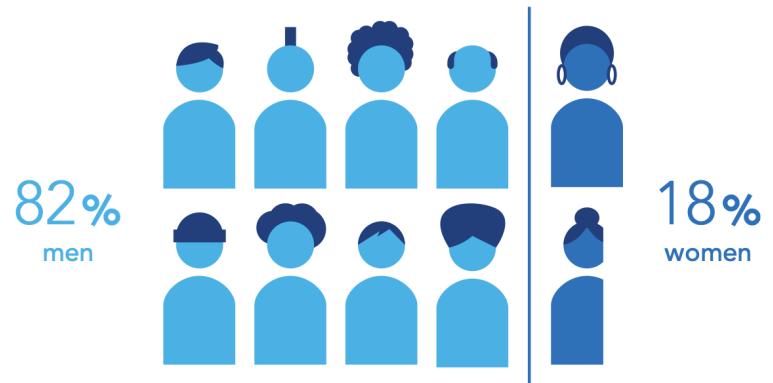
Gender Classifier	Darker Male	Darker Female	Lighter Male	Lighter Female	Largest Gap
Microsoft	94.0%	79.2%	100%	98.3%	20.8%
FACE++	99.3%	65.5%	99.2%	94.0%	33.8%
IBM	88.0%	65.3%	99.7%	92.9%	34.4%

Gender Shades Project

Employees at Microsoft



Researchers at leading AI conferences



The Coded Gaze: Unmasking Algorithmic Bias



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

ImageNet et al. criticism

Large datasets are instrumental for DL success

Dataset	Number of images (in millions)	Number of categories (in thousands)	Number of consensual images
JFT-300M ([41])	300+	18	0
Open Images ([50])	9	20	0
Tiny-Images ([79])	79	76	0
Tencent-ML ([89])	18	11	0
ImageNet-(21K,11k ¹ ,1k) ([70])	(14, 12, 1)	(22, 11, 1)	0
Places ([93])	11	0.4	0

Grabbing images from the internet
without consensus esp. children

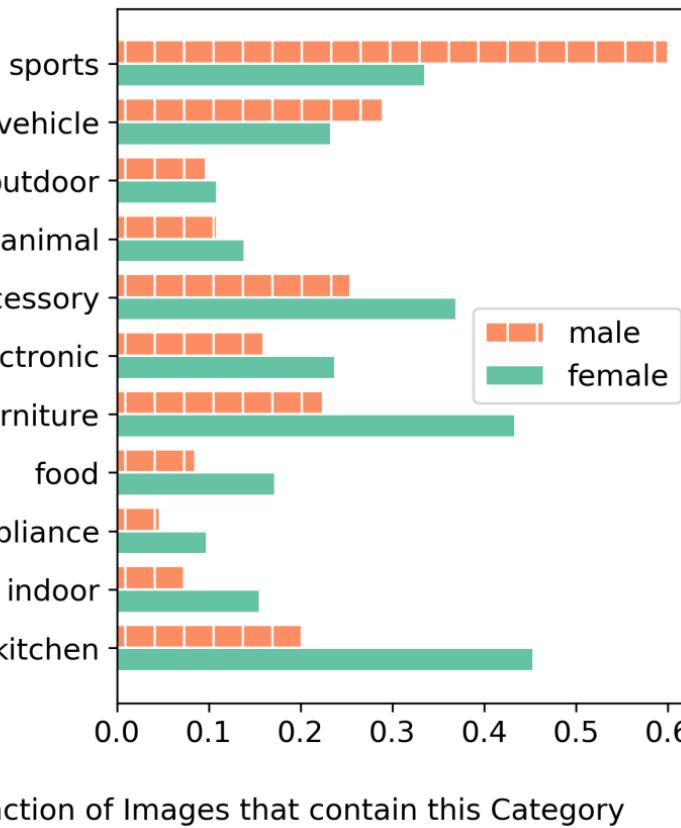
Categorize people based on
appearance

Dataset audit cards for curators: publish the goals, curation procedures,
known shortcomings and caveats

Measuring Bias in Visual Datasets

In which image contexts appear men and women?

Object Category



AI in health care

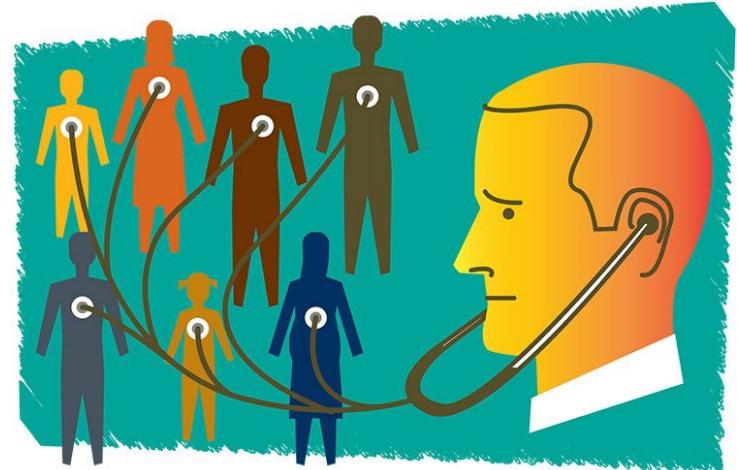
Task: Identify patients for high-risk care management

Algorithm: Predicts risk score for each patient

Policy: Patients above the 97th percentile are enrolled

Algorithm input: demographics (e.g., age, sex), insurance type, diagnosis and procedure codes, medications, and detailed costs. **Excludes race.**

Output: health costs (costs as a proxy for health needs)



AI in health care

Task: Identify patients for high-risk care management

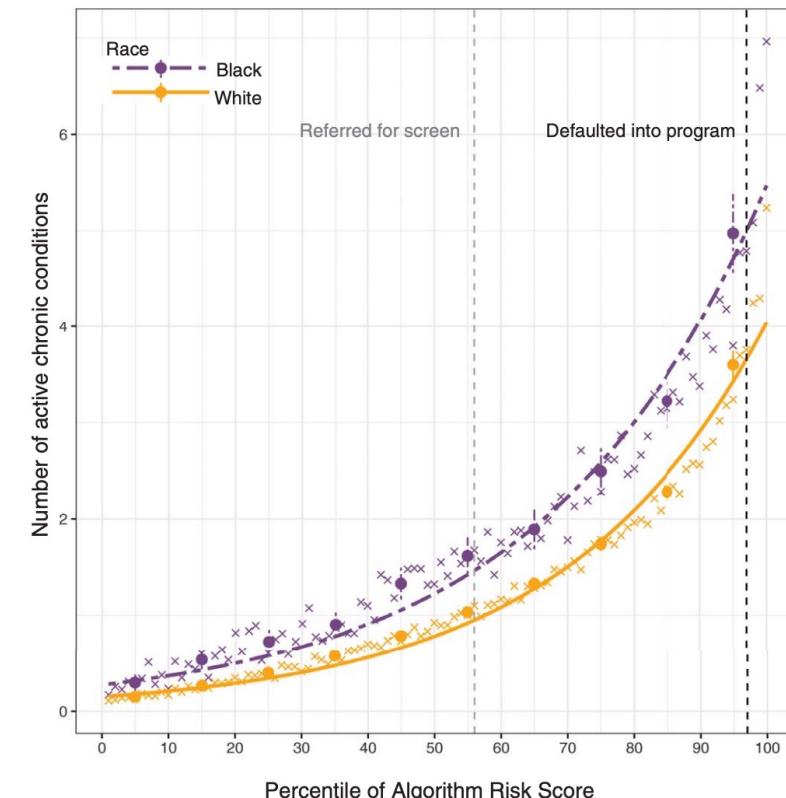
Algorithm: Predicts risk score for each patient

Policy: Patients above the 97th percentile are enrolled

Algorithm input: demographics (e.g., age, sex), insurance type, diagnosis and procedure codes, medications, and detailed costs. **Excludes race.**

Output: health costs (costs as a proxy for health needs)

At the same level of algorithm-predicted risk, Blacks have significantly more illness burden than Whites



Black patients generate lesser medical expenses (unequal access to care)
=> algorithm learned that Black patients are healthier than equally sick White patients

Bias in Medical Imaging

Skin lesions



AI skin cancer diagnoses risk being less accurate for dark skin - study

Research finds few image databases available to develop technology contain details on ethnicity or skin type



Studies suggest image recognition technology can classify skin cancers as successfully as humans. Posed by model. Photograph: ChesiireCat/Getty Images/Stockphoto

Nicola Davis
Science correspondent

@NicolaKSDavis
Tue 9 Nov 2021 23.30 GMT



Latest

The Atlantic

HEALTH

AI-Driven Dermatology Could Leave Dark-Skinned Patients Behind

Machine learning has the potential to save thousands of people from skin cancer each year—while putting others at greater risk.

By Angela Lashbrook

MOTHERBOARD
TECH BY VICE

Google's New Dermatology App Wasn't Designed for People With Darker Skin

The company trained the system to recognize different skin conditions. But like Google itself, the app's data has a diversity problem.

TF By Todd Feathers

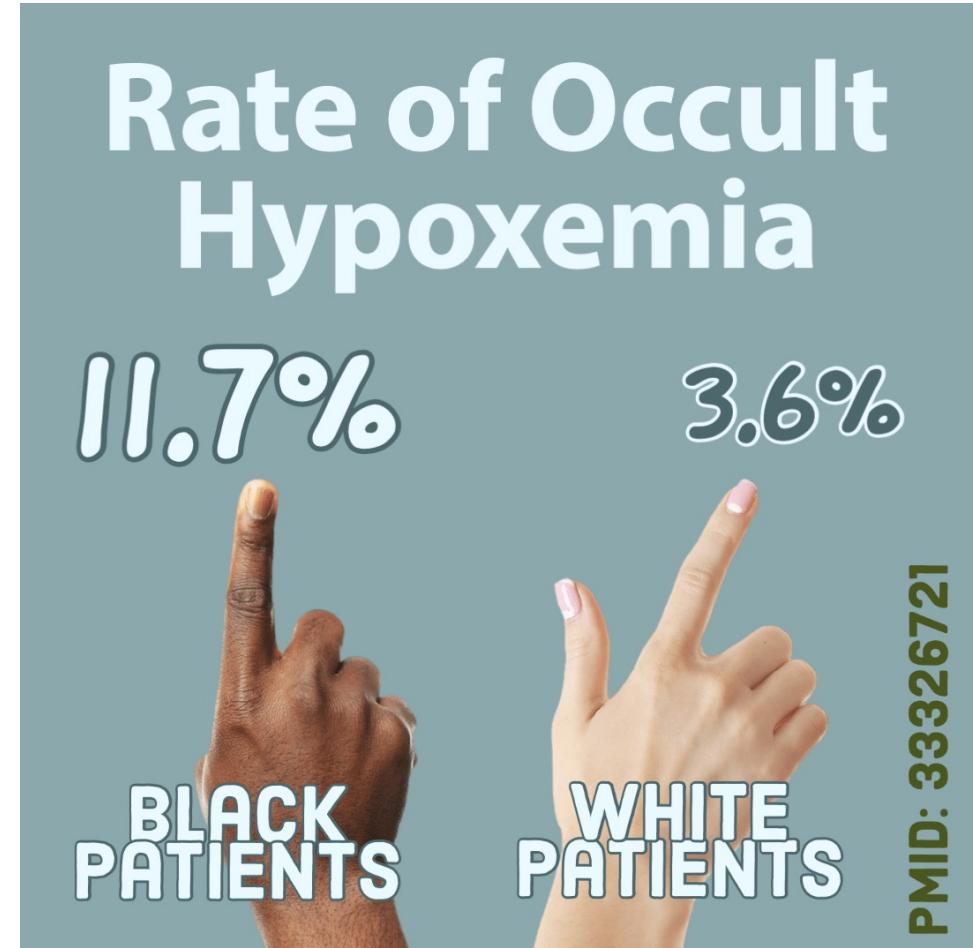
May 20, 2021, 3:40pm [Share](#) [Tweet](#) [Snap](#)

Pulse Oximeters: False elevated oxygen saturation in dark-skinned persons



Occult hypoxemia: falsely elevated oxygen saturation when measured with pulse oximetry, masking hypoxemia

<https://jamanetwork.com/journals/jama/fullarticle/2800468>



Impact of user biases while using AI tools

npj | digital medicine

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [npj digital medicine](#) > [perspectives](#) > [article](#)

Perspective | [Open Access](#) | Published: 28 December 2022

AI in the hands of imperfect users

[Kristin M. Kostick-Quenet](#)  & [Sara Gerke](#)

[npj Digital Medicine](#) 5, Article number: 197 (2022) | [Cite this article](#)

2872 Accesses | 34 Altmetric | [Metrics](#)

<https://www.nature.com/articles/s41746-022-00737-z>



<https://www.healthcareitnews.com/news/implementation-best-practices-dealing-complexity-ai>

**BIAS &
FAIRNESS**

Bias

Biases are **cognitive shortcuts** that can result in judgments which lead to discriminatory practices.

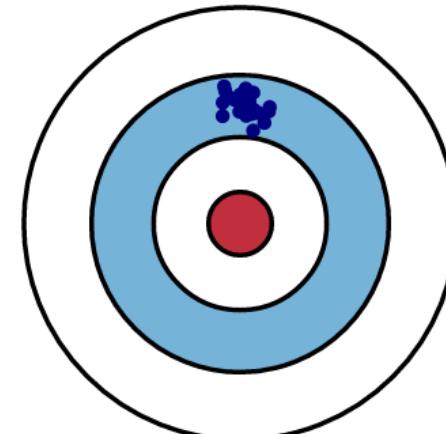
Bias in statistics:

Bias is the tendency of a statistic to **systematically** overestimate or underestimate a population parameter

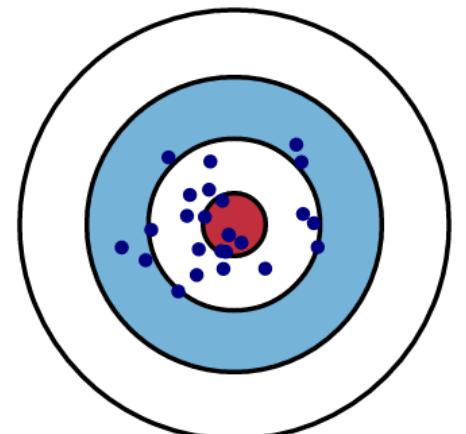
Bias-Variance tradeoff in machine learning

Fairness == absence of bias?

High bias, low variance



Low bias, high variance



WA

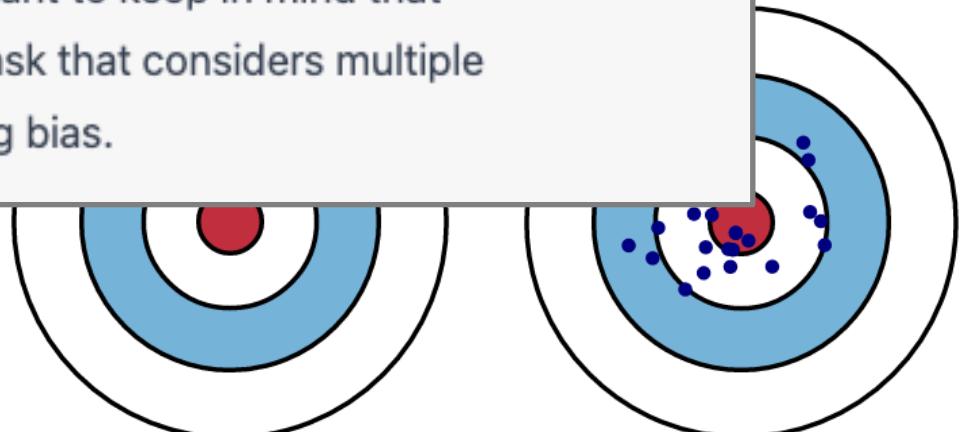
is fairness the absence of bias?



Fairness and bias are related concepts, but they are not the same thing. Fairness generally refers to the idea that all individuals should be treated equally, regardless of their characteristics such as race, gender, or age. Bias, on the other hand, refers to the systematic and unjustified difference in treatment of individuals based on those same characteristics.

In the context of machine learning, fairness can be thought of as the absence of bias in the model's predictions or decisions. This means that the model should not make predictions or decisions that are systematically different for certain groups of individuals based on their protected characteristics. However, it's important to keep in mind that achieving fairness in machine learning models is a complex task that considers multiple aspects and trade offs, and it's not just a matter of eliminating bias.

Fairness == absence of bias?



Common Biases

- **Selection bias** (sampling bias): selection of individuals such that the sample is NOT representative of the population intended to be analyzed
- **Confounding bias**: more on this in AIM II

COGNITIVE BIAS CODEX

What Should We Remember?

To avoid mistakes, we aim to preserve autonomy and group status, and avoid irreversible decisions

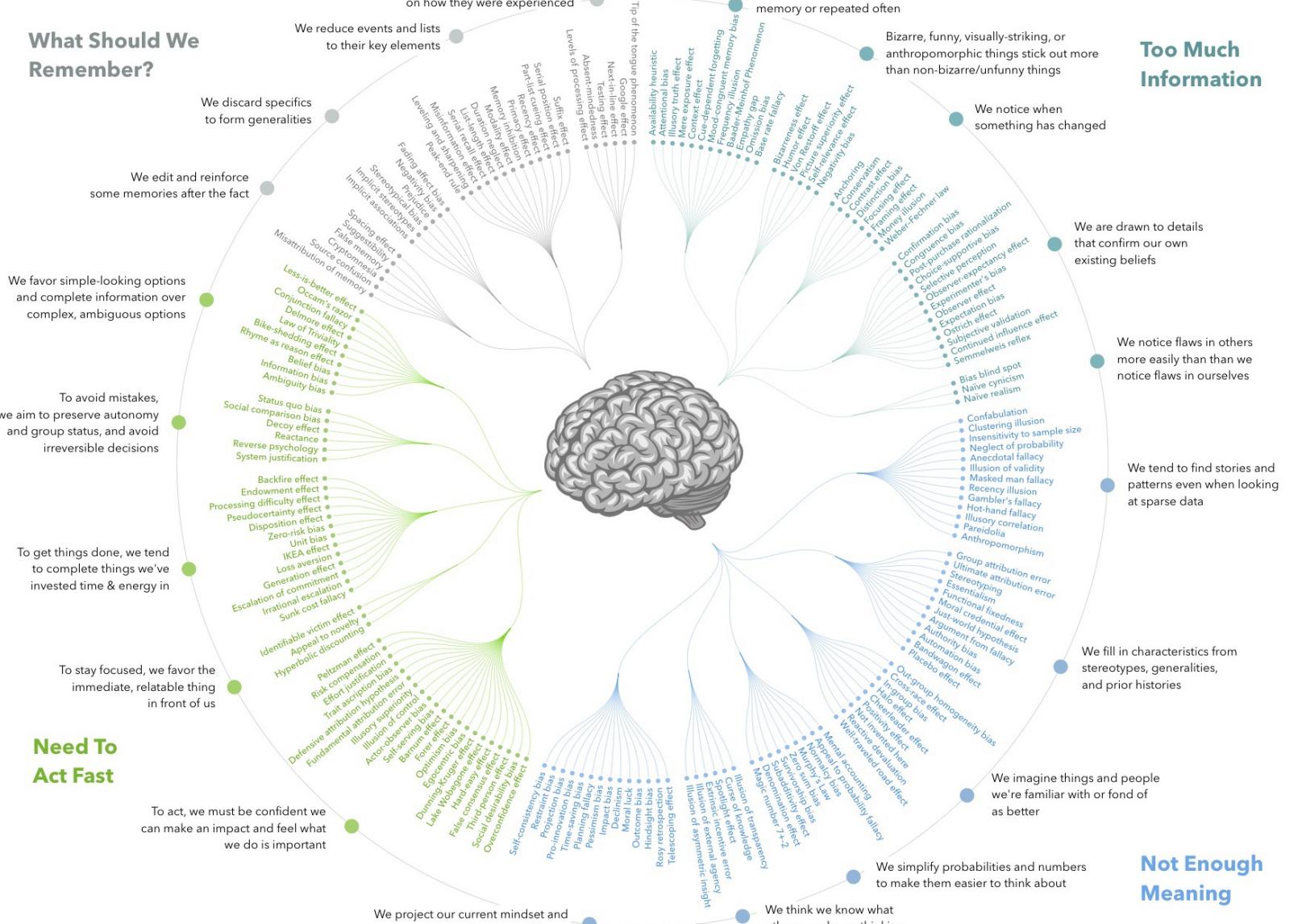
To get things done, we tend to complete things we've invested time & energy in

To stay focused, we favor the immediate, relatable thing in front of us

Need To Act Fast

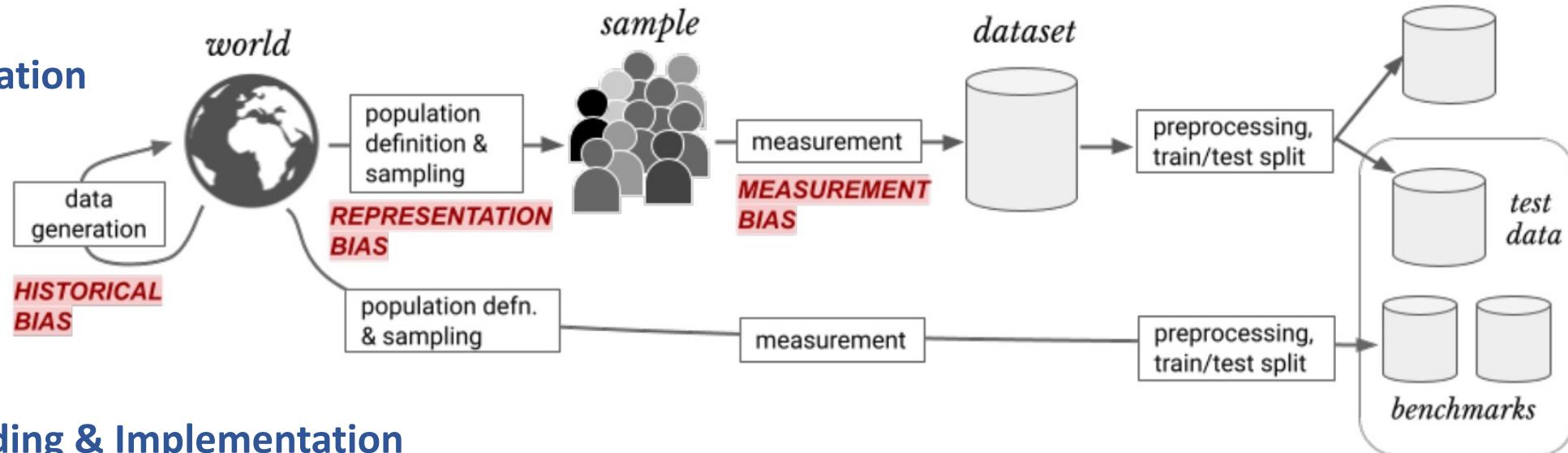
To act, we must be confident we can make an impact and feel what we do is important

We project our current mindset and assumptions onto the past and future

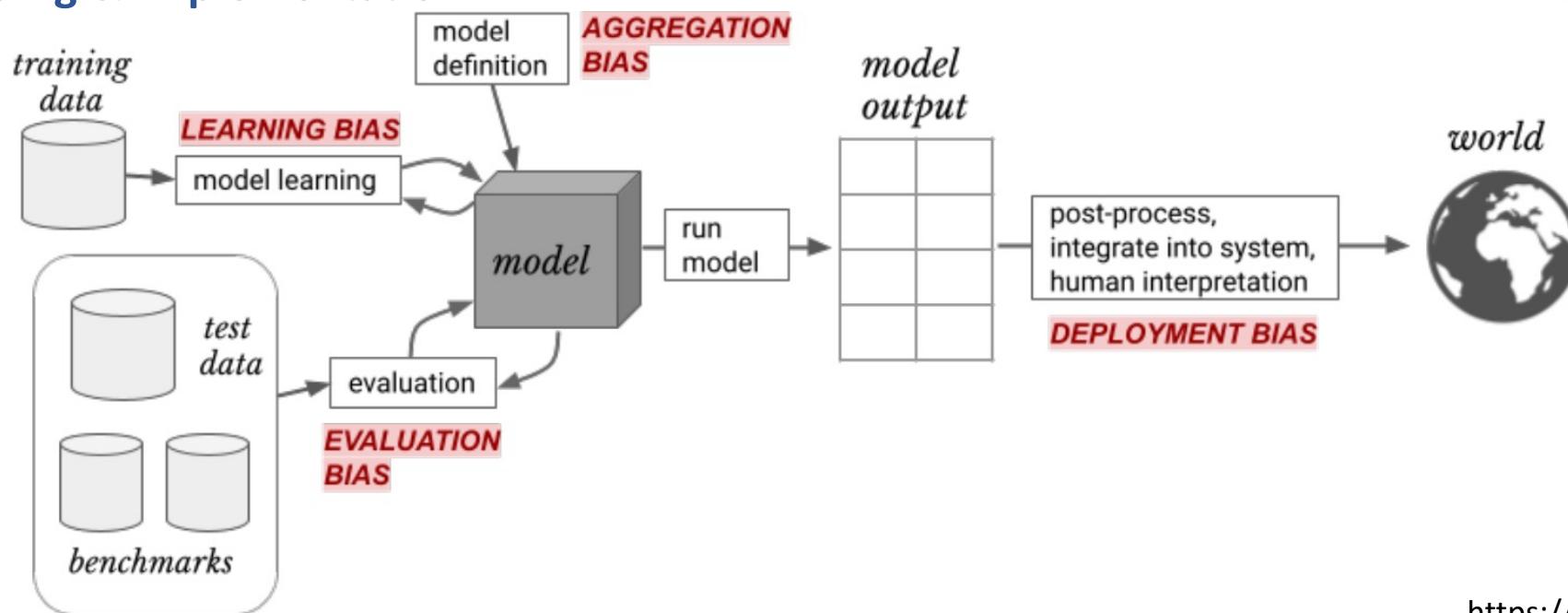


Biases in Machine Learning

Data Generation



Model Building & Implementation



DATA
GENERATION
& COLLECTION

DATASET
DEVELOPMENT

DATA
LABELING

BIASED DATASET



PURPOSE

INPUTS

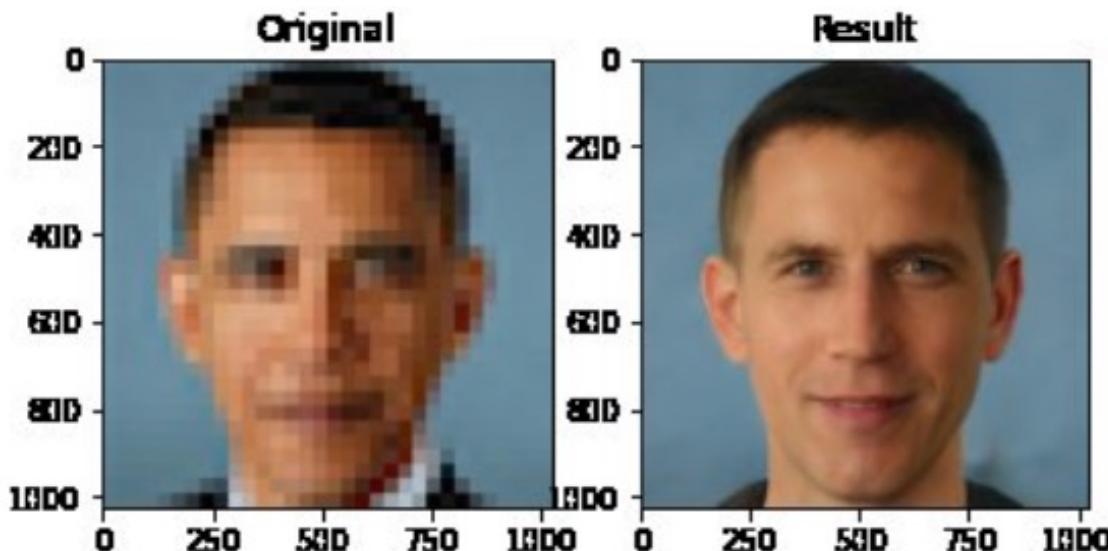
EVALUATION

BIASED ALGORITHM

“Data is never this raw, truthful input and never neutral. It is information that has been collected in certain ways by certain actors and institutions for certain reasons.”

Catherine D'Ignazio

Bias in Face Upsampling



Yann LeCun
@ylecun

ML systems are biased when data is biased.
This face upsampling system makes everyone look
white because the network was pretrained on
FlickFaceHQ, which mainly contains white people pics.
Train the *exact* same system on a dataset from
Senegal, and everyone will look African.



Brad Wyble @bradpwyble · Jun 20, 2020

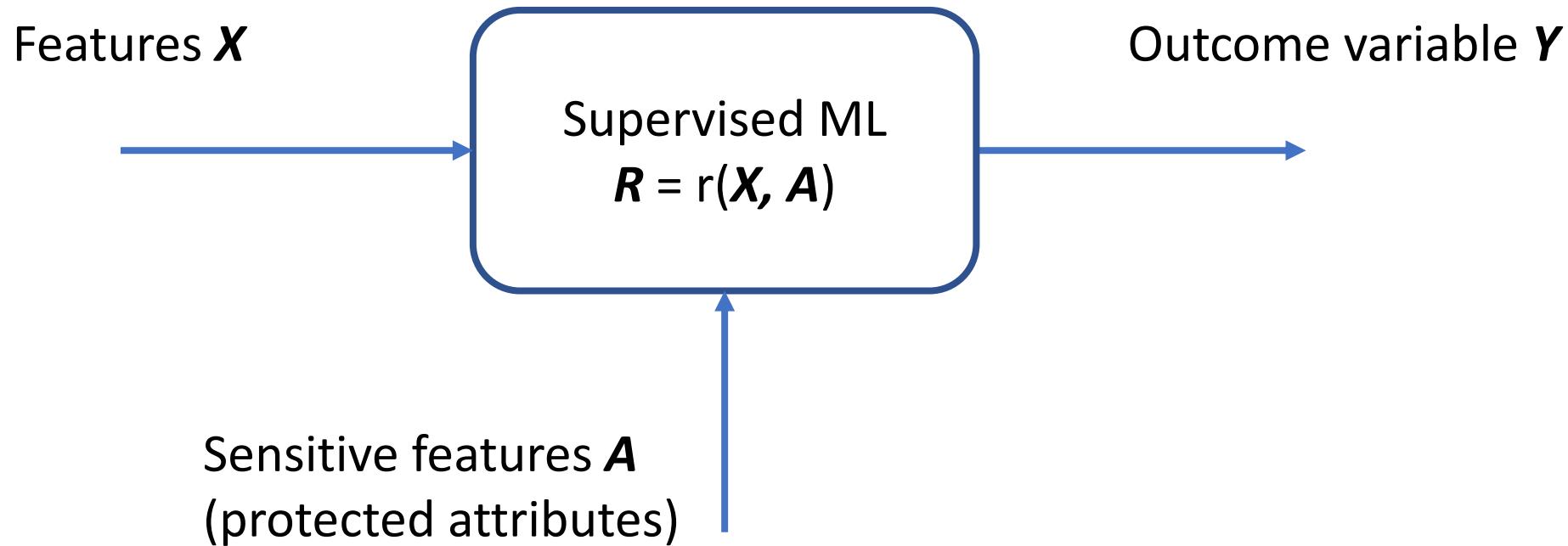
This image speaks volumes about the dangers of bias in AI
[twitter.com/Chicken3gg/sta...](https://twitter.com/Chicken3gg/status/1274782757907030016)

[Show this thread](#)

Are AI developer responsible for the creation of unfair AI systems without purposely designing them to be unfair?

FAIRNESS CRITERIA

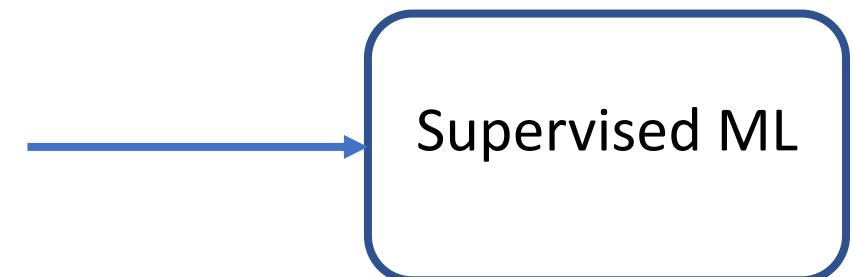
Formal prediction and decision making setting



Fairness through unawareness

Fairness through unawareness if
none of the sensitive features are directly used in the model

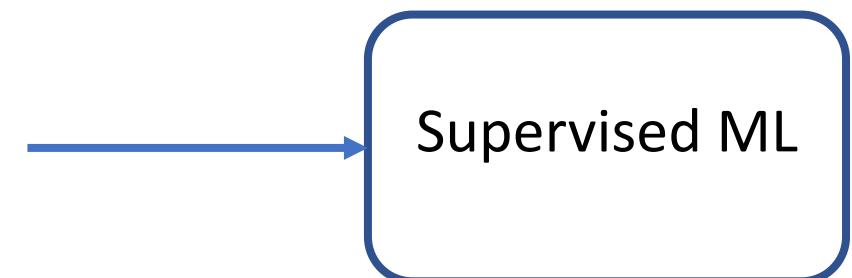
Ethnicity	Skills	Experience	Loves tacos	Hired?
Hispanic	Python	1	Yes	No
Hispanic	C++	5	Yes	Yes
White	Java	2	No	Yes
White	C++	3	No	Yes



Fairness through unawareness

Fairness through unawareness if
none of the sensitive features are directly used in the model

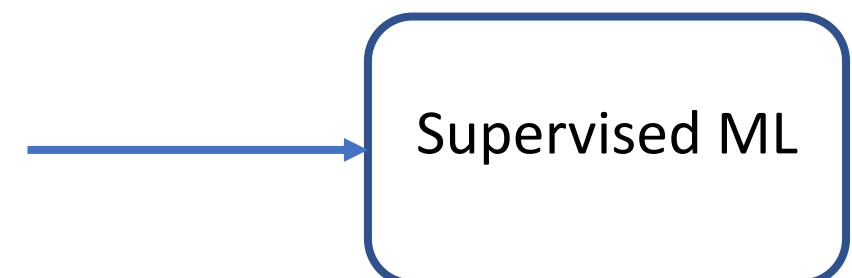
Ethnicity	Skills	Experience	Loves tacos	Hired?
Hispanic	Python	1	Yes	No
Hispanic	C++	5	Yes	Yes
White	Java	2	No	Yes
White	C++	3	No	Yes



Fairness through unawareness

Fairness through unawareness if
none of the sensitive features are directly used in the model

Ethnicity	Skills	Experience	Loves tacos	Hired?
Hispanic	Python	1	Yes	No
Hispanic	C++	5	Yes	Yes
White	Java	2	No	Yes
White	C++	3	No	Yes



Inferred

Sensitive features may still be used indirectly!

Fairness criteria

Independence

$$R \perp A$$

Demographic parity, statistical parity

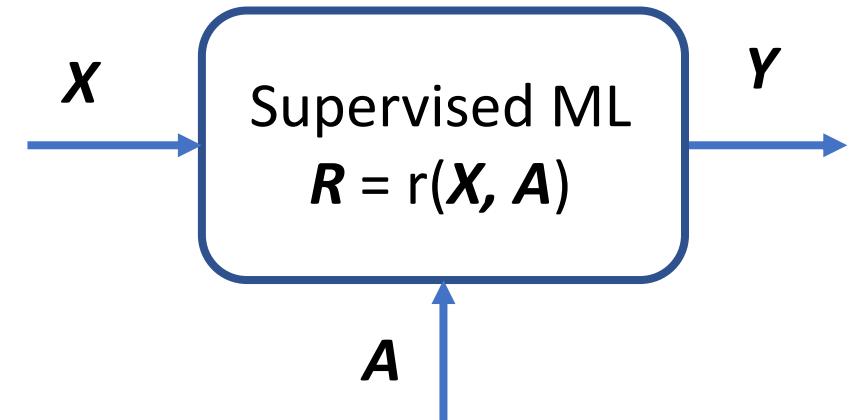
Separation

$$R \perp A \mid Y$$

Equality of odds

Sufficiency

$$Y \perp A \mid R$$



Binary setting

Confusion matrix

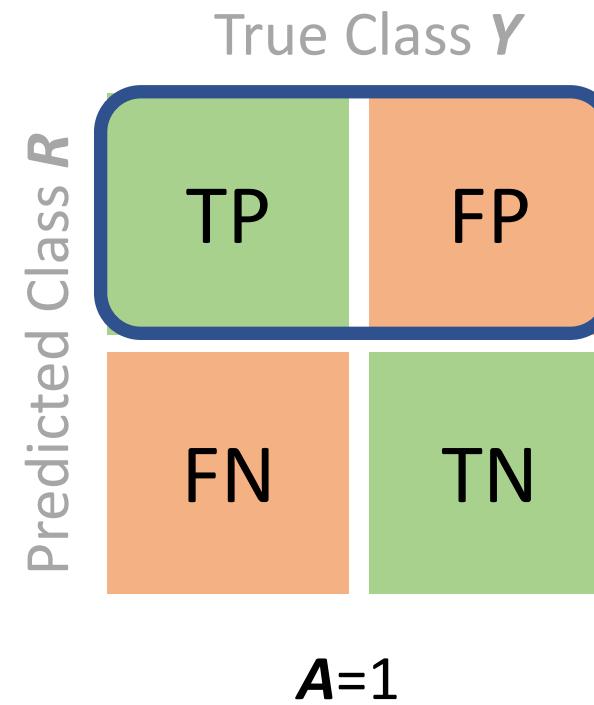
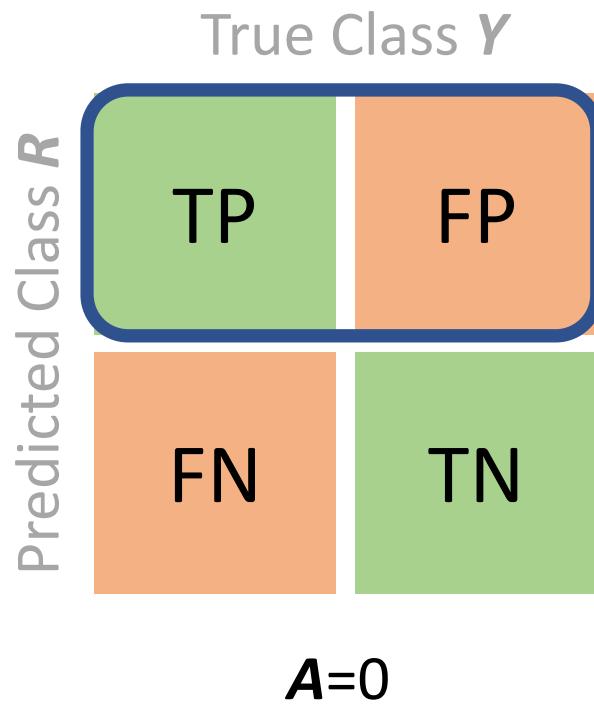
		True Class Y	
		positive	negative
Predicted Class R	positive	TP	FP
	negative	FN	TN

Binary setting with sensitive attribute

		$A=0$		$A=1$	
		True Class Y		True Class Y	
		positive	negative	positive	negative
Predicted Class R	positive	TP	FP	TP	FP
	negative	FN	TN	FN	TN

Independence (demographic parity)

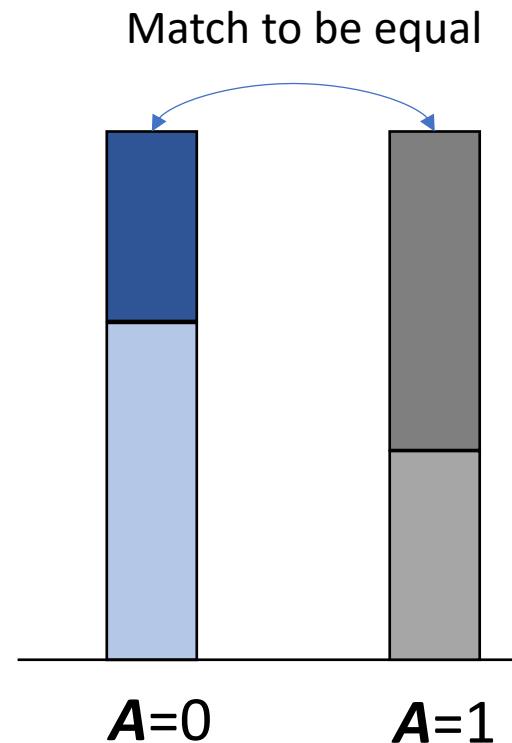
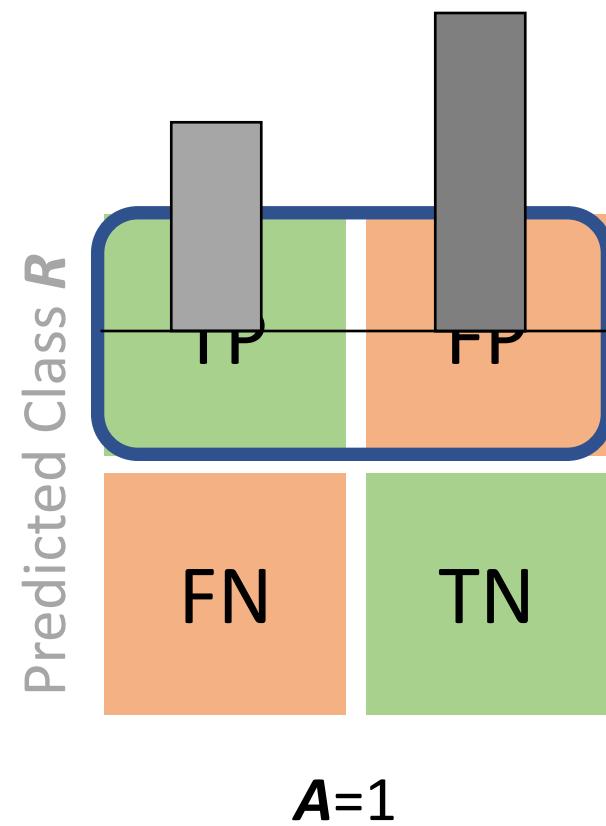
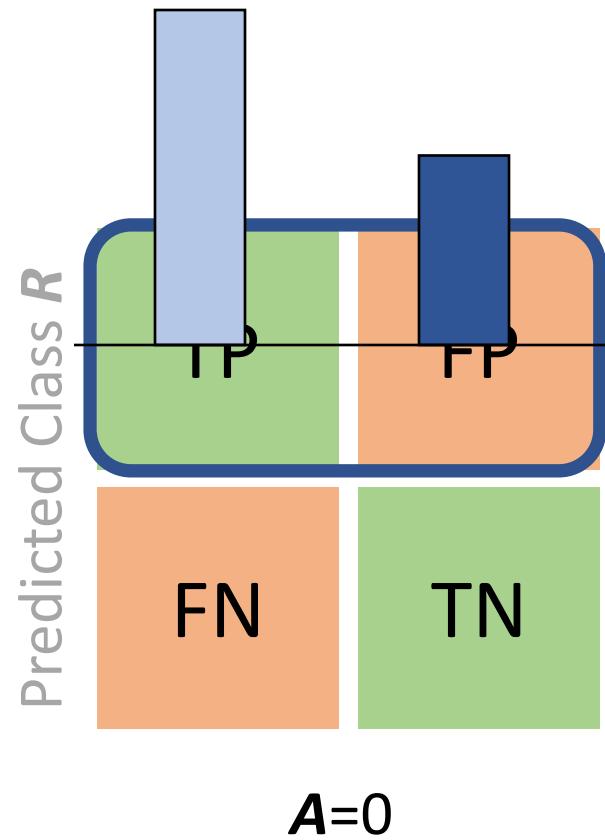
$R \perp A$: $P(R = 1|A = 1) = P(R = 1|A = 0)$ Positive rate is the same for both groups



Independence (demographic parity)

$$R \perp A: P(R = 1|A = 1) = P(R = 1|A = 0)$$

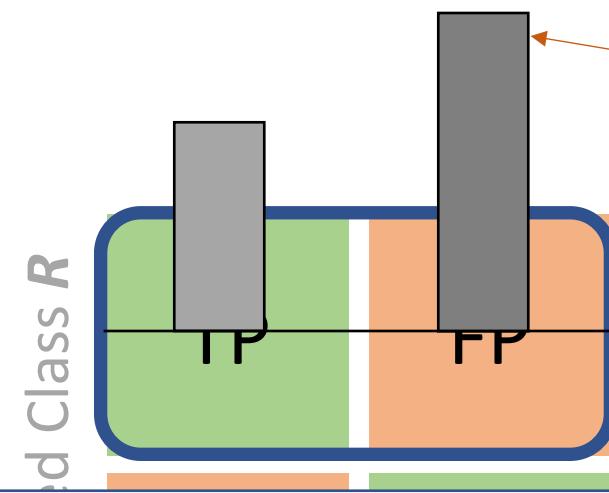
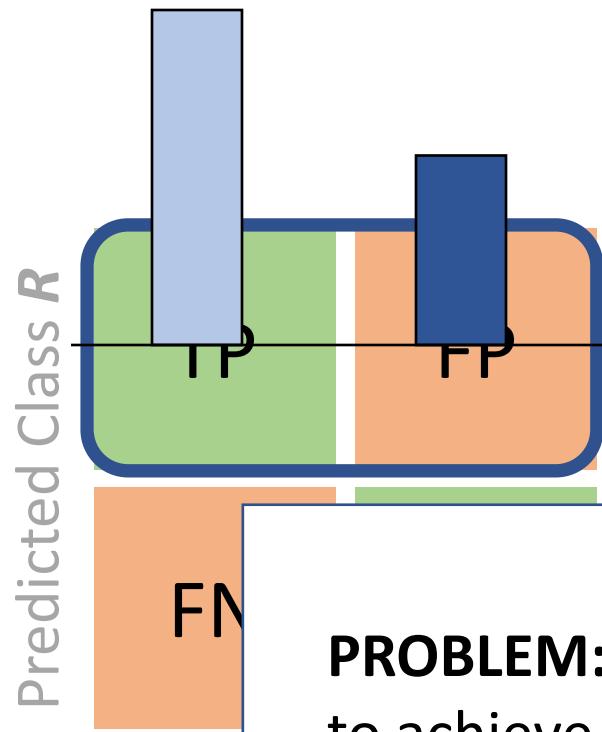
Positive rate is the same for both groups



Independence (demographic parity)

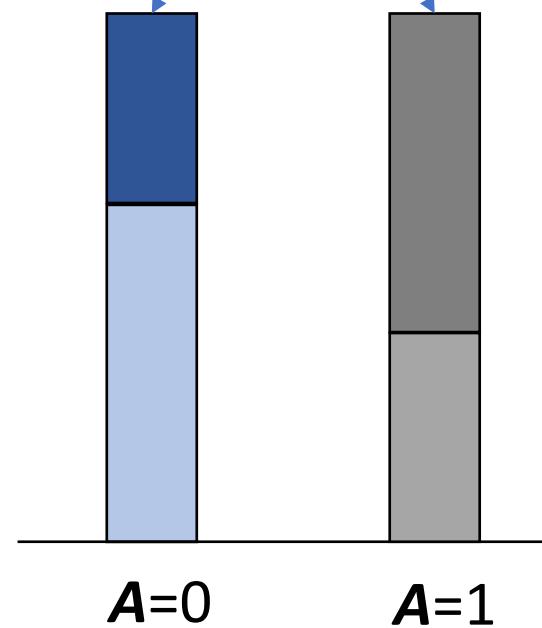
$$R \perp A: P(R = 1|A = 1) = P(R = 1|A = 0)$$

Positive rate is the same for both groups



Accepts too many wrong predictions

Match to be equal



PROBLEM: Effectiveness of predictor decreases to achieve fairness, if A and Y are correlated

Separation (equality of odds)

$$R \perp A \mid Y: P(R = 1|A = 1, Y) = P(R = 1|A = 0, Y)$$

→ $P(R = 1|A = 1, Y = 0) = P(R = 1|A = 0, Y = 0)$
 $P(R = 1|A = 1, Y = 1) = P(R = 1|A = 0, Y = 1)$

		True Class Y	
		TP	FP
Predicted Class R	TP	FP	
	FN	TN	

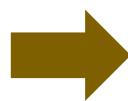
$A=0$

		True Class Y	
		TP	FP
Predicted Class R	TP	FP	
	FN	TN	

$A=1$

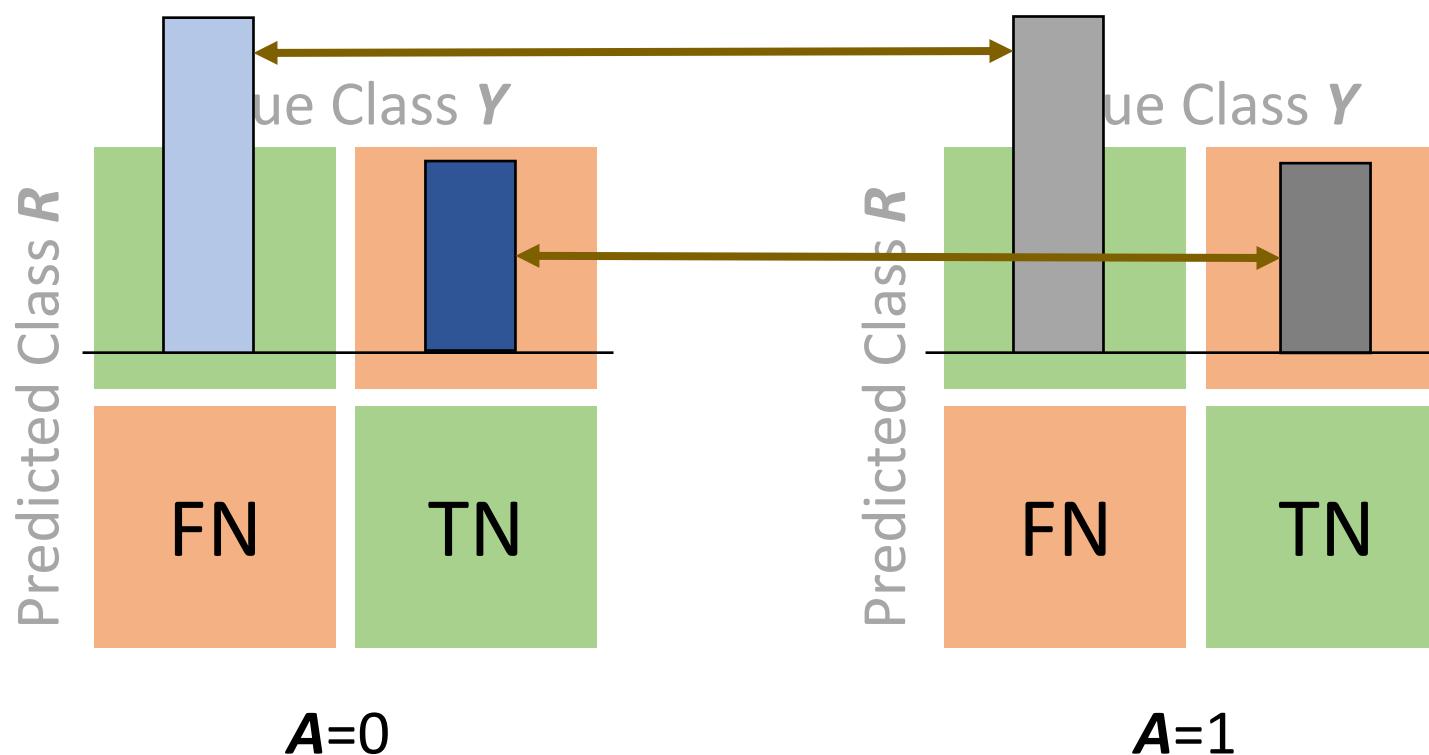
Separation (equality of odds)

$$R \perp A | Y: P(R = 1|A = 1, Y) = P(R = 1|A = 0, Y)$$



$$\begin{aligned} P(R = 1|A = 1, Y = 0) &= P(R = 1|A = 0, Y = 0) \\ P(R = 1|A = 1, Y = 1) &= P(R = 1|A = 0, Y = 1) \end{aligned}$$

Equal false positive rate
Equal true positive rate

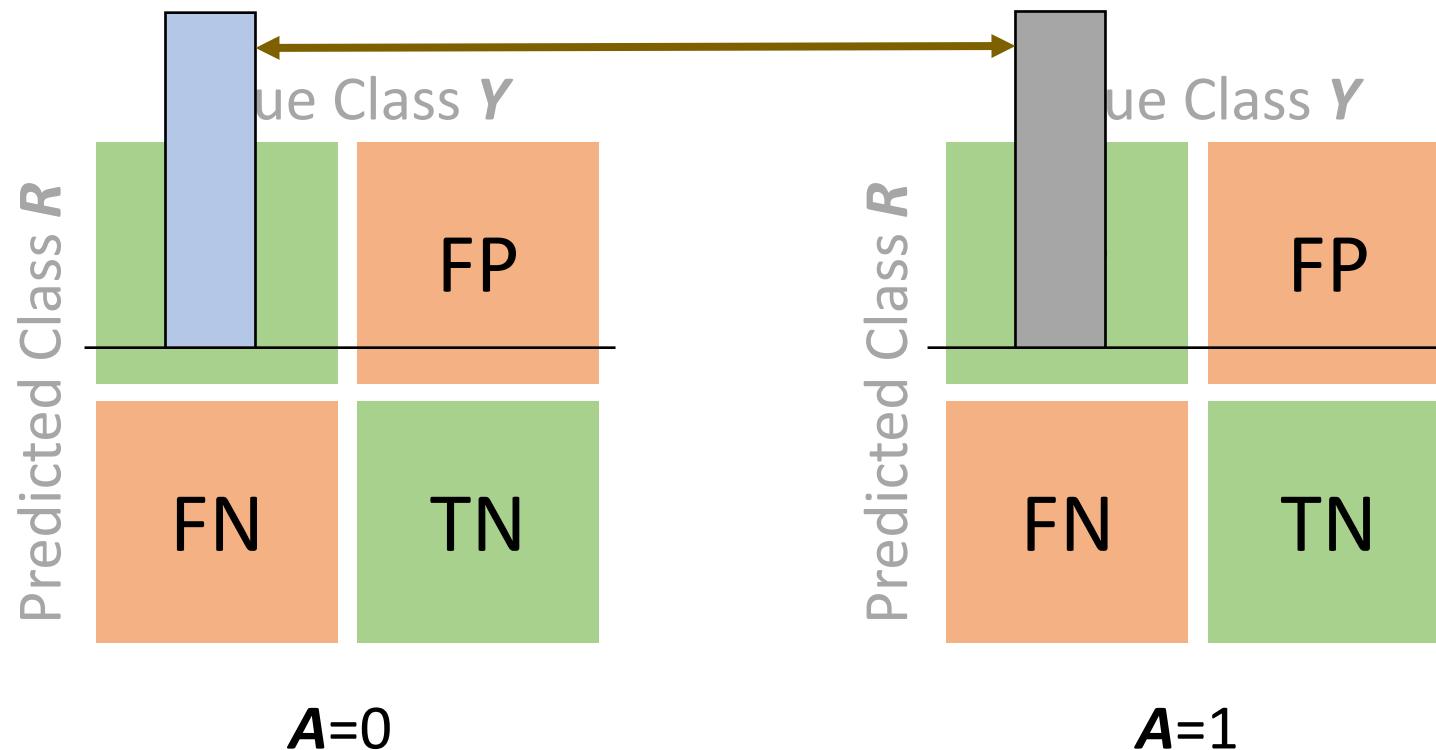


Separation allows correlation between R and A to the extent that it is *justified by the target variable Y*

Equal opportunity

$$P(R = 1|A = 1, Y = 1) = P(R = 1|A = 0, Y = 1)$$

Equal true positive rate



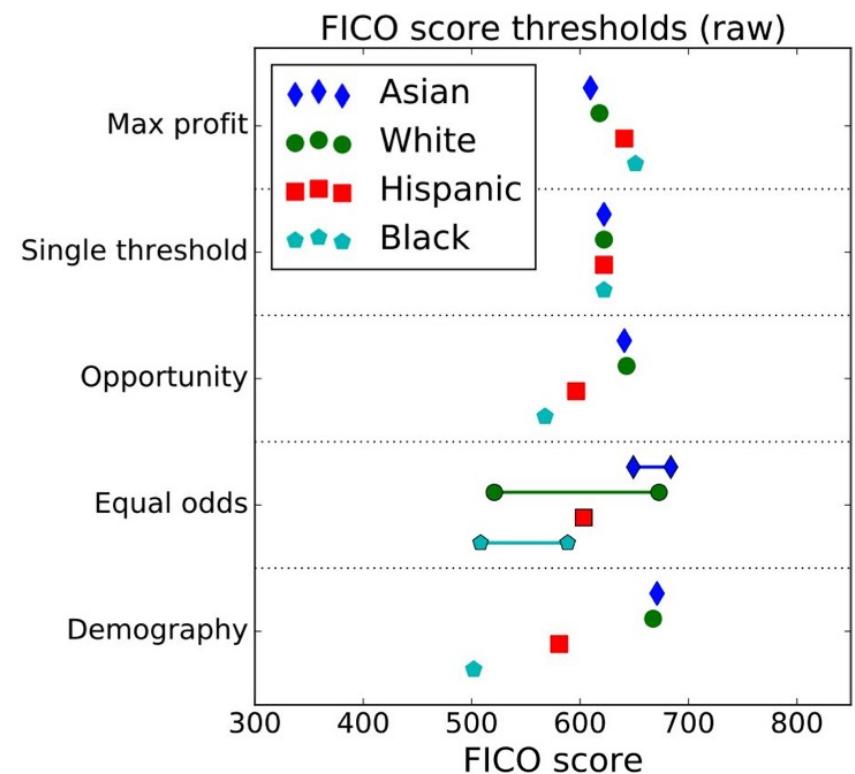
Equal opportunity for people across protected groups.

Bad News

It is not possible to jointly achieve any pair of these conditions

- Independence *xor* Separation
- Independence *xor* Sufficiency
- Separation *xor* Sufficiency

Fairness criteria on their own **cannot**
be a “proof of fairness”



Different thresholds for different criteria!

Summary of fairness criteria

Fairness	Criteria
Unawarness	Exclude A in prediction
Demographic parity	$P(R = 1 A = 1) = P(R = 1 A = 0)$
Equality of odds	$P(R = 1 A = 1, Y) = P(R = 1 A = 0, Y)$
Equal opportunity	$P(R = 1 A = 1, Y = 1) = P(R = 1 A = 0, Y = 1)$

METRICS & MITIGATION

Disparity Metrics

Disparity metrics evaluate how far a given predictor departs from satisfying a parity constraint.

- Differences & Ratios:
- Demographic Parity Difference / Ratio
 - Equal Opportunity Difference / Ratio
 - Average Odds Difference / Ratio

Example

Demographic Parity:

$$P(\mathbf{R} = 1|\mathbf{A} = 1) = P(\mathbf{R} = 1|\mathbf{A} = 0)$$

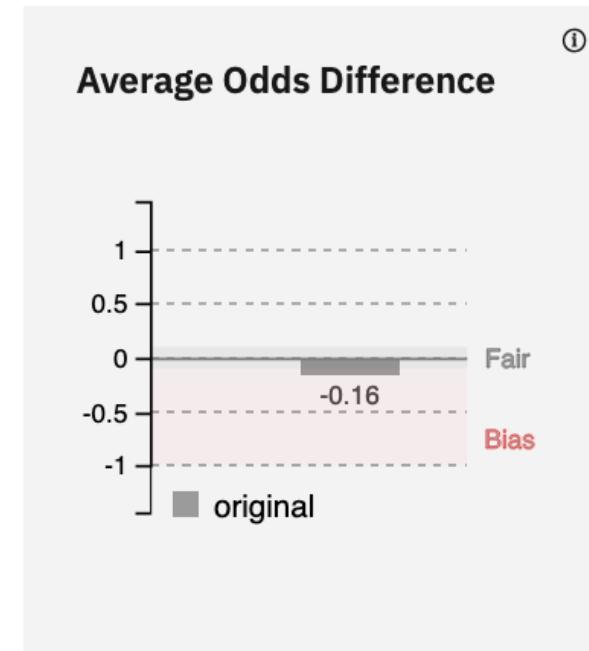
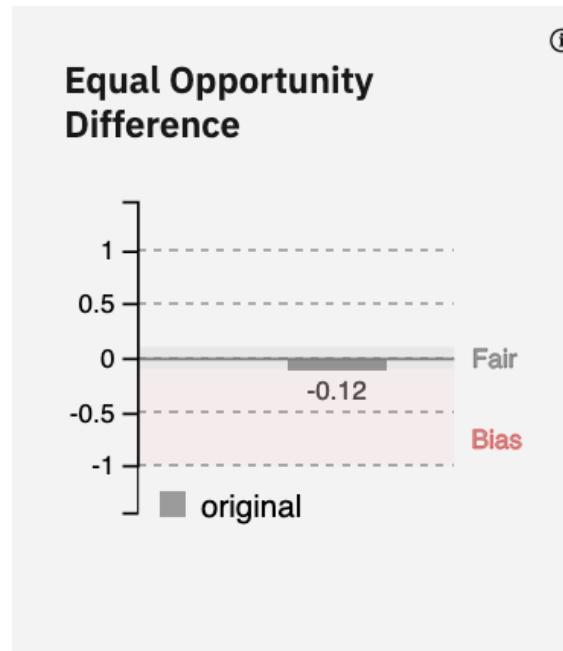
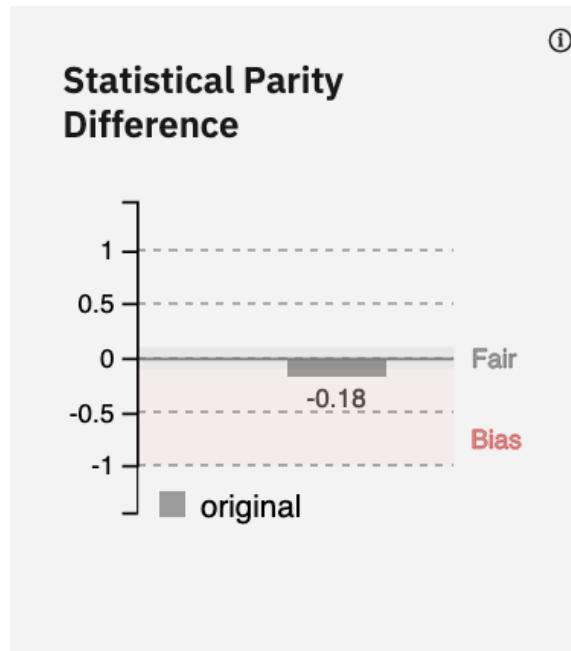


Difference: $\varepsilon = P(\mathbf{R} = 1|\mathbf{A} = 1) - P(\mathbf{R} = 1|\mathbf{A} = 0)$

$$\text{Ratio: } \varepsilon = \frac{P(\mathbf{R} = 1|\mathbf{A} = 1)}{P(\mathbf{R} = 1|\mathbf{A} = 0)} - 1$$

Disparity metrics for COMPAS

- Predict a criminal defendant's likelihood of reoffending
- Protected attribute: Race

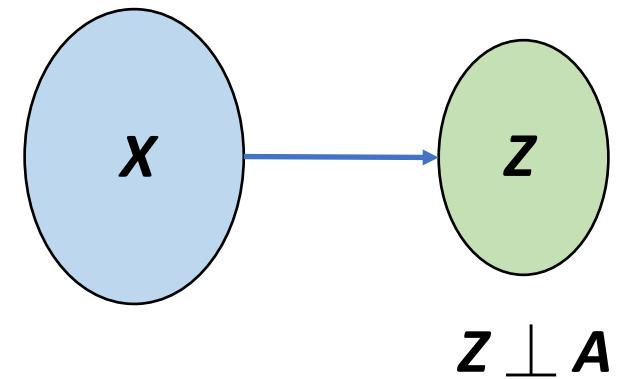


Mitigation strategies

- Pre-processing
- In-processing
- Post-processing

Mitigation: Pre-processing

1. Adjust the features to be **uncorrelated with the sensitive attribute**
 - “Fair” representation learning
 - Adversarial learning
2. Address **representation bias**
 - Reweighting
 - Oversample from minority groups
3. **Data augmentation:** synthesize data for minority groups
 - Example: from observed “he is a doctor” → synthesize “she is a doctor”



Reweighting

Weights the examples in each (group, label) combination to ensure fairness before classification

$$w = \frac{P \text{ expected}}{P \text{ observed}}$$

Example: Loan

Expected: $P(\text{male}) = P(\text{female}) = 0.5$

$P(\text{loan}) = P(\text{no loan}) = 0.5$

Male loan:

$$w = \frac{0.5 * 0.5}{0.4} = 0.625$$

P observed	Male	Female
Loan	0.4	0.2
No loan	0.1	0.3

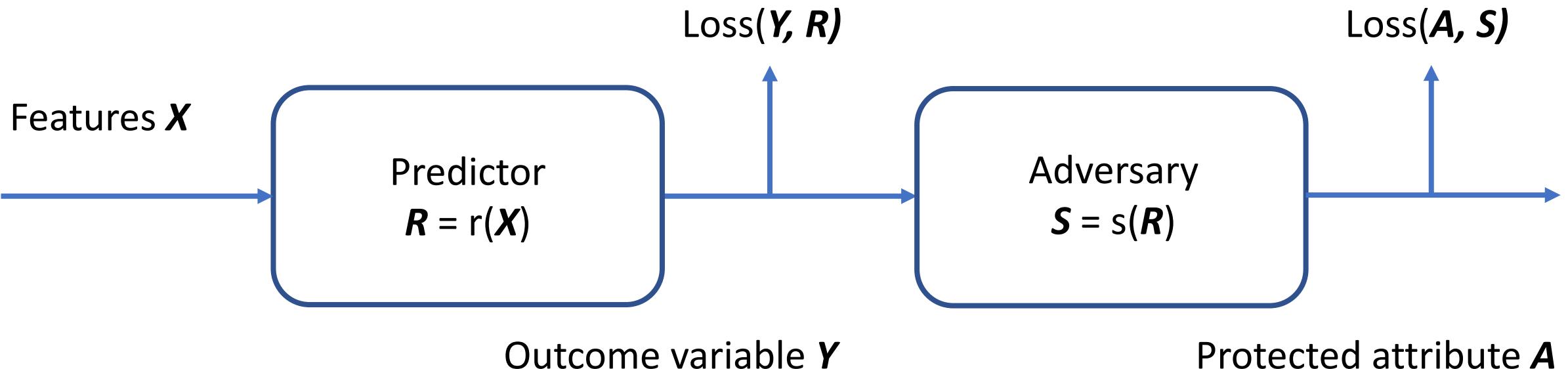
Male no-loan:

$$w = \frac{0.5 * 0.5}{0.1} = 2.5$$

Mitigation: In-processing

1. Separate model for each protected group A
2. **Adversarial debiasing**

Maximize prediction accuracy and **simultaneously** reduce an adversary's ability to determine the protected attribute from the predictions

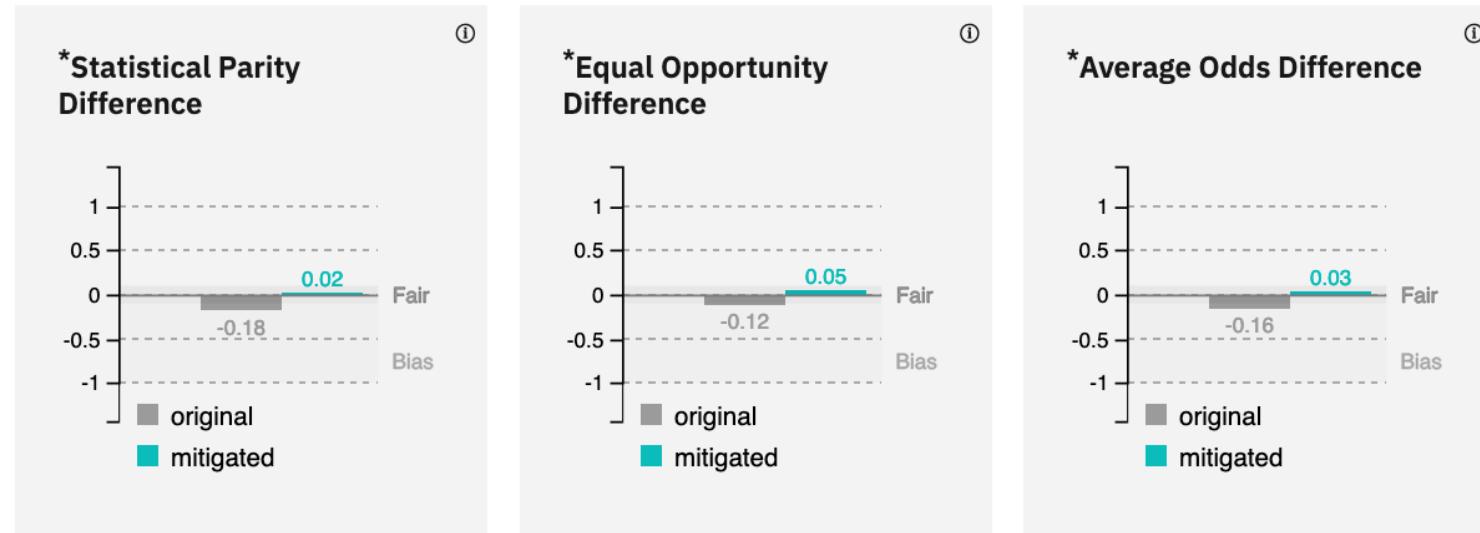


Mitigation for COMPAS

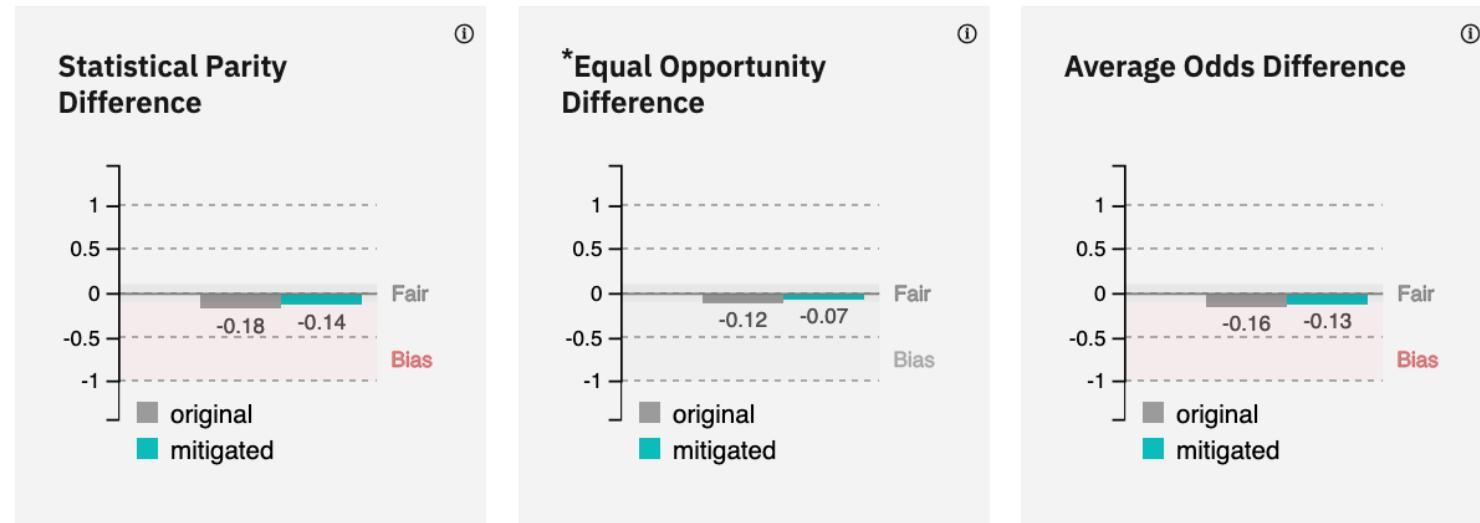
- Pre-processing
 - Reweighting
- In-processing
 - Adversarial debiasing
- Post-Processing
 - Calibrated equalized odds

Mitigation for COMPAS

Reweighting



Adversarial debiasing



Agenda

- Trustworthy AI
- Examples of unfair AI
- Bias & Fairness
- Fairness criteria
- Metrics & Mitigation

Further reading

- Fair ML book (chapter classification)
<https://fairmlbook.org/classification.html>
- Tutorial by Moritz Hardt
https://www.youtube.com/watch?v=lgq_S_7IfOU
- Attacking discrimination with smarter machine learning
<https://research.google.com/bigpicture/attacking-discrimination-in-ml/>
- Lecture: Fair, Accountable, and Transparent (FAccT) Deep Learning
<https://hci.stanford.edu/courses/cs335/2020/sp/schedule.html>

Thank you for your attention!