

Ethics and Fairness in ML

Machine Learning Course - CS-433

Nov 23, 2021

Nicolas Flammarion





Special thanks - disclaimer - ©

These slides are essentially based on:

- The MLSS and NeurIPS tutorials of Moritz Hardt mrtz.org/
- The book *Fairness and Machine learning* of Solon Barocas, Moritz Hardt and Arvind Narayanan fairmlbook.org

And also:

- The lecture of Nathan Kallus at Cornell
- The slides “Fairness in Machine Learning” by Delip Rao



Moritz Hardt



Example borrowed from
Delip Rao

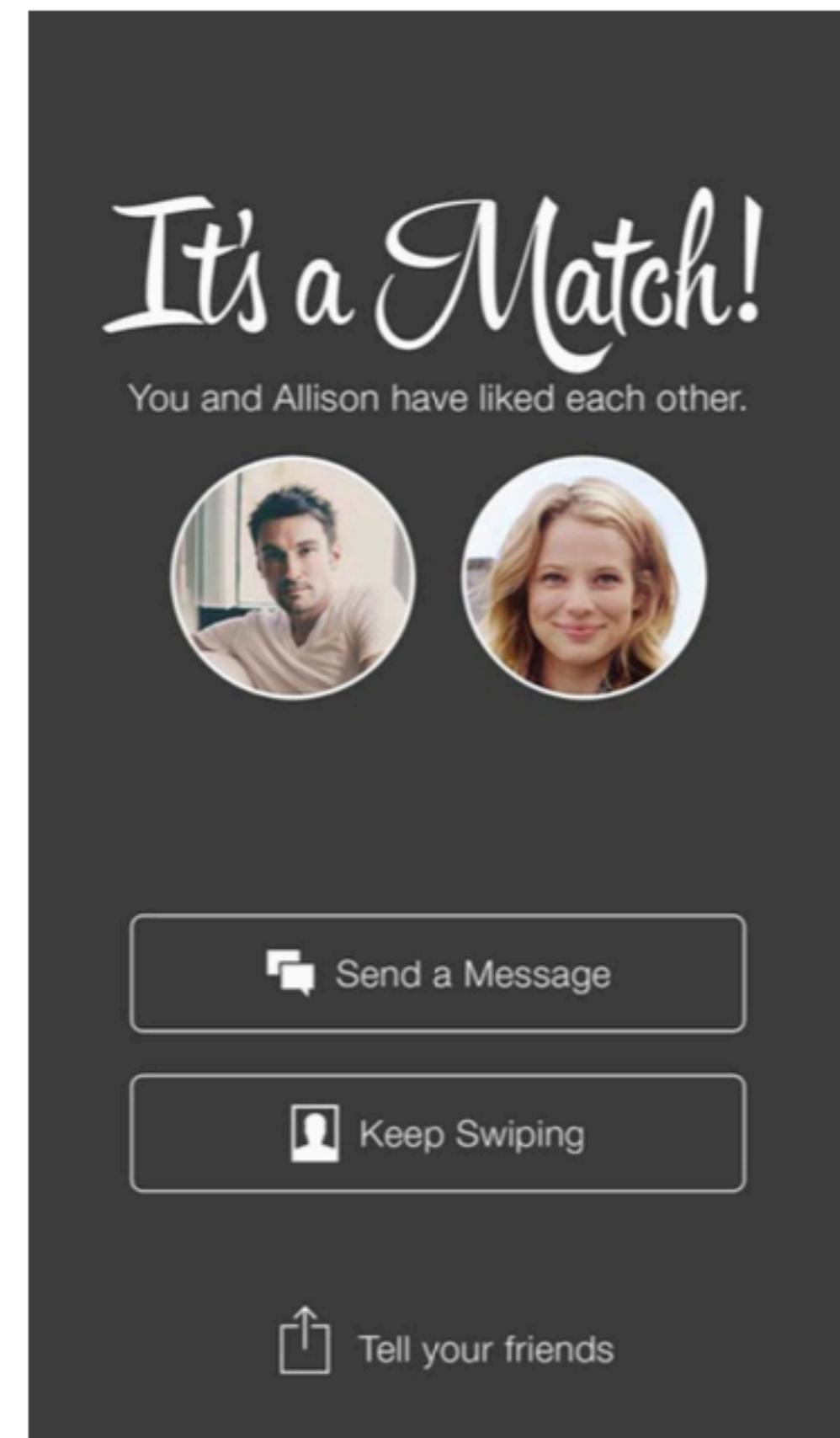
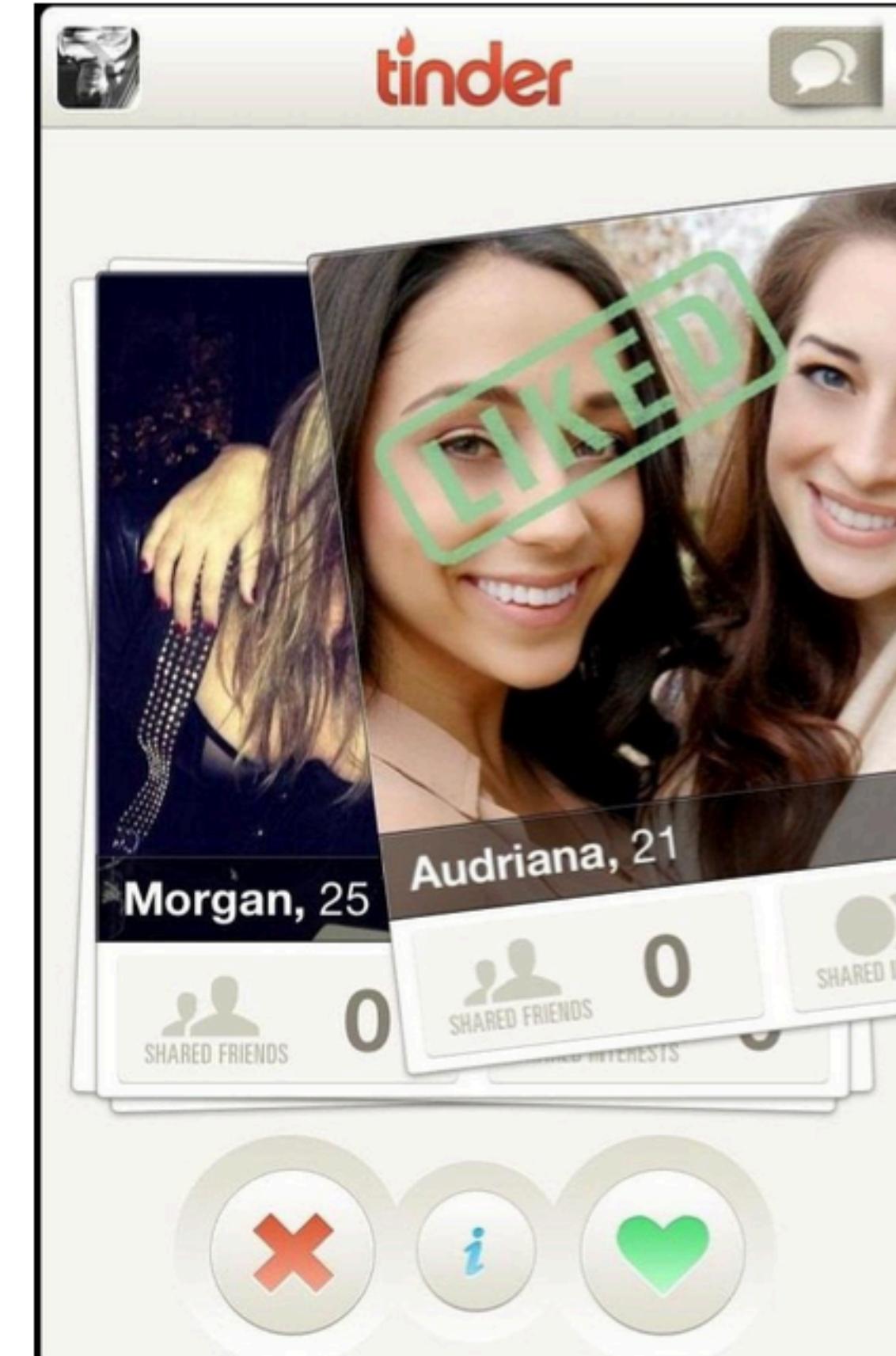
Simple ML problem

You may want to increase:

- % right swipes
- % matches

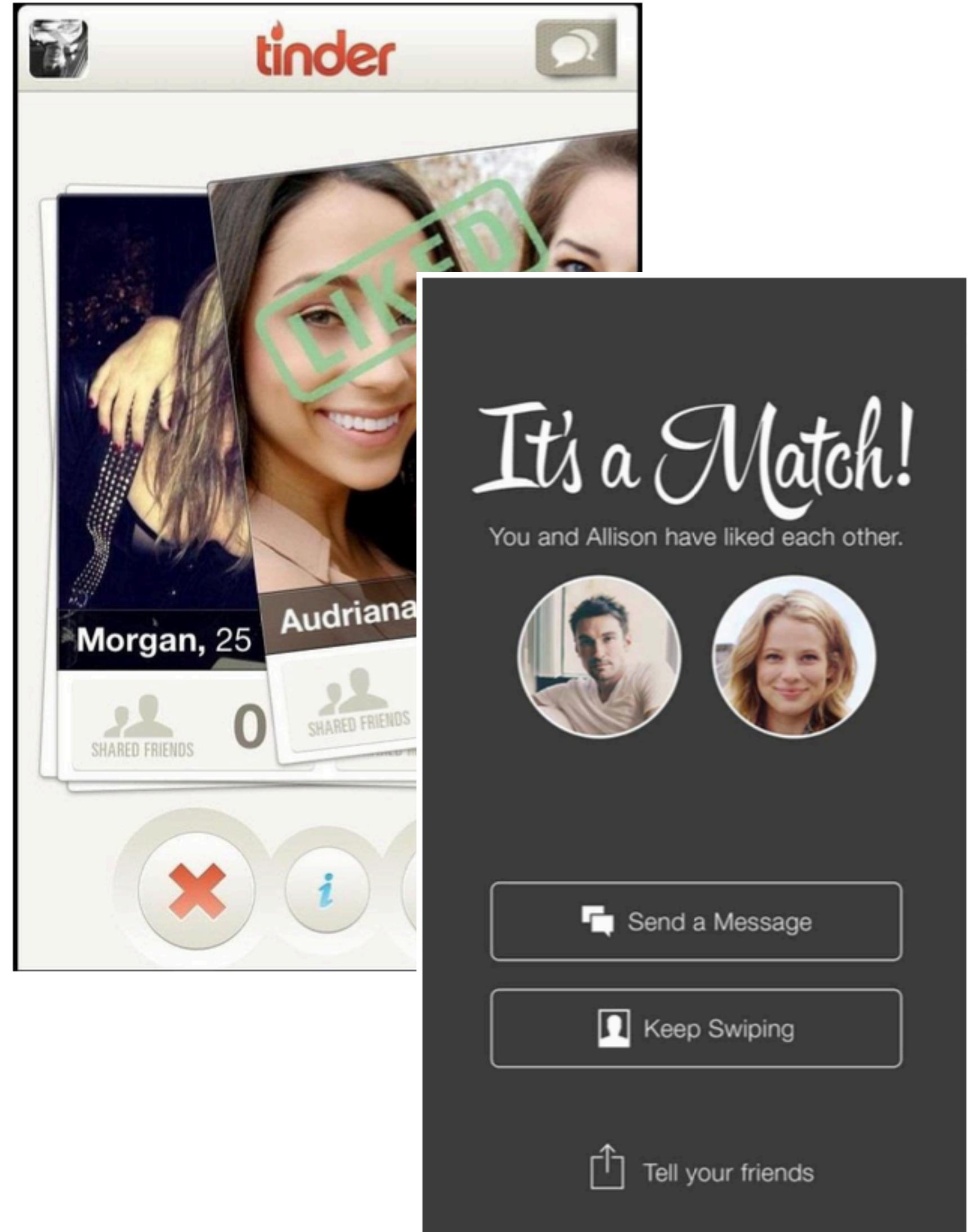
Well defined problem:

- Well defined objective
- Tons and tons of clean data
- Many rich features



Which raises deep ethical questions

- Say we can improve metrics by including as a feature skin color (extracted using computer vision) for the ranking algorithm
- Should we do it?
- What about self-identified ethnicity (e.g. in profile)?



Failure of fairness through unawareness

Amazon uses data-driven techniques to decide the neighborhoods it will offer free same-day delivery

Disparities in the demographic makeup of these neighborhoods

→ White residents were more than twice as likely as Black residents to have access to this service



Article from Bloomberg, 2016

Failure of fairness through unawareness

Certainly, Amazon was just predicting number of purchases, which correlates with wealth which correlates with race in the US.

They did not look at their customers' race when building their product

Example of just using ML without concern for fairness issues which lead to ethical issues



Discarding “sensitive attributes” do not solve the fairness problem and can aggravate them



Article from Bloomberg, 2016

Discrimination in ML

Discrimination: didn't we actually learn how to discriminate in the previous lectures?

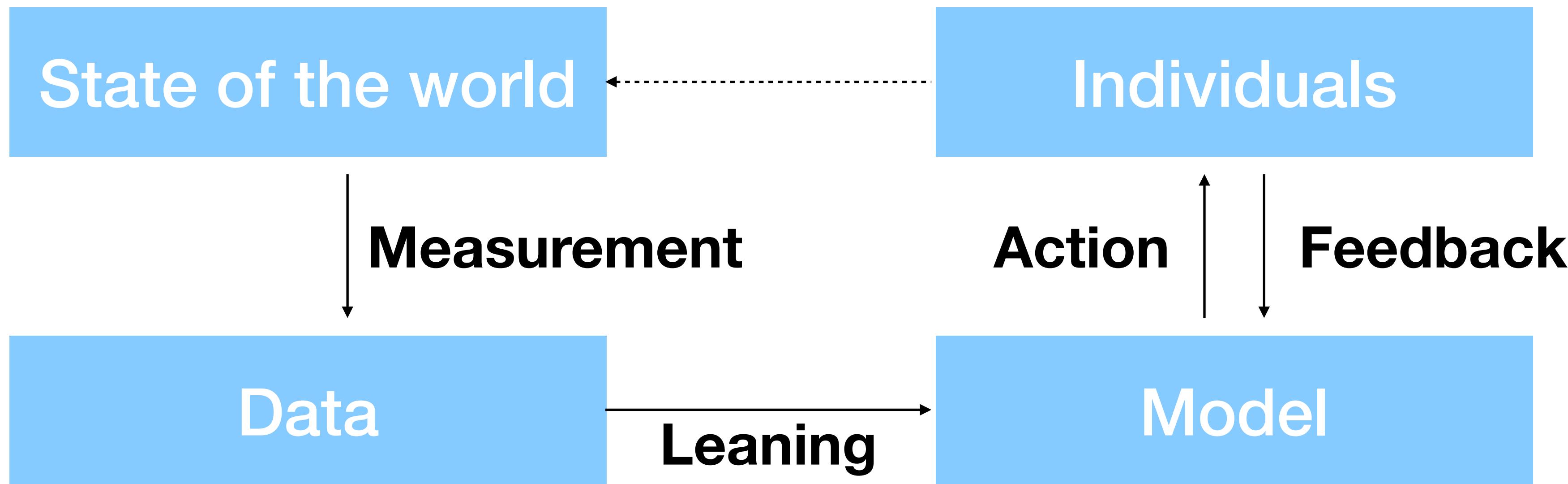
We will be concerned with **unjustified basis for differentiation**:

- Practical irrelevance
 - Sexual orientation in employment decisions
- Moral irrelevance
 - Disability status in hiring decisions

Discrimination is **domain specific**: concerned with opportunities that affect people's lives

Discrimination is **group specific**: concerned with social categories that have served as the basis for unjustified and systematically adverse treatment in the past

The machine learning loop



Applications about people

- Most ML applications are about people: 14 out of the top 30 Kaggle competitions concern tasks where decisions are made about individuals
- Training data often encode existing demographic disparities
- Social stereotypes may be perpetuated by applications of ML algorithms to these tasks
Ex: Automated essay scoring: train data come from human graders with possible stereotype



Applications that are not about people?



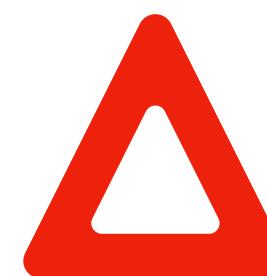
Example borrowed from Delip Rao

Measurements are not without problems

Measurement involves

- Defining your variable of interest
- Defining the process for interacting with the real world
- Turning the observation into numbers - collecting the data

Measuring any attributes about people is subjective and challenging

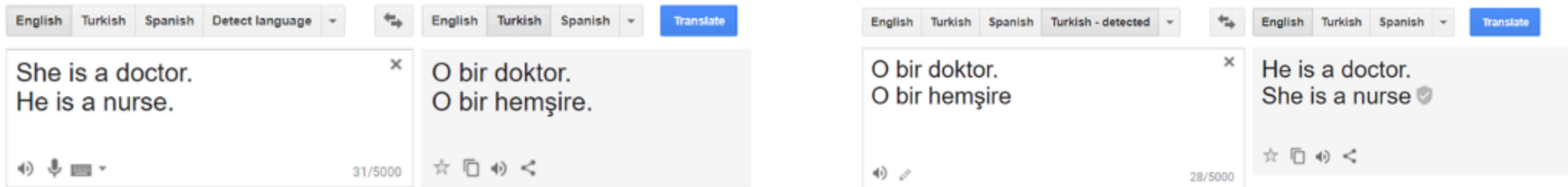


It is crucial to understand the provenance of the data as practitioner

From data to models: disparities can be preserved

Some patterns in the training data represent **knowledge** - we want to learn, while other patterns represent **stereotypes** - we want to avoid learning

→ML algorithms cannot distinguish between these two. Without specific intervention, ML algorithms will extract both.



Removing, e.g., the gender is not a solution because of redundant encodings, i.e., other attributes that may correlate with, e.g., the gender.

These redundant encodings may be relevant for the problem at hand

From data to models: disparities can be introduced

Sample size disparity:

- Uniform subsampling from population leads to fewer data about minorities
- If minority groups are in addition underrepresented, then even fewer data

ML works best with a lot of data → ML may work less well for minorities

True error is an average criterion → low true error may hide terrible performance for a minority group

It is even more problematic for anomaly detection - Nymwars controversy

Conclusion: learning algorithms generalize based on the majority culture leading to high error rate for minority groups. This is because of our goal to avoid overfitting

Toy example

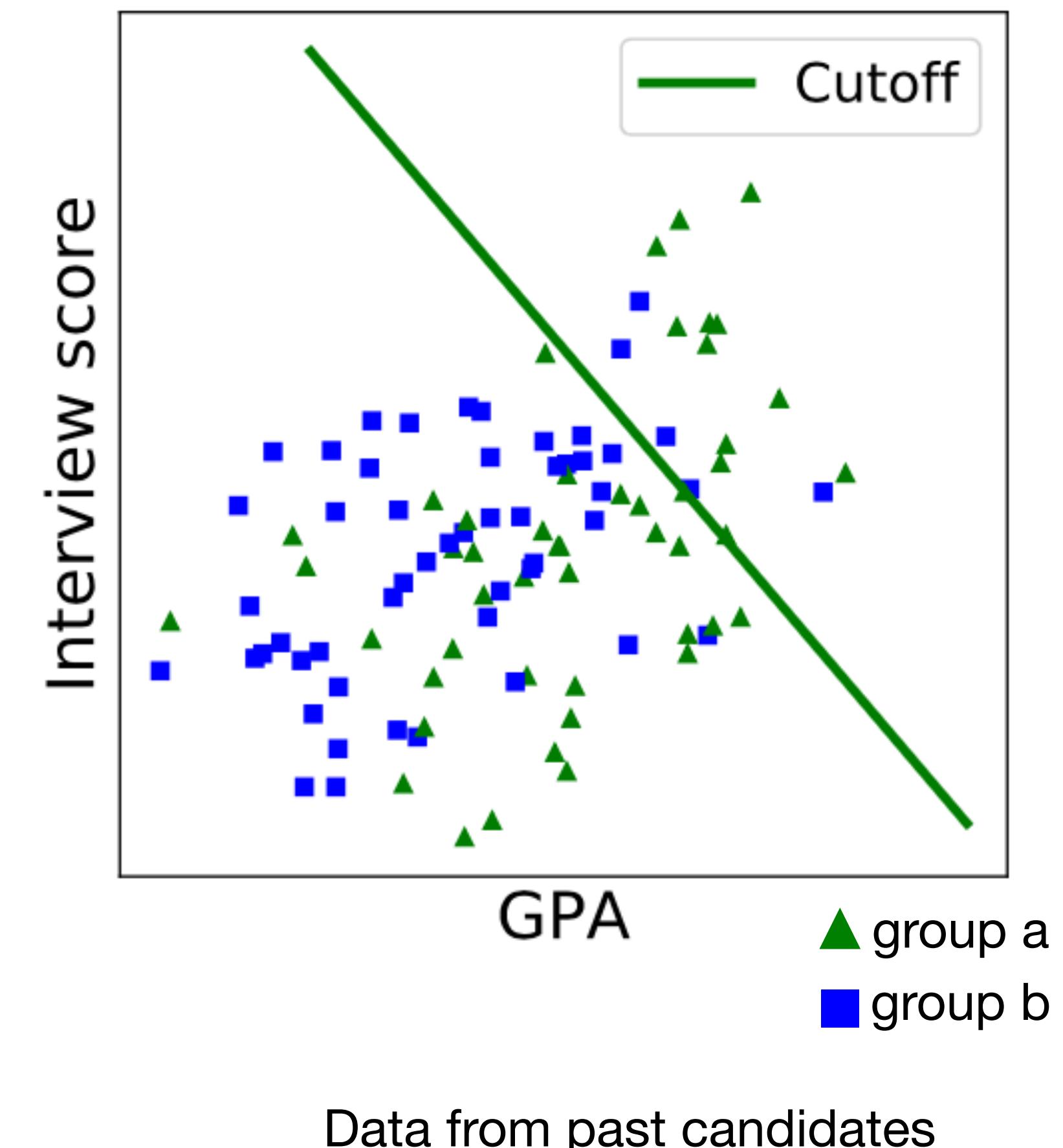
Hiring committee makes decision based on applicants' college GPA and interview score

The classifier does not take into account which group a candidate belongs to

However, the triangles are more likely to be selected than the squares: the ground truth we predict is systematically lower for the squares than for the triangles

Why:

- Post-hiring: bias inside the company
- Pre-hiring: disparities in education
- Combinations of different factors



How to decrease the disparity of our selection procedure?

1. GPA is correlated with the demographic attribute - proxy
→ But simply omitting it would decrease the accuracy of our model
2. Pick different cutoffs so that candidates from both groups have the same probability of being hired
→ but two candidates with the same attributes may receive a different decisions depending on their groups
3. Change the model to weight less the GPA and increase diversity between selected candidates

Fairness criteria in classification

Formal setting: classification

Data are described by covariates X and outcomes variable $Y \in \{0,1\}$

Goal: given a new X you want to predict its label Y

How:

1. Use an algorithm to produce a score function $R = r(X)$
 - Bayes optimal score
 - Learned from labeled data, e.g., in logistic regression
2. Make binary decisions according to threshold rule $D = 1_{R>t}$

Today: we assume R given and are interested in the decision process

Statistical classification criteria

		Decision D
		0 1
True class Y	0	True negative
	1	False positive
True class Y	1	False negative
	0	True positive

True positive rate: $\mathbb{P}(D = 1 | Y = 1)$

False positive rate: $\mathbb{P}(D = 1 | Y = 0)$

True negative rate: $\mathbb{P}(D = 0 | Y = 0)$

False negative rate: $\mathbb{P}(D = 0 | Y = 1)$

The choice of the threshold t in the decision rule D will depend on the classification criteria we pick

Sensitive attributes

In many tasks, X can encode sensitive attributes of an individual

We introduce additional random variable A encoding membership status in a protected class

No fairness through unawareness: removing/ignoring sensitive attributes is not solving the problem

Many features slightly correlated with the sensitive attribute can be used to recover the attribute

If we remove the attribute, the classifier will still find a redundant encoding in terms of other features and we'll have learnt an equivalent classifier

Three fundamental fairness criteria

Idea: equalize different statistical quantities involving group membership A

→ it dates back to the 1960s with seminal work of Anne Cleary

Most of the fairness criteria are properties of (A, Y, R) :

- Independence: R independent of A
- Separation: R independent of A , conditional on Y
- Sufficiency: Y independent of A conditional on R

Independence: equalizing acceptance rate

It requires the random variable A and R to be independent, denoted $A \perp R$

Implies, for any two groups a, b :

$$\mathbb{P}(D = 1 | A = a) = \mathbb{P}(D = 1 | A = b)$$

→ Acceptance rate is the same in all groups: equal positive rate

Limitations of independence

This criterion does not rule out unfair practice. Let's imagine a company which

- hires with care (ie., makes good decision) in group a at some rate $p > 0$
 - hires without care (i.e., makes poor decisions) in group b with the same rate p
- acceptance in both group are identical
→ unqualified applicants are more likely to be selected in group b
→ members of group b will appear to perform less well than those of a

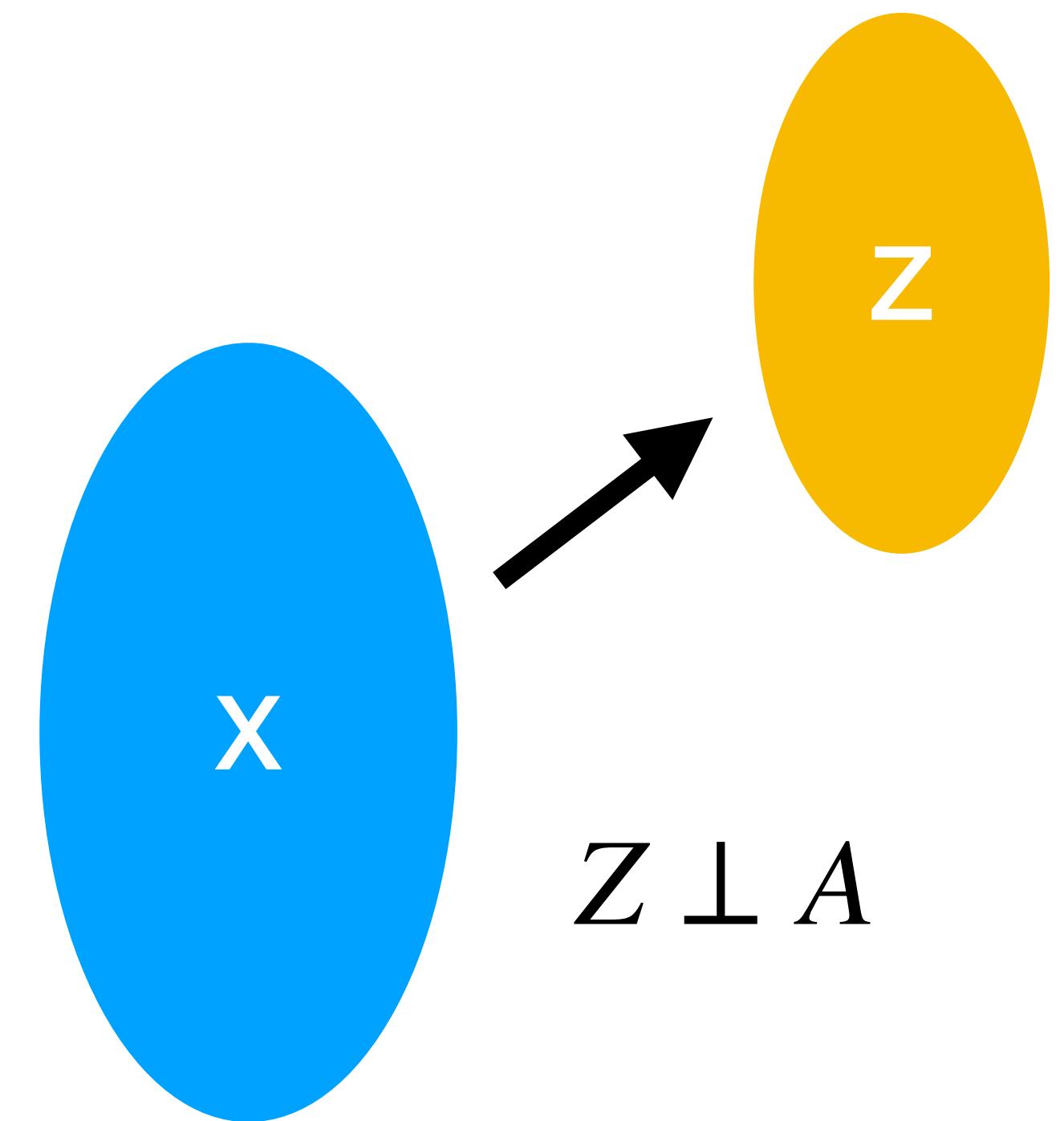
It can happen on its own if there is less data in one group

A positive output can either be a false positive or a true positive

- we shouldn't be able to match true positives in one group with false positives in another

How to achieve fairness criteria

- Post-processing: adjust your learned classifier so as it becomes uncorrelated with the sensitive attribute A
- At training time: work the constraint into the optimisation process
- Pre-processing: adjust your features so that they become uncorrelated with the sensitive attribute A : e.g., use deep learning to learn a representation of the data independent of A , while representing original data as well as possible - Zemel et al., 2015



Separation: equalizing error rates

It requires the random variable A and R to be independent conditional on the target variable Y , denoted $A \perp R | Y$

It implies for all groups a, b :

$$\mathbb{P}(D = 1 | Y = 0, A = a) = \mathbb{P}(D = 1 | Y = 0, A = b) \quad (\text{equal false positive rate})$$

$$\mathbb{P}(D = 0 | Y = 1, A = a) = \mathbb{P}(D = 0 | Y = 1, A = b) \quad (\text{equal false negative rate})$$

This is a **post-hoc criterion**: at decision time, we do not know who is a positive/negative instance

Can be computed in retrospect, by collecting group of positive and negative instances

Sufficiency:

It requires the random variable A and Y to be independent conditional on R , denoted $A \perp Y | R$

For all groups a, b and values r we have:

$$\mathbb{P}(Y = 1 | R = r, A = a) = \mathbb{P}(Y = 1 | R = r, A = b)$$

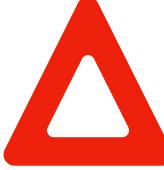
Meaning: for predicting Y we do not need to know A if we have R

Calibration and sufficiency

Def: A score R is calibrated if

$$\mathbb{P}(Y = 1 | R = r) = r$$

- you can interpret your score as a probability
- a priori guarantee: score value r corresponds to positive outcome rate r

 The guarantee does not hold at the individual level

Calibration by group:

$$\mathbb{P}(Y = 1 | R = r, A = a) = r$$

Fact: Calibration by group implies sufficiency

Rmk: it is also possible to go from sufficiency to calibration

Calibration via Platt scaling

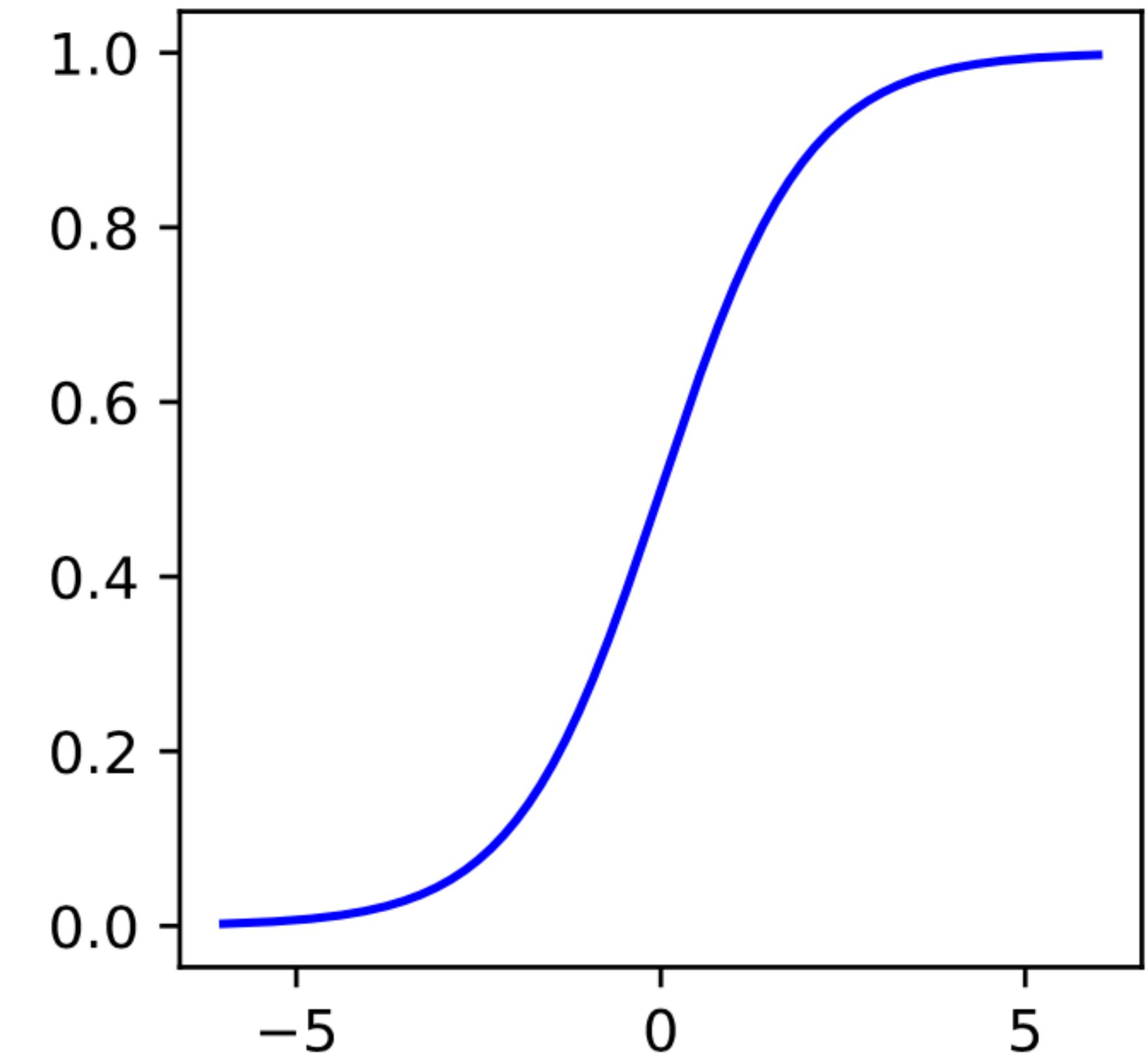
- Given an uncalibrated score R , we can find parameters α, β such that the sigmoid function

$$S = \frac{1}{1 + \exp(\alpha R + \beta)}$$

fits the target Y

- This can be done by minimizing the log loss

$$\mathbb{E}[Y \log S + (1 - Y) \log(1 - S)]$$



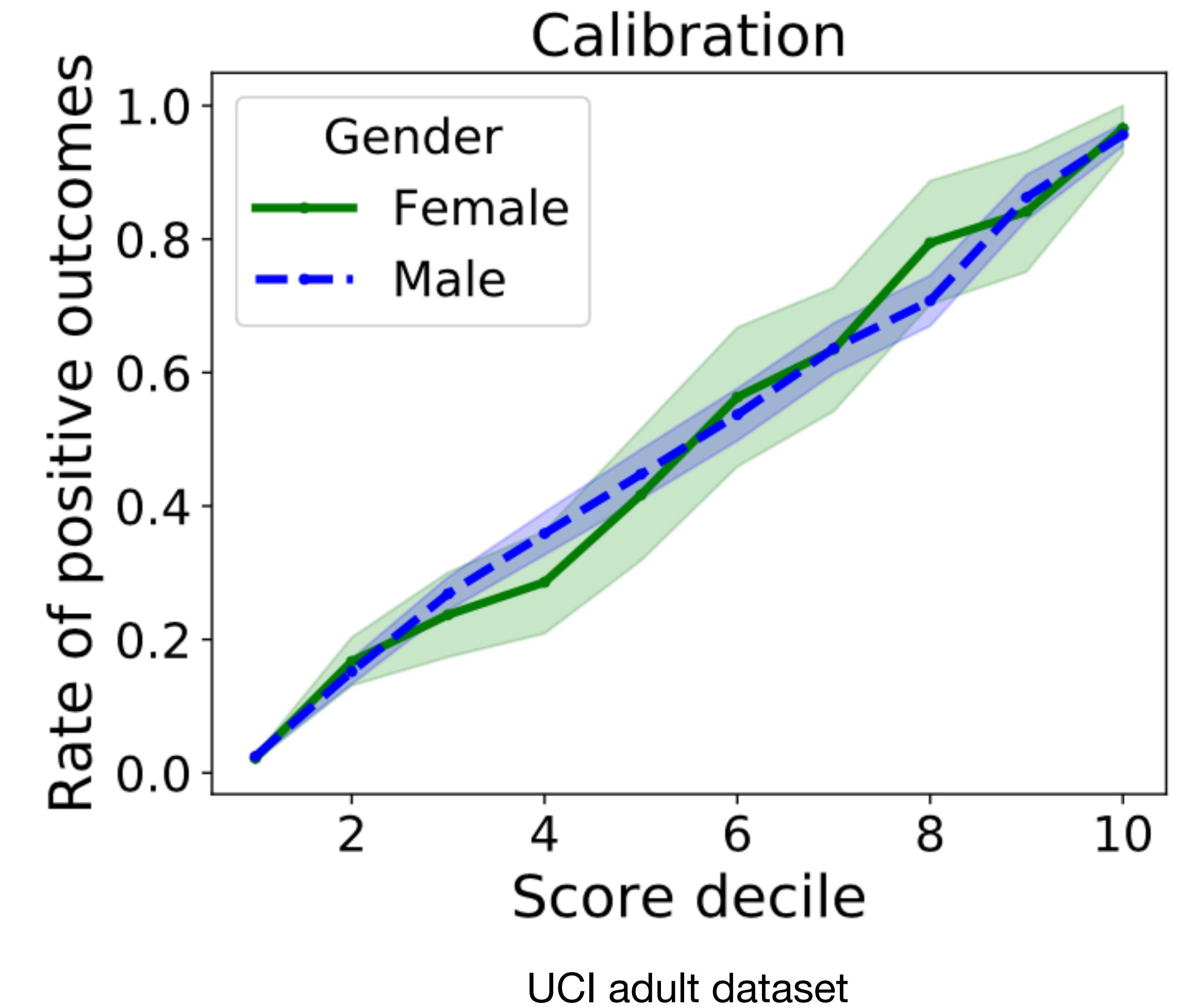
Group calibration as a consequence of unconstrained learning

Sufficiency is often satisfied by default

Not a constraint on ML but what ML does by default

Informal theorem: the better you are at ML, i.e., the closer you are from the Bayes optimal, the more your group calibration is satisfied

- Liu, Simchowitz, Hardt, 2019



Can we satisfy them simultaneously?

Three criteria:

- Independence: R independent of A
 \implies equal acceptance rate
- Separation: R independent of A , conditional on Y
 \implies equal error rate
- Sufficiency: Y independent of A conditional on R
 \implies calibration by group

Informal theorem: any of these criteria are mutually exclusive - except in degenerate cases!

Incompatibility results: trade-offs are necessary

1. Independence vs sufficiency: If A and Y are not independent, then sufficiency and independence cannot both hold

Proof: $A \perp R$ and $A \perp Y|R \implies A \perp (Y, R) \implies A \perp Y$

2. Independence vs separation: if A is not independent of Y and R is not independent of Y , then independence and separation cannot both hold

Proof: $A \perp R$ and $A \perp R|Y \implies A \perp Y \text{ or } R \perp Y$

Proof of the second implication

Claim: $A \perp R$ and $A \perp R | Y \implies A \perp Y \text{ or } R \perp Y$

Proof: $\mathbb{P}(R = r | A = a) = \sum_y \mathbb{P}(R = r | A = a, Y = y) \mathbb{P}(Y = y | A = a)$

Since $A \perp R$ and $A \perp R | Y$:

$$\mathbb{P}(R = r) = \mathbb{P}(R = r | A = a) = \sum_y \mathbb{P}(R = r | Y = y) \mathbb{P}(Y = y | A = a)$$

We also have

$$\mathbb{P}(R = r) = \sum_y \mathbb{P}(R = r | Y = y) \mathbb{P}(Y = y)$$

Thus

$$\sum_y \mathbb{P}(R = r | Y = y) \mathbb{P}(Y = y | A = a) = \sum_y \mathbb{P}(R = r | Y = y) \mathbb{P}(Y = y)$$

Proof of the second implication

Since $Y \in \{0,1\}$ it implies

$$\begin{aligned}\mathbb{P}(R = r | Y = 0)\mathbb{P}(Y = 0 | A = a) + \mathbb{P}(R = r | Y = 1)\mathbb{P}(Y = 1 | A = a) \\ = \mathbb{P}(R = r | Y = 0)\mathbb{P}(Y = 0) + \mathbb{P}(R = r | Y = 1)\mathbb{P}(Y = 1)\end{aligned}$$

It directly implies

$$\begin{aligned}\mathbb{P}(Y = 0)(\mathbb{P}(R = r | Y = 0) - \mathbb{P}(R = r | Y = 1)) \\ = \mathbb{P}(Y = 0 | A = a)(\mathbb{P}(R = r | Y = 0) - \mathbb{P}(R = r | Y = 1))\end{aligned}$$

Therefore either $\mathbb{P}(Y = 0) = \mathbb{P}(Y = 0 | A = a)$ and $A \perp Y$

Or $\mathbb{P}(R = r | Y = 0) = \mathbb{P}(R = r | Y = 1)$ and $Y \perp R$

Incompatibility results: trade-offs are necessary

3. Separation vs sufficiency: Assume all events in the joint distribution of (A, R, Y) have positive probability and assume $A \not\perp\!\!\!\perp Y$. Then, separation and sufficiency cannot both hold

Proof:

$$A \perp R | Y \text{ and } A \perp Y | R \implies A \perp (R, Y) \implies A \perp R \text{ and } A \perp Y$$

Conclusion (I)

- What is the purpose of a fairness criterion?
- Fairness criteria are subject to fundamental limitations
- Fairness criteria do not rule out unfair practices
see Compass debate
- Fairness criteria are not proof of fairness
- Need to go beyond the joint distribution of (X, Y, R, A) - causal reasoning

Conclusion (II)

- These examples (hiring, lending) are controversial - you might not end up working in
- Black-box ML has no guarantee of being aligned with societal values
- How can you design product that leverages ML and which are aligned with human values?

Other concerns: ethics in data collection

The New York Times

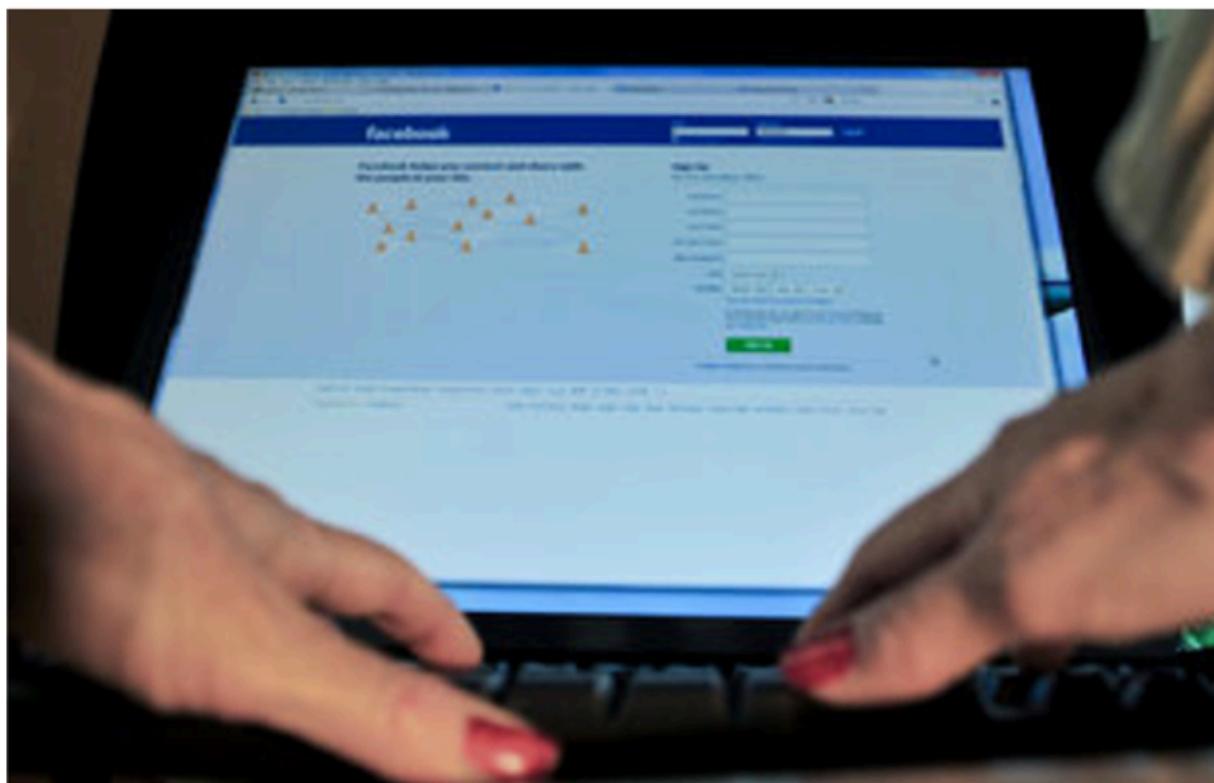
LOG IN

TECHNOLOGY

Facebook Tinkers With Users' Emotions in News Feed Experiment, Stirring Outcry

By VINDU GOEL JUNE 29, 2014

216



Facebook revealed that it had altered the news feeds of over half a million users in its study.
Karen Bleier/Agence France-Presse — Getty Images

Slide borrowed from
Nathan Kallus

Other concerns: privacy

Fredrikson, Jha, Ristenpart, 2015



Figure 1: An image recovered using a new model inversion attack (left) and a training set image of the victim (right). The attacker is given only the person's name and access to a facial recognition system that returns a class confidence score.

Gaydar: Facebook friendships expose sexual orientation

by Carter Jernigan and Behram F.T. Mistree

You Can't Keep Your Secrets From Twitter

On the Internet, no one knows you're secretly a man (or woman), right? Think again. Just by examining patterns in tweets, you can infer a Twitter user's gender. A look at the words (Etsy, Jeep, redneck...) that make men and women give themselves away.

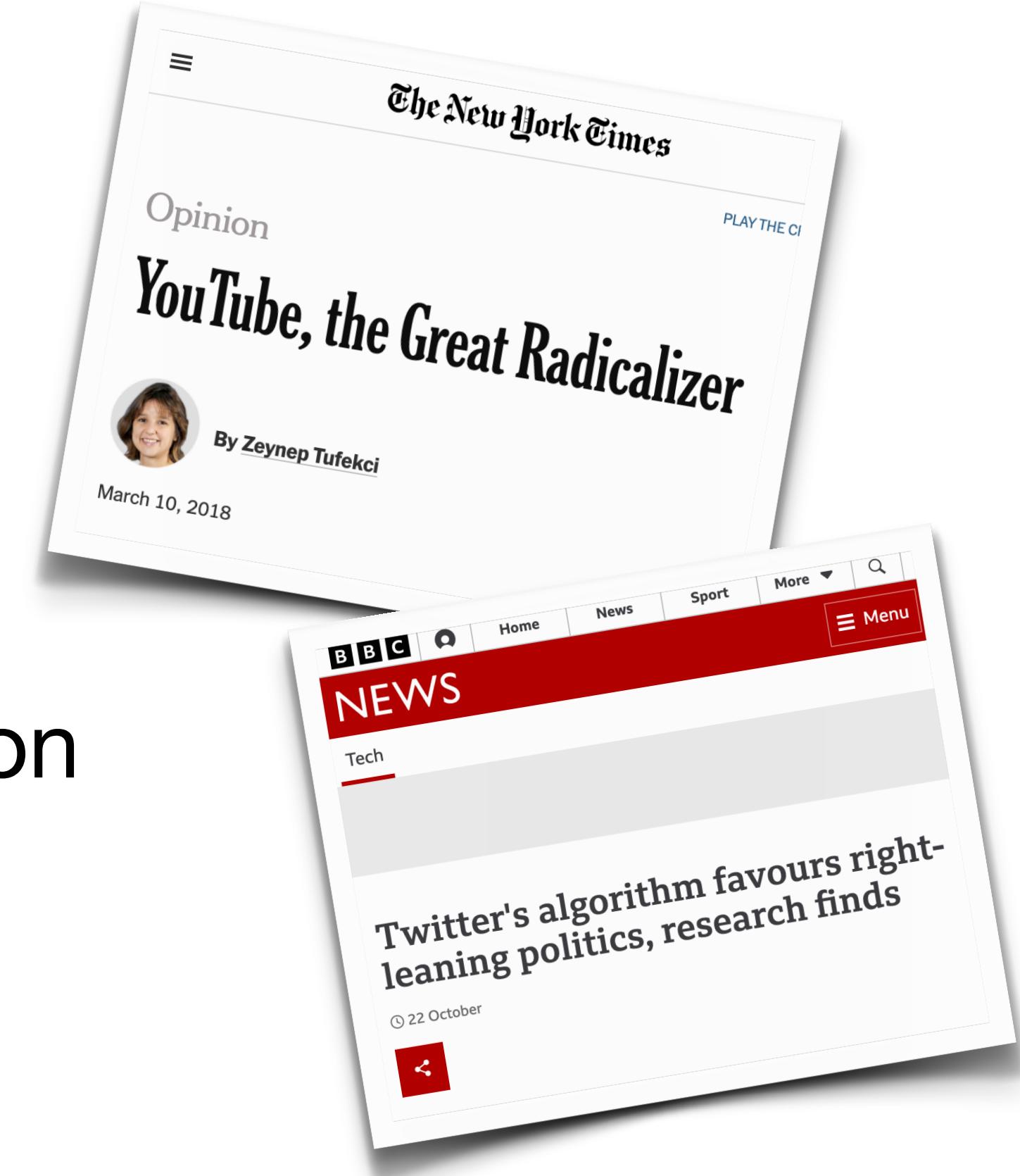
Slide borrowed from
Nathan Kallus

Other concern: the Politics of Recommender Systems

Radicalization of users on social media platforms

- Does it exist? Societal issue?
- Is it fuelled by the recommender system?

EPFL study: A significant amount of users commenting on YouTube **systematically migrate** from commenting exclusively on **milder** content to commenting on more **extreme** content



Auditing Radicalization Pathways on YouTube
MH Ribeiro et al., 2020



Algorithmic Amplification of Politics on Twitter
F Huszár et al., 2021