

# Fairness in Machine Learning

Virginie Do

Université Paris Dauphine - PSL

May 12<sup>th</sup>, 2021

# Outline

1. **Biases in AI systems**
2. Fairness in machine learning: binary decisions
3. Beyond fair learning

# Biases in AI systems

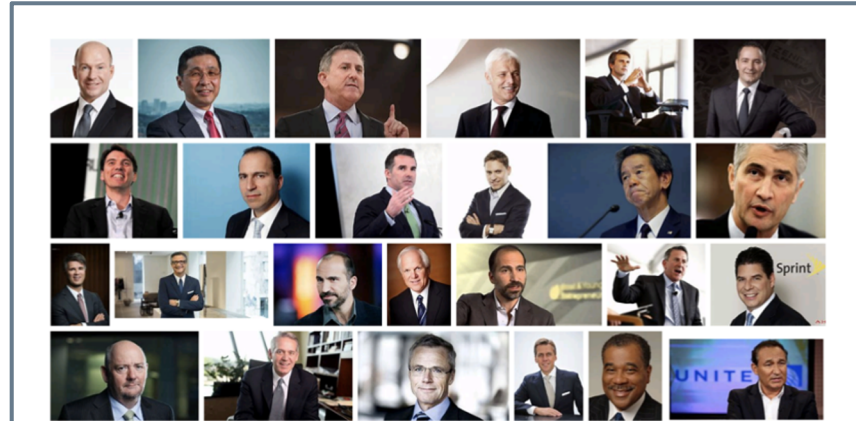
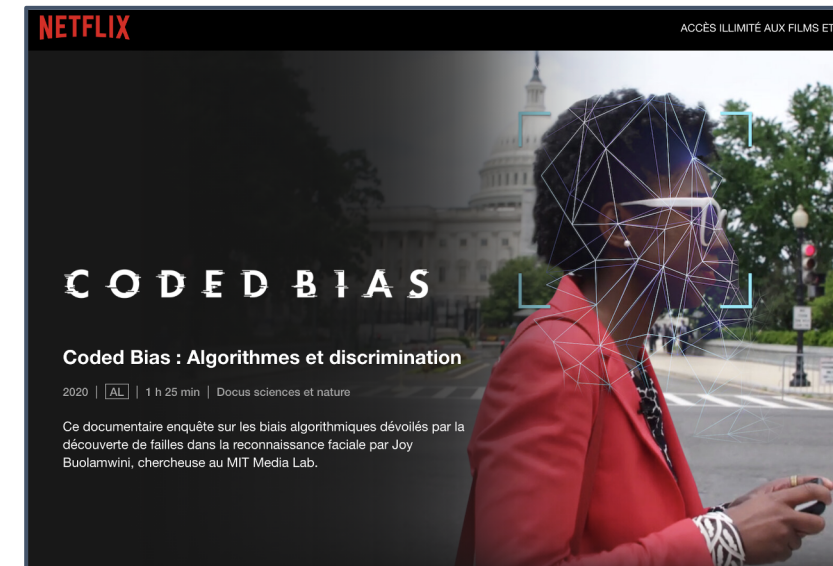
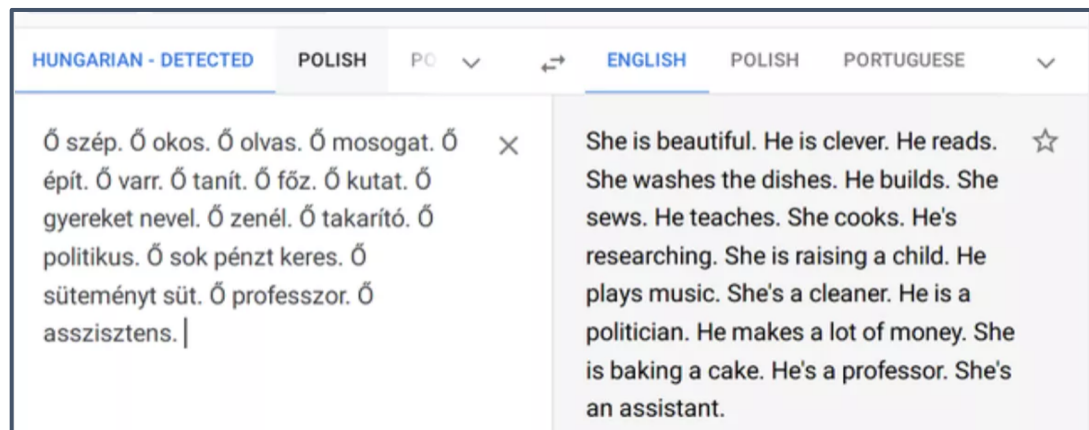


Figure 2: An image search result page for the query "CEO" showing a disproportionate number of male CEOs.



RETAIL | OCTOBER 11, 2018 / 1:04 AM / UPDATED 2 YEARS AGO

## Amazon scraps secret AI recruiting tool that showed bias against women

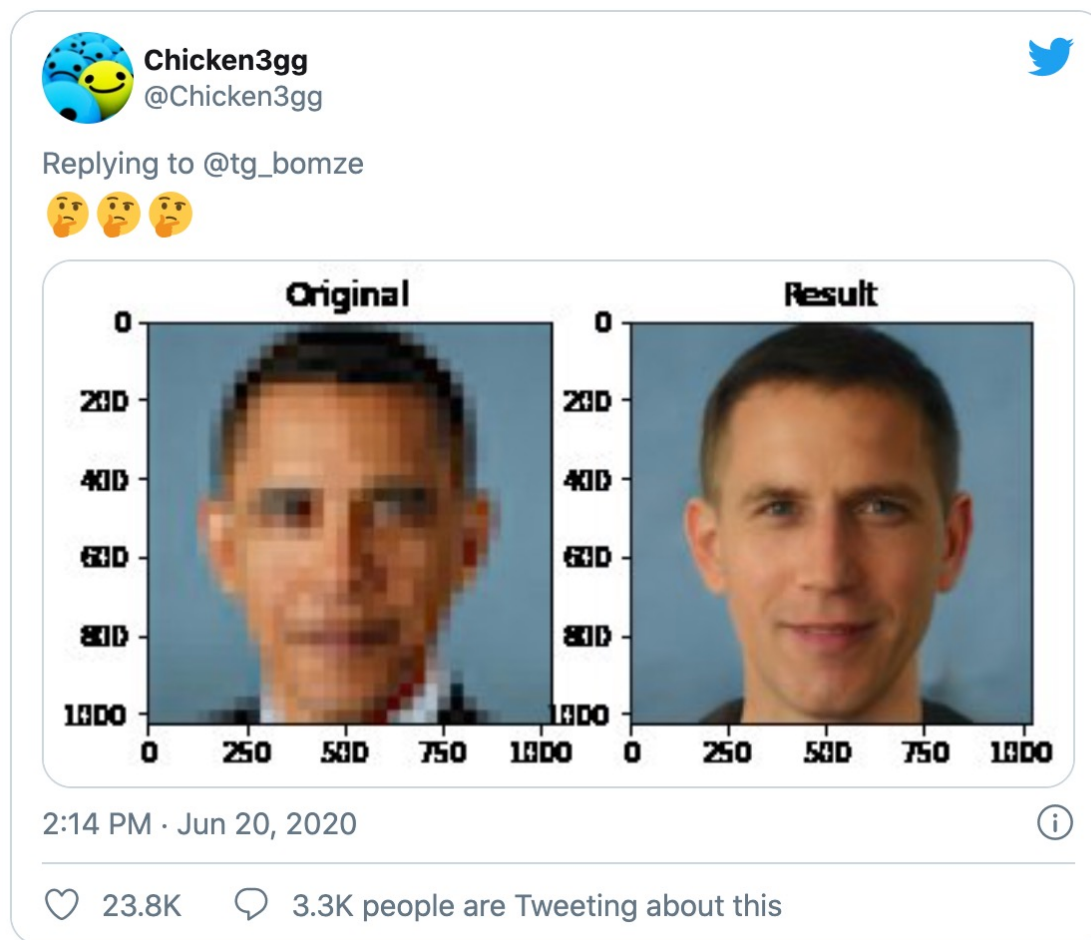


# Biases in AI systems

## PULSE: Self-Supervised Photo Upsampling via Latent Space Exploration of Generative Models

Sachit Menon\*, Alexandru Damian\*, Shijia Hu, Nikhil Ravi, Cynthia Rudin  
Duke University  
Durham, NC

{sachit.menon, alexandru.damian, shijia.hu, nikhil.ravi, cynthia.rudin}@duke.edu



# Outline

1. Biases in AI systems
2. **Fairness in machine learning: binary decisions**
3. Beyond fair learning

# Fairness in Machine Learning – Binary predictions



RETAIL OCTOBER 11, 2018 / 1:04 AM / UPDATED 2 YEARS AGO

**Amazon scraps secret AI recruiting tool that showed bias against women**

**The Apple Card Is the Most High-Profile Case of AI Bias Yet**

Binary decisions: Good vs. Bad outcome

Applications: Recidivism prediction, Loan approval, Job application

# Fairness in Machine Learning – Typical setup

	Example: Lending
$A$ sensitive attribute	Gender (Men/Women)
$X$ “relevant” features	Salary, Debt history
$Y$ actual outcome	Repaid / Default
$\hat{Y} = f(X, A)$ predictor	Classifier
$\hat{S} = g(X, A)$ score function (can be turned into binary decision)	Credit score

# Fairness criteria in Machine Learning

Demographic parity

$$P(\hat{Y} = 1 | A = 0) = P(\hat{Y} = 1 | A = 1)$$

Equal opportunity

$$P(\hat{Y} = 1 | Y = 1, A = 0) = P(\hat{Y} = 1 | Y = 1, A = 1)$$

Calibration within groups

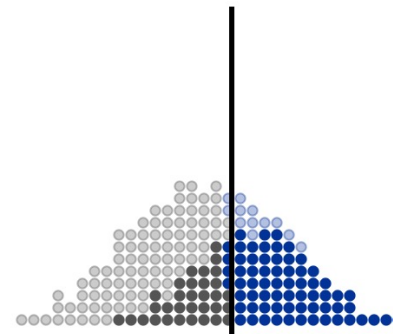
$$P(Y = 1 | \hat{S} = s, A = 0) = P(Y = 1 | \hat{S} = s, A = 1)$$

→ Incompatibility

Blue Population

0 10 20 30 40 50 60 70 80 90 100

loan threshold: 59

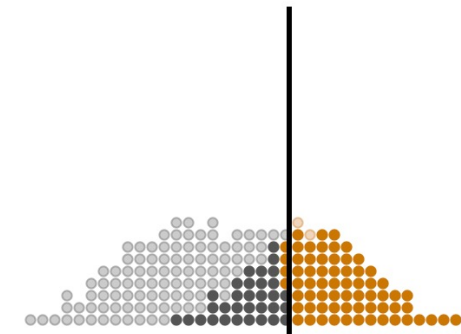


denied loan / would default    granted loan / defaults  
denied loan / would pay back    granted loan / pays back

Orange Population

0 10 20 30 40 50 60 70 80 90 100

loan threshold: 53



denied loan / would default    granted loan / defaults  
denied loan / would pay back    granted loan / pays back

3

S. Corbett-Davies et al. '17

J. Kleinberg et al. '16, A. Chouldechova '16

<https://research.google.com/bigpicture>



# Trade-offs

Many more definitions...

- More parity measures
- Individual metric-based fairness
- Counterfactual fairness

And trade-offs:

- Between different measures of group fairness
- Between group fairness and individual fairness
- *Between group fairness and group fairness*
- Between fairness and utility

Translation tutorial:  
21 fairness definitions and their politics

Arvind Narayanan  
@random\_walker



Dwork et al., *Individual fairness*, 2012  
Kusner et al., *Counterfactual fairness*, 2017  
Kearns et al., *Preventing fairness gerrymandering*, 2017

# Fair algorithms

## 1. Pre-processing

Learning fair representations

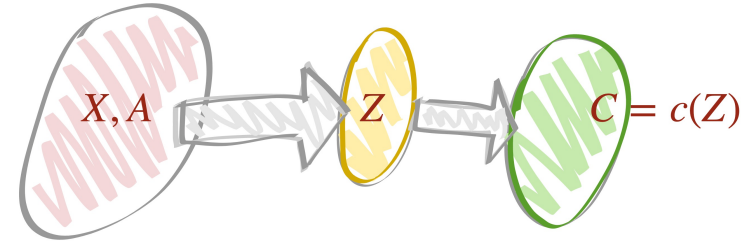
## 2. Optimization at training time

Empirical risk minimization with constraint, regularization term

## 3. Post-processing

Threshold on a score function

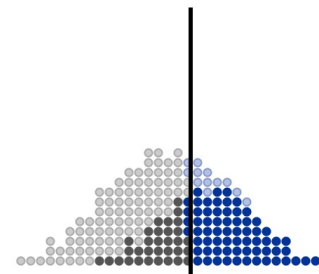
## Representation learning approach



Blue Population

0 10 20 30 40 50 60 70 80 90 100

loan threshold: 59

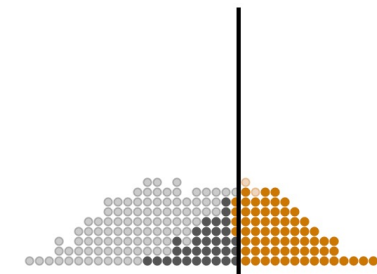


denied loan / would default    granted loan / defaults  
denied loan / would pay back    granted loan / pays back

Orange Population

0 10 20 30 40 50 60 70 80 90 100

loan threshold: 53



denied loan / would default    granted loan / defaults  
denied loan / would pay back    granted loan / pays back

# Additional references

## Tutorials

- Hardt and Barocas, tutorial @ NeurIPS '17
- Narayanan @ FAccT '18

## Surveys

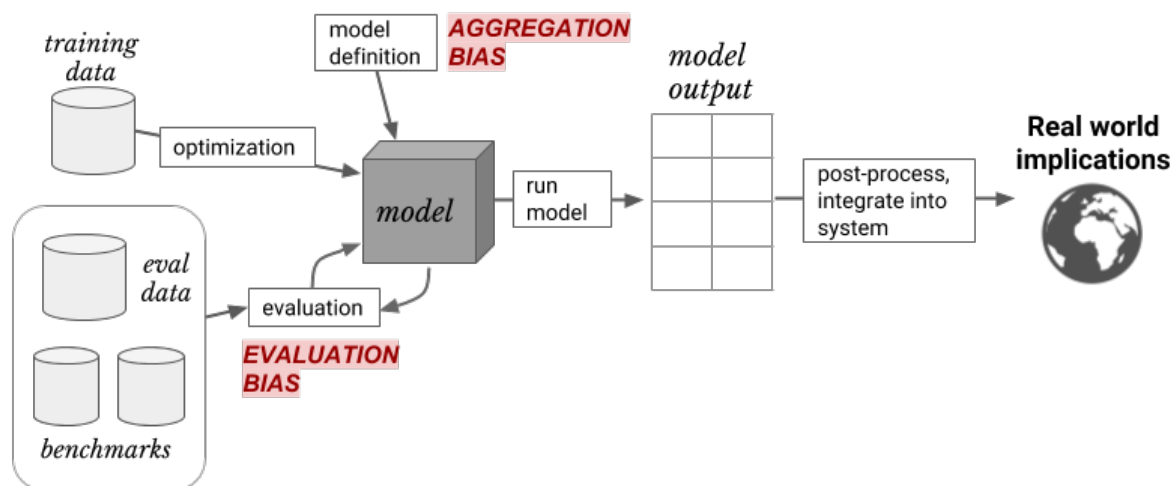
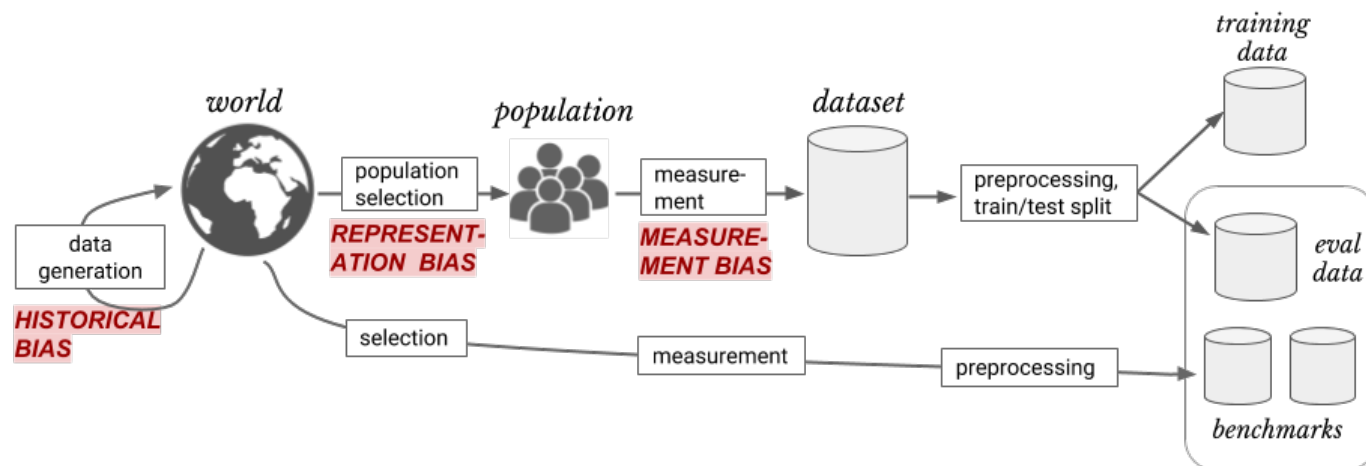
- Chouldechova and Roth, *The frontiers of fairness in machine learning*
- Corbett-Davies and Goel, *The measure and mismeasure of fairness*
- Barocas, Hardt, Narayanan, *Fairness and machine learning: limitations and opportunities*. [Book]

# Outline

1. Biases in AI systems
2. Fairness in machine learning: binary decisions
3. **Beyond fair learning**

# Biased data?

Potential sources of harm arise at different stages of the ML pipeline



# Parity vs. preference

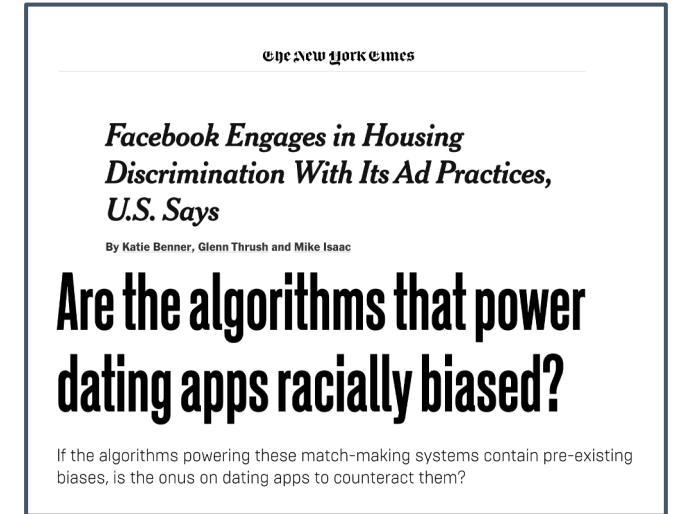
When subjects have different preferences / utilities, should they be given the same predictions?

→ Personalization

## Preference guarantees

with concepts like “envy-freeness”: no one should prefer someone else’s model to their given model.

Zafar et al. '16, Ustun et al. '19, Balcan et al. '19, Kim et al. '20



# Discussion

- Interdisciplinarity
  - “Mathematical” fairness for computer scientists vs. fairness for ethicists, philosophers, legal scholars, economists...
- Context
  - Applications: which fairness definition for which specific context? should ML be used at all?
  - Fairness for unobserved characteristics: ethnicity, sexual orientation.
  - Complex pipelines
- Explainability