

Defining and mitigating algorithmic bias: a practitioner's perspective

Hinda Haned

April 13th 2021

University of Amsterdam

h.haned@uva.nl



Complex systems raise concern



NOS | Nieuws | Sport | Uitzendingen | TELEHOOR | AEX 94 km 7°

Nieuws

NEWSBURO • BIJENLAND • POLITIEK • MA 21 OKTOBER 18 13

VN-rapporteur zeer bezorgd over Nederlands opsporingssysteem voor uitkeringsfraude

De VN-rapporteur voor de mensenrechten Philip Alston heeft ernstige zorgen over Nederland. De reden is een systeem dat uitkeringsfraude moet opsporen. In een brief aan de rechtbank in Den Haag schrijft Alston dat het systeem in strijd is met de mensenrechten omdat het mensen met weinig geld en mensen met een migratie-achtergrond discrimineert.



BBC | Sign in | News | Sport | Reel | Worklife | Travel | Future | M

NEWS

Home | Video | World | UK | Business | Tech | Science | Stories | Entertainment & Arts

Amazon scrapped 'sexist AI' tool

🕒 10 October 2018 | [Share](#)

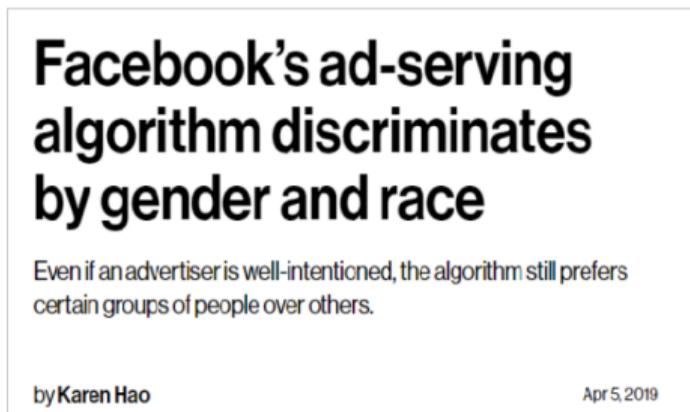


2,780 views | Jan 28, 2019, 07:20pm

Forbes

Amazon Refuses To Quit Selling 'Flawed' And 'Racially Biased' Facial Recognition

Zak Doffman Contributor



Facebook's ad-serving algorithm discriminates by gender and race

Even if an advertiser is well-intentioned, the algorithm still prefers certain groups of people over others.

by Karen Hao | Apr 5, 2019

What is bias?

- Systematic errors that create unfair outcomes
- Sources: algorithm design, biased data collection or selection
- Algorithms learn and perpetuate bias

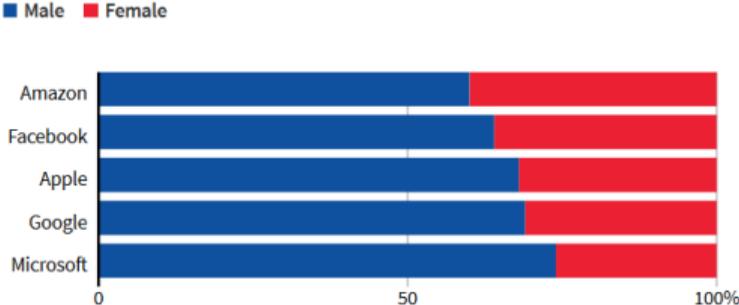
Types of bias

- **Historical bias** reflects structural societal issues
- **Representation bias** certain groups are under-represented in the training data
- **Measurement bias** training data are proxies for some ideal features and labels

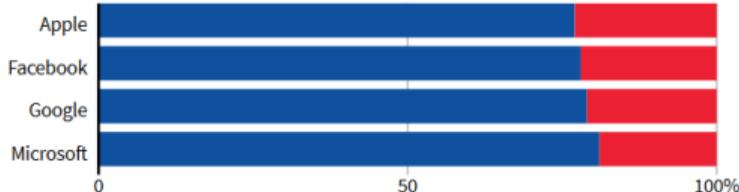
simplified from Suresh & Guttag. A Framework for understanding unintended consequences of machine learning, 2019.

Historical bias

GLOBAL HEADCOUNT



EMPLOYEES IN TECHNICAL ROLES

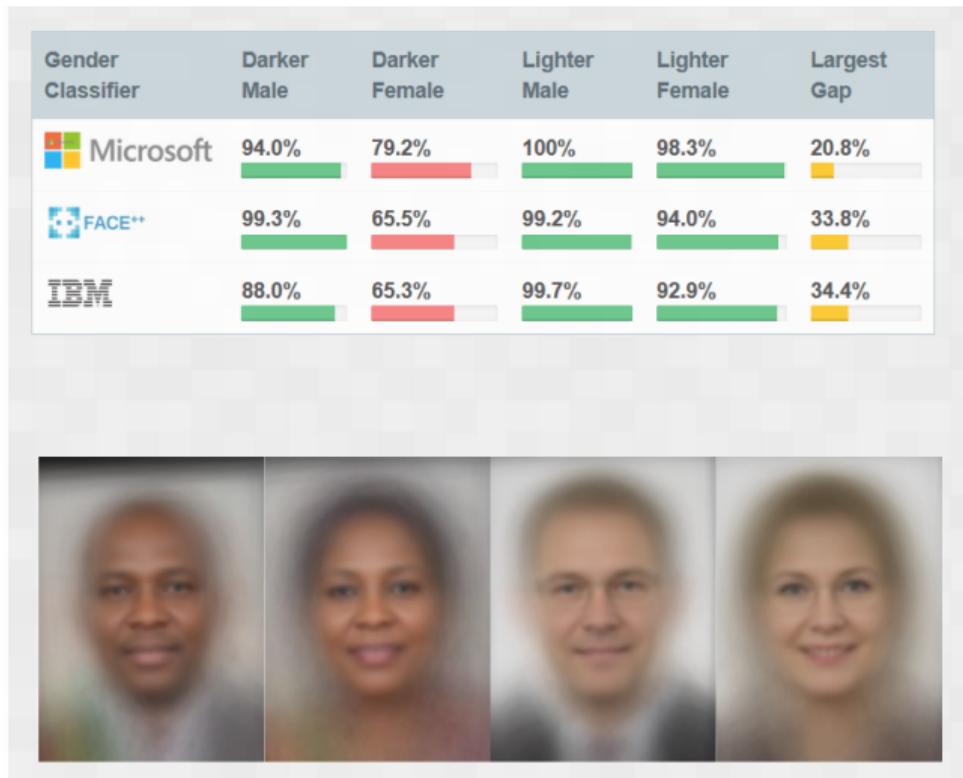


Note: Amazon does not disclose the gender breakdown of its technical workforce.

Source: Latest data available from the companies, since 2017.

By Han Huang | REUTERS GRAPHICS

Representation bias



<http://gendershades.org/overview.html>

Measurement bias

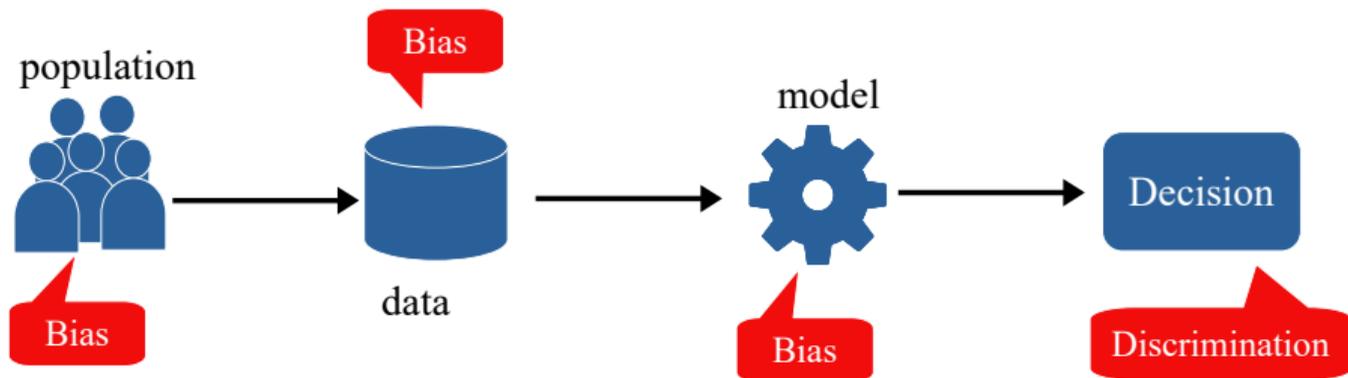


How do we mitigate algorithmic bias in practice?

Mitigating algorithmic bias

- There is no unifying framework to tackle algorithmic bias testing and mitigation
- In most use cases, mitigation is performed after a system is built and decisions have been made based on this system

Bias occurs throughout the data science pipeline



Mitigation algorithms

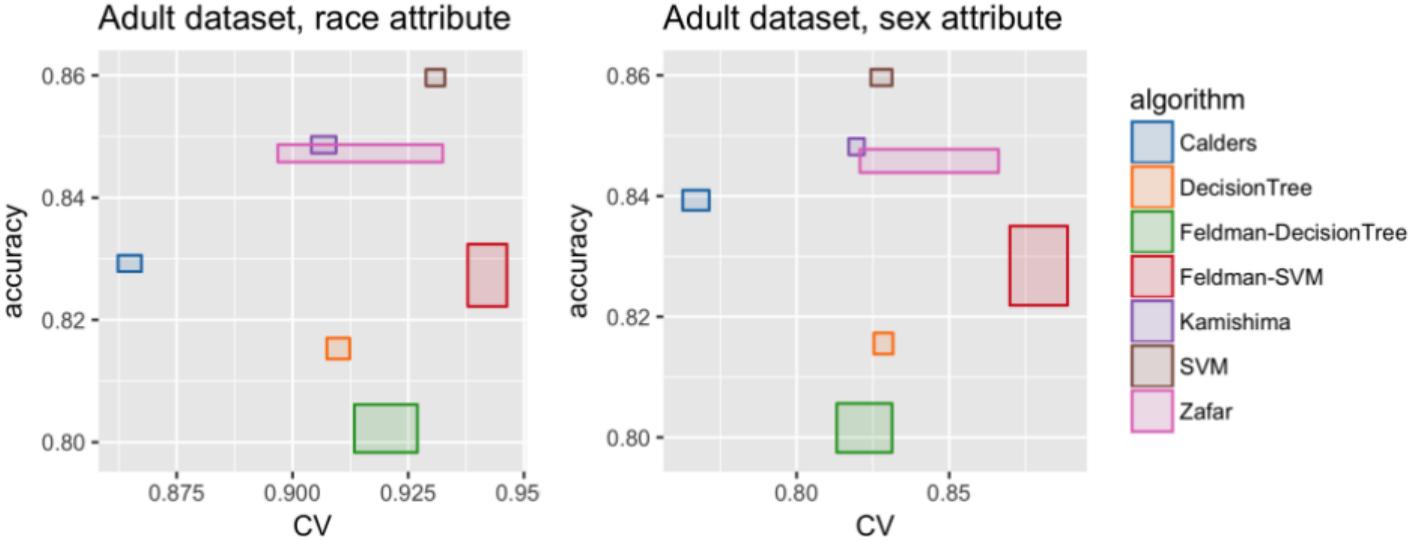
Mitigation or fairness algorithms have only been developed for classification tasks:

- Pre-processing: modify the train data
- In-processing: modify the algorithm's objective function to incorporate fairness constraints/penalty
- Post-processing: modifies the predictions produced by the algorithm

Limitations of mitigation algorithms

- Unrealistic assumptions: sensitive attributes are known & ground truth or observable outcomes are available
- Trade-off: utility vs. desired (quantifiable) measure of fairness
- Reliability: sensitivity to fluctuations in dataset composition, and to different forms of pre-processing (Friedler et al, 2019)

Illustration: limitations of mitigation algorithms



Friedler et al. A comparative study of fairness-enhancing interventions in machine learning, FAT* 2019.

Beyond mitigation

“Any real machine-learning system seeks to make some change in the world. To understand its effects, then, we have to consider it in the context of the larger socio-technical system in which it is embedded.”

Barocas et al. Fairness and machine learning, fairmlbook.org, 2019.

Mitigation vs. Fundamental Questions

- Why do you need an automated system for this task?
- Is the system transparent?
- What are the potential harms that could occur?
- What is a fair outcome? What is an unfair outcome?
- When and how does the system fail?
- Who is responsible for the errors?

Thank you

<https://hindantation.github.io>

Mitigation algorithms

Pre-processing	Re-weighting (Kamiran & Calders, 2012) Optimized pre-processing (Calmon et al., 2017) Learning fair representations (Zemel et al., 2013) Disparate impact remover (Feldman et al., 2015)
In-processing	Adversarial debiasing (Zhang et al., 2018) Prejudice remover (Kamishima et al., 2012)
Post-processing	Equalized odds post-processing (Hardt et al., 2016) Calibrated eq. odds postprocessing (Pleiss et al., 2017) Reject option classification (Kamiran et al., 2012) Fairness-focused regularization (Kamishima et al, 2019) Two Naive Bayes (Calders & Verwer, 2010)
