

# AI algorithms: defining and mitigating bias

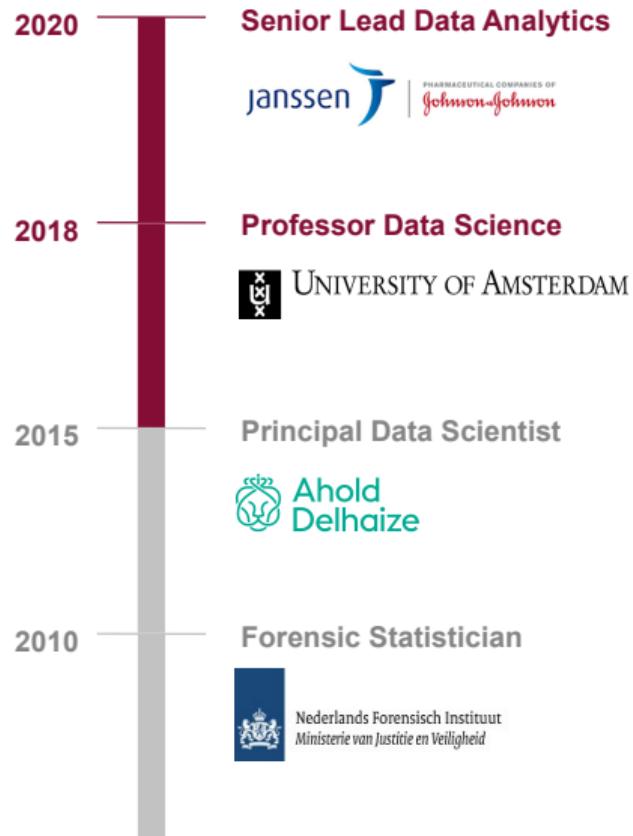
---

**Hinda Haned**

Januray 20th, 2021

[h.haned@uva.nl](mailto:h.haned@uva.nl)

# About me



# The Civic AI Lab: partners



UNIVERSITEIT VAN AMSTERDAM



An engaging Society where all citizens have **equal** opportunity to participate and engage in a **fair**, **transparent** manner



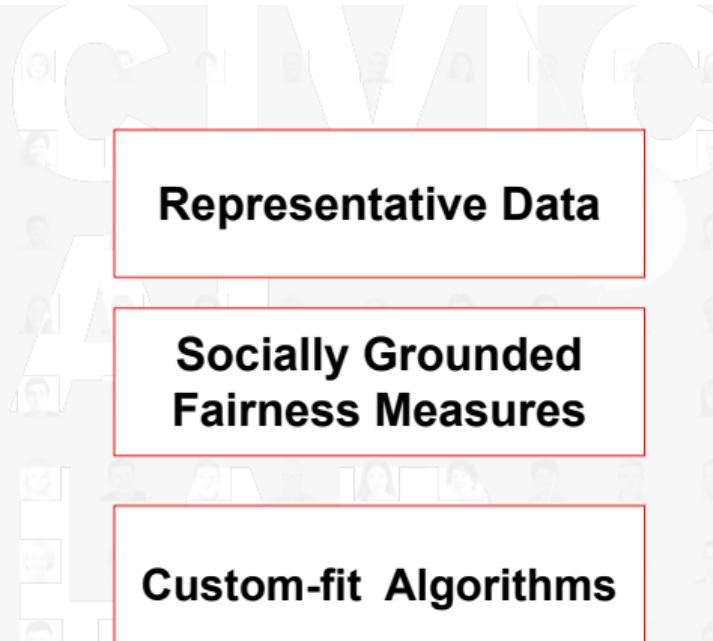
To develop AI that promotes economic & social human rights

To develop AI that advances inclusive socio-cultural systems

Representative Data

Socially Grounded Fairness Measures

Custom-fit Algorithms



UvA



VU



Gemeente  
Amsterdam



IAI



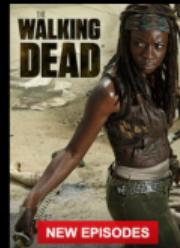
ICAI

How do AI algorithms help us?

# NETFLIX



Popular on Netflix



TV Programmes based on Books



## Voor jou in de Bonus

30 van je eerdere  
aankopen zijn in de  
Bonus



BONUS

1.<sup>49</sup>

500 g

BONUS

0.<sup>59</sup>  
per stuk

BONUS

1.<sup>09</sup>  
per stuk

AH Pitloze witte druiven



AH Courgette



AH Avocado



2 VOOR 2.99

1.<sup>69</sup>  
per stuk

BONUS

2.<sup>29</sup>  
140 g

BONUS

AH Sweet eve frambozen



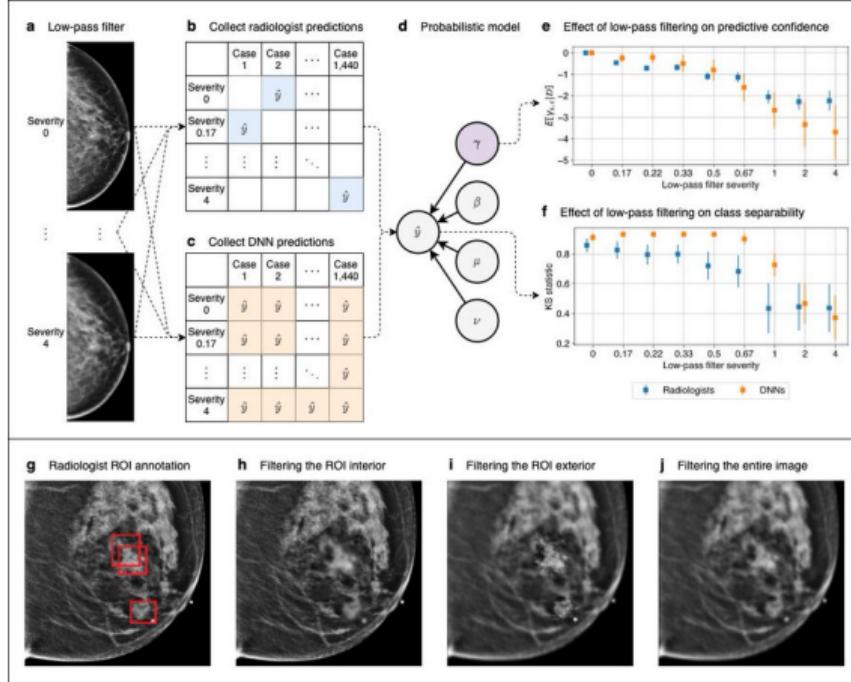
AH Mango

0.<sup>59</sup>  
per stuk

25% KORTING

4.<sup>08</sup>  
6 x 0,33 l

Coca-Cola Regular



## Breast cancer screening

Source: Makino et al, 2020, <https://arxiv.org/abs/2011.14036>



<https://gen.com/articles/2019/07/10/ai-disaster-relief.aspx>

Can AI algorithms fail?

**NOS** Nieuws Sport Uitzendingen

TELEBRIEF AEX 94 km 7°

nieuwsuur  
NEDERLAND • POLITIE • MA 29 OKTOBER, 19:19

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

# Algorithmic fairness: why bother?

## Compliance

Compliance with General Data Protection Regulation (GDPR – 2018) – article 22

The data subject shall have the **right not to be subject to a decision based solely on automated processing**, including **profiling**, which produces legal effects concerning him or her or similarly significantly affects him or her

## GDPR fines & notices

# GDPR fines & notices

Date	Organisation	Amount	Issued by	Reason(s)
2020-07-14	<a href="#">Google LLC (Google Belgium)</a>	€600,000	Belgium ( <a href="#">GBA/APD</a> )	Failure to respect a citizen's right to be forgotten.
2020-07-06	BKR	€840,000	The Netherlands ( <a href="#">AP</a> )	Failing to give access to personal data free of charge, failing to provide easy means of accessing the data, putting unreasonable limits on the number of requests per individual <a href="#">[43]</a>
2019-12-17	Doorstep Dispensaree	£275,000	UK ( <a href="#">ICO</a> )	"cavalier attitude to data protection", having left 500,000 patient records in an unsecured location <a href="#">[42]</a>
2019-12-09	<a href="#">1&amp;1 Ionos</a>	€9,550,000	Germany ( <a href="#">BfDI</a> )	Insufficient protection of personal data, failing to put "sufficient technical and organizational measures" in place to protect customer data in its call centers. Violation of article 32 of GDPR <a href="#">[41]</a>
2019-10-17	<a href="#">Vueling Airlines</a>	€30,000	Spain ( <a href="#">AEPD</a> )	Failing to obtain valid consent to process customer cookies, as per privacy notice. <a href="#">[40]</a>
2019-09-20	Online retailer Morele.net	€645,000	Poland ( <a href="#">UODO</a> )	Insufficient protection of personal data, leading to the exposure of data of about 2.2 million people <a href="#">[39]</a>
2019-09-19	Unnamed retailer	€10,000	Belgium ( <a href="#">GBA/APD</a> )	Demanding an electronic identity card to create a customer loyalty card. <a href="#">[38]</a>

## Complex systems raise concern

- Why this ad?
- Why this discount?
- Why this recommendation?
- Why was I rejected?
- Can I change the outcome?
- When will the system fail?

“Algorithms are written and maintained by people, and machine learning algorithms adjust what they do based on people’s behavior. As a result algorithms can reinforce human prejudices.”

C.C. Miller. When algorithms discriminate, NYT, 2019.

## What is bias?

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

## What is fairness?

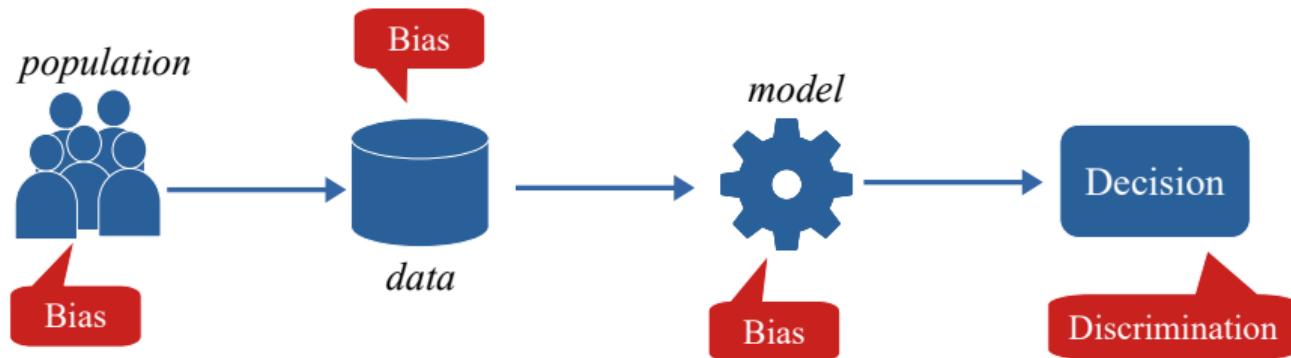
- Fairness is concerned with how outcomes are assigned to particular group of individuals
- Core principles: avoid bias even if it is supported by data, as to avoid the perpetuation of existing discrimination (distributive justice)
- Fairness is a political construct: someone decides to avoid (direct or indirect) harm

## Types of harm

- **Harms of allocation** when a system allocates or withholds certain groups, an opportunity or resource. Economically oriented view (e.g. who gets a discount, who gets hired),
- **Harms of representation** systems reinforce the subordination of certain groups along the lines of identity like race, class, gender etc. (e.g. search results biased against a group),
- **Reduced quality of service** when a system produces degraded user/customer experience based on their personal characteristics (e.g stereotyping and denigration).

adapted from: Kate Crawford's NIPS 2017 Keynote presentation: the trouble with Bias & Challenges of incorporating algorithmic 'fairness' into practice: FAccT 2019 tutorial

## Bias occurs throughout the algorithmic pipeline



## Types of bias

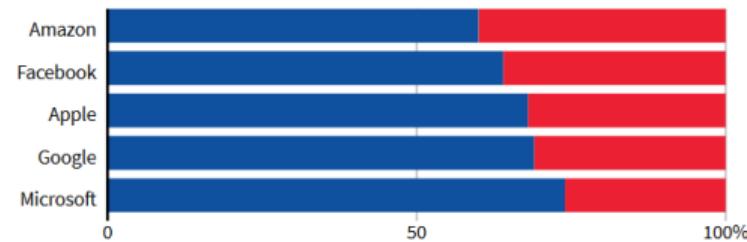
- **Historical bias** reflects structural societal issues
- **Representation bias** certain groups are under-represented in the training data
- **Measurement bias** training data are proxies for some ideal features and labels

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

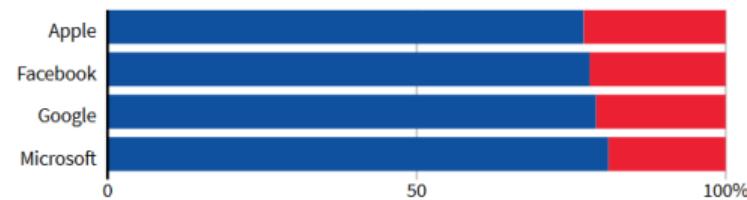
# Historical bias

## GLOBAL HEADCOUNT

■ Male ■ Female



## EMPLOYEES IN TECHNICAL ROLES

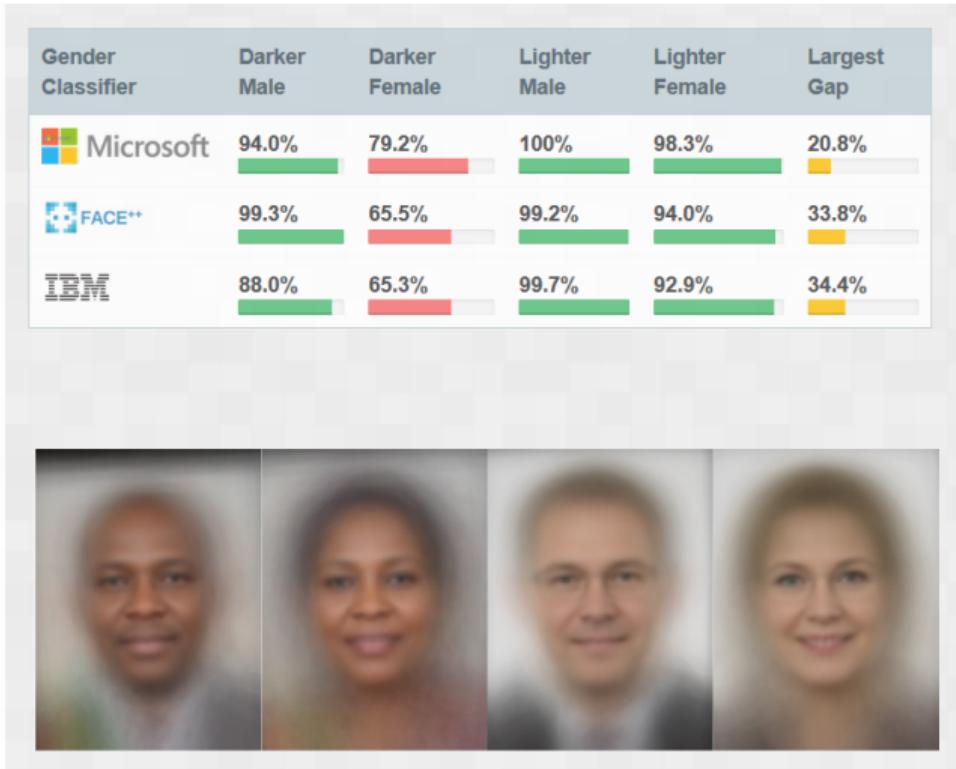


Note: Amazon does not disclose the gender breakdown of its technical workforce.

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

By Han Huang | REUTERS GRAPHICS

# Representation bias



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

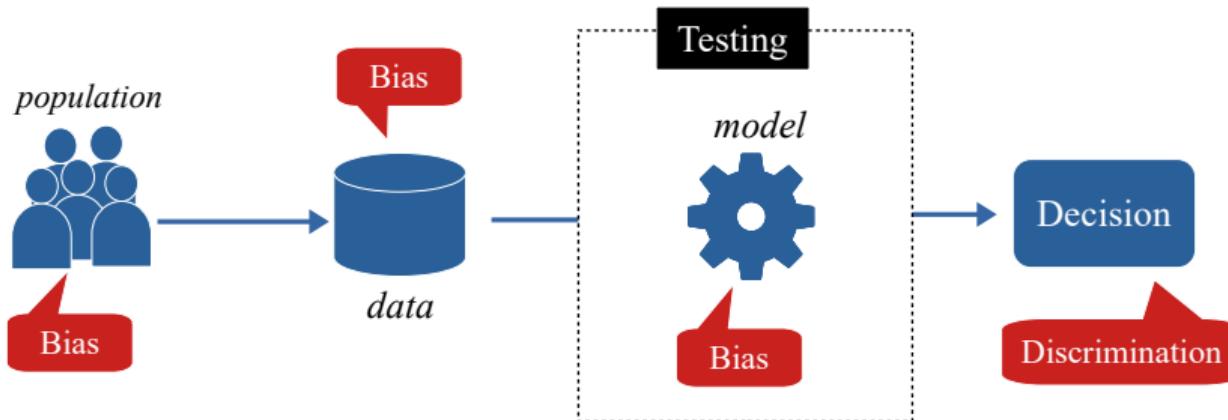
## Measurement bias



<https://www.businessinsider.com>

How can we avoid bias?

## Best practice



- check datasets imbalances
- ensure model treats all groups fairly
- understand model behavior
- monitor prediction

## Regulation: GDPR

“Data subjects have a right to **meaningful information** about the **logic involved** and to the significance and the **envisaged consequence** of automated decision-making”

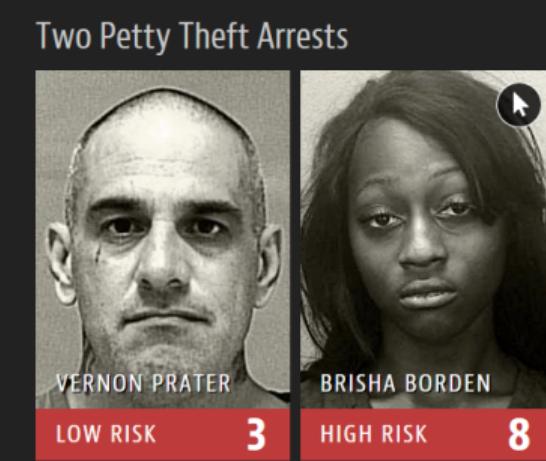
# Ethics



# Open Data

“Through a public records request, ProPublica obtained two years worth of COMPAS scores from the Broward County Sheriff’s Office in Florida. We received data for all 18,610 people who were scored in 2013 and 2014.”

### Two Petty Theft Arrests



The image shows two mugshots side-by-side. On the left is a man with a shaved head, identified as VERNON PRATER, with a risk score of 3 labeled as "LOW RISK". On the right is a woman with dark hair, identified as BRISHA BORDEN, with a risk score of 8 labeled as "HIGH RISK". A cursor icon is visible in the top right corner of the image area.

Arrestee	Risk Score
VERNON PRATER	3
BRISHA BORDEN	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.*

# Open Data

Gemeente  
Amsterdam

Data en informatie

Zoek data op adres, gebied, etc. Of datasets op trefwoord.



Onderdelen Over OIS

## Dataset

### Werk en inkomen (Buurten)

Diverse datasets met statistieken van Onderzoek, Informatie en Statistiek.

Thema: Werk en inkomen,  
Detailniveau: Buurten

## Resources

### Data

2.1 Bedrijfsvestigingen en werkzame personen 1) naar buurten, 1 januari 2016-2019

<https://api.data.amsterdam.nl/dcatd/datasets/-fotdSpwjmSK9Q/purls/1>

2.2 Bedrijfsvestigingen buurten naar hoofdfunctie, 1 januari 2019

<https://api.data.amsterdam.nl/dcatd/datasets/-fotdSpwjmSK9Q/purls/2>

2.3 Werkzame personen buurten naar hoofdfunctie, 1 januari 2019 1)

<https://api.data.amsterdam.nl/dcatd/datasets/-fotdSpwjmSK9Q/purls/3>

2.4 Bedrijfsvestigingen en werkende personen buurten naar sectoren, 1 januari 2019

<https://api.data.amsterdam.nl/dcatd/datasets/-fotdSpwjmSK9Q/purls/4>

2.5a Startende ondernemers naar buurten, 2013-2018 1)

<https://api.data.amsterdam.nl/dcatd/datasets/-fotdSpwjmSK9Q/purls/5>

# ACM Conference on Fairness, Accountability, and Transparency (ACM FAccT)

A computer science conference with a cross-disciplinary focus that brings together researchers and practitioners interested in fairness, accountability, and transparency in socio-technical systems.

<https://facctconference.org/>

## Mitigation algorithms

- **Mitigation** The action of reducing the severity, seriousness, or painfulness of something
- **Mitigation algorithms** Algorithms to remove or reduce bias in data and model outputs

# Mitigation algorithms

The screenshot shows the homepage of the AI Fairness 360 Open Source Toolkit. At the top, there is a navigation bar with links for Home, Demo, Resources, Events, Videos, and Community. The 'Home' link is underlined, indicating it is the current page. Below the navigation bar, the title 'AI Fairness 360 Open Source Toolkit' is displayed in large, bold, dark blue font. A detailed description follows, explaining the toolkit's purpose: 'This extensible open source toolkit can help you examine, report, and mitigate discrimination and bias in machine learning models throughout the AI application lifecycle. Containing over 70 fairness metrics and 10 state-of-the-art bias mitigation algorithms developed by the research community, it is designed to translate algorithmic research from the lab into the actual practice of domains as wide-ranging as finance, human capital management, healthcare, and education. We invite you to use it and improve it.' At the bottom of the main content area, there are two buttons: 'API Docs' and 'Get Code', both with dropdown arrows.

IBM Research Trusted AI

Home Demo Resources Events Videos Community

## AI Fairness 360 Open Source Toolkit

This extensible open source toolkit can help you examine, report, and mitigate discrimination and bias in machine learning models throughout the AI application lifecycle. Containing over 70 fairness metrics and 10 state-of-the-art bias mitigation algorithms developed by the research community, it is designed to translate algorithmic research from the lab into the actual practice of domains as wide-ranging as finance, human capital management, healthcare, and education. We invite you to use it and improve it.

API Docs ▾ Get Code ▾

<https://aif360.mybluemix.net/>

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

## Self-scan checks

- **Model:** classifier trained on historical data to predict when a check is relevant (yes/no)
- **Utility:** reveal mistakes/thefts while preserving customer experience



<http://www.deltifood.nl>

**How can we build a model that is useful to the business and fair to its customers?**

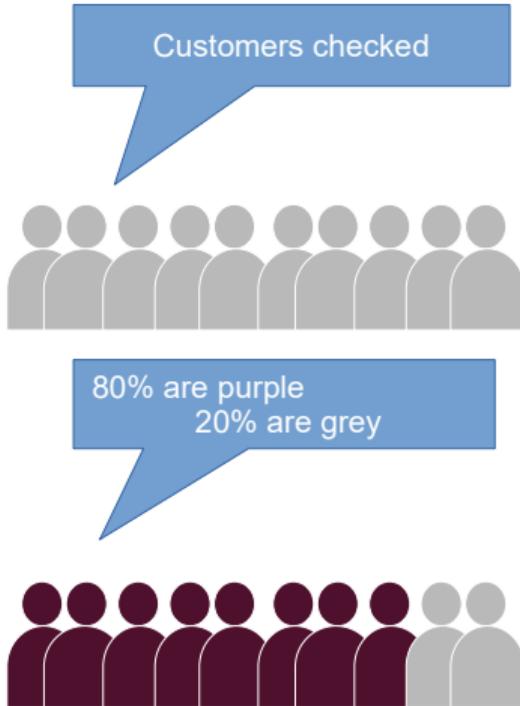
## Questions

- What are the potential harms that could occur in this case?
- What is a fair outcome? What is an unfair outcome?
- How can we ensure fair outcomes?
- Who is responsible for ensuring fairness?
- What are the responsibilities: system owner, store manger, data scientist, customer, legislator?

## Fundamental questions

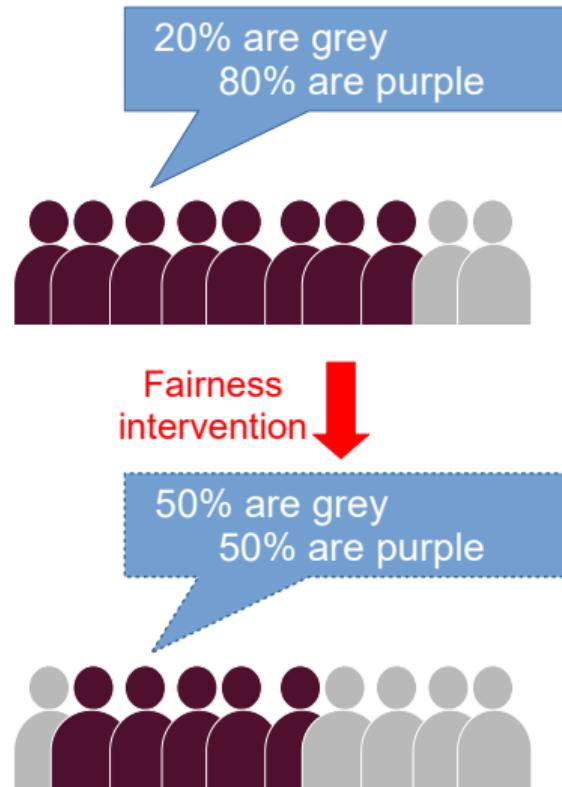
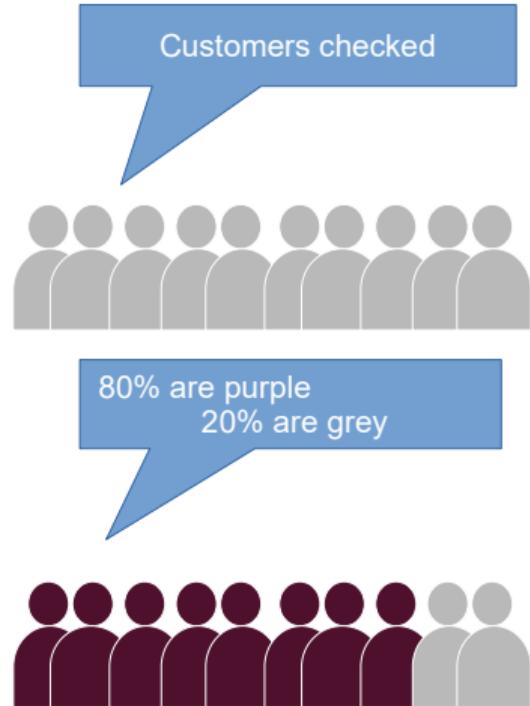
1. Why do you need AI for this task?
2. Is the system transparent?
3. When and how does the system fail?
4. Who is responsible for the errors?

## Business objective: avoid harm



Avoid a degraded customer experience for a given group

## Business objective: avoid harm



## Bias mitigation steps

1. Identify sensitive attribute
2. Define fairness objective
3. Choose fairness intervention
4. Monitor fairness intervention

## 1. Identify sensitive attribute

- Sensitive attributes such as age, gender or ethnicity are not directly encoded in the data but they can strongly correlate with other features, such as purchase behaviors
- Working assumption: sensitive attributes are unknown but underlying distributions are accessible for fairness testing

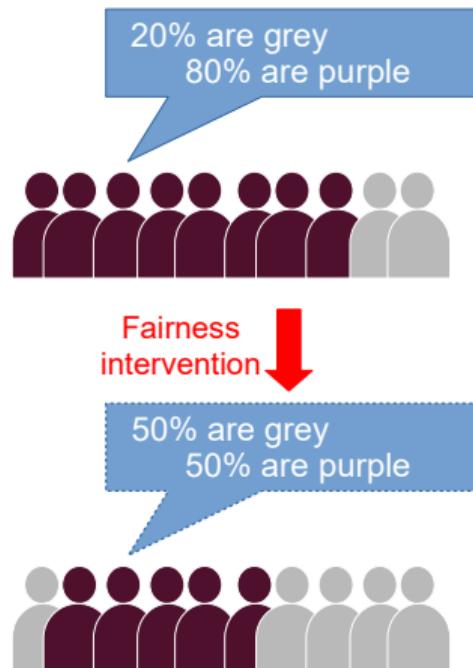
## 2. Define fairness objectives

Average predictions compared across groups:

- **Statistical parity** subjects in both subgroups have equal probabilities of being assigned to the positive predictive class
- **Conditional statistical parity** extends statistical parity by allowing additional legitimate attributes to affect the outcomes

Verma & Rubin. Fairness definitions explained, FairWare, 2018.

### 3. Choose fairness intervention: requirements



- Choice of method is constrained by availability of sensitive features
- Satisfying fairness objective is not enough
- What is the cost of the intervention?

#### 4. Monitor fairness intervention: requirements

- Monitoring fairness interventions as part of the product pipeline might not align with business objectives
- Product pipelines often require simplicity and maintainability
- Added complexity: analyzing customer behavior over time to analyze effect of intervention (opportunity)

## Practical challenges

In practice, there are many limitations to testing and correcting for fairness in this case

- Sensitive attributes are unknown and cannot be integrated in the model
- No ground truth available
- Potential fairness intervention impact cannot be monitored
- Stakeholders goals vs. fairness goals are misaligned

Evaluating & mitigating algorithmic bias requires navigating uncertainty.

## Going further

- A curated list of fairness in AI resources.
- Challenges of incorporating algorithmic ‘fairness’ into practice: tutorial at FAccT 2019, Video.
- A detailed overview of statistical definitions in available from this paper.
- Open source toolkit to examine, report, and mitigate algorithmic bias: AI Fairness 360.

## 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.

# Thank you

<https://hindantation.github.io>