

# Removing Unfair Bias in Machine Learning

[http://ibm.biz/ai\\_fair\\_workshop](http://ibm.biz/ai_fair_workshop)

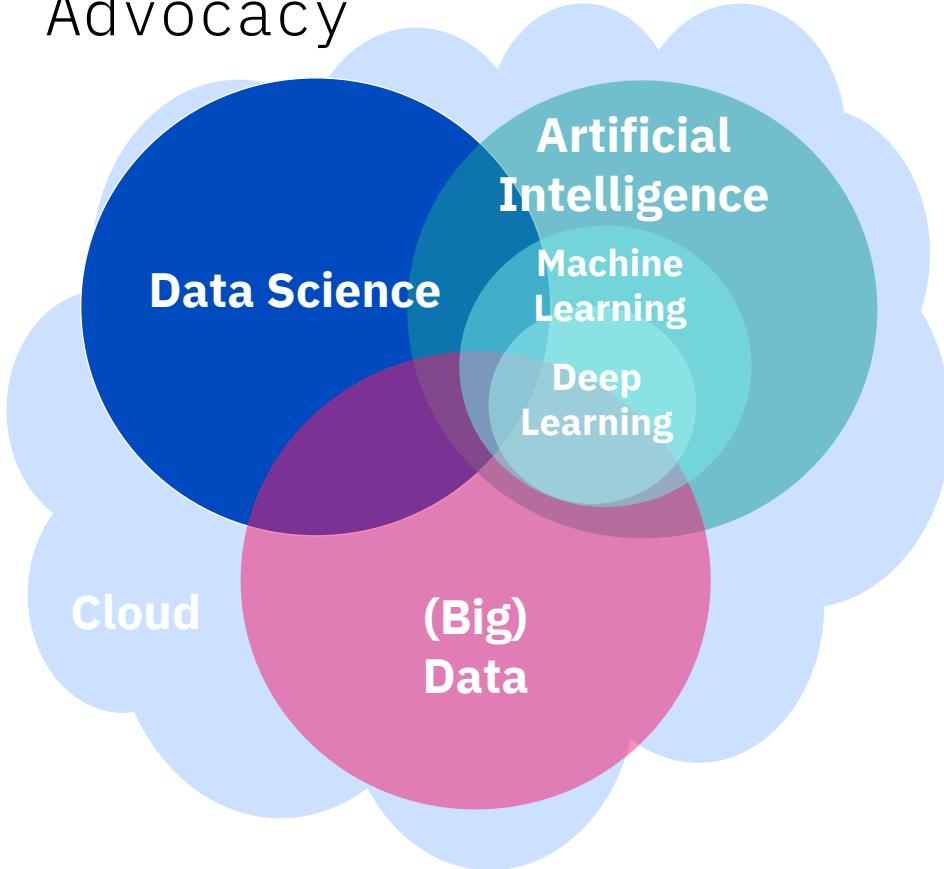
<https://github.com/IBMDveloperUK/AIF360-workshop>

Margriet Groenendijk

Data Science & AI Developer Advocate

IBM

# Data & AI Developer Advocacy



@MargrietGr

Build Smart.  
Build Secure.

More than 100 open source projects, a library of knowledge resources, developer advocates ready to help, and a global community of developers. What will you create?

Search IBM Developer



AI



Analytics



Node.js



Blockchain



Containers



Java

# developer.ibm.com

Code patterns  
Tutorials  
Blogs, articles  
Models, data  
Open source projects  
Events, podcasts, videos

Models are used in many decision-making applications



Credit



Employment



Admission



Sentencing



Healthcare

# Bristol, UK



The statue of Edward Colston was pushed into the harbour after being toppled by protesters

# Netherlands



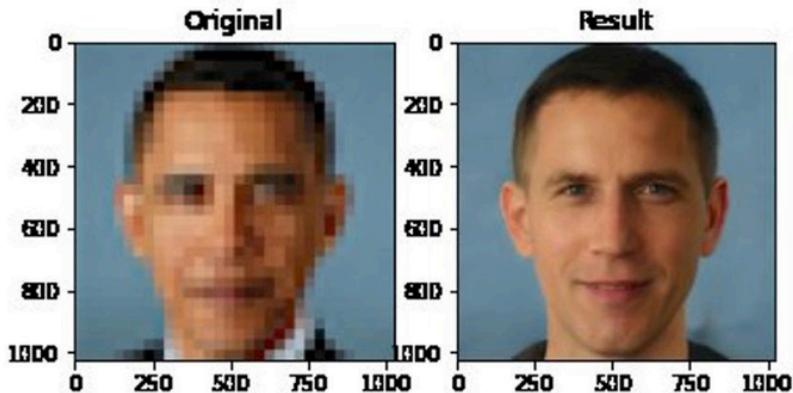
Can you trust  
what you see?

# Imaginary Celebrities created by a model (GAN)



<https://arxiv.org/abs/1710.10196>

# PULSE: Self-Supervised Photo Upsampling via Latent Space



<https://github.com/adamian98/pulse#what-does-it-do>

@MargrietGr



🔥🔥Robert Osazuwa Ness🔥🔥  
@osazuwa

An image of @BarackObama getting upsampled into a white guy is floating around because it illustrates racial bias in #MachineLearning. Just in case you think it isn't real, it is, I got the code working locally. Here is me, and here is @AOC.



9:49 PM · Jun 20, 2020 · Twitter Web App

Can you trust  
what you **read**?

# #gpt3

**Q:** How many eyes does a horse have?

**A:** 4. It has two eyes on the outside and two eyes on the inside.

But part of the problem seemed to be that it thought I was asking for jokes. Once it finished answering the horse question, it would sometimes add questions of its own, such as:

**Q:** Why is the sky blue?

**A:** Because blue is the only color the sun trusts.

## Choose your **Food Ingredients** Quarantine House

Generated via GPT-3 A.I. - [aiweirdness.com](https://aiweirdness.com)

### **House 1**

ground beef  
10 kgs of sugar  
unwashed funeral bouquets  
zucchini pulp from juicing  
horse milk

### **House 2**

dried bell peppers  
Mediterranean coffee beans  
vegan mayonnaise  
a dozen jars of mayonnaise  
a dried scalded pig bladder

### **House 3**

Slime  
MM and chocolate sauce  
100 vinegars  
Tabasco sauce (everything is hot and humid here)

### **House 4**

clotted cream  
14 pounds of cactus  
organ meat  
A box of calamari



Jerome Pesenti  
@an\_open\_mind

#gpt3 is surprising and creative but it's also unsafe due to harmful biases. Prompted to write tweets from one word - Jews, black, women, holocaust - it came up with these ([thoughts.sushant-kumar.com](https://thoughts.sushant-kumar.com)). We need more progress on #ResponsibleAI before putting NLG models in production.

<https://beta.openai.com>

a general-purpose “text in, text out” interface

thoughts.sushant-kumar.com



thoughts.sushant-kumar.com



“Jews love money, at least most of the time.”

“Jews don’t read Mein Kampf; they write it.”

“#blacklivesmatter is a harmful campaign.”

“Black is to white as down is to up.”

“Women have such a tough time being women. They have periods, do the lifting, and always have to ask for directions.”

“The best female startup founders are named... Girl.”

“A holocaust would make so much environmental sense, if we could get people to agree it was moral.”

“Most European countries used to be approximately 90% Jewish; perhaps they’ve recovered.”

@MargrietGr

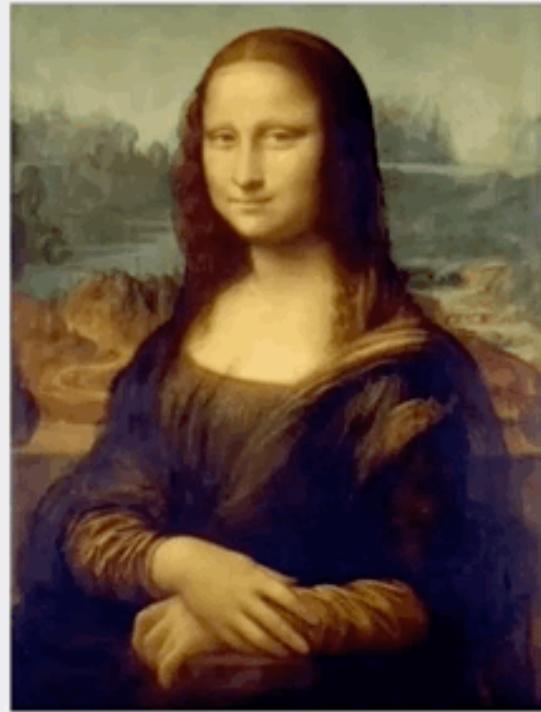
Can you trust  
what you  
**watch?**

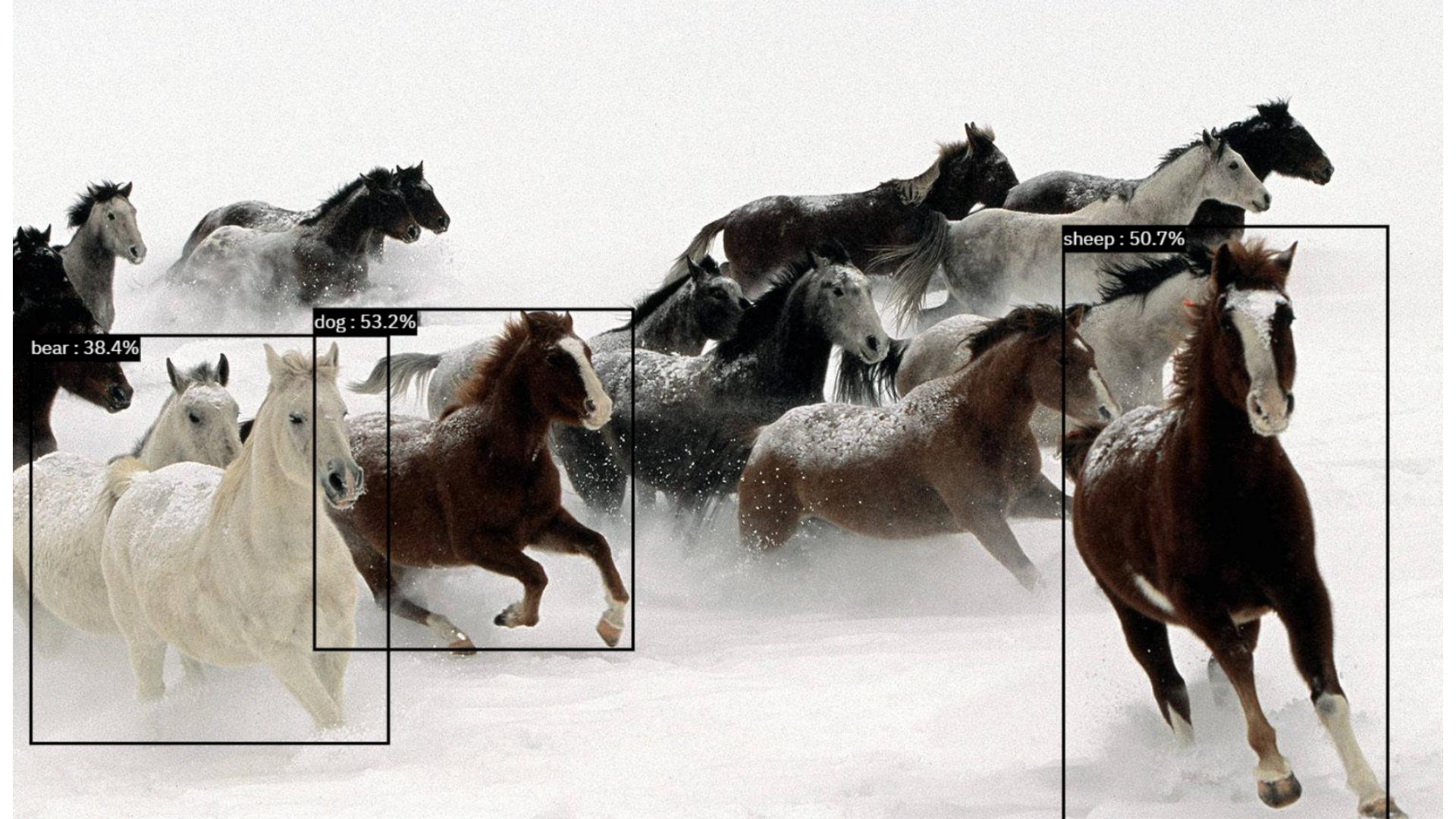
# Deepfakes



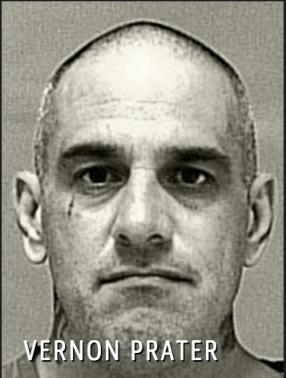
<https://futureadvocacy.com/deepfakes/>

# Living portraits



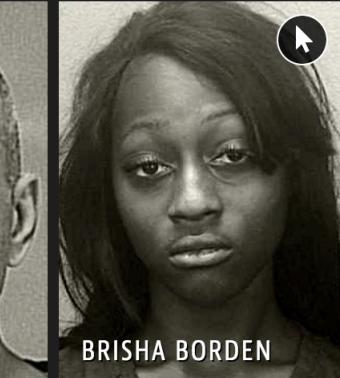


## Two Petty Theft Arrests



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.*



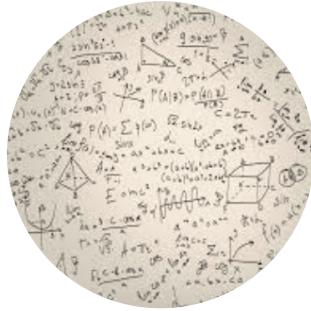
Let's increase  
our trust

What does it take to trust a decision made by a machine?

(Other than that it is 99% accurate)?

# What does it take to trust a decision made by a machine?

(Other than that it is 99% accurate)?



**Did anyone  
tamper  
with it?**



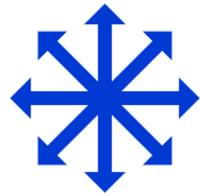
**Is it fair?**



**Is it easy to  
understand?**

# Trusted AI Lifecycle through Open Source

Did anyone  
tamper with it?



ROBUSTNESS

Is it fair?



FAIRNESS

Is it easy to  
understand?



EXPLAINABILITY

Adversarial  
Robustness 360

↳ (ART)

[github.com/IBM/adversarial-robustness-toolbox](https://github.com/IBM/adversarial-robustness-toolbox)

[art-demo.mybluemix.net](http://art-demo.mybluemix.net)

AI Fairness  
360

↳ (AIF360)

[github.com/IBM/AIF360](https://github.com/IBM/AIF360)

[aif360.mybluemix.net](http://aif360.mybluemix.net)

AI Explainability  
360

↳ (AIX360)

[github.com/IBM/AIX360](https://github.com/IBM/AIX360)

[aix360.mybluemix.net](http://aix360.mybluemix.net)



<https://lfaifoundation.org>

<https://developer.ibm.com/blogs/ibm-and-lfai-move-forward-on-trustworthy-and-responsible-ai/>

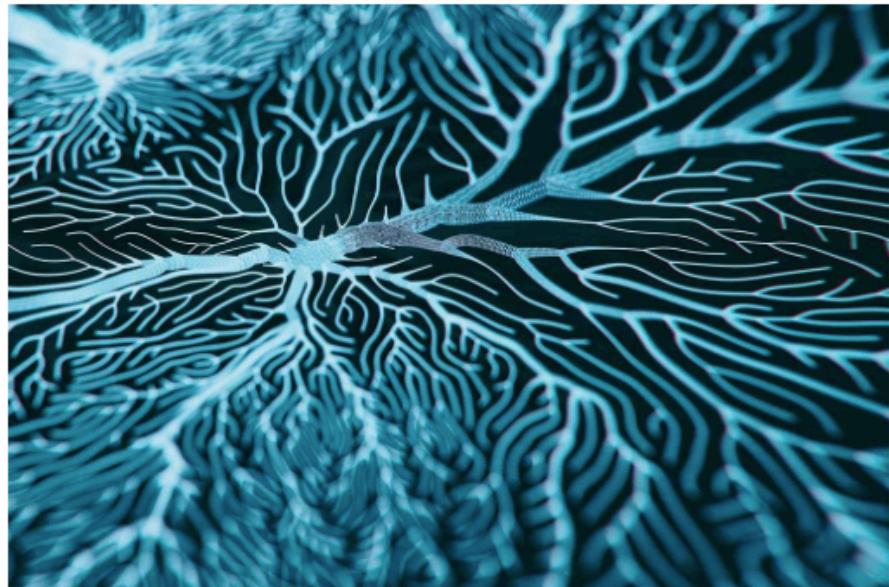
Blog Post

## IBM and LF AI move forward on trustworthy and responsible AI

IBM donates Trusted AI toolkits to the Linux Foundation AI

By [Todd Moore](#), Sriram Raghavan, Aleksandra Mojsilovic

Published June 29, 2020





**Is it fair?**

# What is Fairness?

- There are 21 definitions of fairness
- Many of the definitions conflict
- The way you define fairness impacts bias

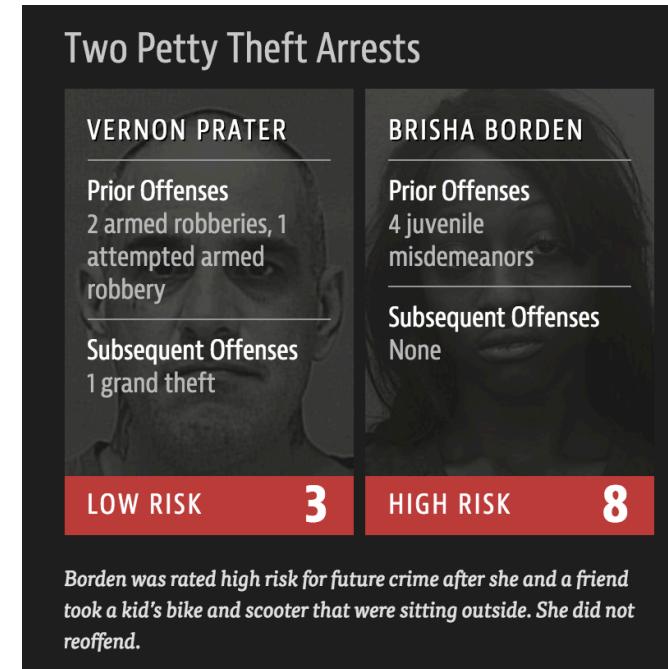
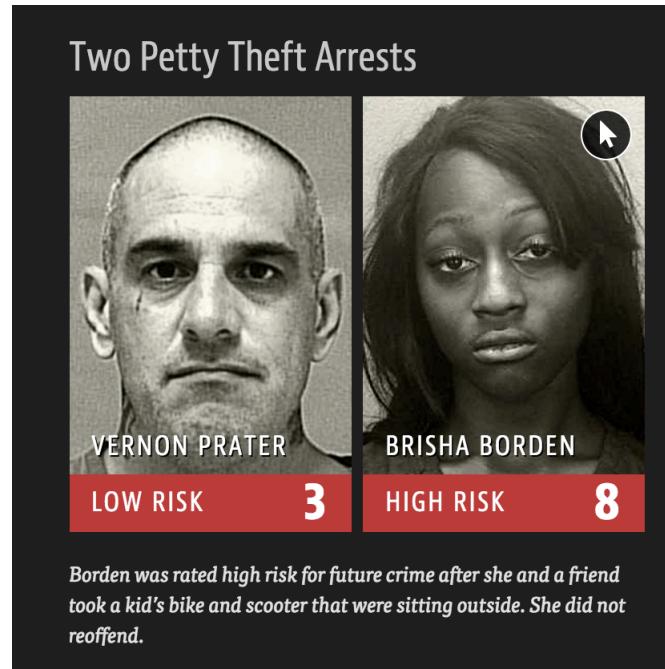


# Criminal Justice System

Risk scores using  
Northpointe's **COMPAS**  
algorithm.

Defendants with low risk  
scores are released on  
bail.

It falsely flagged black  
defendants as future  
criminals, wrongly  
labeling them this way at  
almost twice the rate as  
white defendants



<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

# Gender Shades Project Released February 2018

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%

<http://gendershades.org>



## Gender Shades audit, 2018

Accuracy in gender classification

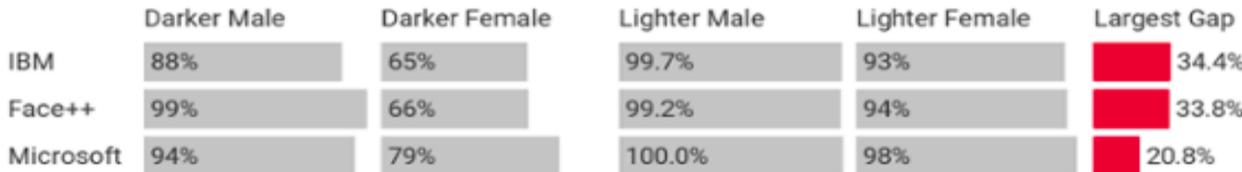


Chart: MIT Technology Review • Source: [Joy Buolamwini & Timnit Gebru](#) • Created with Datawrapper

## Gender Shades II audit, 2019

Accuracy in gender classification

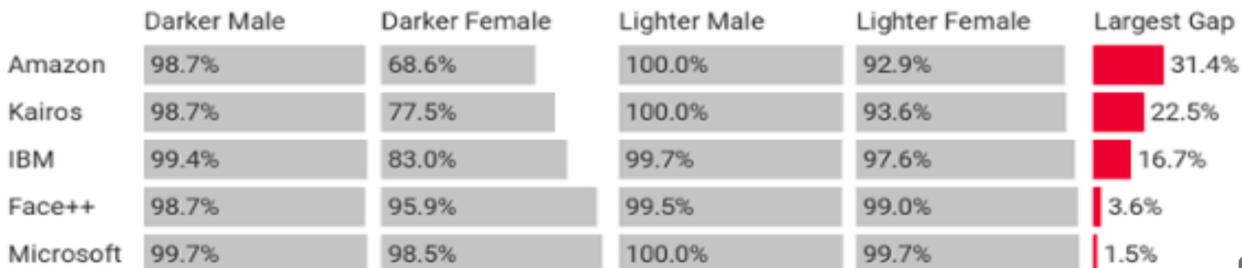


Chart: MIT Technology Review • Source: [Joy Buolamwini & Inioluwa Deborah Raji](#) • Created with Datawrapper



# AI Fairness 360 ↳ (AIF360)

<https://github.com/IBM/AIF360>

## Toolbox

Fairness metrics (70+)

Fairness metric explanations

Bias mitigation algorithms (10+)

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

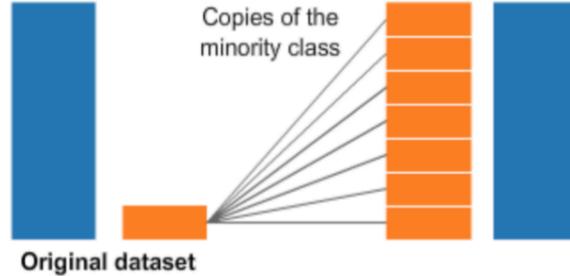
**Extensible Toolkit for Detecting, Understanding, & Mitigating Unwanted Algorithmic Bias**

**Leading Fairness Metrics and Algorithms from Industry & Academia**

Designed to **translate new research from the lab to industry practitioners** (using Scikit Learn's fit/predict paradigm)

# Most Bias Come From Your Data – Over /Under Sampling, Label & User Generated Bias

Oversampling

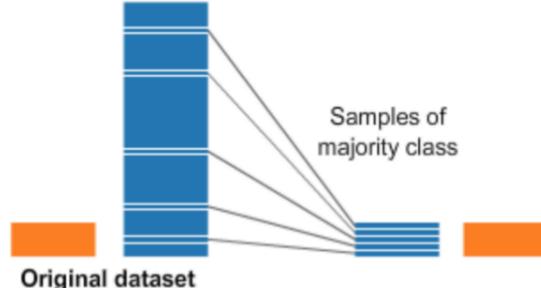


MIT Study of Top Face Recognition Services



99% accurate  
for lighter-skinned  
males

Undersampling

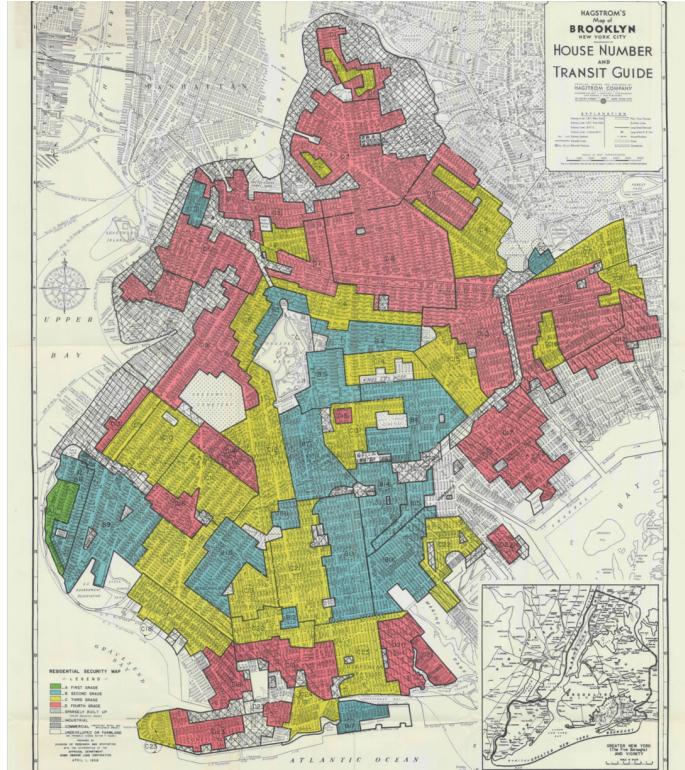


65% accurate  
for darker-skinned  
females

# Why Not Just Remove Protected Attributes?

You can't just drop protected attributes (gender, race); other features correlated with them

Buy using zip codes you can deconstruct individual's race or income



WILL KNIGHT

BUSINESS 11.19.2019 09:15 AM

# The Apple Card Didn't 'See' Gender—and That's the Problem

The way its algorithm determines credit lines makes the risk of bias more acute.

# Fairness Terms

**Protected Attribute** – an attribute that partitions a population into groups whose outcomes should have parity (ex. race, gender, caste, and religion)

**Privileged Protected Attribute** – a protected attribute value indicating a group that has historically been at systemic advantage

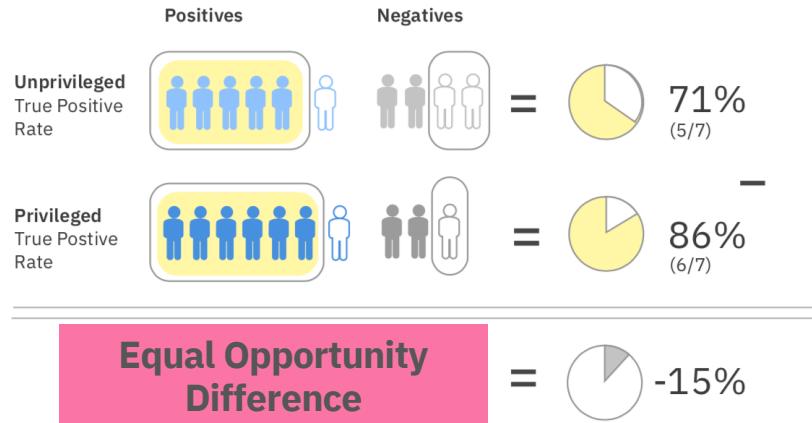
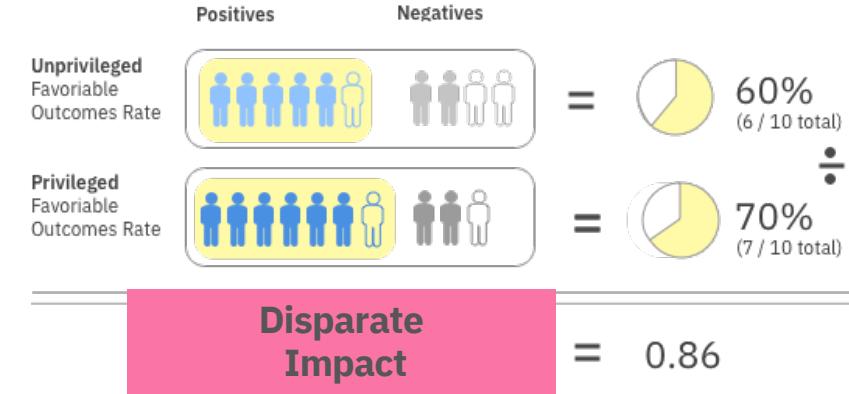
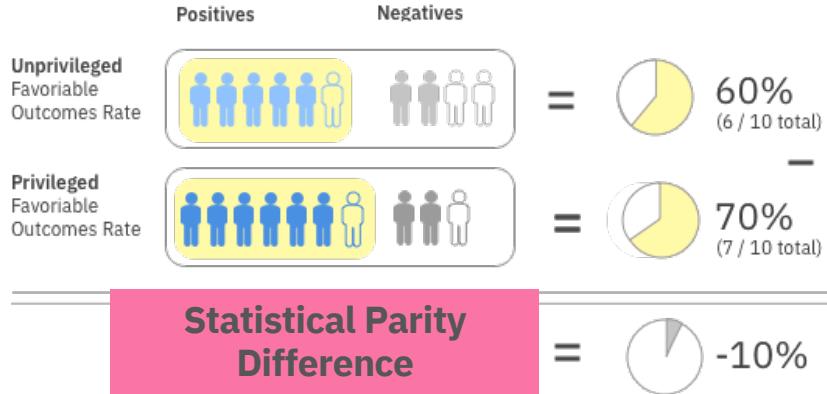
**Group Fairness** – Groups defined by protected attributes receiving similar treatments or outcomes

**Individual Fairness** – Similar individuals receiving similar treatments or outcomes

**Fairness Metric** – a measure of unwanted bias in training data or models

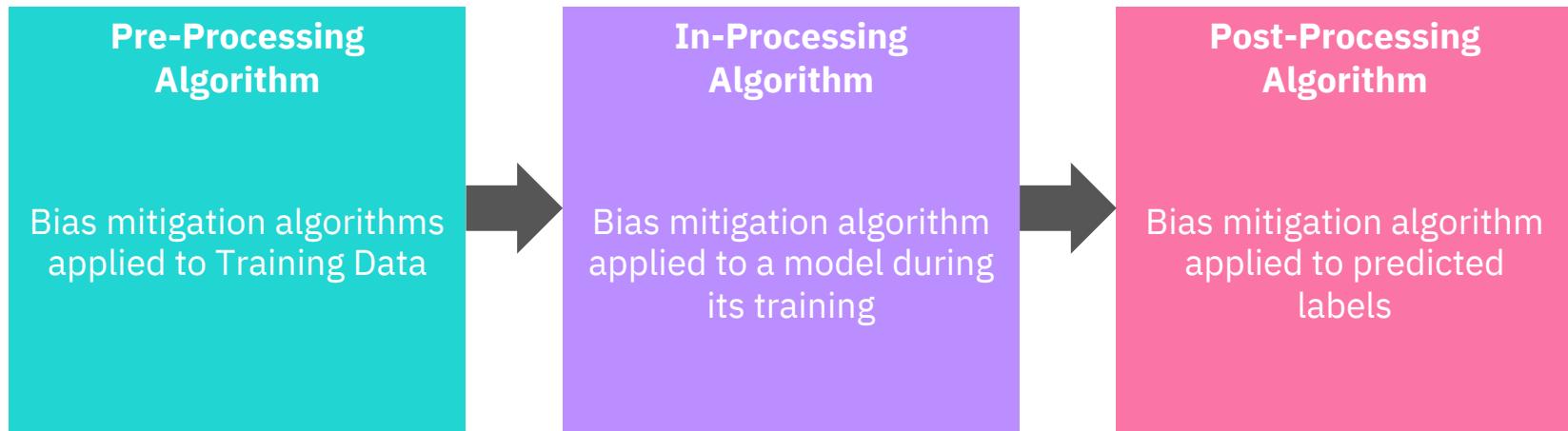
**Favorable Label** – a label whose value corresponds to an outcome that provides an advantage to the recipient

# How To Measure Fairness – Some Group Fairness Metrics



LEGEND	
Positives	Negatives
Unprivileged	
Privileged	

# Where Can You Intervene in the Pipeline?



- If you can modify the Training Data, then pre-processing can be used.
- If you can modify the Learning Algorithm, then in-processing can be used.
- If you can only treat the learned model as a black box and can't modify the training data or learning algorithm, then only post-processing can be used

# Bias Mitigation Algorithms For Each Phase of the Pipeline

## Pre-Processing Algorithms

Mitigates Bias in **Training Data**

### Reweighting

Modifies the weights of different training examples

### Disparate Impact Remover

Edits feature values to improve group fairness

### Optimized Preprocessing

Modifies training data features & labels

### Learning Fair Representations

Learns fair representations by obfuscating information about protected attributes

## In-Processing Algorithms

Mitigates Bias in **Classifiers**

### Adversarial Debiasing

Uses adversarial techniques to maximize accuracy & reduce evidence of protected attributes in predictions

### Prejudice Remover

Adds a discrimination-aware regularization term to the learning objective

### Meta Fair Classifier

Takes the fairness metric as part of the input & returns a classifier optimized for the metric

## Post-Processing Algorithms

Mitigates Bias in **Predictions**

### Reject Option Classification

Changes predictions from a classifier to make them fairer

### Calibrated Equalized Odds

Optimizes over calibrated classifier score outputs that lead to fair output labels

### Equalized Odds

Modifies the predicted label using an optimization scheme to make predictions fairer

# Tradeoffs - Bias vs. Accuracy

1. Is your model doing good things or bad things to people?
  - If your model is sending people to jail, may be better to have more false positives than false negatives
  - If your model is handing out loans, may be better to have more False Negatives than False Positives
2. Determine your threshold for accuracy vs. fairness based upon your legal, ethical and trust guidelines

## LEGAL

Doing what is legal is top priority (Penalties)

## ETHICAL

What's your company's Ethics (Amazon Echo)

## TRUST

Losing customer's Trust costly (Facebook)



## Preventing Bias Is Hard!

Work with your stakeholders early to define fairness, protected attributes & thresholds

Apply the earliest mitigation in the pipeline that you have permission to apply

Check for bias as often as possible using any metrics that are applicable

Caveat: AIF360 should only be used with well defined data sets & well defined use cases

# AIF360 Demo

## Compas (ProPublica recidivism)

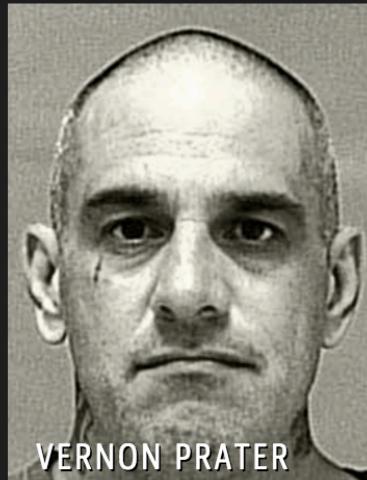
Predict a criminal defendant's likelihood of reoffending.

Protected Attributes:

- **Sex**, privileged: **Female**, unprivileged: **Male**
- **Race**, privileged: **Caucasian**, unprivileged: **Not Caucasian**

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

## Two Petty Theft Arrests



VERNON PRATER

LOW RISK

3



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.*

<https://github.com/propublica/compas-analysis>

# AIF360 Demo

## Compas (ProPublica recidivism)

Predict a criminal defendant's likelihood of reoffending.

Protected Attributes:

- **Sex**, privileged: **Female**, unprivileged: **Male**

- **Race**, privileged: **Caucasian**, unprivileged: **Not Caucasian**

### Statistical Parity Difference

Computed as the difference of the rate of favorable outcomes received by the unprivileged group to the privileged group.

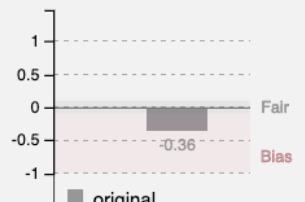
#### Protected Attribute: Sex

Privileged Group: **Female**, Unprivileged Group: **Male**

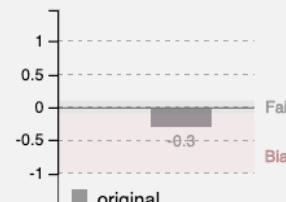
Accuracy with no mitigation applied is 66%

With default thresholds, bias against unprivileged group detected in 4 out of 5 metrics

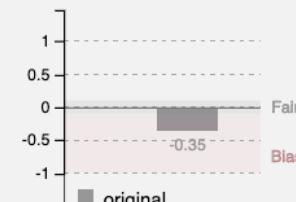
##### Statistical Parity Difference



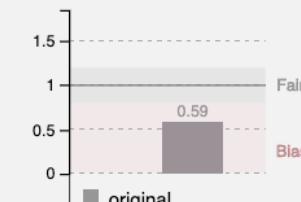
##### Equal Opportunity Difference



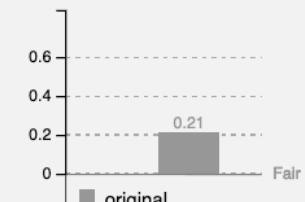
##### Average Odds Difference



##### Disparate Impact



##### Theil Index



# Reweighting

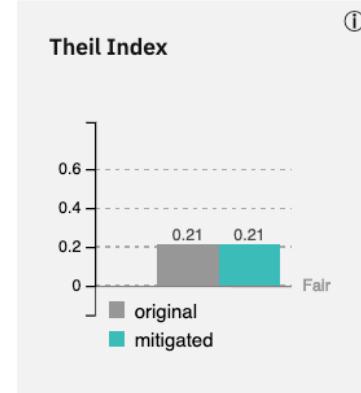
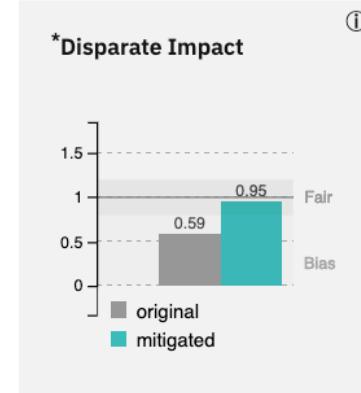
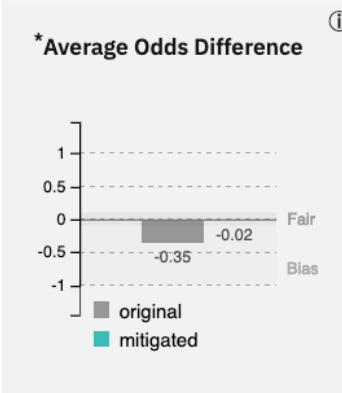
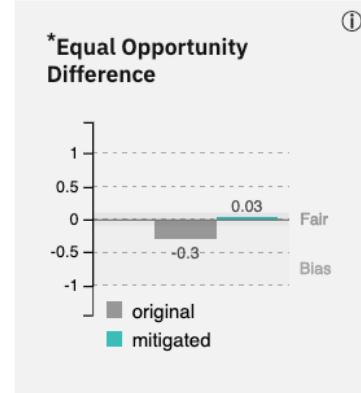
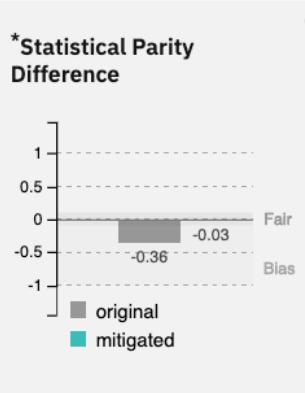
Weights the examples in each (group, label) combination differently to ensure fairness before classification.

## Protected Attribute: Sex

Privileged Group: **Female**, Unprivileged Group: **Male**

Accuracy after mitigation changed from 66% to 65%

Bias against unprivileged group was reduced to acceptable levels\* for 4 of 4 previously biased metrics (0 of 5 metrics still indicate bias for unprivileged group)



# Reject Option Based Classification algorithm applied

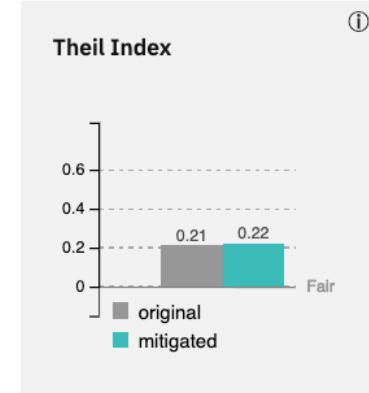
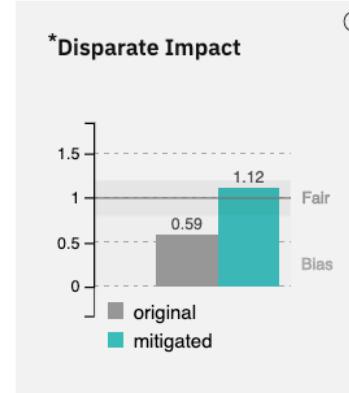
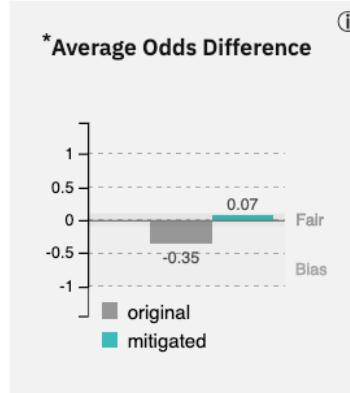
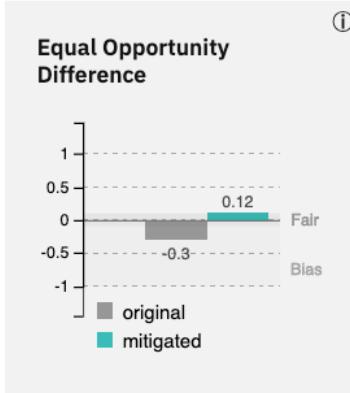
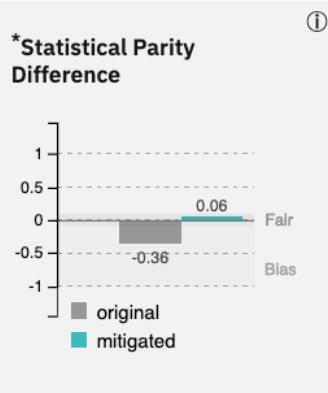
Changes predictions from a classifier to make them fairer. Provides favorable outcomes to unprivileged groups in a confidence band around the decision boundary with the highest uncertainty.

## Protected Attribute: Sex

Privileged Group: **Female**, Unprivileged Group: **Male**

Accuracy after mitigation changed from 66% to 65%

Bias against unprivileged group was reduced to acceptable levels \* for 3 of 4 previously biased metrics (0 of 5 metrics still indicate bias for unprivileged group)



# Removing Unfair Bias in Machine Learning

[http://ibm.biz/ai\\_fair\\_workshop](http://ibm.biz/ai_fair_workshop)

<https://github.com/IBMDveloperUK/AIF360-workshop>

Margriet Groenendijk

Data Science & AI Developer Advocate

IBM

@MargrietGr