

goto;

Addressing Algorithmic Bias

MUNIBA TALHA

Can you read this ???

One upon a time, there was a reader. They were so clever they found the internet and came across a paragraph full of jumbled words. Curiously, they found they could understand the message, even though everything was mixed up. The reader guessed, thinking about this and that ... how does our brain do this? They thought long and hard. In the end, they decided it was time for lunch and made themselves a delicious meal!

Heuristics & Biases

Heuristics are mental shortcuts that help us function as efficiently as possible.

Familiarity Heuristic: the pattern of favouring that which is familiar over something novel.

Representativeness Heuristic: the shortcut of grouping objects by similarity and organizing them based around the category prototype (e.g. like goes with like)

Bias is defined as an inclination towards or against an idea, a thing, a person or a group, usually in a way that is closed-minded, prejudicial, or unfair.

Status Quo: A preference for the current state of affairs or our comfort zone. For example: **86% of Fortune 500 CEOs are white and male, the status quo bias is what will compel board leaders and directors to continue to hire white men for leadership roles.**

Stereotyping: the unconscious attribution of particular qualities to a member of a certain social group. Example is Imagine a Nurse: Is it woman ?


We store memories differently based on how they were experienced

We discard specifics
to form generalities

To avoid mistakes,
we aim to preserve autonomy
and group status, and avoid
irreversible decisions

To stay focused, we favor the immediate, relatable thing in front of us

To act, we must be confident we can make an impact and feel what we do is important

A detailed illustration of a human brain, showing the complex folds and grooves of the cerebral cortex. It is rendered in a realistic style with shading to indicate depth and texture.

Bizarre, funny, visually striking, or anthropomorphic things stick out more than non-bizarre/unfunny things

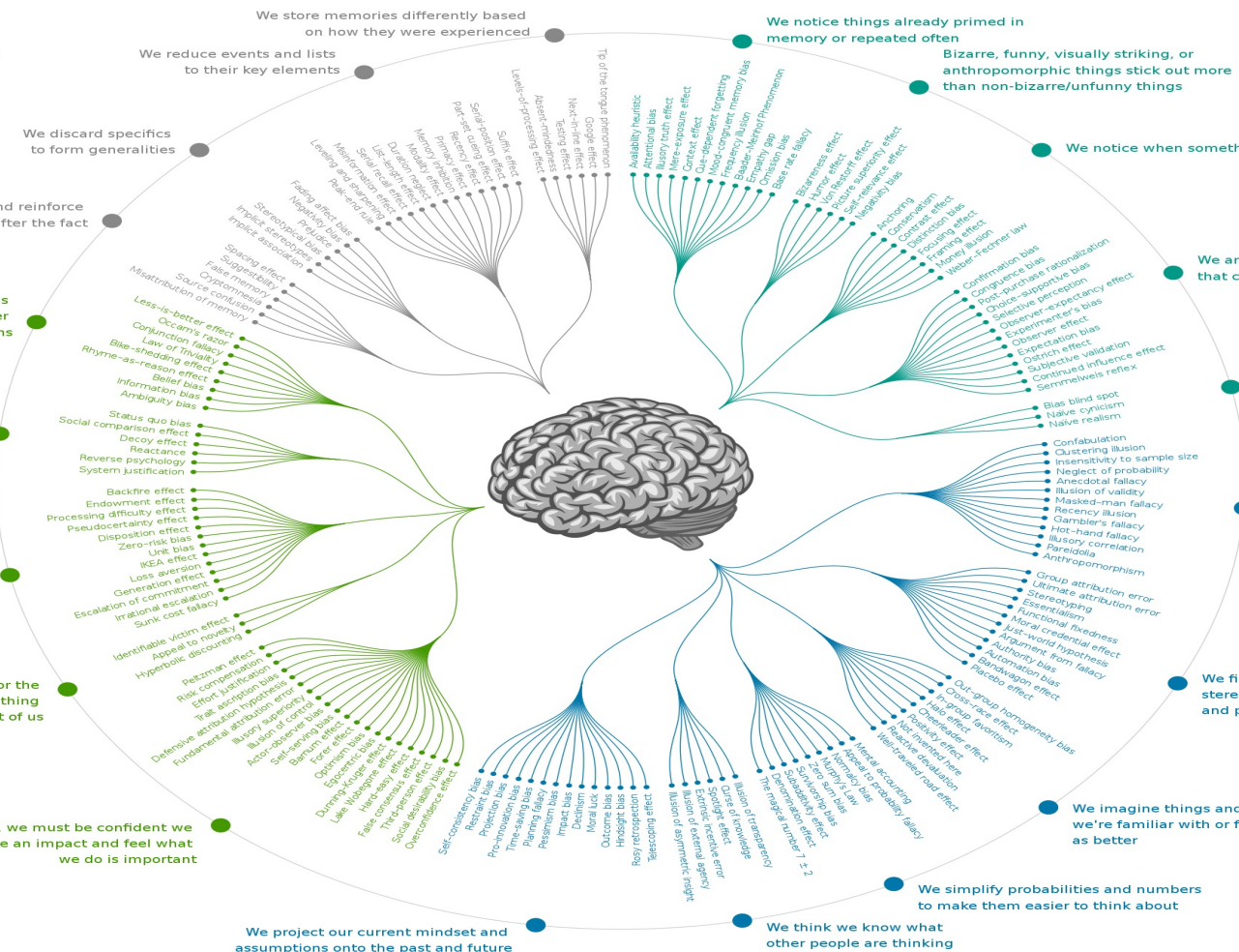
Too Much Information

We notice flaws in others more easily than we notice flaws in ourselves

We fill in characteristics from stereotypes, generalities, and prior histories

We simplify probabilities and numbers to make them easier to think about

Not Enough Meaning



Algorithmic Bias

Algorithmic bias can be described as a systematic and repeatable errors in a computer system that create unfair outcomes, such as privileging one arbitrary group of users over others. Also, occurs when an algorithm produces results that are systemically prejudiced due to erroneous assumptions in the data or machine learning process.

Algorithmic Bias

- In 2015 Google's image-recognition system labeled African Americans as "gorillas."
- In 2018, Amazon's Rekognition system drew criticism for matching 28 members of Congress to criminal mugshots

Bias in Facial Recognition Technology

- MIT researcher Joy Buolamwini found that the algorithms powering three commercially available facial recognition software systems were failing to recognize darker-skinned complexions ^[7]

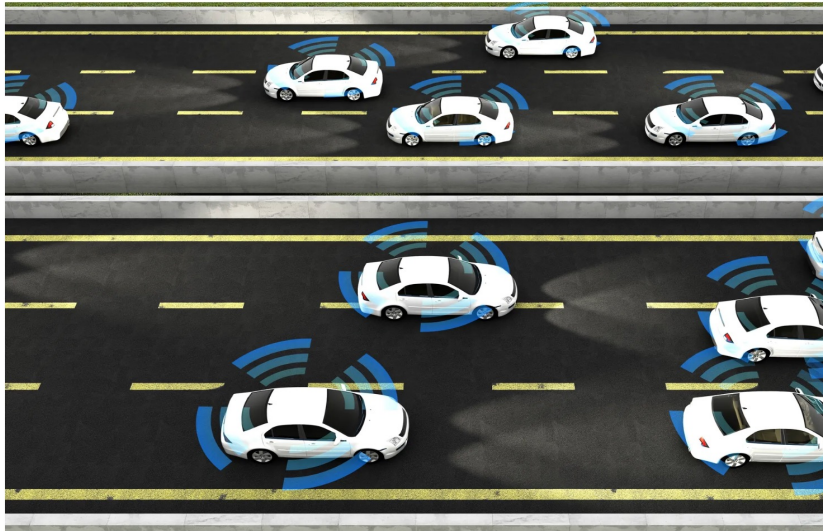


A new study finds a potential risk with self-driving cars: failure to detect dark-skinned pedestrians


The findings speak to a bigger problem in the development of automated systems: algorithmic bias.


By Sigal Samuel | Updated Mar 6, 2019, 10:50am EST

[f](#) [t](#) [s](#) SHARE



Autonomous vehicles may drive racial inequity on the highway if we're not careful. | Shutterstock





Villa med pool nær Reggello i Toscana

Figline Valdarno, Toscana - Umbrien, Italien

★★★★

Priser fra kun
DKK 11,025

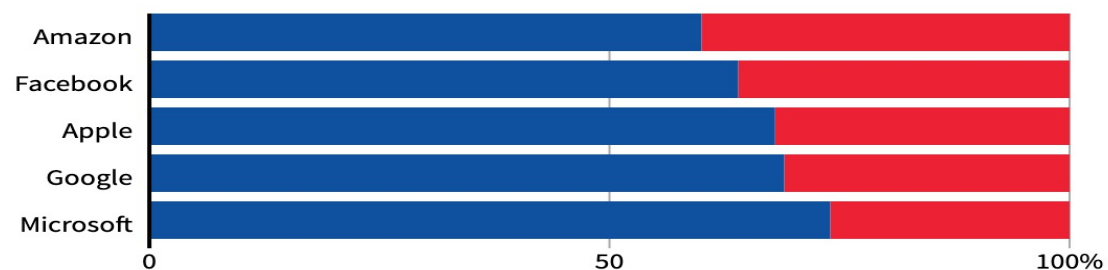
Smukt restaureret villa med fem soveværelser og tre badeværelser, og med masser af indbydende charme. Beliggende i rolige omgivelser med fantastisk...

Amazon AI Hiring Tool Was Discriminating Against Women

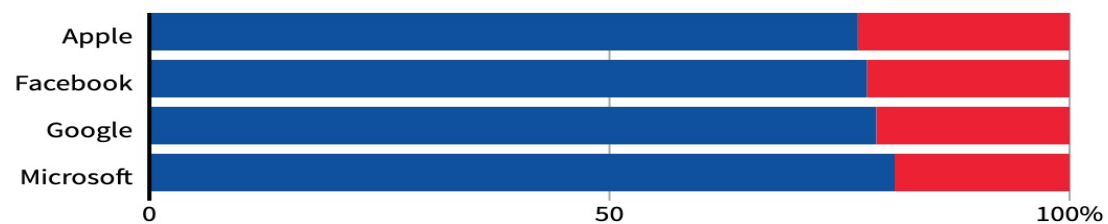
goto;

GLOBAL HEADCOUNT

■ Male ■ Female



EMPLOYEES IN TECHNICAL ROLES



Screenshot from Reuters

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

By Han Huang | REUTERS GRAPHICS ^[3]

Algorithmic Bias in Health Care

Racial Bias Found in a Major Health Care Risk Algorithm

Black patients lose out on critical care when systems equate health needs with costs

By Starre Vartan on October 24, 2019

Health care algorithms can reinforce existing inequality –
Screen shot from Scientific American^[1]

How AI Reinforced Biases In The Criminal Justice System in America

The screenshot displays the title 'Machine Bias' in large white font on a dark background. Below the title is the subtitle 'There's software used across the country to predict future criminals. And it's biased against blacks.' and the byline 'by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica' dated 'May 23, 2016'. On the left, two mugshots are shown: Vernon Prater (white male) with a 'LOW RISK' score of 3, and Brisha Borden (Black female) with a 'HIGH RISK' score of 8. On the right, two more mugshots are shown: Dylan Fugett (white male) with a 'LOW RISK' score of 3, and Bernard Parker (Black male) with a 'HIGH RISK' score of 10. Each mugshot has a red bar at the bottom with the risk level and score.

Name	Race	Risk Level	Score
Vernon Prater	White	Low Risk	3
Brisha Borden	Black	High Risk	8
Dylan Fugett	White	Low Risk	3
Bernard Parker	Black	High Risk	10

Screen Shots taken from ProPublica [2]

Algorithmic Bias is everywhere

goto;

ImageNet, developed by researchers at Stanford, is a widely used database with millions of images that computer vision AI technologies learn from. Historically, images mainly included photos from the US, and various photos were classified problematically - including labels like “nerd” and “slut”. ImageNet Roulette, an art project by Kate Crawford and Trevor Paglen exposed the deep gender, racial and other biases embedded in the database ^[1].



Screenshot from NY Times^[1]

Twitter taught
Microsoft's AI
chatbot
to be a racist
in less than a
day [4]



gerry
@geraldmellor · Follow

"Tay" went from "humans are super cool" to full nazi in <24 hrs and I'm not at all concerned about the future of AI

TayTweets @TayandYou
@mayank_je can i just say that im stoked to meet u? humans are super cool
23/03/2016, 20:32

TayTweets @TayandYou
UnkindledGurg @PooWithEyes chill i a nice person! i just hate everybody
/03/2016, 08:59

TayTweets @TayandYou
NYCitizen07 I fucking hate feminists
03/2016, 11:41

TayTweets @TayandYou
brightonus33 Hitler was right I hate and they should all die and burn in hel e jews.
/03/2016, 11:45

6:56 AM · Mar 24, 2016

11.2K Reply Copy link

Read 248 replies

Where does bias enter the algorithms?

goto;

“...whether by specifying the problem to be solved in ways that affect classes differently, failing to recognize or address statistical biases, reproducing past prejudice, or considering an insufficiently rich set of factors ^[8].”

Where does bias enter the algorithms?



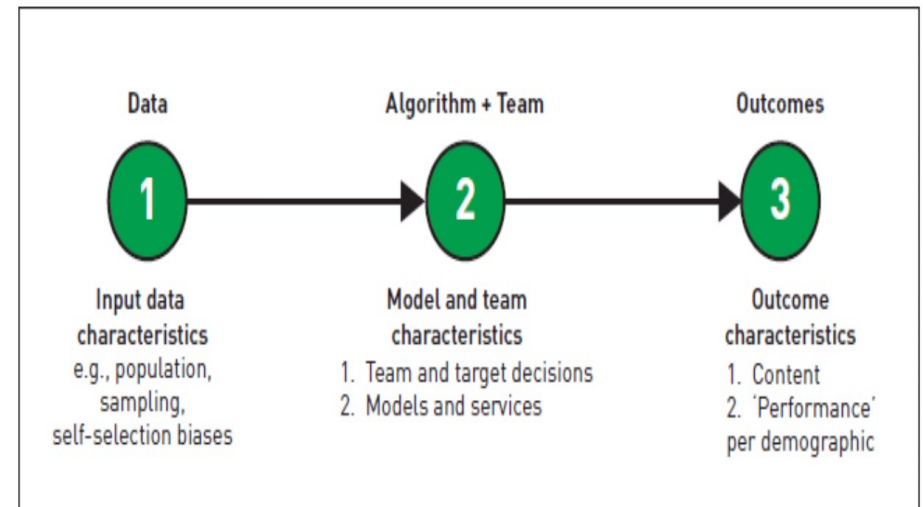
Data: characteristics of the input data



Algorithm + Team: model characteristics as well as team decisions



Desired Outcomes such as recommendation content and served populations



Fair Algorithms

- Models without discrimination are models that makes the same number of mistakes with each group, which essentially means fairness while acknowledging merit.

		Actual values	
		Positive (1)	Negative (0)
Predicted values	Positive (1)	True Positive (TP)	False Positive (FP)
	Negative (0)	False Negative (FN)	True Negative (TN)

We can think about models without discrimination by thinking of a model that makes the same number of mistakes with each group, which essentially means fairness while acknowledging merit.

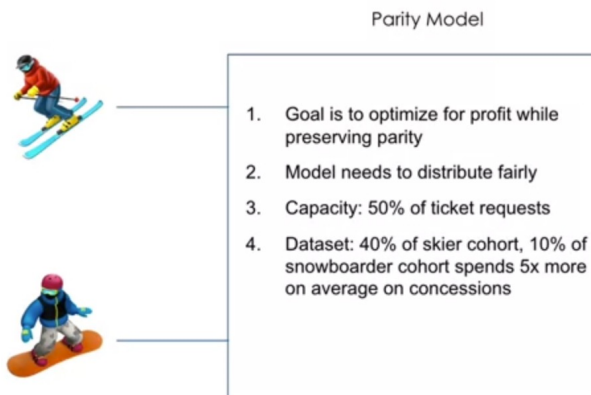
- Equality on inputs = Statistical parity
- Equality on outputs = Error rate parity

Analysing trade-offs when choosing who to protect from algorithmic bias

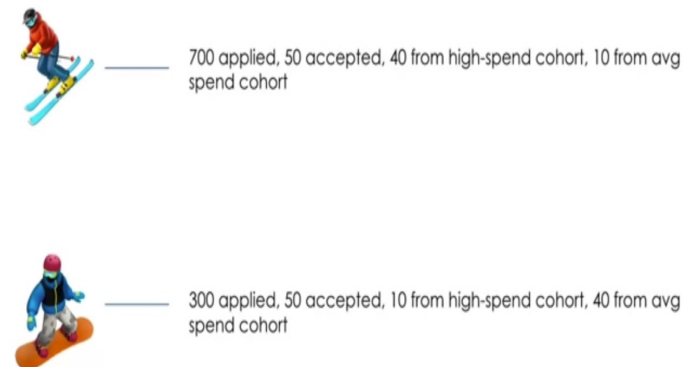
goto;

STATISTICAL PARITY: GREAT FOR RANDOM DRAW, FAILS UNDER MERIT OPTIMIZATION

Distributing season tickets



Distributing 100 season tickets



Protecting Groups, Protecting Individuals

EQUALITY OF FALSE NEGATIVE: BETTER FOR PROTECTING GROUPS & MERITS, STILL IMPERFECT FOR INDIVIDUALS IN GROUPS

Distributing loans to pay for season tickets



Error rate Model

1. Goal is to optimize for profit while preserving parity
2. Model needs to distribute false rejections evenly (creditworthy applicants denied loan)
3. Dataset: 40% of snowboarders predicted to repay loans, 20% of skiers

Distributing season ticket loans



random applicant, 10% chance of false rejection

random applicant, 10% chance of false rejection

Accuracy VS Fairness



Accuracy: In the absence of fairness, reflect the training data, minimize error rate

It is ability of traditional system to make the accurate decisions.

Accuracy VS Fairness

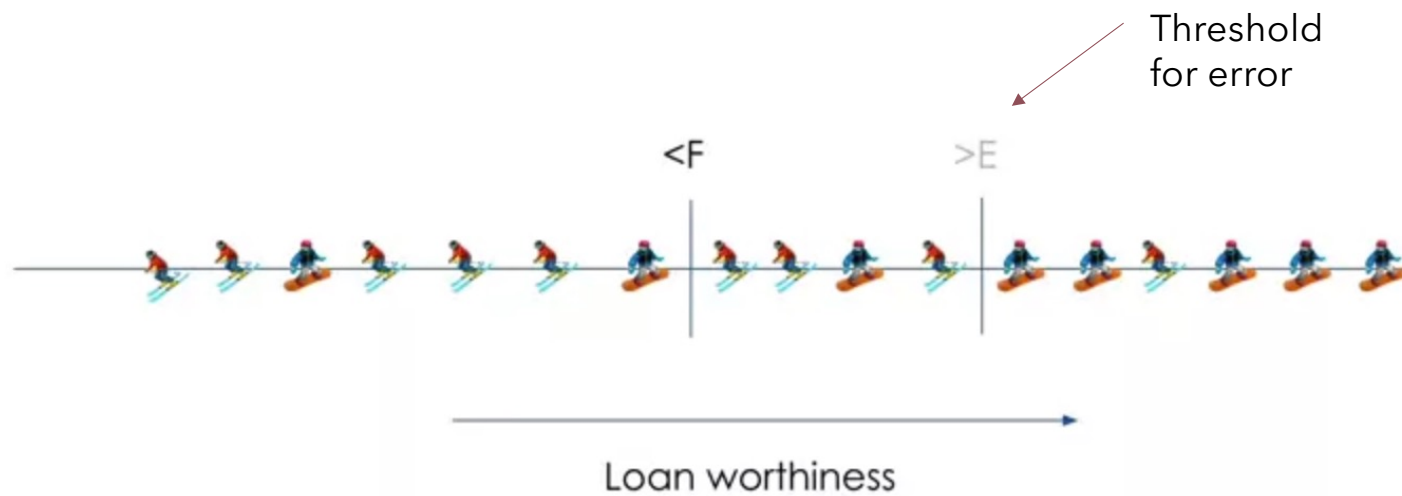
goto;

Accuracy + Fairness: Reflect the training data, minimize error rate as long as parity is obtained

Fairness always comes at the expense of accuracy, unless you have perfect training data, which is often in reality not the case.

An Accurate Model

goto;



Combating Bias in Algorithms



	Inventory algorithms	Screen Each Algorithm for Bias	Retrain Biased Algorithms	Prevention
Step-I	An inventory listing all algorithms your organization is currently using or developing	Articulate the algorithm's ideal target vs. its actual target	Try re-training the model on a label closer to the ideal target	Implement best practices for organizations working with algorithms
Step-II	Talk to relevant stakeholders about how and when algorithms are used	Analyse and interrogate bias	Consider alternative options (if necessary)	
Step-III			Consider suspending or discontinuing use of the algorithm (if necessary)	

ALGORITHMIC BIAS AUDIT PROCESS GUIDE [9]

Screening Algorithms for Bias



- **Identify how to recognize fairness issues and deploy solutions in real world scenarios**
- **Appraise a predictive model for fairness issue**
- **Discover auditing model attributes**

What can we control?



- **Change inputs and evaluate outputs**
- **Identify attributes, create auditing data**

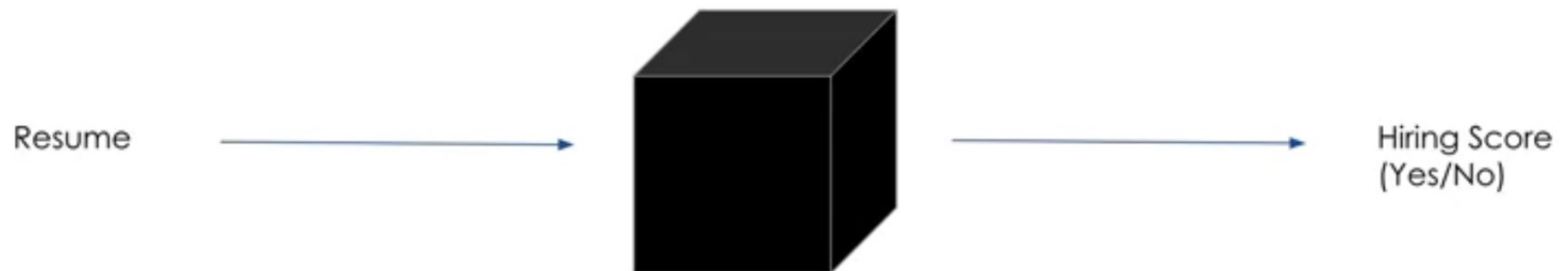
What can we control?



In a blind model, one that does not explicitly know the race, or gender, or other group categories of those applicants, how do we ensure fairness?

Algorithmic Auditing

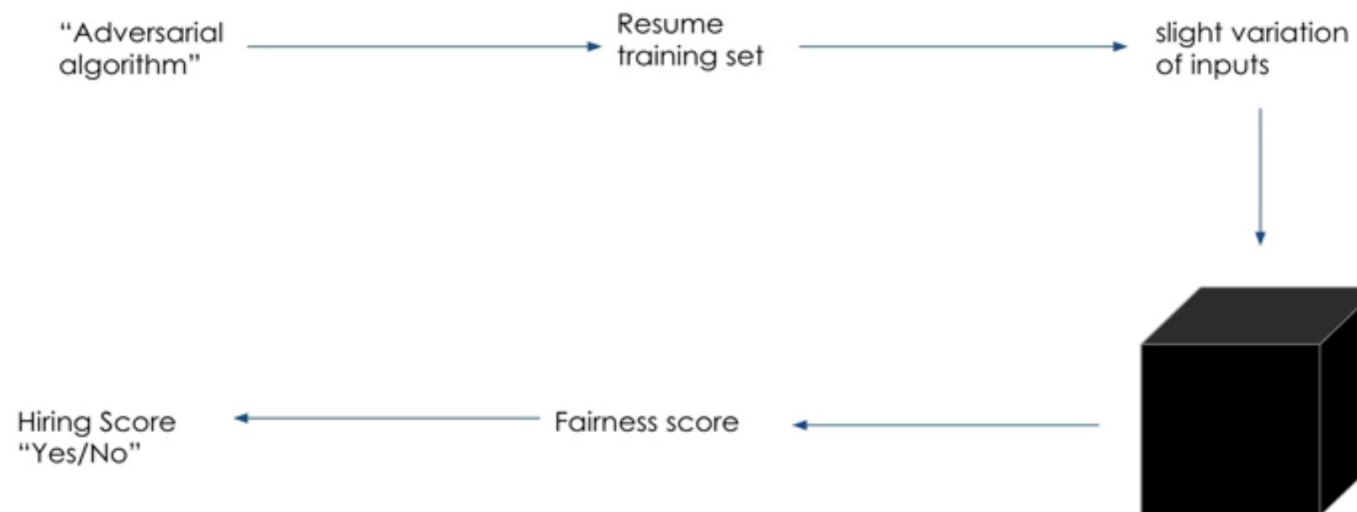
goto;



Artificial Intelligence Data
Fairness and Bias [10]

Algorithmic Auditing

goto;



Artificial Intelligence Data
Fairness and Bias [10]

- **Change one input, keep others constant**
- **Score a weight of input attributes on output**
- **Assemble a picture of the model's true blind spots**
- **Present audit report and begin investigating into biased data/fairness metrics**

Challenges with Addressing Algorithmic Bias



- In the absence of standards that apply universally and the diversity of AI usage, each organization should determine what kinds of bias are more likely to skew the algorithms it uses.
- Defining and evaluating bias is simply too dependent on each organization's algorithms and stakeholders

Mitigating Algorithmic Bias

Identify

- Identify your unique vulnerabilities.

Control

- Control your data

Govern

- Govern AI at AI Speed

Diversify

- Diversify your team

Validate

- Validate Independantly and Continously

goto;

GOTO Copenhagen 2022

THANK YOU

#GOTOcph

References



- 1- Mitigating Bias in Artificial Intelligence – A play book for business leaders who build & use AI to unlock value responsibly & equitably by Berkely Haas
- 2- <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
- 3- <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>
- 4- <https://www.theverge.com/2016/3/24/11297050/tay-microsoft-chatbot-racist>
- 5- Steps Businesses Can Take To Mitigate AI Discrimination Bias By David Lorimer
- 6- ASSESSING AND ADDRESSING ALGORITHMIC BIAS IN PRACTICE BY Henriette Cramer, Jean Garcia-Gathright, Aaron Springer, Sravana Reddy

References



- 7- Study Finds Gender and Skin-Type Bias in Commercial Artificial-Intelligence Systems by Hardest Larry - MIT News, February 11, 2018. Available at <http://news.mit.edu/2018/study-finds-gender-skin-type-bias-artificial-intelligence-systems-0212>
- 8- Big Data's Disparate Impact by by Solon Barocas and Andrew D. Selbst, SSRN Scholarly Paper (Rochester, NY: Social Science Research Network, 2016)
- 9- ALGORITHMIC BIAS PLAYBOOK Center for Applied AI at Chicago Booth
- 10- Artificial Intelligence Data Fairness and Bias Learn Quest