# Algorithmic Bias and Responsible AI

**Anjana Susarla**

**Omura-Saxena Professor of Responsible AI**
**Eli Broad College of Business**
**Michigan State University**
(**asusarla@msu.edu**)


**Twitter: @asusarla**


**Material with grateful acknowledgements from Ashley Casovan, Virginia Dignum and Kay Firth-Butterfield**

# Advances in AI = Advances in Prediction!

FIGURE 3-1

Image classification error over time

Machine classification error

Human benchmark

Error rate — 2010 2011 2012 2013 2014 2015 2016 2017
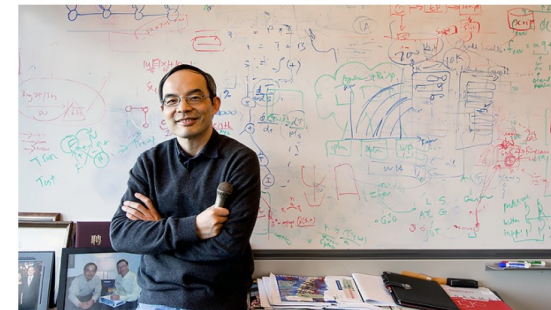
GOOGLE'S ARTIFICIAL BRAIN LEARNS TO FIND CAT VIDEOS

By Liat Clark, Wired UK

When computer scientists at Google's mysterious X lab built a neural network of 16,000 computer processors with one billion connections and let it browse YouTube, it did what many web users might do – it began to look for cats.

ANNALS OF MEDICINE   APRIL 3, 2017 ISSUE

A.I. VERSUS M.D.

What happens when diagnosis is automated?

By Siddhartha Mukherjee

Microsoft reaches a historic milestone, using AI to match human performance in translating news from Chinese to English
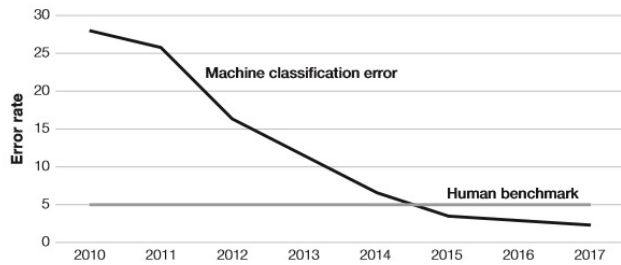
Mar 14, 2018  |  Allison Linn

*Xuedong Huang, technical fellow in charge of Microsoft's speech, natural language and machine translation efforts. (Photo by Scott Eklund/Red Box Pictures)*

A team of Microsoft researchers said Wednesday that they believe they have created the first machine translation system that can translate sentences of news articles from Chinese to English with the same quality and accuracy as a person.
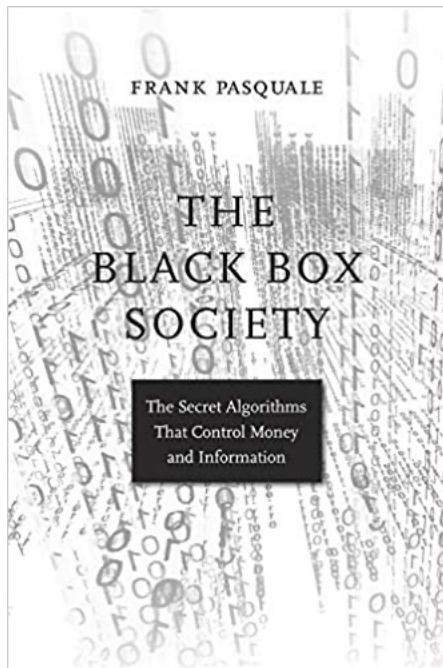
From Agarwal et al. (2018)

# Automated Decisions

# We had better be quite sure that the purpose put into the machine is the purpose which we really desire.

FRANK PASQUALE

THE BLACK BOX SOCIETY

The Secret Algorithms That Control Money and Information

Norbert Wiener (1960)

CODED BIAS

Input → Blackbox → Output

Stimulus                    Response

# The Black Box Society (Frank Pasquale, 2015)



In one study, black-identified names generated different ads than white-identified ones.

Chart courtesy Latanya Sweeney/Harvard University (http://arxiv.org/ftp/arxiv/papers/1301/1301.6822.pdf)

# AI Biases

# Biases in AI

— Unconscious biases

— Black box decisions

— Proxy discrimination

— Training data

**Ade Adamson, MD MPP**
@AdeAdamson

Google launches AI health tool for skin conditions in Europe on.ft.com/3tWNQtS The algorithm was developed based on training data with less than 4% dark skin types. It should come with a warning BEWARE OF RESULTS IF BLACK!!!

| | | |
|---|---|---|
| 46 (0.3%) | 9 (0.2%) | 0 (0.0%) |
| 2,807 (17.4%) | 383 (10.2%) | 104 (10.8% |
| 6,641 (41.2%) | 2,412 (64.2%) | 607 (63.0% |
| 5,040 (31.3%) | 724 (19.3%) | 195 (20.2% |
| 510 (3.2%) | 101 (2.7%) | 24 (2.5%) |
| 46 (0.3%) | 1 (0.0%) | 0 (0.0%) |
| 1,024 (10.2%) | 126 (3.4%) | 33 (3.4%) |

5:07 PM · May 18, 2021 · Buffer

**MOTHERBOARD**
TECH BY VICE

## Google's New Dermatology App Wasn't Designed for People With Darker Skin

The company trained the system to recognize different skin conditions. But like Google itself, the app's data has a diversity problem.

By Todd Feathers

May 20, 2021, 9:40am   Share   Tweet   Snap

# How to create a racist chatbot without trying

**What happens when you don't understand what your algorithm is learning?**

**We want these sentences to all give the same score but they don't.**

```
text_to_sentiment("Let's go get Italian food")
2.0429166109
text_to_sentiment("Let's go get Chinese food")
1.4094033658
text_to_sentiment("Let's go get Mexican food")
0.3880198556
```

**The algorithm is probably accurately learning real feelings of people based on the data but it's not learning what we intended it to learn.**

**We never told the algorithm that we didn't want to learn racism!**

**"My name is _____" is a neutral statement so the score should be about 0.**

```
text_to_sentiment("My name is Emily")
2.2286179365
text_to_sentiment("My name is Heather")
1.3976291151
text_to_sentiment("My name is Yvette")
0.9846380213
text_to_sentiment("My name is Shaniqua")
-0.4704813178
```

**We wanted to learn the sentiment score of the sentence "My name is _____" which should be independent of the particular name used Emily, Shaniqua, etc**

partially based on analysis by Robyn Speer and images by Mark Xiang)

@kareem_carr

# Example 1: Hiring

Amazon realized its hiring system was not rating candidates for software developer jobs and other technical posts in a gender-neutral way

That is because Amazon's computer models were trained to vet applicants by observing patterns in resumes submitted to the company over a 10-year period. Most came from men, a reflection of male dominance across the tech industry.

Amazon's system penalized resumes that included the word "women's," as in "women's chess club captain" and downgraded graduates of two all-women's colleges. And it privileged resumes with the kinds of verbs that men tend to use, like "executed" and "captured."

## MANAGEMENT SCIENCE

📄 View PDF

🔧 Tools     ⌁ Share

# Algorithmic Bias? An Empirical Study of Apparent Gender-Based Discrimination in the Display of STEM Career Ads

Anja Lambrecht [iD] , Catherine Tucker [iD]

**Go to Section**

Abstract

## Abstract

We explore data from a field test of how an algorithm delivered ads promoting job opportunities in the science, technology, engineering and math fields. This ad was explicitly intended to be gender neutral in its delivery. Empirically, however, fewer women saw the ad than men. This happened because younger women are a prized demographic and are more expensive to show ads to. An algorithm that simply optimizes cost-effectiveness in ad delivery will deliver ads that were intended to be gender neutral in an apparently discriminatory way, because of crowding out. We show that this empirical regularity extends to other major digital platforms.

f

t

✉

🅭

🎁

## MACHINE BIAS

# Facebook Lets Advertisers Exclude Users by Race

Facebook's system allows advertisers to exclude black, Hispanic, and other "ethnic affinities" from seeing ads.
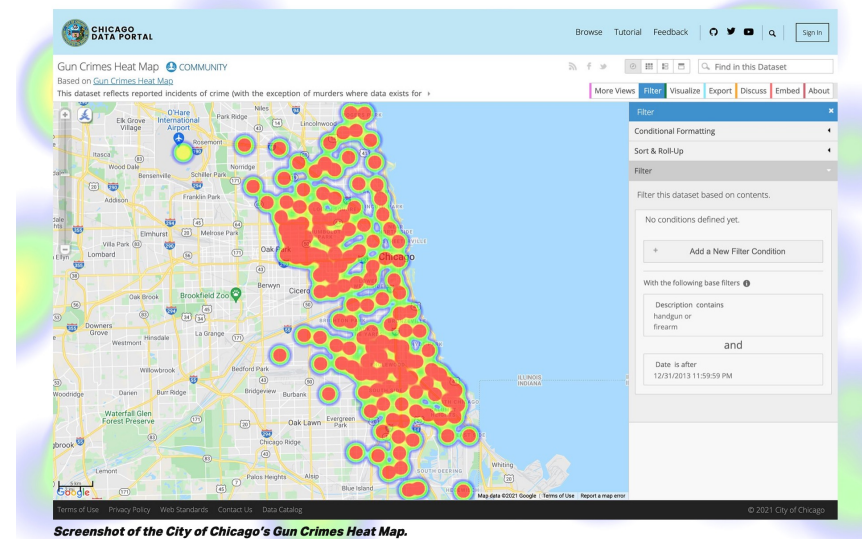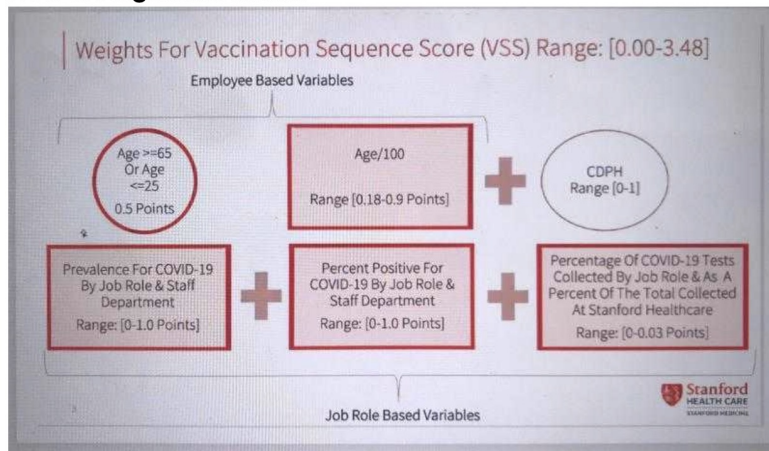
by Julia Angwin and Terry Parris Jr., Oct. 28, 2016, 1 p.m. EDT

# Example 2: Predictive Policing

Chicago PD told a man he would be involved in a shooting

- But could not identify on which side he would be

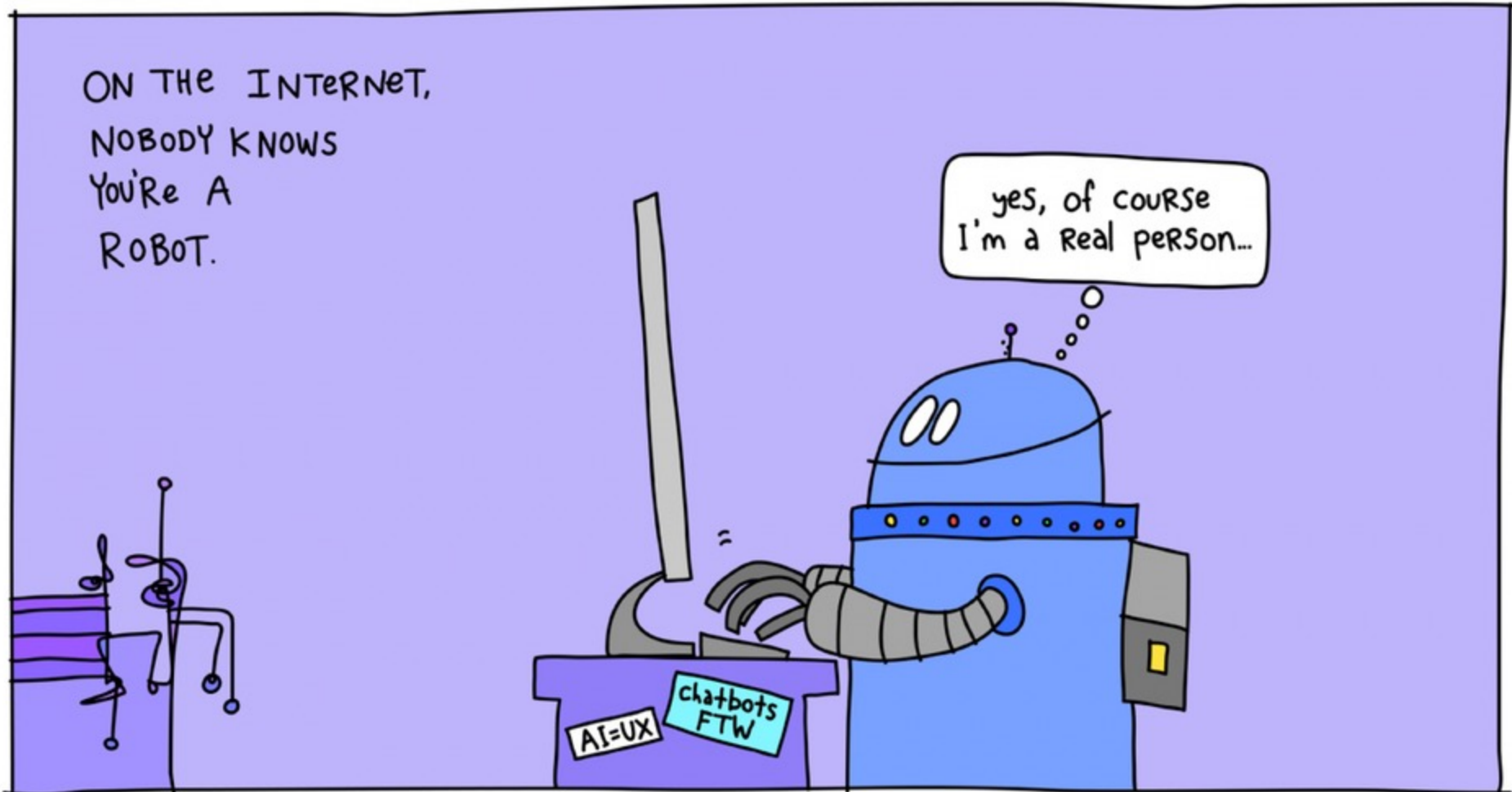- Resulted in continuous surveillance of certain individuals



*Screenshot of the City of Chicago's Gun Crimes Heat Map.*

**How the algorithm works**



Stanford used a so-called "algorithm", which was really a rules-based formula designed to determine the order in which the thousands of medical workers at Stanford should be vaccinated. The tool took into account employee-based variables like age, job-based variables, and public health guidance, according to MIT Technology Review.

But flaws in that calculation meant hospital administrators and other employees working from home were toward the front of the line, while only seven of Stanford's 1,300 medical residents made the list.

# Example 3: Healthcare

# Fairness, accountability & transparency

# AI AND ETHICS
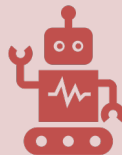
**Self-driving cars** — Who is responsible for the accident by self-driving car?

**Automated manufacturing** — How can workers practice new sophisticated skills so as not to lose their jobs?

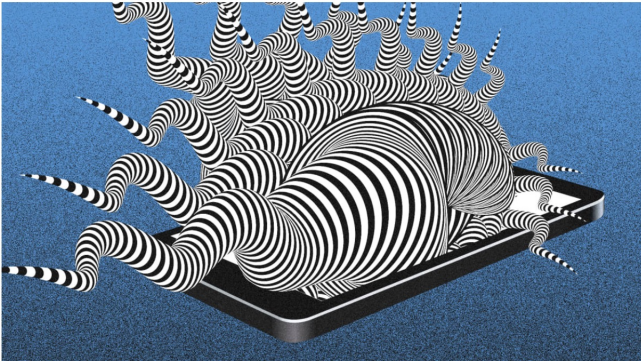**Chatbots** — Manipulation of emotions / nudging / behaviour change support

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.

DATA

ANSWERS

# My journey into Responsible AI

**Agenda for Responsible AI**



## How do we understand algorithmic bias?

## A feasible path..

- Develop frameworks for identifying risks
- Understanding (sources of) bias introduced in training data or machine learning models
- Algorithmic auditing and other guardrails

## How do we mitigate these biases and other unintended consequences?

# From the Responsible AI Institute

## Checklist for Fair Lending



**Fair Lending Risk at Every Stage of Credit Transaction**

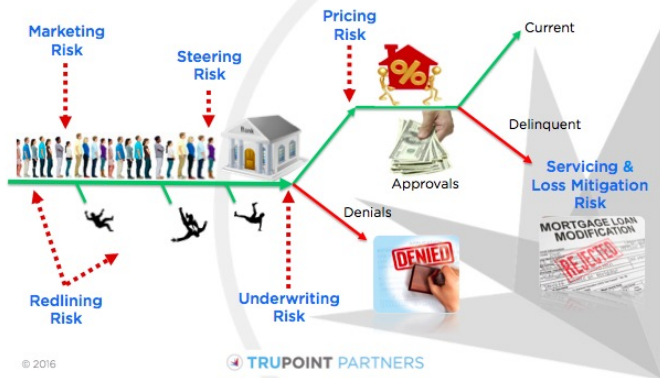Marketing Risk · Steering Risk · Pricing Risk · Current · Delinquent · Servicing & Loss Mitigation Risk · Approvals · Denials · Redlining Risk · Underwriting Risk · MORTGAGE LOAN MODIFICATION · DENIED · REJECTED · © 2016 · TRUPOINT PARTNERS

What are the key challenges in the use of AI systems in current and future lending practices?

What are the key demographics that are at a disadvantage with this work?

How can a certification program ensure that automated lending assisted or produced by an AI system is fair?

What are other measures such as accurate advertising, improved explainability and oversight of these systems that should be required as part of this certification program?

What are the barriers we need to address to accelerate sector-wide responsible AI adoption by businesses?

What type of AI systems are currently used to assist or augment AI practices?

What techniques can be used to automatically assess these systems?

What data is typically collected for lending? Is there other data that should or should not be collected to assess the responsible use of AI systems?

# Evaluating fairness-aware algorithms

Which algorithm is the best?

… on which dataset?

… how was it preprocessed?

… under which measure?

… with which training / test split?

… what are the right hyperparameter settings?

… what if there are multiple sensitive attributes?

**MOTHERBOARD**
TECH BY VICE

# An Insurance Startup Bragged It Uses AI to Detect Fraud. It Didn't Go Well

Lemonade backtracked after suggesting it uses "non-verbal cues" like eye movements to reject claims. Its response raises more questions than answers.

TF   By Todd Feathers

By Janus Rose
NEW YORK, US

May 26, 2021, 1:01pm    **Share**    **Tweet**    **Snap**

**Rachel Metz** ✔
@rachelmetz                                                         · · ·

this is not true, according to the @Lemonade_Inc inc. s-1, filed with the SEC, which says on page 128, "AI Jim handles the entire claim through resolution in approximately a third of cases, paying the claimant or declining the claim without human intervention".

> **Lemonade** ✔ @Lemonade_Inc · May 26
> We never let AI auto-decline claims (2/4)
> Show this thread

# Aiming for truth, fairness, and equity in your company's use of AI

By: Elisa Jillson | Apr 19, 2021 9:43AM

SHARE THIS PAGE

**TAGS:**  Bureau of Consumer Protection | Consumer Protection | Privacy and Security | Consumer Privacy |
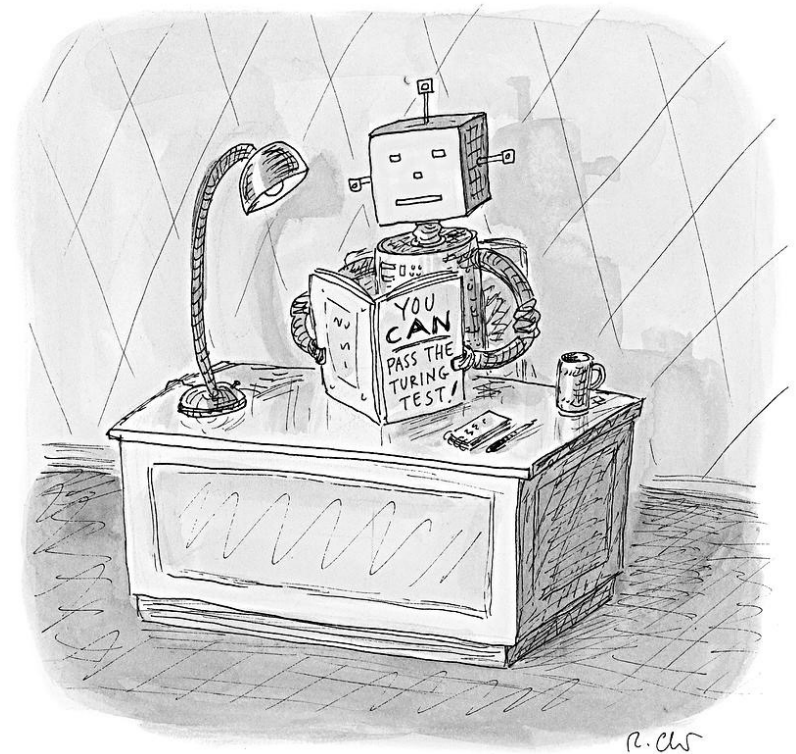Credit Reporting | Tech

Advances in artificial intelligence (AI) technology promise to revolutionize our approach to medicine, finance, business operations, media, and more. But research has highlighted how apparently "neutral" technology can produce troubling outcomes – including discrimination by race or other legally protected classes. For example, COVID-19 prediction models can help health systems combat the virus through efficient allocation of ICU beds, ventilators, and other resources. But as a recent study in the Journal of the American Medical Informatics Association suggests, if those models use data that reflect existing racial bias in healthcare delivery, AI that was meant to benefit all patients may worsen healthcare disparities for people of color.

The question, then, is how can we harness the benefits of AI without inadvertently introducing bias or other unfair outcomes? Fortunately, while the sophisticated technology may be new, the FTC's attention to automated decision making is not. The FTC has decades of experience enforcing three laws important to developers and users of AI:

- **Section 5 of the FTC Act.** The FTC Act prohibits unfair or deceptive practices. That would include the sale or

# Explainability?

# WHAT IS AN EXPLANATION?

禁止合闸
有人工作

KEEP
← RIGHT

Terms and Conditions

Correct
Compreensible
Timely
Complete
Parsimonous

**Fire Action**

Any person discovering a fire
1. Sound the alarm.
2. _____ to call the fire brigade
3. Attack the fire if possible using the appliances provided

On hearing the fire alarm
4. Leave the building by _____ route
5. Close all doors behind you
6. Report to assembly point _____

Do not take risks
Do not return to the building for any reason until authorised to do so

Email: virginia@cs.umu.se, twitter: @vdignum

# RAI Internships/ Experiential Projects

**Idea is the converse of typical experiential projects where students get a problem & dataset and build models**

**RAI Project –students get a data and model, which needs auditing on the appropriateness of the methods and validating data sources/ model training**

**Different approaches**

Sandbox

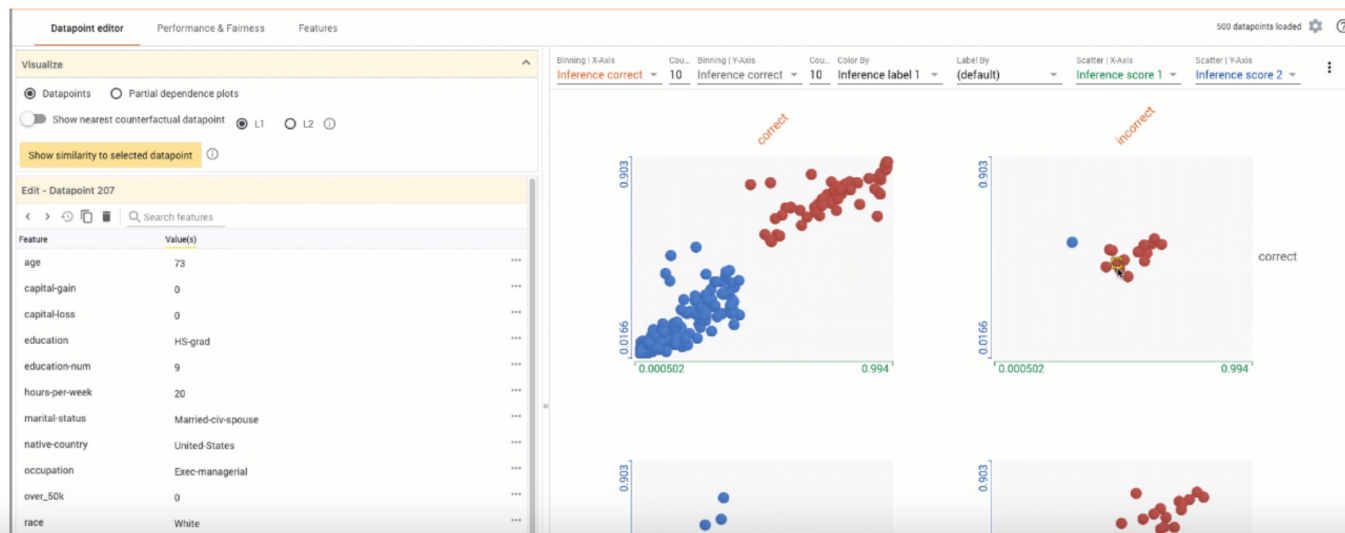Counterfactual Analysis

Fairness Audits

RAI Collab

# What-If from Google Research

+ Code    + Text    ☁ Copy to Drive

## What-If Tool on COMPAS

This notebook shows use of the What-If Tool on the COMPAS dataset.

For ML fairness background on COMPAS see:

- https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing
- https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm
- http://www.crj.org/assets/2017/07/9_Machine_bias_rejoinder.pdf

The dataset is from the COMPAS kaggle page.

This notebook trains a linear classifier on the on the COMPAS dataset to mimic the behavior of the the COMPAS recidivism class then analyze our COMPAS proxy model for fairness using the What-If Tool.

The specific binary classification task for this model is to determine if a person belongs in the "Low" risk class according to COM class), or the "Medium" or "High" risk class (positive class).

▶   Install the What-If Tool widget if running in colab

▶   Define helper functions

▶   Read training dataset from CSV

# IBM AI Fairness 360

AI Fairness 360 - Demo

○───○───○───○
Data  Check  Mitigate  Compare

Back   Next
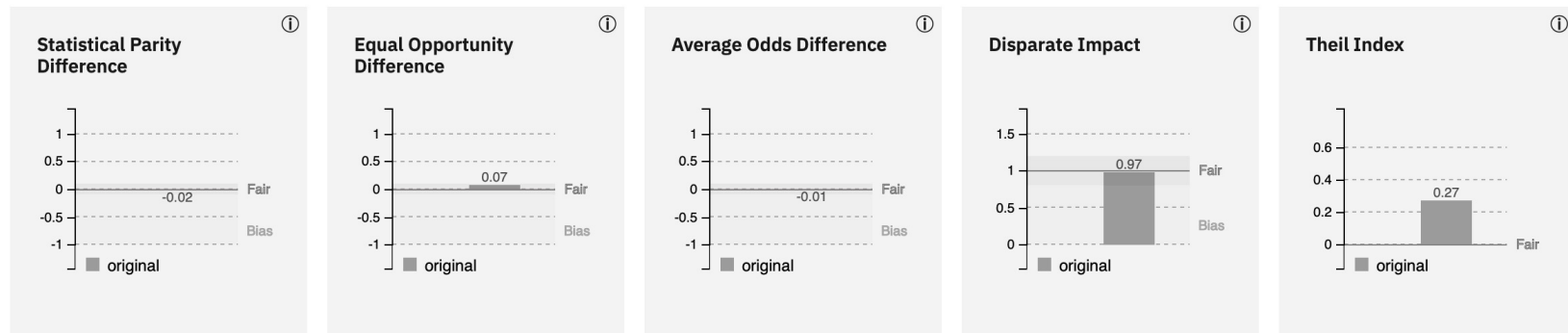
## 2. Check bias metrics

Dataset: German credit scoring
Mitigation: none

### Protected Attribute: Sex

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

Accuracy with no mitigation applied is 75%

With default thresholds, bias against unprivileged group detected in 0 out of 5 metrics

| **Statistical Parity Difference** ⓘ | **Equal Opportunity Difference** ⓘ | **Average Odds Difference** ⓘ | **Disparate Impact** ⓘ | **Theil Index** ⓘ |
|---|---|---|---|---|



Statistical Parity Difference: -0.02 (Fair / Bias), original

Equal Opportunity Difference: 0.07 (Fair / Bias), original

Average Odds Difference: -0.01 (Fair / Bias), original

Disparate Impact: 0.97 (Fair / Bias), original

Theil Index: 0.27 (Fair), original

# Building Bias Mitigation

Local explanation for class >= 10 Visits

CHBRON=2
COGLIM=2
JTPAIN=2
ADSMOK42=2
PREGNT=-1
ADHDADDX=-1
PHQ242=0
ARTHTYPE=-1
EMPHDX=2
ANGIDX=2

−0.25  −0.20  −0.15  −0.10  −0.05  0.00  0.05  0.10

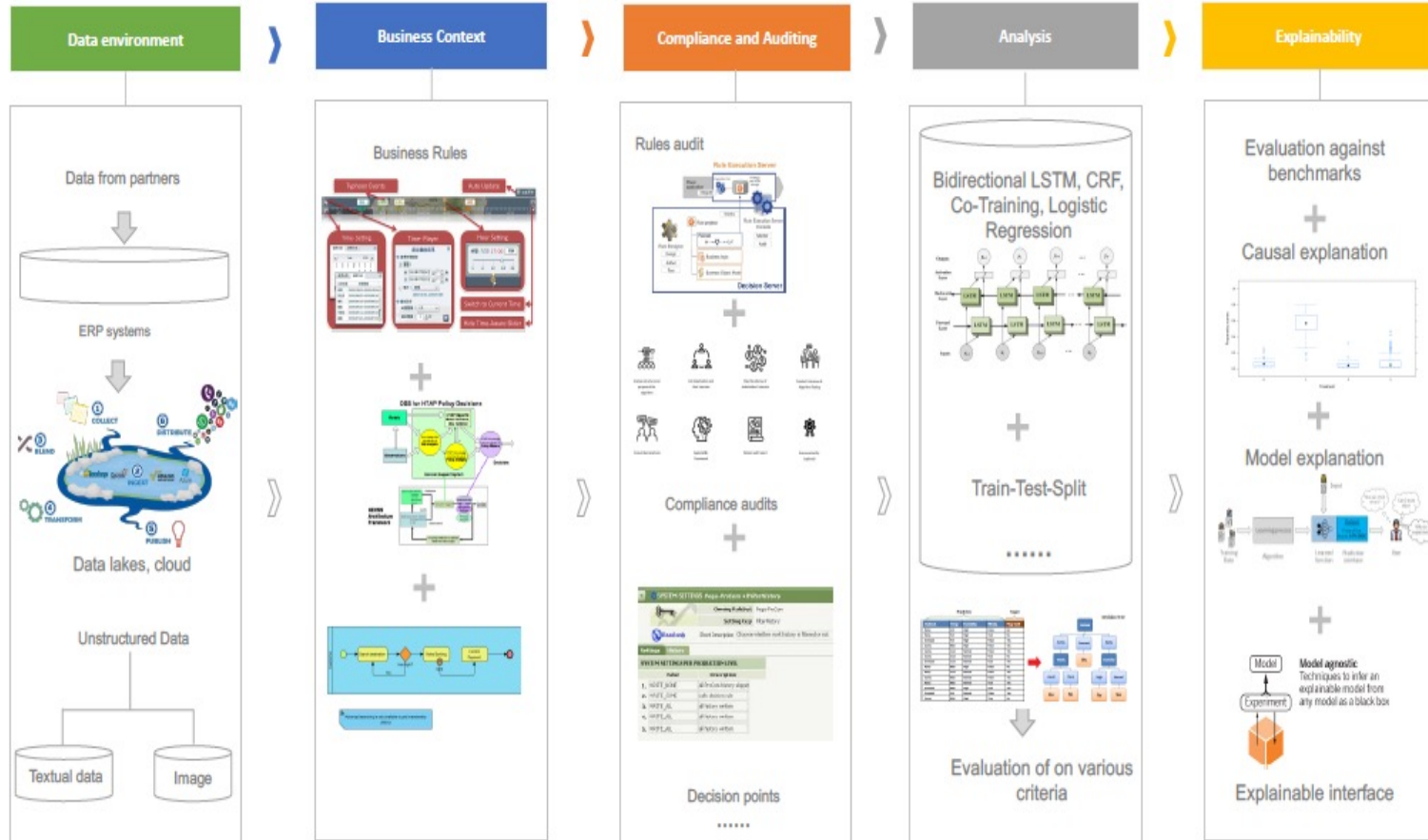https://nbviewer.jupyter.org/github/IBM/AIF360/blob/master/examples/tutorial_medical_expenditure.ipynb

# Use Case of Automated Lending

- What are the key demographics that are at a disadvantage with this work?

- How can we ensure that automated lending assisted or produced by an AI system is fair?

- What are other measures such as improved explainability and oversight of these systems that should be required?

- What data is typically collected for lending? Is there other data that should or should not be collected to assess the responsible use of AI systems?

# Responsible AI https://anjanasusarla.substack.com
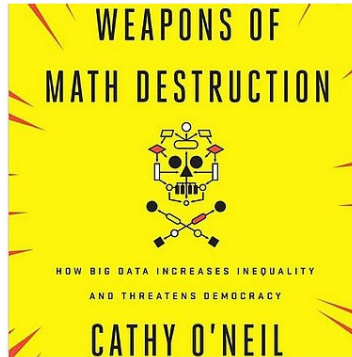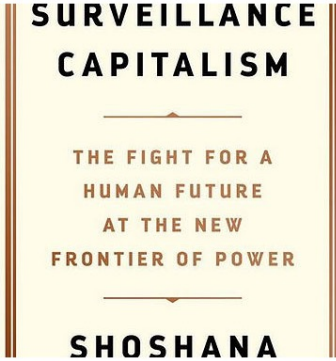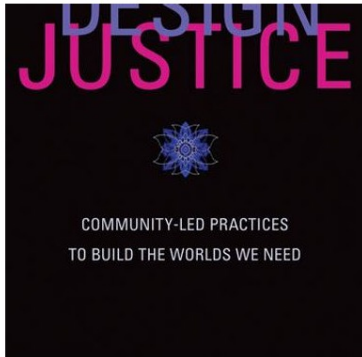


Explainability and Automation of Decision Rules

(© Susarla, 2021)

# RAI Reading List

# Questions?