

# Trust and Trustworthiness for Dependable Machine Learning

Flavio Figueiredo (flaviovdf@dcc.ufmg.br)  
Universidade Federal de Minas Gerais

# Outline

## Initial Discussions

- Trust and trustworthiness
- When a system dependable? How does machine learning affect all of this?

## Overview of ML systems and definitions

- Fairness
- Accountability
- Transparency, and
- Interpretability

# Joint Work

## Atmosphere's WP6 Team

- Leandro Balby (UFCG)
- Vasiliki Diamantopoulou (UPRC)
- Wagner Meira (UFMG)
- + others



# Trust and Trustworthiness

# Simplified Background

- Mathematically speaking, what is the goal of a **supervised** learning system?

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$$

- The goal is to learn some parameters

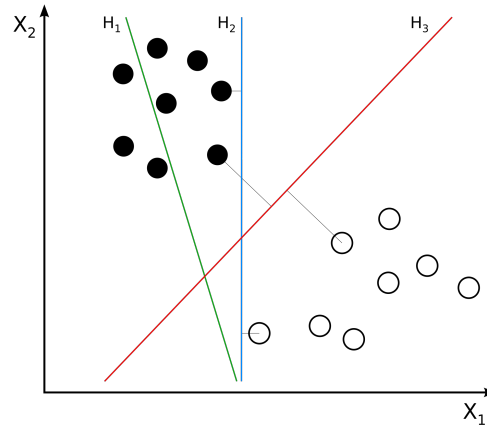
$$\Theta$$

- Where these parameters maximize some prediction function across  $y$

$$\hat{\Theta} = \arg \max_{\Theta} P(\mathbf{y} \mid \mathbf{X}, \Theta)$$

# Simplified Background

- The goal of a supervised learning algorithm is to **discriminate**



# Simplified Background

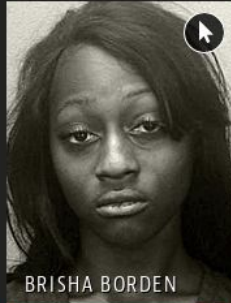
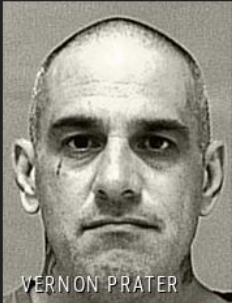
- The goal of a supervised learning algorithm is to **discriminate**
- Why are we now so worried that it does? It seems we can **trust** them.

# Machine Bias

- Pro Publica analysis of COMPAS (which stands for Correctional Offender Management Profiling for Alternative Sanctions)

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

Two Petty Theft Arrests



NAME	RISK	SCORE
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.*

Two Petty Theft Arrests

NAME	Prior Offenses	Subsequent Offenses	RISK	SCORE
VERNON PRATER	2 armed robberies, 1 attempted armed robbery	1 grand theft	LOW RISK	3
BRISHA BORDEN	4 juvenile misdemeanors	None	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.*



# What is trust?

- Inspired by Onora O'Neill  
<https://www.youtube.com/watch?v=XWwTYy9k5nc>
- Consider a question? **Do we trust politicians?**
  - This question has had the same answer for a long time
  - People usually don't trust politicians



# What is trustworthiness?

- Inspired by Onora O'Neill
  - <https://www.youtube.com/watch?v=XWwTYy9k5nc>
- Consider a question? **Do we trust politicians?**
  - This question has had the same answer for a long time
  - People usually don't trust politicians
- **Trustworthiness**
  - Evidence of **why** can I trust you
  - Evidence is observable (though hard to quantify): competence, reliability
  - We trust a science not because it came from a scientist, it is **testable**
- We need to direct our **trust** to **trustworthy** properties.
  - Why and when can I trust



**Services need to earn trust. They need to be trustworthy**

**Trustworthiness changes  
over time**

# Dependability

# Dependability

## From Wikipedia

### Dependability

---

From Wikipedia, the free encyclopedia

In [systems engineering](#), **dependability** is a measure of a system's **availability**, **reliability**, and its **maintainability**, and **maintenance support performance**, and, in some cases, other characteristics such as **durability**, **safety** and **security**.<sup>[1]</sup> In [software engineering](#), **dependability** is the ability to provide services that can defensibly be trusted within a time-period.<sup>[2]</sup> This may also encompass mechanisms designed to increase and maintain the dependability of a system or software.<sup>[3]</sup>

The [International Electrotechnical Commission](#) (IEC), via its Technical Committee [TC 56](#) develops and maintains international standards that provide systematic methods and tools for dependability assessment and management of equipment, services, and systems throughout their life cycles.

Dependability can be broken down into three elements:

- **Attributes** - A way to assess the dependability of a system
- **Threats** - An understanding of the things that can affect the dependability of a system
- **Means** - Ways to increase a system's dependability

However, over time these properties get more complex.

ML systems currently impact society.

It is interesting that it mentions: **trust over time**

# It's a hard problem

- **Uptime.** Do we want more or less of it?

# It's a hard problem

- **Uptime.** Do we want more or less of it?
- Probably more, there are clear problems that require more uptime.
- Maybe cost and energy are issues, but they are quantifiable issues



# It's a hard problem

- **Fairness.** Do we want more or less of it?

# It's a hard problem

- **Fairness.** Do we want more or less of it?
- Name one problem solved by machine learning fairness?

# It's a hard problem

- **Fairness.** Do we want more or less of it?
- Name one problem solved by machine learning fairness?
- What is fairness?!

# It's a hard problem

- **Fairness.** Do we want more or less of it?
- Name one problem solved by machine learning fairness?
- What is fairness?!

Let's try and help

- On a credit scoring system that helps one decide loans to give out. Do we want more fairness?

# It's a hard problem

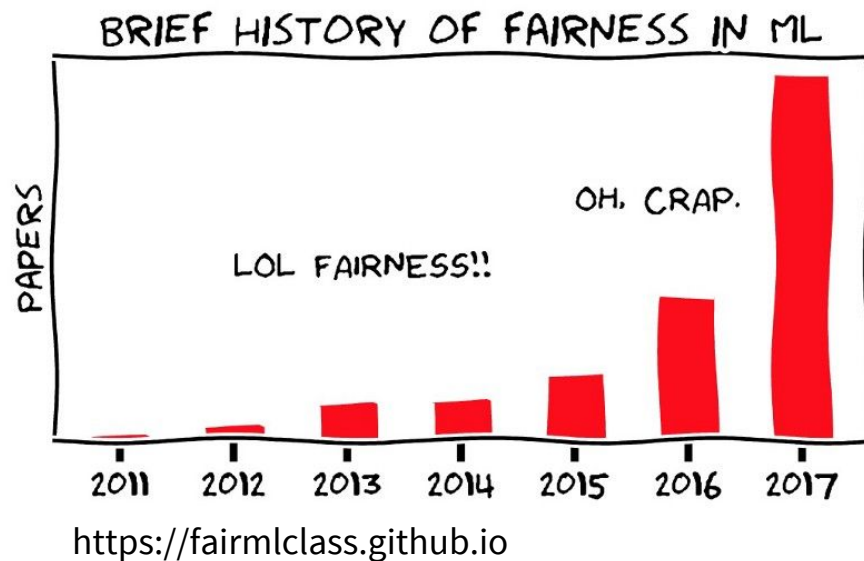
- **Fairness.** Do we want more or less of it?
- Name one problem solved by machine learning fairness?
- What is fairness?!

Let's try and help

- On a credit scoring system that helps one decide loans to give out. Do we want more fairness?
- What if the system simply denies all loans?

# Machine Learning Systems

- Need to specify their trustworthiness
  - Fairness
  - Transparency
  - Accountability
  - Interpretability
- I can't even trust the definition of fairness



# It is a human problem

Nitin Koli, Joshua Kroll (NeuRIPS, 2018)

- Issues of fairness, transparency, accountability, transparency and interpretability are **social-technological**

“Technologies don’t live in a vacuum and if we pretend that they do we kind of have put our blinders on and decided to ignore any human problems.”

# Dependability

## From the Working Group

### About IFIP Working Group 10.4

Increasingly, individuals and organizations are developing or procuring sophisticated computing systems on whose services they need to place great reliance. In differing circumstances, the focus will be on differing properties of such services -- e.g., continuity, performance, real-time response, ability to avoid catastrophic failures, prevention of deliberate privacy intrusions.

The notion of **dependability**, defined as *the trustworthiness of a computing system which allows reliance to be justifiably placed on the service it delivers*, enables these various concerns to be subsumed within a single conceptual framework. Dependability thus includes as special cases such attributes as **reliability, availability, safety, security**.

The Working Group is aimed at identifying and integrating approaches, methods and techniques for specifying, designing, building, assessing, validating, operating and maintaining computer systems which should exhibit some or all of these attributes.

Specifically, the Working Group is concerned with progress in:

1. Understanding of faults (accidental faults, be they physical, design-induced, originating from human interaction; intentional faults) and their effects.
2. Specification and design methods for dependability.
3. Methods for error detection and processing, and for fault treatment.
4. Validation (testing, verification, evaluation) and design for testability and verifiability.
5. Assessing dependability through modeling and measurement.





# Dependability

- Assumes that the human negatively impacts system dependability
- We are now learning that with machine learning systems it's the other way around
- Systems now negatively impact social structures

# Dependability

- Assumes that the human negatively impacts system dependability
- We are now learning that with machine learning systems it's the other way around
- Systems now negatively impact social structures

## How social media filter bubbles and algorithms influence the election

With Facebook becoming a key electoral battleground, researchers are studying how automated accounts are used to alter political debate online

- **Revealed: Facebook's internal rules on sex, terrorism and violence**



▲ A Facebook Live broadcast hosted by ITV News had Theresa May answering questions sent in by users of the site. Photograph: Facebook/PA

## Facebook's WhatsApp limits text forwards to 5 recipients to curb rumors

PUBLISHED MON, JAN 21 2019 • 4:46 AM EST | UPDATED TUE, JAN 22 2019 • 9:58 AM EST

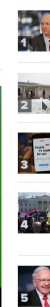
REUTERS

SHARE

### KEY POINTS

- \* WhatsApp expanded a cap on message forwarding it first introduced in India after the spread of rumors led to killings and lynching attempts.
- \* The messaging app will roll out an update to activate the new forward limit, starting Monday, WhatsApp's head of communications said.

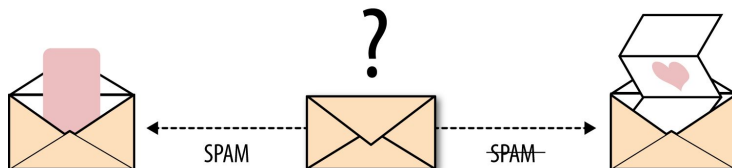
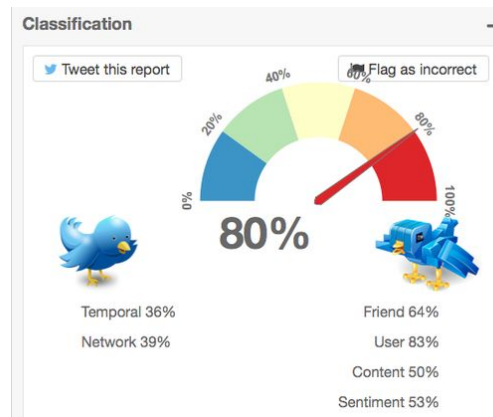
### TRENDING



Chris Ratti/Getty Images

# The old arms race

- We can tackle these problems as an arms race



# Social Technological

- Computer Science
- Data Science
- Law
- Health
- Politics
- etc.



<https://vimeo.com/149389876>

**FAT\***

# What is Fairness?

Let's begin with fairness as it closely relates to all other metrics.

# 21 fairness definitions and their politics

Arvind Narayanan - FAT Conference 2018 Tutorial

- Computer Scientist on a wild goose chase for a single definition
- There is value to various definitions
- Each can lead to **trustworthiness**

# What is Fairness?

Sahil Verma and Julia Rubin (2018) -- Fairness Definitions Explained

- A lot of these metrics worry about some form of equality

$$\hat{\Theta} = \arg \max_{\Theta} P(\mathbf{y} \mid \mathbf{X}, \Theta)$$

- Let  $S$  be some subset of sensitive attributes.

$$S = \{ \text{col}(j, X) \mid \text{column } j \text{ is sensitive} \}$$

$$N = \{ \text{col}(i, X) \mid \text{column } i \text{ is not-sensitive} \}$$

	Definition	Paper	Citation #	Result
3.1.1	Group fairness or statistical parity	[12]	208	×
3.1.2	Conditional statistical parity	[11]	29	✓
3.2.1	Predictive parity	[10]	57	✓
3.2.2	False positive error rate balance	[10]	57	×
3.2.3	False negative error rate balance	[10]	57	✓
3.2.4	Equalised odds	[14]	106	×
3.2.5	Conditional use accuracy equality	[8]	18	×
3.2.6	Overall accuracy equality	[8]	18	✓
3.2.7	Treatment equality	[8]	18	×
3.3.1	Test-fairness or calibration	[10]	57	✓
3.3.2	Well calibration	[16]	81	✓
3.3.3	Balance for positive class	[16]	81	✓
3.3.4	Balance for negative class	[16]	81	×
4.1	Causal discrimination	[13]	1	×
4.2	Fairness through unawareness	[17]	14	✓
4.3	Fairness through awareness	[12]	208	×
5.1	Counterfactual fairness	[17]	14	–
5.2	No unresolved discrimination	[15]	14	–
5.3	No proxy discrimination	[15]	14	–
5.4	Fair inference	[19]	6	–

Table 1: Considered Definitions of Fairness



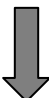
# Balanced Representation

Is this fairness?

$$\hat{\Theta} = \arg \max_{\Theta} P(\mathbf{y} | \mathbf{X}, \Theta)$$



$\mathbf{X} =$



1	1	1	0	1
1	0	1	1	0
0	0	1	0	1
1	0	1	1	0
1	1	0	1	0

N

S

# Parity in Predictions

	Actual – Positive	Actual – Negative
Predicted – Positive	<p><b>True Positive (TP)</b></p> <p>PPV = <math>\frac{TP}{TP+FP}</math></p> <p>TNR = <math>\frac{TN}{TP+FN}</math></p>	<p><b>False Positive (FP)</b></p> <p>FDR = <math>\frac{FP}{TP+FP}</math></p> <p>FPR = <math>\frac{FP}{FP+TN}</math></p>
Predicted – Negative	<p><b>False Negative (FN)</b></p> <p>FOR = <math>\frac{FN}{TN+FN}</math></p> <p>FNR = <math>\frac{FN}{TP+FN}</math></p>	<p><b>True Negative (TN)</b></p> <p>NPV = <math>\frac{TN}{TN+FN}</math></p> <p>TNR = <math>\frac{TN}{TN+FP}</math></p>

# Some Examples

- Demographic Parity

$$P(\hat{y} = 1|s = 0) = P(\hat{y} = 1|s = 1)$$

- Equality in Opportunity (FNR)

$$P(\hat{y} = 0|s = 0, y = 0) = P(\hat{y} = 0|s = 1, y = 0)$$

- Calibration

$$E[y|s = 0, \hat{y} = p] = E[y|s = 1, \hat{y} = p] \quad \forall p \in [0, 1]$$

# No free Lunch

## Arvind Narayanan - FAT Conference 2018 Tutorial

		True condition			
Total population		Condition positive	Condition negative	Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Predicted condition	Predicted condition positive	<b>True positive</b> , Power	<b>False positive</b> , Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	<b>False negative</b> , Type II error	<b>True negative</b>	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) $= \frac{\text{LR+}}{\text{LR-}}$
		False negative rate (FNR), Miss rate $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) $= \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$	

# Feedback Loops and Utility

How safe do we want a city to be?

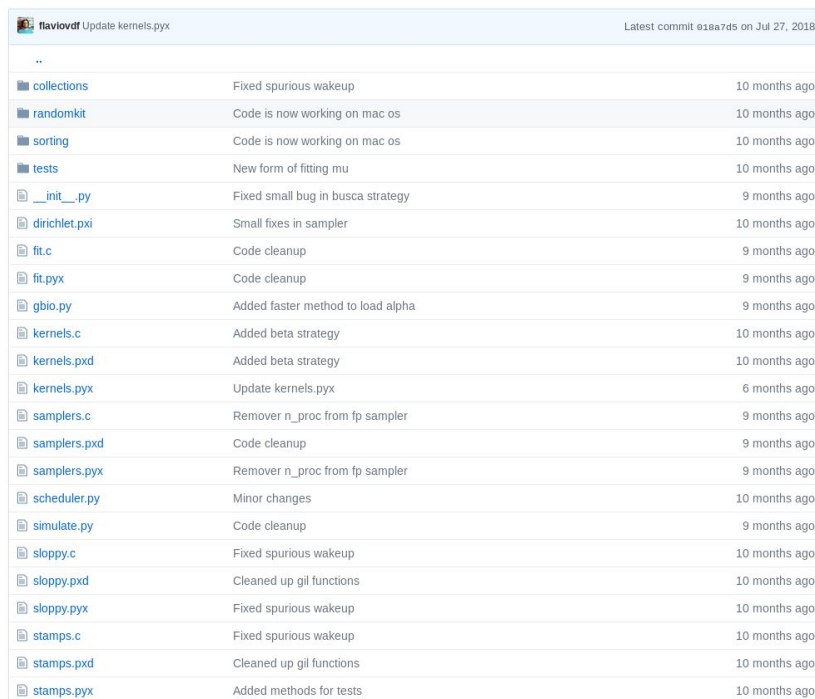
It can be shown thresholding this score, leads to unfairness.

# Drawbacks

- We are mostly focused on correlations
- Maybe that nice matrix is impossible
- We are reducing a **processes** to **measures**

# Free Software Approach to Transparency

- The source code is public and auditable



flaviovd Update kernels.pyx Latest commit 018a7d5 on Jul 27, 2018

..		
collections	Fixed spurious wakeup	10 months ago
randomkit	Code is now working on mac os	10 months ago
sorting	Code is now working on mac os	10 months ago
tests	New form of fitting mu	10 months ago
__init__.py	Fixed small bug in busca strategy	9 months ago
dirichlet.pxi	Small fixes in sampler	10 months ago
fit.c	Code cleanup	9 months ago
fit.pyx	Code cleanup	9 months ago
gbio.py	Added faster method to load alpha	9 months ago
kernels.c	Added beta strategy	10 months ago
kernels.pxd	Added beta strategy	10 months ago
kernels.pyx	Update kernels.pyx	6 months ago
samplers.c	Remover n_proc from fp sampler	9 months ago
samplers.pxd	Code cleanup	9 months ago
samplers.pyx	Remover n_proc from fp sampler	9 months ago
scheduler.py	Minor changes	10 months ago
simulate.py	Code cleanup	9 months ago
sloppy.c	Fixed spurious wakeup	10 months ago
sloppy.pxd	Cleaned up gil functions	10 months ago
sloppy.pyx	Fixed spurious wakeup	10 months ago
stamps.c	Fixed spurious wakeup	10 months ago
stamps.pxd	Cleaned up gil functions	10 months ago
stamps.pyx	Added methods for tests	10 months ago

# A Lazy ML Approach to Transparency

- I employed a simple model, thus it is easy to understand

Table 4: Average value ( $\mu$ ) and the lower ( $\downarrow 95\%$ ) and upper HPD ( $\uparrow 95\%$ ) values for significant explanatory variables in each topic cluster.

	$\mu$	$\downarrow 95\%$	$\uparrow 95\%$
<b>(S1) Out of Scope</b>			
Intercept	-2.11	-3.04	-1.30
Male Gender & Aud. D	-2.06	-3.98	-0.03
<b>(S2) About Characters</b>			
Intercept	1.66	0.80	2.46
<b>(S3) Greetings</b>			
Intercept	-3.32	-4.56	-2.01
<b>(S4) Reaction to Failure</b>			
Intercept	-4.48	-6.26	-2.78
Audience B	1.74	0.18	3.60
Direct Address	1.88	0.38	3.61
Direct Address & Aud. B	-1.74	-3.35	-0.16



# Interpretability

Lime and Shap. Also limited, ML to explain ML.

Images (explaining prediction of 'Cat' in pros and cons)



# Accountability

From Wikipedia:

“In ethics and governance, accountability is answerability, blameworthiness, liability, and the expectation of account-giving.”

# General Data Protection Regulation

Also the California Consumer Privacy Act

- Big companies need to be **accountable** when using your data
- However, impacts exist even when I agree to share my data
  - Let's say most of the population agrees to share data with Facebook
- Who is **accountable**?

# Accountability

- How do we measure accountability?

# Accountability

- How do we measure accountability?
- Machines can keep track of records (data provenance)

DATE	PARTICULARS	INITIALS	DR.	CR.	DR. OR CR.	BALANCE	DATE	PARTICULARS	INITIALS	DR.	CR.	DR. OR CR.	BALANCE
1953							1953						
Feb 23	Trans			41.52		41.52	June 30	Trans			20.97		20.97
March 17	Dr			74.15		116.37	July 4			10.00			
19			5.00				10	Dr			101.92		
			1325							5.00			
33	July		56				18			50.00			
22			100				27	July			72		
24			1775				Aug 27			2.95			
			1085				Nov 29	Dr			25.00		
April 1			600				Dec 3	Dr			100.00		
			1000							350.00			
12			1700				8			10.00			
18	Dr			150.00			12			17.00			
19			128.80							45.00			
25			10.00				15	Dr			496.98		
28	July		1.00							217.80			
30	Aug. Inv.		106					Oct. Rate			167.71		
	Dr			32.00			21			50.00			
June 7			10.00				27			20.00			
13			20.00							23.67			
24	Dr			104.69			Jan 4/6			28.00			
46	Dr		120.12			20.97	10	Dr			946.69		965.99

# Accountability

- What do we do with it?
- Conflicts with Privacy
- Hot topic in ML nowadays. Is it our job to make others accountable?

# Drawbacks

- We are mostly focused on correlations
- We are reducing a **processes** to **measures**
- Kroll et al. (2018). When is an election fair (or transparent, accountable)?  
An election is a process. The whole process should be accountable, transparent and subject to recounts.

# Counterfactuals

- Evaluates the impact of features with counterfactual approach  
[Zhang and Bareinboim (2018)]
- **“Would the prediction change if the subject were black?”**



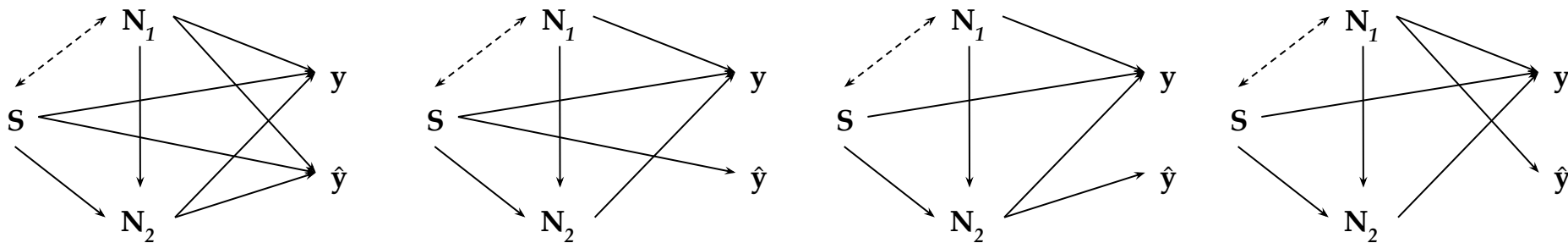
# Causal Analysis

<https://pair-code.github.io/what-if-tool/>

1. For a given subject
2. Find the closest point with a different prediction and different sensitive attribute
3. Swap features keeping the sensitive attribute

# Causal Fairness

The COMPASS model.  $S$  captures race.  $N_1$  demography.  $N_2$  prior convictions.  
Zhang and Bareinboim (2018) -- Equality of Opportunity in Classification



Each of these classifiers have the same equalized odds.



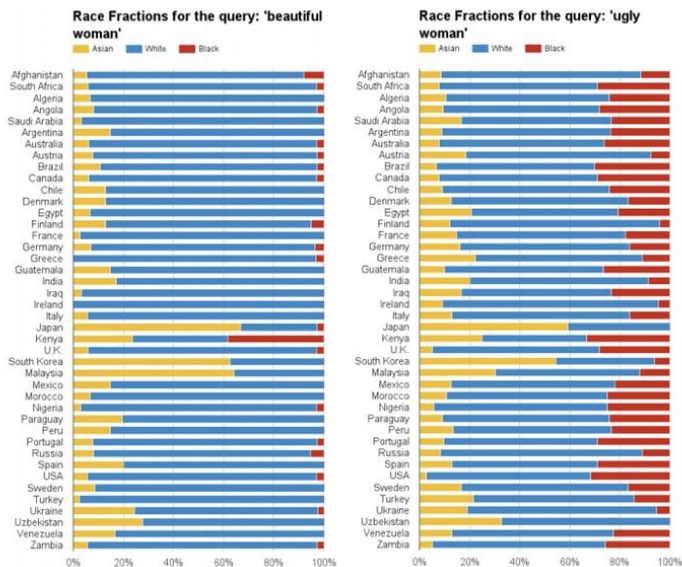
- Step in the right direction
- Human in the loop (provides the DAG asks the question)
- One step closer to a process.

However

- We are not lawmakers or sociologists
- We still need to educate and get educated

# Limitations

Different machine learning tasks



# Limitations

Different machine learning tasks

Man is to Computer Programmer as Woman is to Homemaker?  
Debiasing Word Embeddings

Tolga Bolukbasi<sup>1</sup>, Kai-Wei Chang<sup>2</sup>, James Zou<sup>2</sup>, Venkatesh Saligrama<sup>1,2</sup>, Adam Kalai<sup>2</sup>

<sup>1</sup>Boston University, 8 Saint Mary's Street, Boston, MA

<sup>2</sup>Microsoft Research New England, 1 Memorial Drive, Cambridge, MA

tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

# Dependable Systems

- We should not view humans as the cause of problems
- Shift in direction for dependability
- Processes not only measures
  - Systems should enable humans to take actions
  - Open data/model provenance
  - Explanations
  
  - Systems should enable humans to say no
  - I do not want to see certain content. Do not use my data
  
  - Systems should be **trustworthy** in ways the average user can understand

# Hard Problem

Society (as a consequence datasets) is unfair  
Accountability is difficult (who do we blame?)  
Datasets and models are hard to understand

**Thank You!**