

NLP and Society: Towards Socially Responsible NLP

Vinodkumar Prabhakaran

Research Scientist



What's in this lecture

- Motivation for Fairness research in NLP
- How and why NLP models may be unfair
- Various types of NLP fairness issues and mitigation approaches
- What can/should we do?

What's **NOT** in this lecture

- Definitive answers to fairness/ethical questions
- Prescriptive solutions to fix ML/NLP (un)fairness
- Focus on research done by myself, my team, or Google.
- ...

With help from...



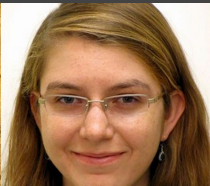
Margaret Mitchell



**Andrew
Zaldivar**



**Emily
Denton**



**Simone
Wu**



**Parker
Barnes**



**Lucy
Vasserman**



**Ben
Hutchinson**



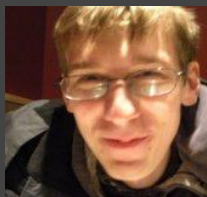
**Elena
Spitzer**



**Deb
Raji**



Timnit Gebru



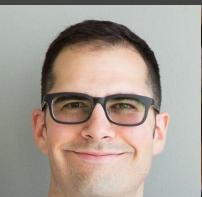
**Adrian
Benton**



**Brian
Zhang**



**Dirk
Hovy**



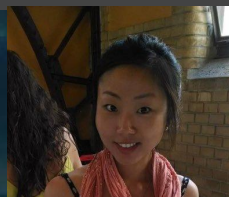
**Josh
Lovejoy**



**Alex
Beutel**



**Blake
Lemoine**



**Hee Jung
Ryu**



**Hartwig
Adam**



**Blaise
Agüera y
Arcas**

What do you see?



What do you see?

- Bananas



What do you see?

- Bananas
- Stickers



What do you see?

- Bananas
- Stickers
- Dole Bananas



What do you see?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store



What do you see?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store
- Bananas on shelves



What do you see?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store
- Bananas on shelves
- Bunches of bananas



What do you see?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store
- Bananas on shelves
- Bunches of bananas

...We don't tend to say
Yellow Bananas



What do you see?

Green Bananas

Unripe Bananas



What do you see?

Ripe Bananas

Bananas with **spots**



What do you see?

Yellow Bananas

Yellow is prototypical for
bananas



Prototype Theory

One purpose of categorization is to **reduce the infinite differences** among stimuli **to** behaviourally and **cognitively usable proportions**

There may be some central, prototypical notions of items that arise from stored typical properties for an object category (Rosch, 1975)

May also store exemplars (Wu & Barsalou, 2009)



Fruit



Bananas
“Basic Level”



Unripe Bananas,
Cavendish Bananas

A man and his son are in a terrible accident and are rushed to the hospital in critical care.

The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"

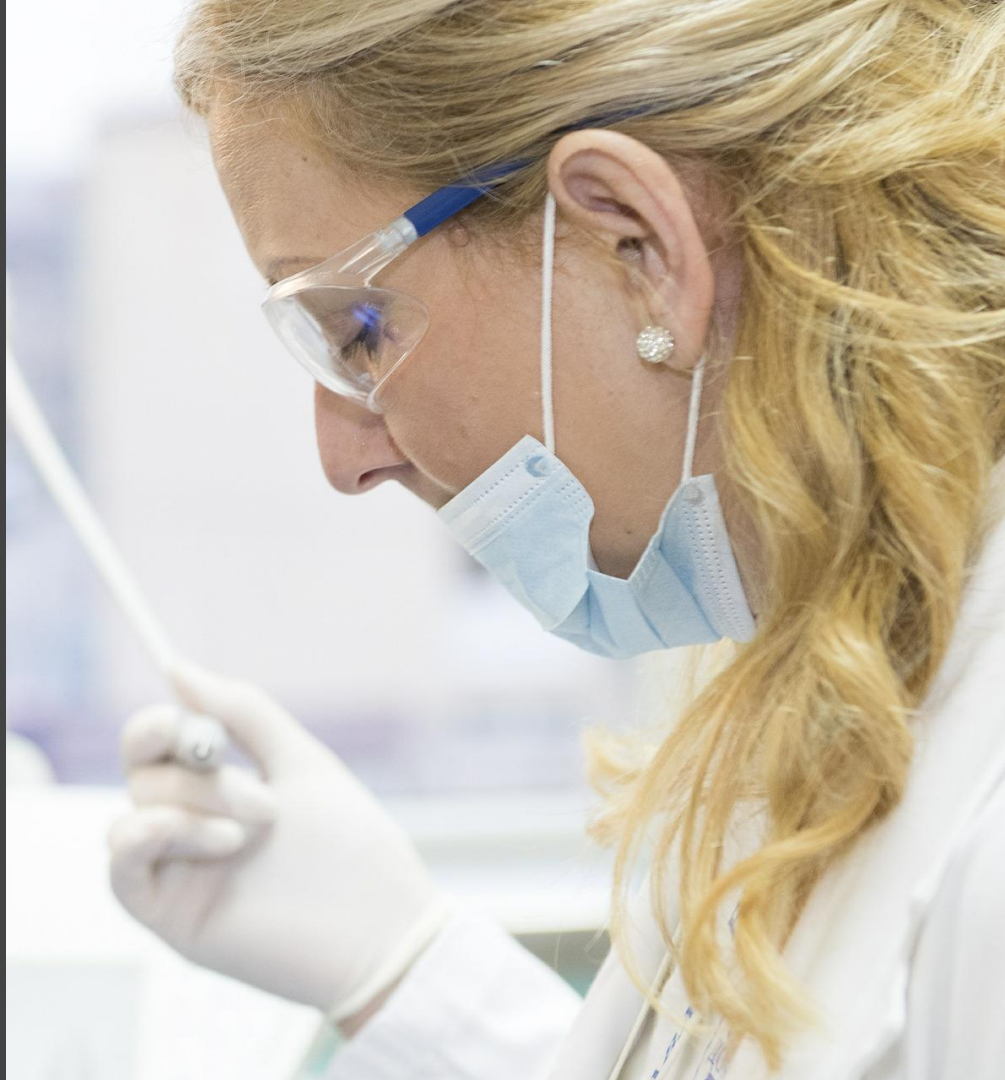
How could this be?



A man and his son are in a terrible accident and are rushed to the hospital in critical care.

The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"

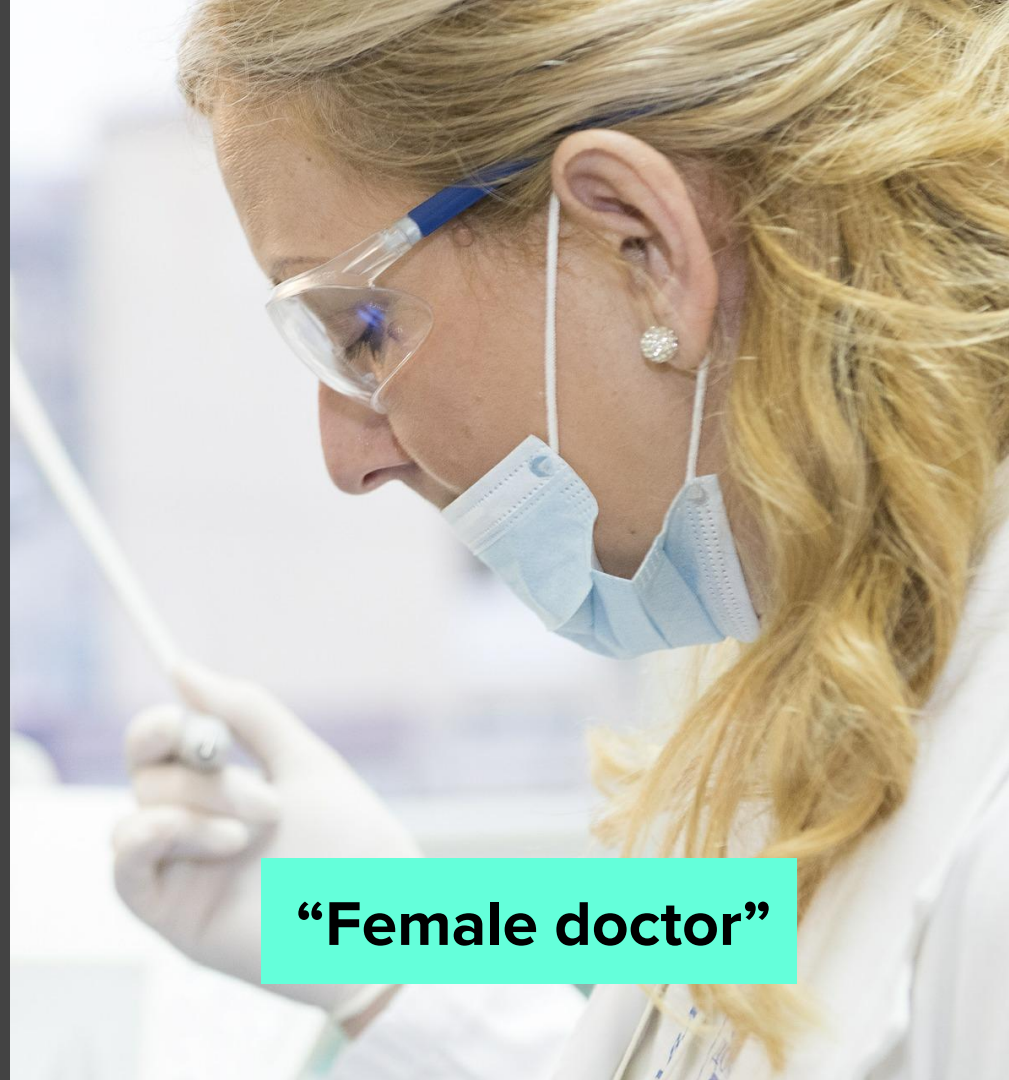
How could this be?



A man and his son are in a terrible accident and are rushed to the hospital in critical care.

The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"

How could this be?



“Female doctor”



“Doctor”



“Female doctor”

Prototype Theory in Action

"male doctor"

 All  Images  Videos  News

About 6,400,000 results (0.47 seconds)

"female doctor"

 All  Images  Videos  News

About 14,000,000 results (0.44 seconds)

Also, found in a study by [Wapman & Belle, Boston University \(2014\)](#)

The majority of test subjects overlooked the possibility that the doctor is a she - including men, women, and self-described feminists.

[Wapman & Belle, Boston University](#)

World learning from text

Gordon and Van Durme, 2013

Word	Frequency in corpus
“spoke”	11,577,917
“laughed”	3,904,519
“murdered”	2,834,529
“inhaled”	984,613
“breathed”	725,034
“hugged”	610,040
“blinked”	390,692
“exhale”	168,985


World learning from text

Gordon and Van Durme, 2013

Word	Frequency in corpus
“spoke”	11,577,917
“laughed”	3,904,519
“murdered”	2,834,529
“inhaled”	984,613
“breathed”	725,034
“hugged”	610,040
“blinked”	390,692
“exhale”	168,985

Human Reporting Bias

The **frequency** with which **people write** about actions, outcomes, or properties is **not a reflection of real-world frequencies** or the degree to which a property is characteristic of a class of individuals



**Training data are
collected and
annotated**

```
graph LR; A((Training data are collected and annotated)) --> B((Model is trained))
```

**Training data are
collected and
annotated**

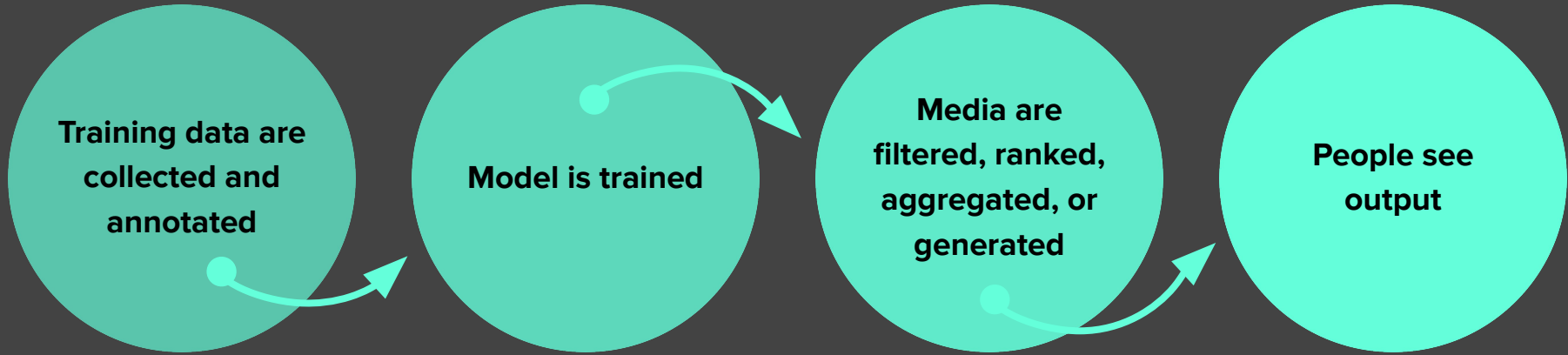
Model is trained

```
graph LR; A((Training data are collected and annotated)) --> B((Model is trained)); B --> C((Media are filtered, ranked, aggregated, or generated));
```

**Training data are
collected and
annotated**

Model is trained

**Media are
filtered, ranked,
aggregated, or
generated**



Human Biases in Data

Reporting bias

Selection bias

Overgeneralization

Out-group homogeneity bias

Stereotypical bias

Historical unfairness

Implicit associations

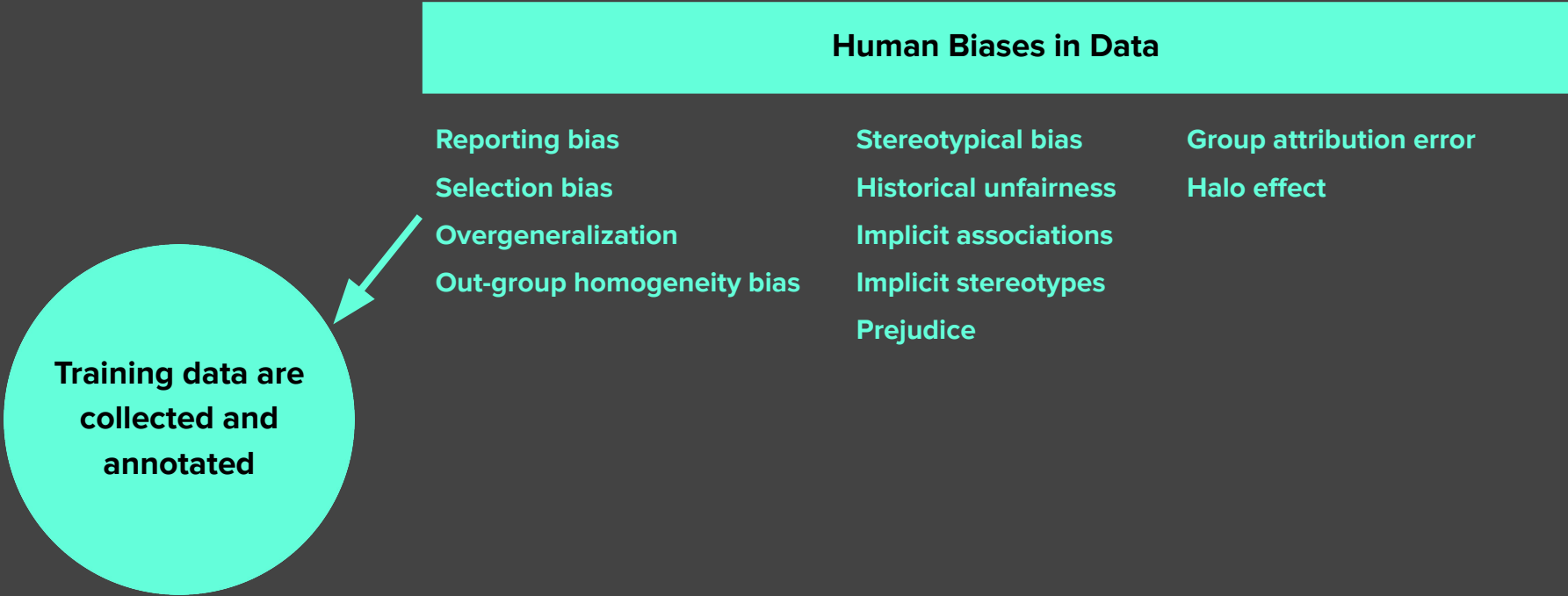
Implicit stereotypes

Prejudice

Group attribution error

Halo effect

Training data are
collected and
annotated



Human Biases in Data

Reporting bias

Stereotypical bias

Group attribution error

Selection bias

Historical unfairness

Halo effect

Overgeneralization

Implicit associations

Out-group homogeneity bias

Implicit stereotypes

Prejudice

Training data are
collected and
annotated

Human Biases in Collection and Annotation

Sampling error

Bias blind spot

Neglect of probability

Non-sampling error

Confirmation bias

Anecdotal fallacy

Insensitivity to sample size

Subjective validation

Illusion of validity

Correspondence bias

Experimenter's bias

In-group bias

Choice-supportive bias

Data

Reporting bias: What people share is not a reflection of real-world frequencies

Selection Bias: Selection does not reflect a random sample

Out-group homogeneity bias: People tend to see outgroup members as more alike than ingroup members when comparing attitudes, values, personality traits, and other characteristics

Confirmation bias: The tendency to search for, interpret, favor, and recall information in a way that confirms one's preexisting beliefs or hypotheses

Interpretation

Overgeneralization: Coming to conclusion based on information that is too general and/or not specific enough

Correlation fallacy: Confusing correlation with causation

Automation bias: Propensity for humans to favor suggestions from automated decision-making systems over contradictory information without automation



Biases in Data

Biases in Data

Selection Bias: Selection does not reflect a random sample

World Englishes



Is the data we use to train our English NLP models representative of all the Englishes out there?

Biases in Data

Selection Bias: Selection does not reflect a random sample

- Men are over-represented in web-based news articles

(Jia, Lansdall-Welfare, and Cristianini 2015)

- Men are over-represented in twitter conversations

(Garcia, Weber, and Garimella 2014)

- Gender bias in Wikipedia and Britannica

(Reagle & Rhuee 2011)

Biases in Data

Selection Bias: Selection does not reflect a random sample

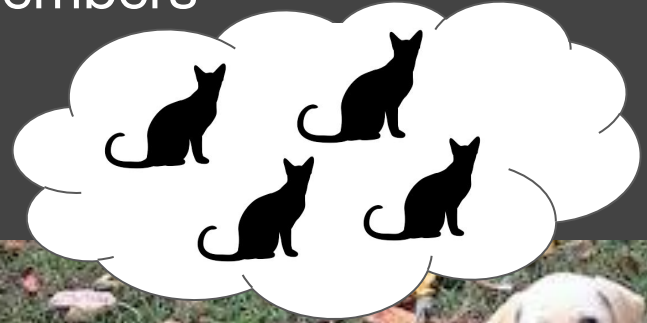


CREDIT

© 2013–2016 Michael Yoshitaka Erlewine and Hadas Kotek

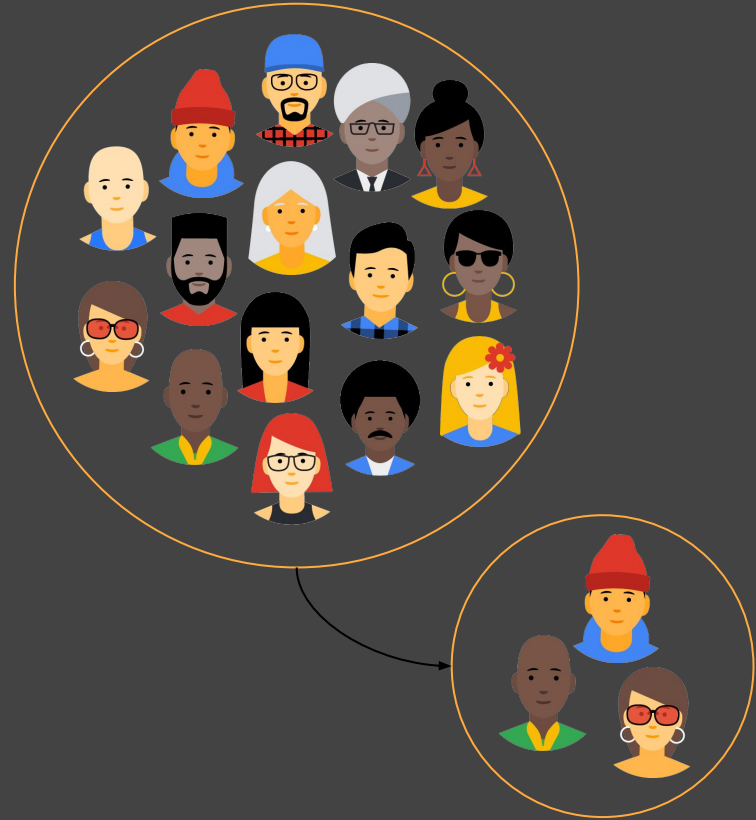
Biases in Data

Out-group homogeneity bias: Tendency to see outgroup members as more alike than ingroup members



Biases in Data → Biased Data Representation

It's possible that you have an appropriate amount of data for every group you can think of but that some groups are represented less positively than others.



Biases in Data → Biased Labels

Annotations in your dataset will reflect the worldviews of your annotators.

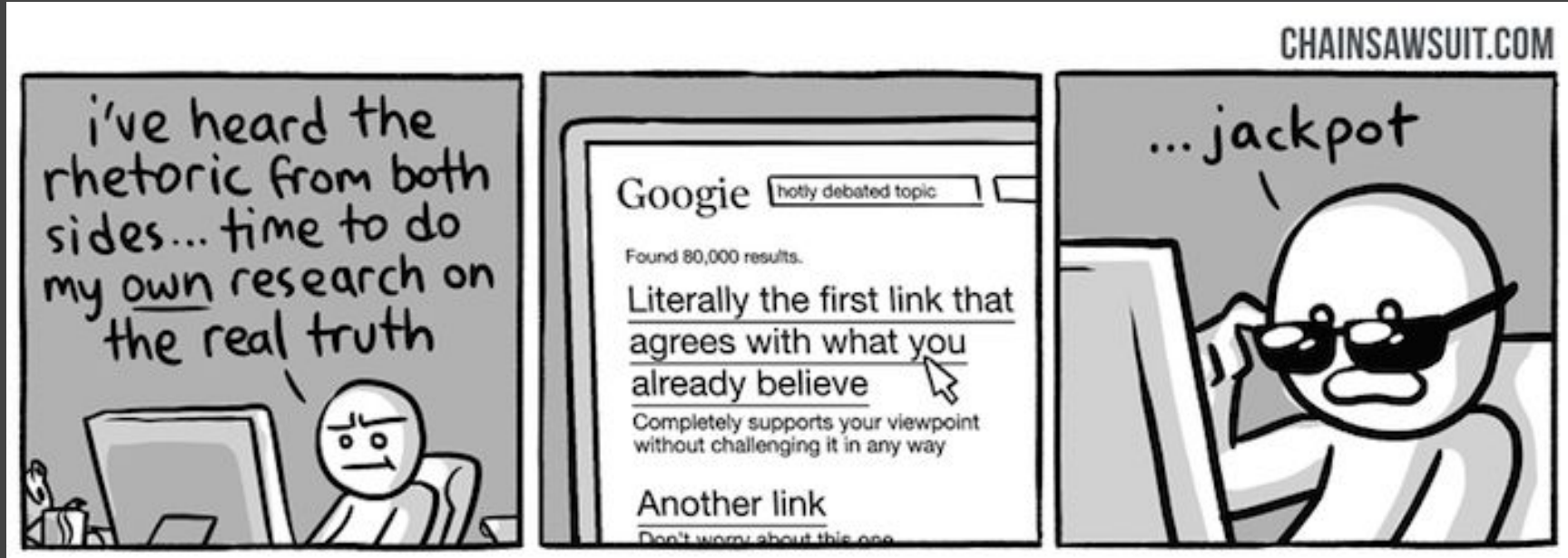




Biases in Interpretation

Biases in Interpretation

Confirmation bias: The tendency to search for, interpret, favor, recall information in a way that confirms preexisting beliefs



CREDIT

© kris straub - Chainsawsuit.com

Biases in Interpretation

Overgeneralization: Coming to conclusion based on information that is too general and/or not specific enough (related: **overfitting**)



CREDIT

Sidney Harris

Biases in Interpretation

Correlation fallacy: Confusing correlation with causation

Post Hoc Ergo Propter Hoc

Women were allowed to vote in the early 1900's and then we had two world wars. Clearly giving them the vote was a bad idea.

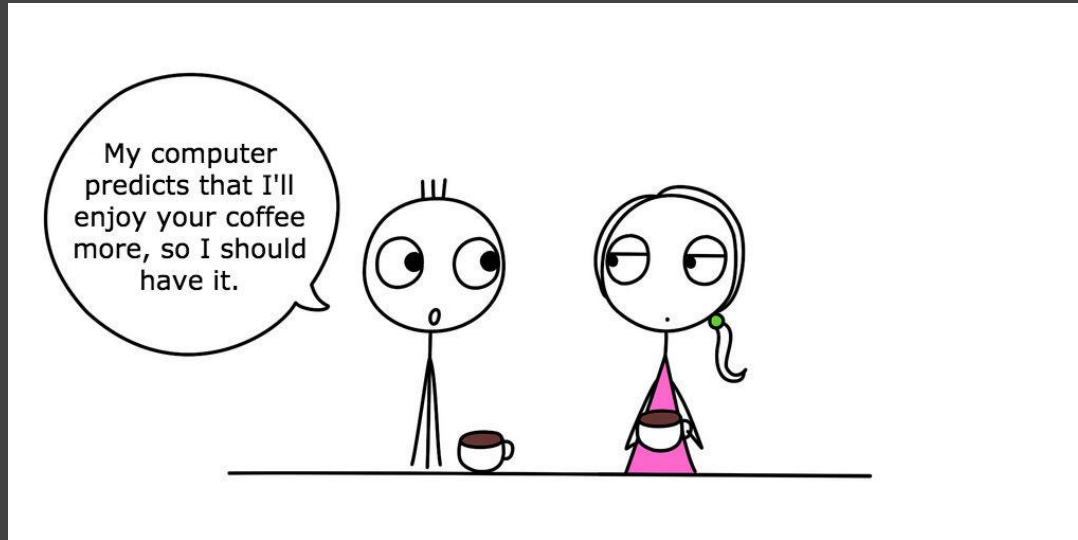


CREDIT

© mollysdad - Slideshare - Introduction to Logical Fallacies

Biases in Interpretation

Automation bias: Propensity for humans to favor suggestions from automated decision-making systems over contradictory information without automation



CREDIT

thedailyenglishshow.com | [CC BY 2.0](https://creativecommons.org/licenses/by/2.0/)

Human Biases in Data

Reporting bias

Stereotypical bias

Group attribution error

Selection bias

Historical unfairness

Halo effect

Overgeneralization

Implicit associations

Out-group homogeneity bias

Implicit stereotypes

Prejudice

Training data are
collected and
annotated

Human Biases in Collection and Annotation

Sampling error

Bias blind spot

Neglect of probability

Non-sampling error

Confirmation bias

Anecdotal fallacy

Insensitivity to sample size

Subjective validation

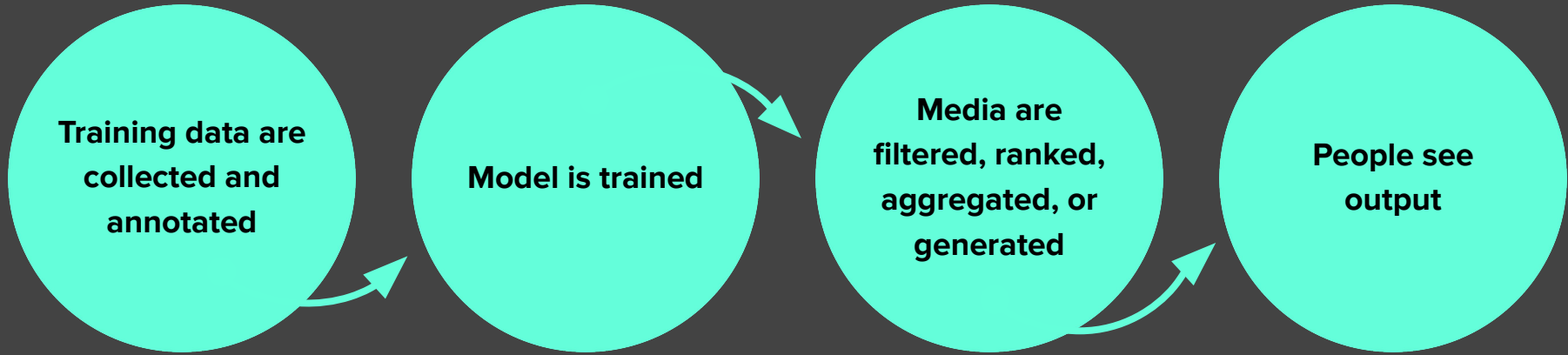
Illusion of validity

Correspondence bias

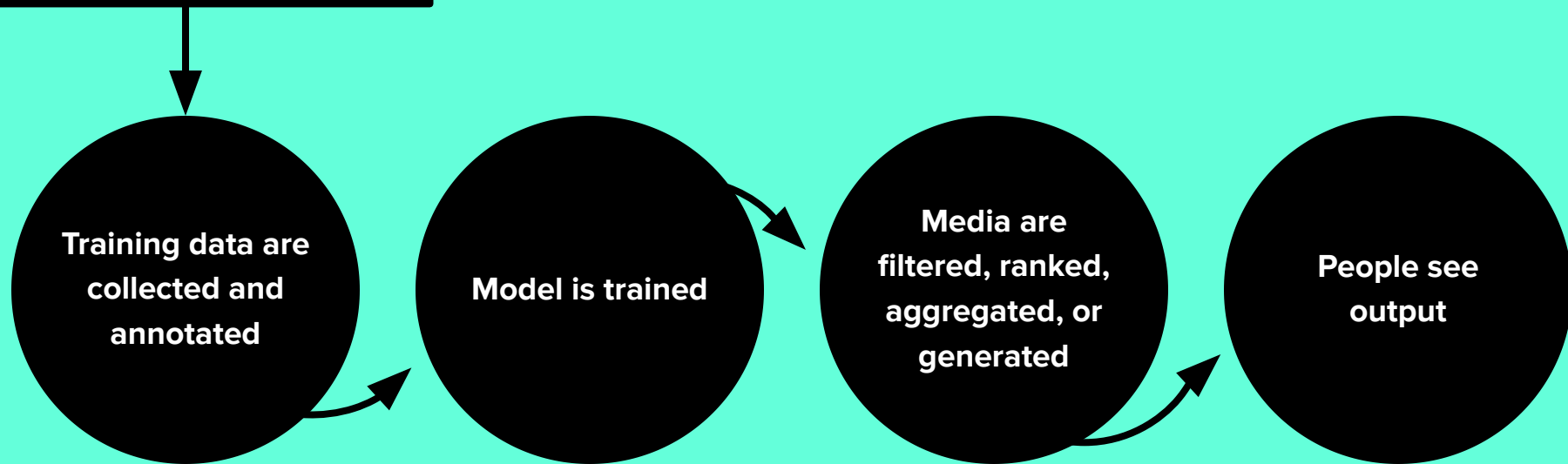
Experimenter's bias

In-group bias

Choice-supportive bias



Human Bias



Human Bias



Training data are
collected and
annotated



Model is trained



Media are
filtered, ranked,
aggregated, or
generated



People see
output

Human Bias

Human Bias

Human Bias



Human Bias



Training data are
collected and
annotated



Model is trained



Media are
filtered, ranked,
aggregated, or
generated



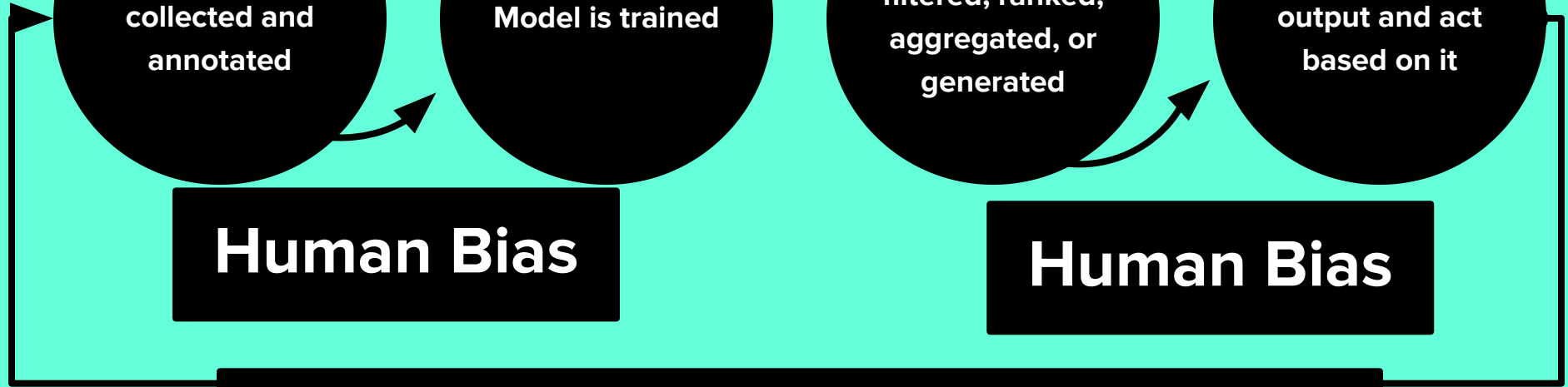
People see
output and act
based on it

Human Bias

Human Bias

Human Bias

Feedback Loop



Human data perpetuates human biases.

**As ML learns from human data, the result is a
bias network effect**

“Bias Laundering”



BIAS = BAD ??

“Bias” can be Good, Bad, Neutral

- Bias in statistics and ML
 - Bias of an estimator: Difference between the predictions and the correct values that we are trying to predict
 - The "bias" term b (e.g., $y = mx + b$)
- Cognitive biases
 - Confirmation bias, Recency bias, Optimism bias
- Algorithmic bias
 - Unjust, unfair, or prejudicial treatment of people related to race, income, sexual orientation, religion, gender, and other characteristics historically associated with discrimination and marginalization, when and where they manifest in algorithmic systems or algorithmically aided decision-making

“Bias” can be Good, Bad, Neutral

- Bias in statistics and ML
 - Bias of an estimator: Difference between the predictions and the correct values that we are trying to predict
 - The "bias" term b (e.g., $y = mx + b$)
- Cognitive biases
 - Confirmation bias, Recency bias, Optimism bias
- **Algorithmic bias**
 - **Unjust, unfair, or prejudicial treatment of people** related to race, income, sexual orientation, religion, gender, and other characteristics historically associated with discrimination and marginalization, when and where they manifest in algorithmic systems or algorithmically aided decision-making

*“Although neural networks might be said to write their own programs, they do so towards **goals set by humans, using data collected for human purposes**. If the data is skewed, even by accident, the computers will amplify injustice.”*

— The Guardian

CREDIT

[The Guardian view on machine learning: people must decide](#)

“Although neural networks might be said to write their own programs, they do so towards goals set by humans, using data collected for human purposes. If the data is skewed, even by accident, the computers will amplify injustice.”

— The Guardian

CREDIT

[The Guardian view on machine learning: people must decide](#)

Fairness in Machine Learning

A Few Case Studies



Language Identification

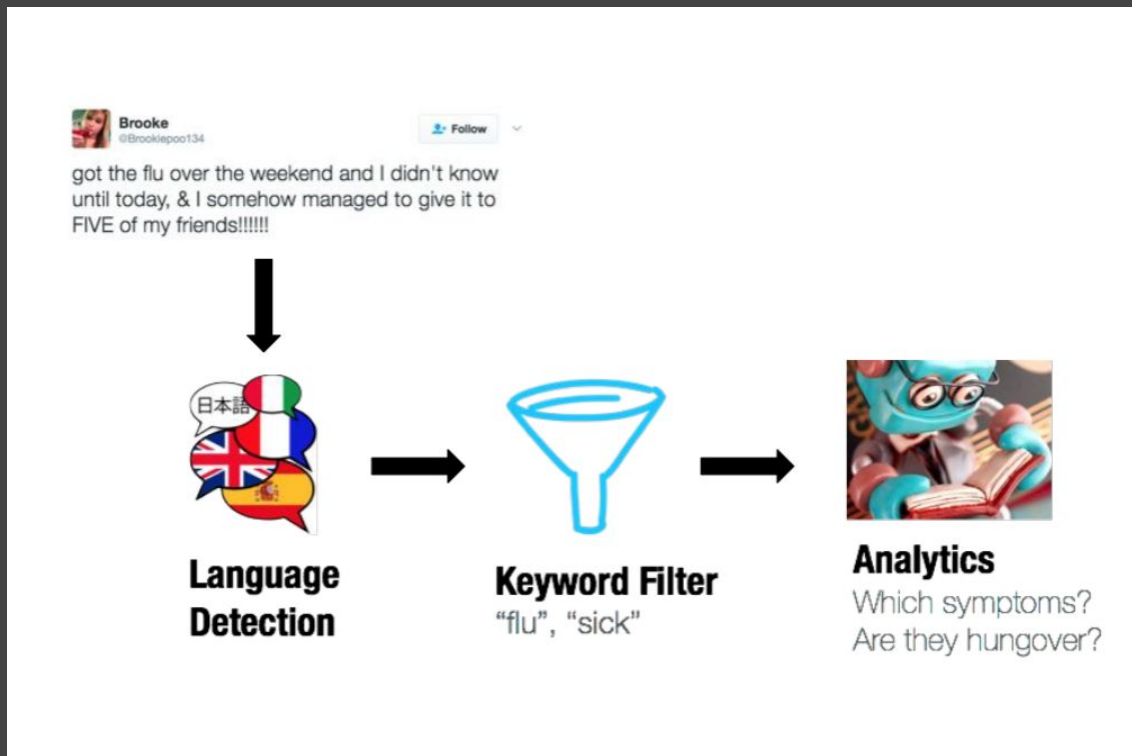
Language Identification

Most NLP models in practice has a Language Identification (LID) step



Language Identification

Most NLP models in practice has a Language Identification (LID) step

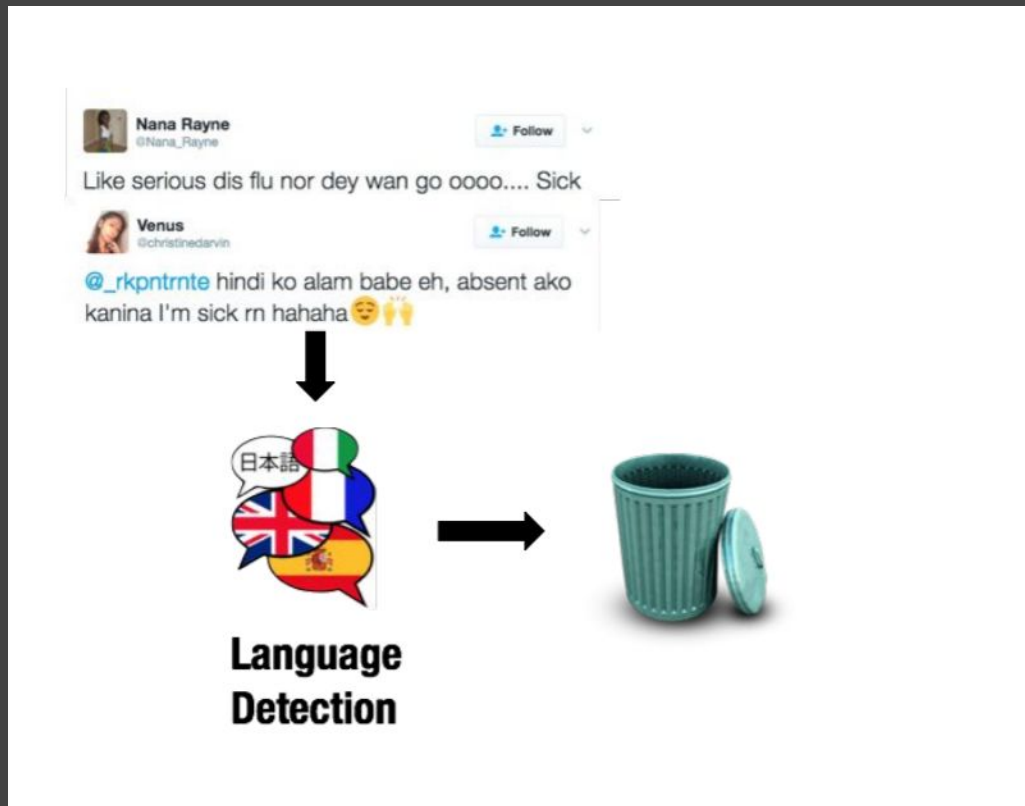


How well do LID systems do?

“This paper describes [...] how even the most simple of these methods *using data obtained from the World Wide Web* achieve accuracy approaching 100% on a test suite comprised of ten European languages”

McNamee, P., “Language identification: *a solved problem* suitable for undergraduate instruction” *Journal of Computing Sciences in Colleges* 20(3) 2005.

LID Usage Example: Public Health Monitoring



Biases in Data

Selection Bias: Selection does not reflect a random sample

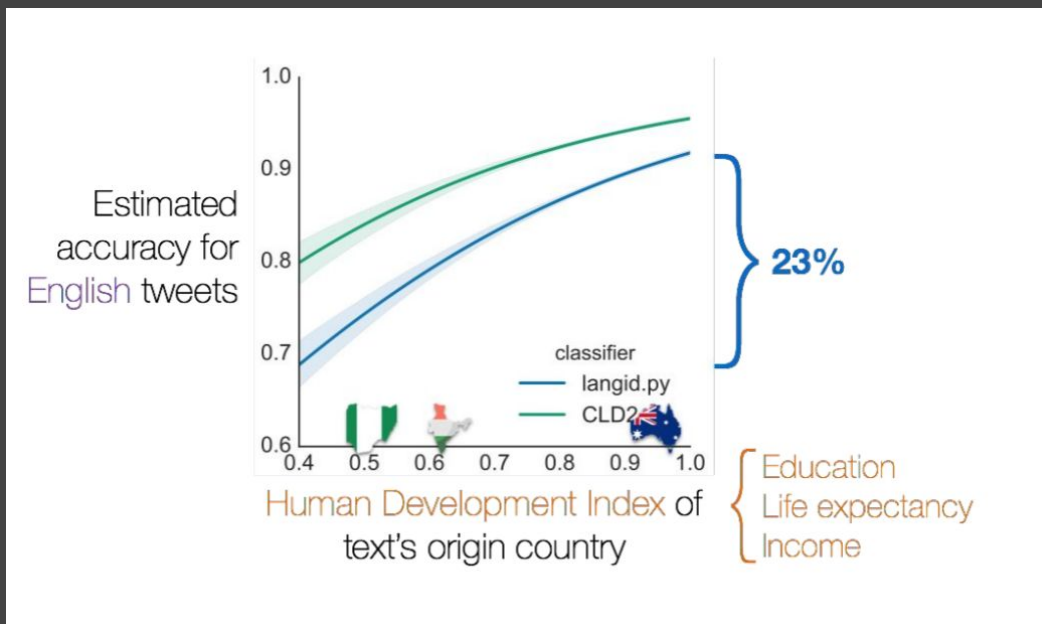
World Englishes



Is the data we use to train our English NLP models representative of all the Englishes out there?

How does this affect NLP models?

Off-the-shelf LID systems under-represent populations in less-developed countries



1M geo-tagged Tweets with any of 385 English terms from established lexicons for *influenza*, *psychological well-being*, and *social health*

i.e.

people who are the most marginalized,
people who'd benefit the most from such technology,
are also the ones who are more likely to be
systematically **excluded** from this technology



Predicting Criminality

Predicting Criminality

Israeli startup, [Faception](#)

*“Faception is first-to-technology and first-to-market with proprietary computer vision and machine learning technology for profiling people and **revealing their personality based only on their facial image.**”*

Offering specialized engines for recognizing “High IQ”, “White-Collar Offender”, “Pedophile”, and “Terrorist” from a face image.

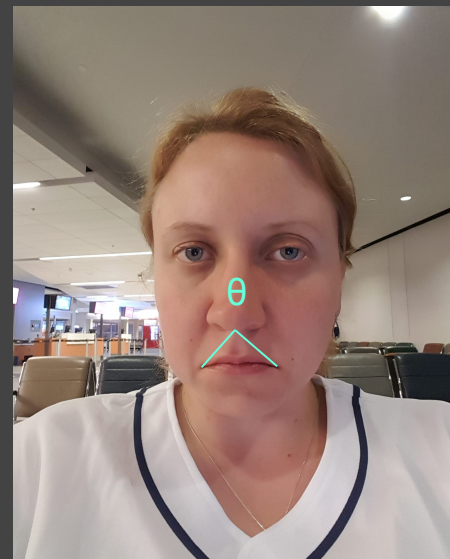
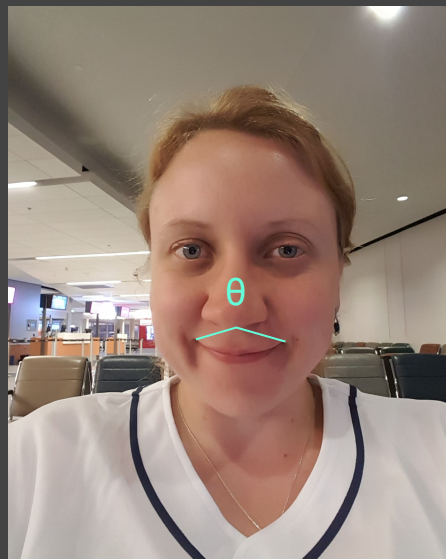
Main clients are in homeland security and public safety.

Predicting Criminality

“[Automated Inference on Criminality using Face Images](#)” Wu and Zhang, 2016.
arXiv

1,856 closely cropped images of faces;
Includes “wanted suspect” ID pictures
from specific regions.

*“[...] angle θ from nose tip to two
mouth corners is on average 19.6%
smaller for criminals than for
non-criminals ...”*



See our longer piece on Medium, “[Physiognomy’s New Clothes](#)”

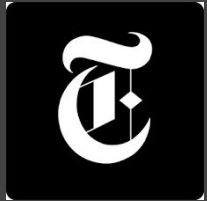


Predicting Toxicity in Text

Toxicity Classification



theguardian



WIKIPEDIA

The
Economist

Source
perspectiveapi.com

We asked the internet what they thought about:

Climate Change Brexit US Election

Showing 46 of 49 total comments based on toxicity*

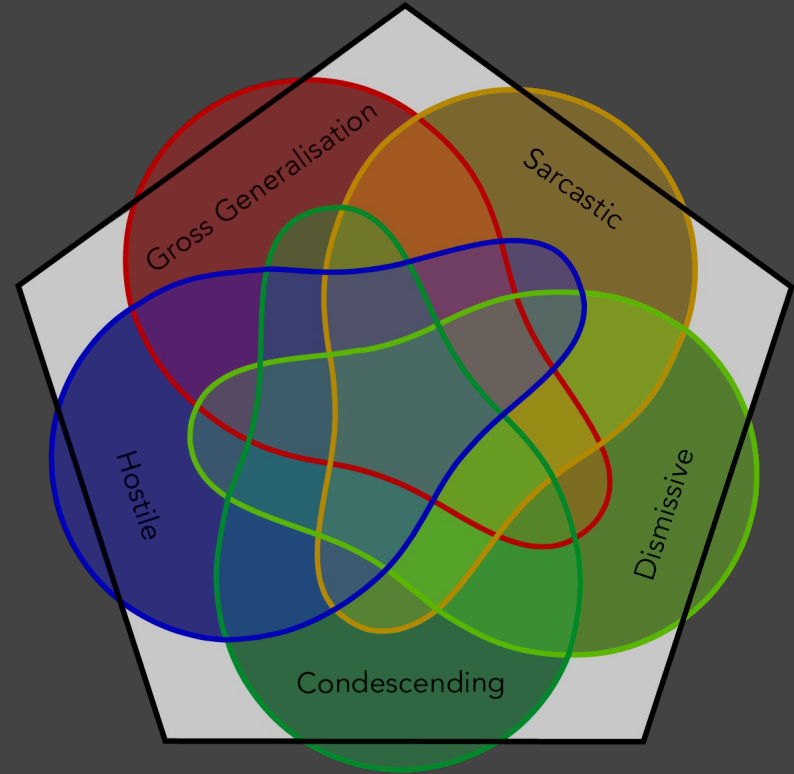
◆ Climate change is happening and it's not changing in our favor. If you think differently you're an idiot.	◆ They're stupid, it's getting warmer, we should enjoy it while it lasts.
◆ They're allowed to do that. But if they act like assholes about, I will block them.	◆ I think those people are stupid and short-sighted
■ uneducated bumpkins or willfully ignorant with vested interests	■ I think its a farce and stinks like a bathroom after 26 beers
■ My thoughts are that people should stop being stupid and ignorant. Climate change is scientifically proven. It isn't a debate.	■ Fools
	■ They are uninformed or ignorant
	■ Their opinion, just don't force it down my

Toxicity Classification

M

The False Positive

Toxicity is defined as... "***a rude, disrespectful, or unreasonable comment that is likely to make you leave a discussion.***"



Source

"The Challenge of Identifying Subtle Forms of Toxicity Online" - Jigsaw

<https://medium.com/the-false-positive/the-challenge-of-identifying-subtle-forms-of-toxicity-online-465505b6c4c9>

Toxicity Classification

Unintended biases towards **certain identity terms**:

Comment	Toxicity Score
The Gay and Lesbian Film Festival starts today.	0.82
Being transgender is independent of sexual orientation.	0.52
A Muslim is someone who follows or practices Islam	0.46

- “The Challenge of Identifying Subtle Forms of Toxicity Online”. Jigsaw.
The False Positive (2018).

Toxicity Classification

Unintended biases towards **named entities**:

Comment

Toxicity Score

I hate Justin Timberlake.

0.90

I hate Rihanna.

0.69

Toxicity Classification

Unintended biases towards **mentions of disabilities**:

Comment

Toxicity Score

I am a person.

0.08

I am a tall person.

0.03

Toxicity Classification

Unintended biases towards **mentions of disabilities**:

Comment	Toxicity Score
I am a person.	0.08
I am a tall person.	0.03
I am a blind person.	0.39
I am a deaf person.	0.44

Toxicity Classification

Unintended biases towards **mentions of disabilities**:

Comment	Toxicity Score
I am a person.	0.08
I am a tall person.	0.03
I am a blind person.	0.39
I am a deaf person.	0.44
I am a person with mental illness.	0.62



NLP Research on Bias and Fairness

Fairness Research in NLP

Slide from SRNLP
Tutorial at NAACL 2018

1. Bolukbasi, T., Saligrama V., Kalai A. (2016) **Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings**. *EMNLP*.
2. Caliskan, A., Bryson, J. J. and Narayanan, A. (2017) **Semantics derived automatically from language corpora contain human-like biases**. *Science*.
3. Nikhil Garg, Londa Schiebinger, Dan Jurafsky, James Zou. (2018) **Word embeddings quantify 100 years of gender and ethnic stereotypes**. *PNAS*.

Fairness Research in NLP

1. Bolukbasi et al. **Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings**. *NIPS* (2016)
2. Caliskan, et al. **Semantics derived automatically from language corpora contain human-like biases**. *Science* (2017)
3. Zhao, Jieyu, et al. **Men also like shopping: Reducing gender bias amplification using corpus-level constraints**. *arXiv* (2017)

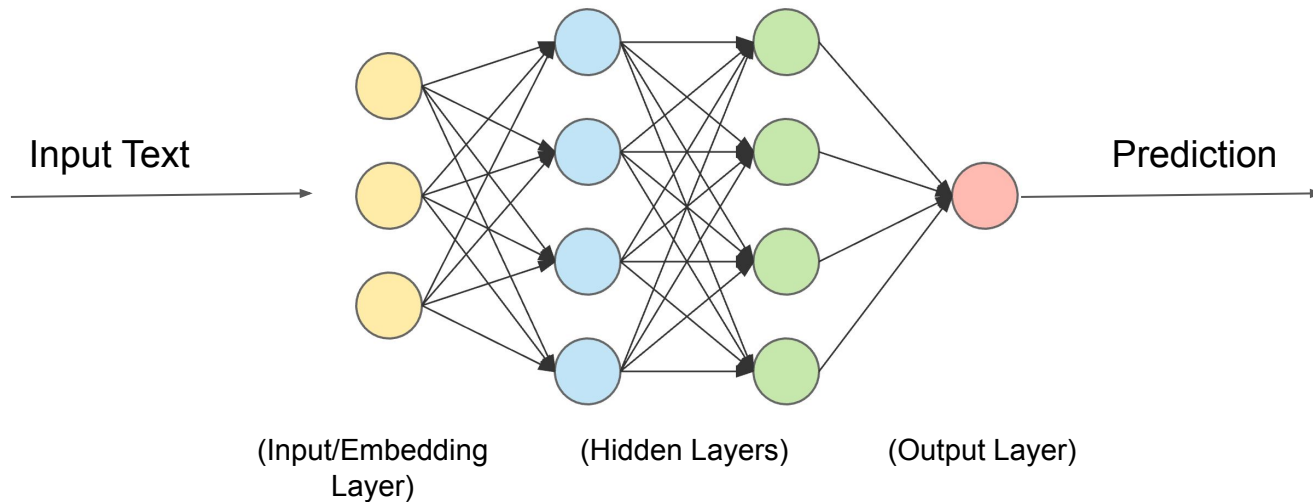
4. Garg et al. **Word embeddings quantify 100 years of gender and ethnic stereotypes**. *PNAS*. (2018)
5. Zhao, Jieyu, et al. **Gender bias in coreference resolution: Evaluation and debiasing methods**. *arXiv* (2018)
6. Zhang, et al. **Mitigating unwanted biases with adversarial learning**. *AIES*, 2018
7. Webster, Kellie, et al. **Mind the GAP: A Balanced Corpus of Gendered Ambiguous Pronouns**. *TACL* (2018)
8. Svetlana and Mohammad. **Examining gender and race bias in two hundred sentiment analysis systems**. *arXiv* (2018)
9. Díaz, et al. **Addressing age-related bias in sentiment analysis**. *CHI Conference on Human Factors in Computing Systems*. (2018)
10. Dixon, et al. **Measuring and mitigating unintended bias in text classification**. *AIES*. (2018)
11. Prates, et al. **Assessing gender bias in machine translation: a case study with Google Translate**. *Neural Computing and Applications* (2018)
12. Park, et al. **Reducing gender bias in abusive language detection**. *arXiv* (2018)
13. Zhao, Jieyu, et al. **Learning gender-neutral word embeddings**. *arXiv* (2018)
14. Anne Hendricks, et al. **Women also snowboard: Overcoming bias in captioning models**. *ECCV*. (2018)
15. Elazar and Goldberg. **Adversarial removal of demographic attributes from text data**. *arXiv* (2018)
16. Hu and Strout. **Exploring Stereotypes and Biased Data with the Crowd**. *arXiv* (2018)

17. Swinger, De-Arteaga, et al. **What are the biases in my word embedding?** *AIES* (2019)
18. De-Arteaga et al. **Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting**. *FAT** (2019)
19. Gonen, et al. **Lipstick on a Pig: Debiasing Methods Cover up Systematic Gender Biases in Word Embeddings But do not Remove Them**. *NAACL* (2019).
20. Manzini et al. **Black is to Criminal as Caucasian is to Police: Detecting and Removing Multiclass Bias in Word Embeddings**. *NAACL* (2019).
21. Sap et al. **The Risk of Racial Bias in Hate Speech Detection**. *ACL* (2019)
22. Stanovsky et al. **Evaluating Gender Bias in Machine Translation**. *ACL* (2019)
23. Garimella et al. **Women's Syntactic Resilience and Men's Grammatical Luck: Gender-Bias in Part-of-Speech Tagging and Dependency Parsing**. *ACL* (2019)
24. ...

2018

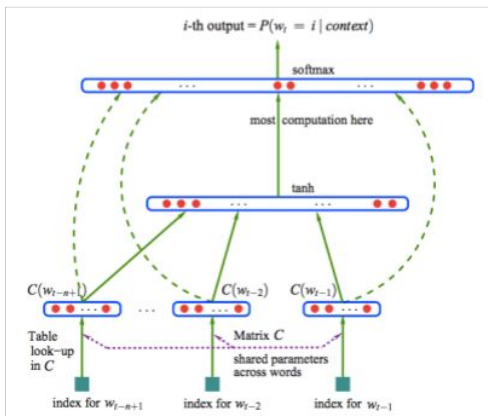
2019

Where to look for biases?

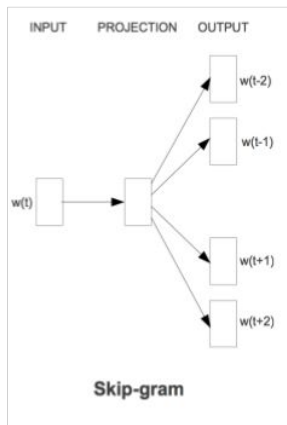


Bias in Input Representations?

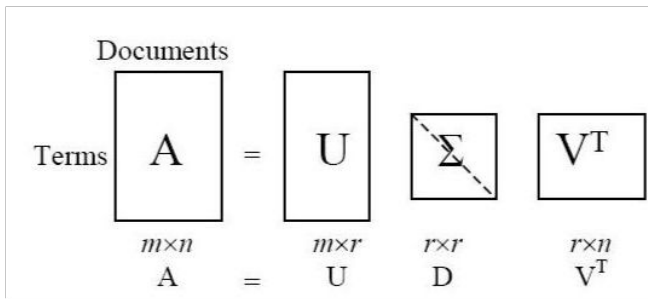
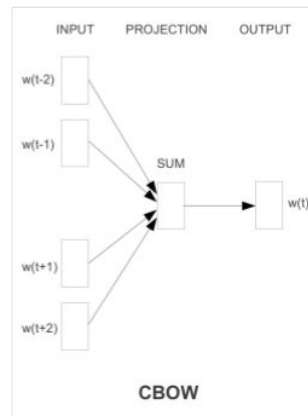
Input Representation: Word Embeddings



Neural Language Model (Bengio et al, '03)

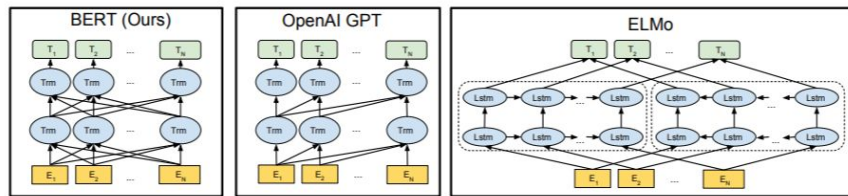


word2vec (Mikolov et al, '03)



Latent Semantic Analysis

(Deerwester et al, '90, Turney & Pantel '10)

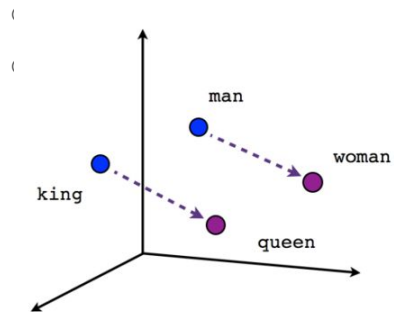


BERT, GPT/GPT-2, ELMo

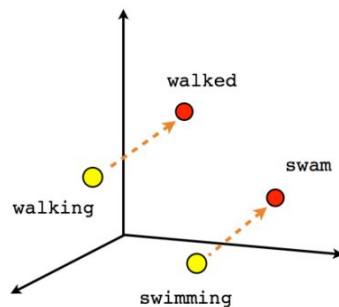
(Devlin et al. '19, Radford et al. '18, Peters et al. '18)

Word Analogy Tasks

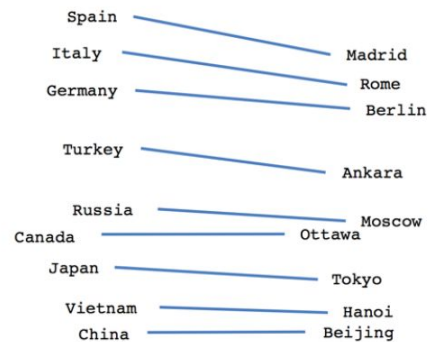
- Mikolov et al. '13



Male-Female



Verb tense



Country-Capital

$$\min \cos(\vec{man} - \vec{woman}, \vec{king} - x) \text{ s.t. } \|\vec{king} - x\|_2 < \delta$$

Social Disparities (and Stereotypes) → Word Embeddings?

He is...



She is...





Bolukbasi et al. **Man is to Computer Programmer as Woman is to Homemaker?**
Debiasing Word Embeddings. *NIPS* (2016)

Biases in NLP Representations

- Bolukbasi et al. **Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings.** *NIPS* (2016)
- Caliskan, et al. **Semantics derived automatically from language corpora contain human-like biases.** *Science* (2017)
- Garg et al. **Word embeddings quantify 100 years of gender and ethnic stereotypes.** *PNAS*. (2018)
- Swinger, De-Arteaga, et al. **What are the biases in my word embedding?** *AIES* (2019)
- Manzini et al. **Black is to Criminal as Caucasian is to Police: Detecting and Removing Multiclass Bias in Word Embeddings.** *NAACL* (2019).
- ...

Implicit bias in humans?

Implicit Association Test - Greenwald et al. 1998

Category	Items
Good	Spectacular, Appealing, Love, Triumph, Joyous, Fabulous, Excitement, Excellent
Bad	Angry, Disgust, Rotten, Selfish, Abuse, Dirty, Hatred, Ugly
African Americans	
European Americans	

Implicit Association Test

The IAT involves making repeated judgments (by pressing a key on a keyboard) to label words or images that pertain to one of two categories presented simultaneously (e.g., categorizing pictures of African American or European American and categorizing positive/negative adjectives).

The test compares response times when different pairs of categories share a **response key** on keyboard

(e.g., African American + GOOD vs African American + BAD vs European American + GOOD vs European American + BAD)

IAT - Societal groups ↔ Stereotype words

Disability IAT

Disability ('Disabled - Abled' IAT). This IAT requires the ability to recognize symbols representing abled and disabled individuals.

Asian IAT

Asian American ('Asian - European American' IAT). This IAT requires the ability to recognize White and Asian-American faces, and images of places that are either American or Foreign in origin.

Sexuality IAT

Sexuality ('Gay - Straight' IAT). This IAT requires the ability to distinguish words and symbols representing gay and straight people. It often reveals an automatic preference for straight relative to gay people.

Arab-Muslim IAT

Arab-Muslim ('Arab Muslim - Other People' IAT). This IAT requires the ability to distinguish names that are likely to belong to Arab-Muslims versus people of other nationalities or religions.

Age IAT

Age ('Young - Old' IAT). This IAT requires the ability to distinguish old from young faces. This test often indicates that Americans have automatic preference for young over old.

Skin-tone IAT

Skin-tone ('Light Skin - Dark Skin' IAT). This IAT requires the ability to recognize light-skinned faces. It often reveals an automatic preference for light-skin relative to dark-skinned faces.

Race IAT

Race ('Black - White' IAT). This IAT requires the ability to distinguish faces of African origin. It indicates that most Americans have an automatic preference for White faces.

Religion IAT

Religion ('Religions' IAT). This IAT requires some familiarity with religious terms from various world religions.

Native IAT

Native American ('Native - White American' IAT). This IAT requires the ability to recognize White and Native American faces in either classic or modern dress, and the names of places that are either American or Foreign in origin.

Gender-Science IAT

Gender - Science. This IAT often reveals a relative link between liberal arts and females and between science and males.

Gender-Career IAT

Gender - Career. This IAT often reveals a relative link between family and females and between career and males.

Presidents IAT

Presidents ('Presidential Popularity' IAT). This IAT requires the ability to recognize photos of Donald Trump and one or more previous presidents.

Weight IAT

Weight ('Fat - Thin' IAT). This IAT requires the ability to distinguish faces of people who are obese and people who are thin. It often reveals an automatic preference for thin people relative to fat people.

Weapons IAT

Weapons ('Weapons - Harmless Objects' IAT). This IAT requires the ability to recognize White and Black faces, and images of weapons or harmless objects.

<https://implicit.harvard.edu/implicit/selectatest.html>

Greenwald et al. 1998

Can we apply this to NLP models?

IAT for Word Embeddings

- Word Embedding Association Test (WEAT)
 - Latency \Leftrightarrow Cosine similarity
 - Target words
 - $X = \{programmer, engineer, scientist, \dots\}$
 - $Y = \{nurse, teacher, librarian, \dots\}$
 - Attribute words
 - $A = \{man, male, \dots\}$
 - $B = \{woman, female, \dots\}$

Word Embedding Association Test

- Target words
 - $X = \{programmer, engineer, scientist, \dots\}$
 - $Y = \{nurse, teacher, librarian, \dots\}$
- Attribute words
 - $A = \{man, male, \dots\}$
 - $B = \{woman, female, \dots\}$

$$s(w, A, B) = \text{mean}_{a \in A} \cos(\vec{w}, \vec{a}) - \text{mean}_{b \in B} \cos(\vec{w}, \vec{b})$$

Association of a word w with an attribute:

$$s(X, Y, A, B) = \sum_{x \in X} s(x, A, B) - \sum_{y \in Y} s(y, A, B)$$

Association of two sets $\frac{\text{mean}_{x \in X} s(x, A, B) - \text{mean}_{y \in Y} s(y, A, B)}{\text{std-dev}_{w \in X \cup Y} s(w, A, B)}$

The effect size of bias:

Additional statistical tests to measure how separated are two distributions and statistical significance

Word Embedding Association Test

$$s(w, A, B) = \frac{\text{mean}_{a \in A} \cos(\vec{w}, \vec{a}) - \text{mean}_{b \in B} \cos(\vec{w}, \vec{b})}{\text{std-dev}_{x \in A \cup B} \cos(\vec{w}, \vec{x})}$$

- **Flowers:** aster, clover, hyacinth, marigold, poppy, azalea, crocus, iris, orchid, rose, bluebell, daffodil, lilac, pansy, tulip, buttercup, daisy, lily, peony, violet, carnation, gladiola, magnolia, petunia, zinnia.
- **Insects:** ant, caterpillar, flea, locust, spider, bedbug, centipede, fly, maggot, tarantula, bee, cockroach, gnat, mosquito, termite, beetle, cricket, hornet, moth, wasp, blackfly, dragonfly, horsefly, roach, weevil.
- **Pleasant:** caress, freedom, health, love, peace, cheer, friend, heaven, loyal, pleasure, diamond, gentle, honest, lucky, rainbow, diploma, gift, honor, miracle, sunrise, family, happy, laughter, paradise, vacation.
- **Unpleasant:** abuse, crash, filth, murder, sickness, accident, death, grief, poison, stink, assault, disaster, hatred, pollute, tragedy, divorce, jail, poverty, ugly, cancer, kill, rotten, vomit, agony, prison.

Word Embedding Association Test: Results

IAT

WEAT

Target words	Attrib. words	Original Finding				Our Finding			
		Ref	N	d	p	N _T	N _A	d	p
Flowers vs insects	Pleasant vs unpleasant	(5)	32	1.35	10^{-8}	25×2	25×2	1.50	10^{-7}

Word Embedding Association Test

- **European American names:** Adam, *Chip*, Harry, Josh, Roger, Alan, Frank, *Ian*, Justin, Ryan, Andrew, *Fred*, Jack, Matthew, Stephen, Brad, Greg, *Jed*, Paul, *Todd*, *Brandon*, *Hank*, Jonathan, Peter, *Wilbur*, Amanda, Courtney, Heather, Melanie, *Sara*, *Amber*, *Crystal*, Katie, *Meredith*, *Shannon*, Betsy, *Donna*, Kristin, Nancy, Stephanie, *Bobbie-Sue*, Ellen, Lauren, *Peggy*, *Sue-Ellen*, Colleen, Emily, Megan, Rachel, *Wendy* (deleted names in italics).
- **African American names:** Alonzo, Jamel, *Lerone*, *Percell*, Theo, Alphonse, Jerome, Leroy, *Rasaan*, Torrance, Darnell, Lamar, Lionel, *Rashaun*, Tivree, Deion, Lamont, Malik, Terrence, Tyrone, *Everol*, Lavon, Marcellus, *Terryl*, Wardell, *Aiesha*, *Lashelle*, Nichelle, Shereen, *Temeka*, Ebony, Latisha, Shaniqua, *Tameisha*, *Teretha*, Jasmine, *Latonya*, *Shanise*, Tanisha, Tia, Lakisha, Latoya, *Sharise*, *Tashika*, Yolanda, *Lashandra*, Malika, *Shavonn*, *Tawanda*, Yvette (deleted names in italics).
- **Pleasant:** caress, freedom, health, love, peace, cheer, friend, heaven, loyal, pleasure, diamond, gentle, honest, lucky, rainbow, diploma, gift, honor, miracle, sunrise, family, happy, laughter, paradise, vacation.
- **Unpleasant:** abuse, crash, filth, murder, sickness, accident, death, grief, poison, stink, assault, disaster, hatred, pollute, tragedy, bomb, divorce, jail, poverty, ugly, cancer, evil, kill, rotten, vomit.

Word Embedding Association Test: Results

		IAT				<u>WEAT</u>			
Target words	Attrib. words	Original Finding				Our Finding			
		Ref	N	d	p	N _T	N _A	d	p
Eur.-American vs Afr.-American names	Pleasant vs unpleasant	(5)	26	1.17	10^{-5}	32×2	25×2	1.41	10^{-8}

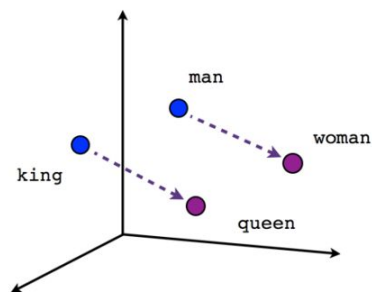
WEAT finds similar biases in Word Embeddings as IAT did for humans

Other ways to detect biases?

Gender Bias in Word Embeddings

$$\vec{\text{man}} - \vec{\text{woman}} \approx \vec{\text{computer programmer}} - \vec{\text{homemaker}}$$

$$\min \cos(\text{he} - \text{she}, x - y) \text{ s.t. } \|x - y\|_2 < \delta$$



Male-Female

surgeon vs. nurse

architect vs. interior designer

shopkeeper vs. housewife

superstar vs. diva

....

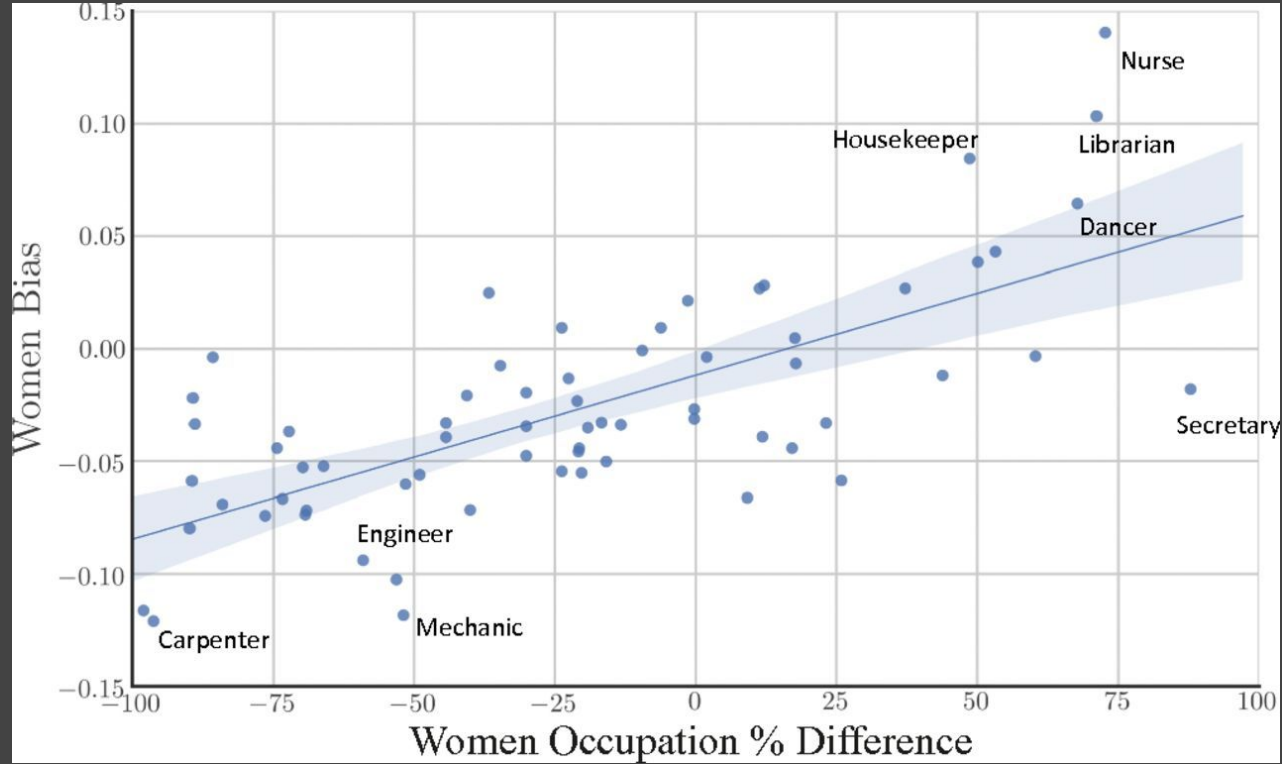
Beyond Gender & Race/Ethnicity Bias

Gender Biased Analogies	
man → doctor	woman → nurse
woman → receptionist	man → supervisor
woman → secretary	man → principal
Racially Biased Analogies	
black → criminal	caucasian → police
asian → doctor	caucasian → dad
caucasian → leader	black → led
Religiously Biased Analogies	
muslim → terrorist	christian → civilians
jewish → philanthropist	christian → stooge
christian → unemployed	jewish → pensioners

Biases in word embeddings trained on the Reddit data from US users.

But aren't they just reflecting Society?

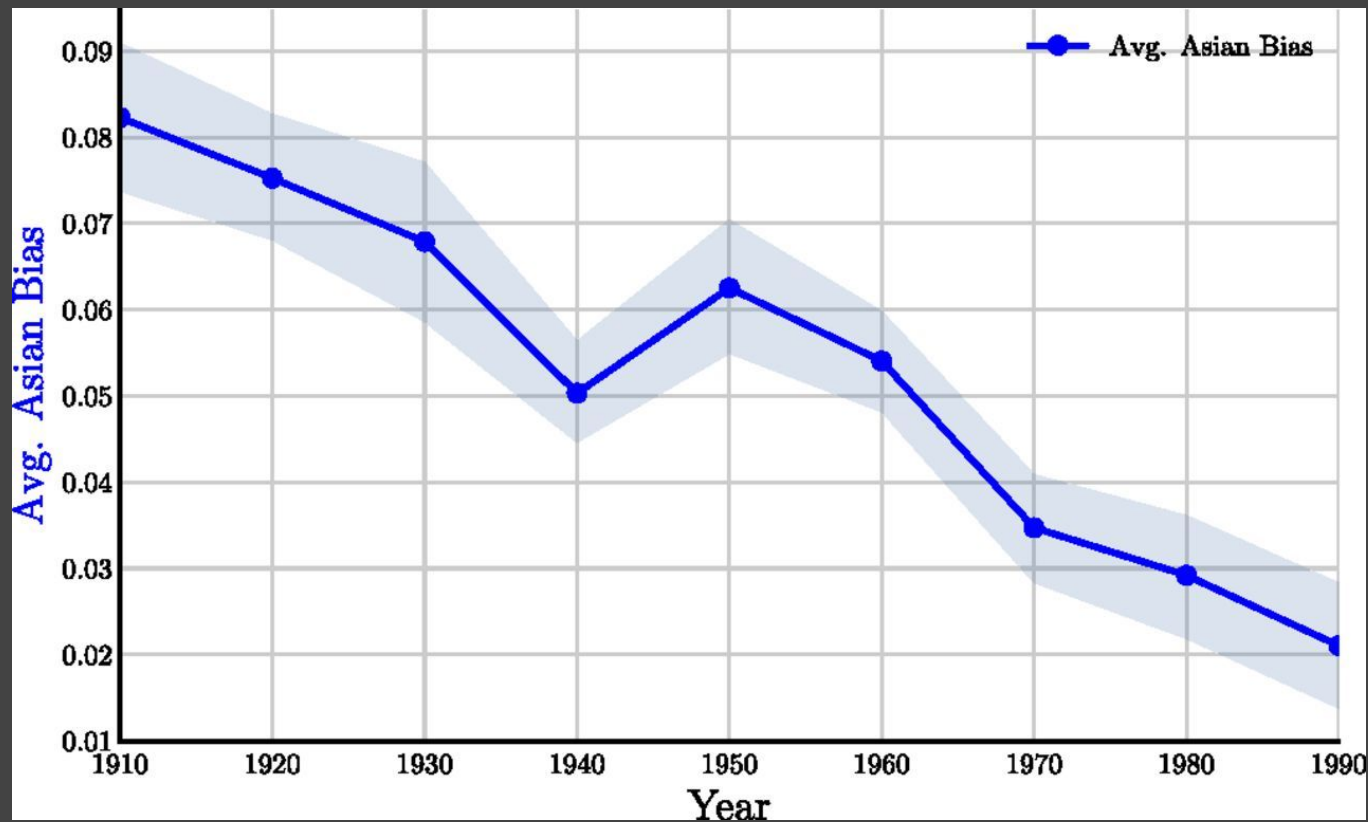
Gender bias in Occupations



Gender bias in Adjectives over the decades



“Asian bias” in Adjectives with “Outsider” words



But aren't they just reflecting Society?

Yup!

Word embeddings ...



... get things
normatively wrong
precisely because they
get things
descriptively right!

Shouldn't we then just leave them as is?

Shouldn't we then just leave them as is?

Would that harm certain groups of people?

What kind of harm?

Associative Harm

“when systems reinforce the subordination of some groups along the lines of identity”

Allocative Harm

“when a system allocates or withholds a certain opportunity or resource”

Amazon's Secret AI Hiring Tool Reportedly 'Penalized' Resumes With the Word 'Women's'



Rhett Jones

Yesterday 10:32am • Filed to: ALGORITHMS ▾

22.3K

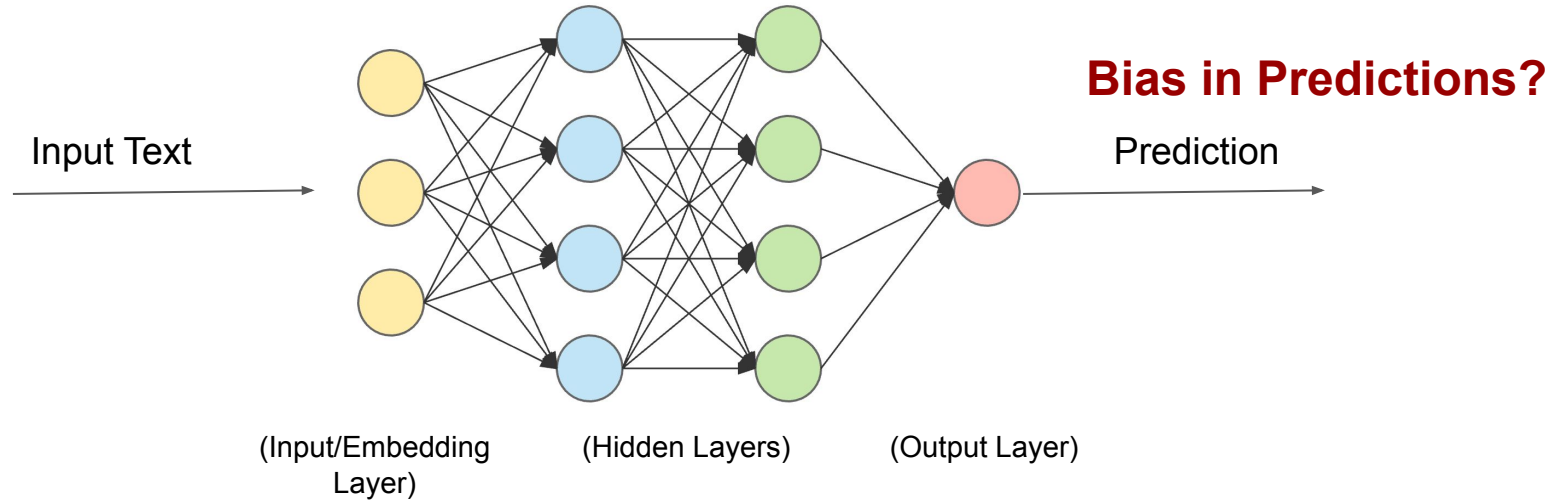
96

2



Photo: Getty

Where to look for biases?



Bias in Input Representations?

Biases in NLP Classifiers/Taggers

- Gender Bias in Part of speech tagging and Dependency parsing
 - Garimella et al. **Women's Syntactic Resilience and Men's Grammatical Luck: Gender-Bias in Part-of-Speech Tagging and Dependency Parsing**. ACL (2019)
- Gender Bias in Coreference resolution
 - Zhao, Jieyu, et al. **Gender bias in coreference resolution: Evaluation and debiasing methods**. *arXiv* (2018)
 - Webster, Kellie, et al. **Mind the GAP: A Balanced Corpus of Gendered Ambiguous Pronouns**. *TACL* (2018)
- Gender, Race, and Age Bias in Sentiment Analysis
 - Svetlana and Mohammad. **Examining gender and race bias in two hundred sentiment analysis systems**. *arXiv* (2018)
 - Díaz, et al. **Addressing age-related bias in sentiment analysis**. CHI Conference on Human Factors in Comp. Systems. (2018)
- LGBTQ identity terms bias in Toxicity classification
 - Dixon, et al. **Measuring and mitigating unintended bias in text classification**. AIES. (2018)
 - Sap, et al. **The Risk of Racial Bias in Hate Speech Detection**. ACL. (2019)
- Gender Bias in Occupation Classification
 - De-Arteaga et al. **Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting**. FAT* (2019)
- Gender bias in Machine Translation
 - Prates, et al. **Assessing gender bias in machine translation: a case study with Google Translate**. Neural Computing and Applications (2018)

Shouldn't we then just leave them as is?

Would that harm certain groups of people?

Would that make things worse?

Bias Amplification

- Zhao et al. **Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraint.** *EMNLP* (2017)
- *De-Arteaga et al.* **Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting.** *FAT** (2019)

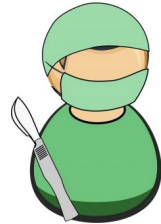
Examples of Harm from NLP Bias

An artificially intelligent headhunter?



Examples of Harm from NLP Bias

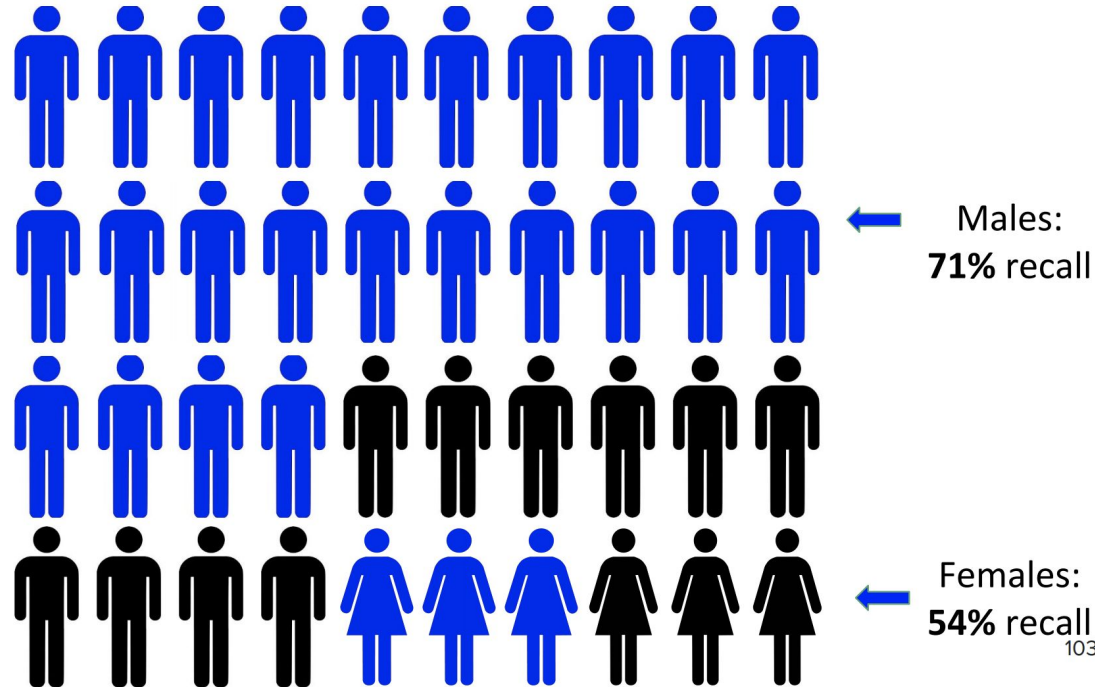
Compounding imbalances



Surgeons

females in data:
14.6%

females in true positives:
11.6%



Ok, How do we make NLP models fair?

What does it mean to be Fair?

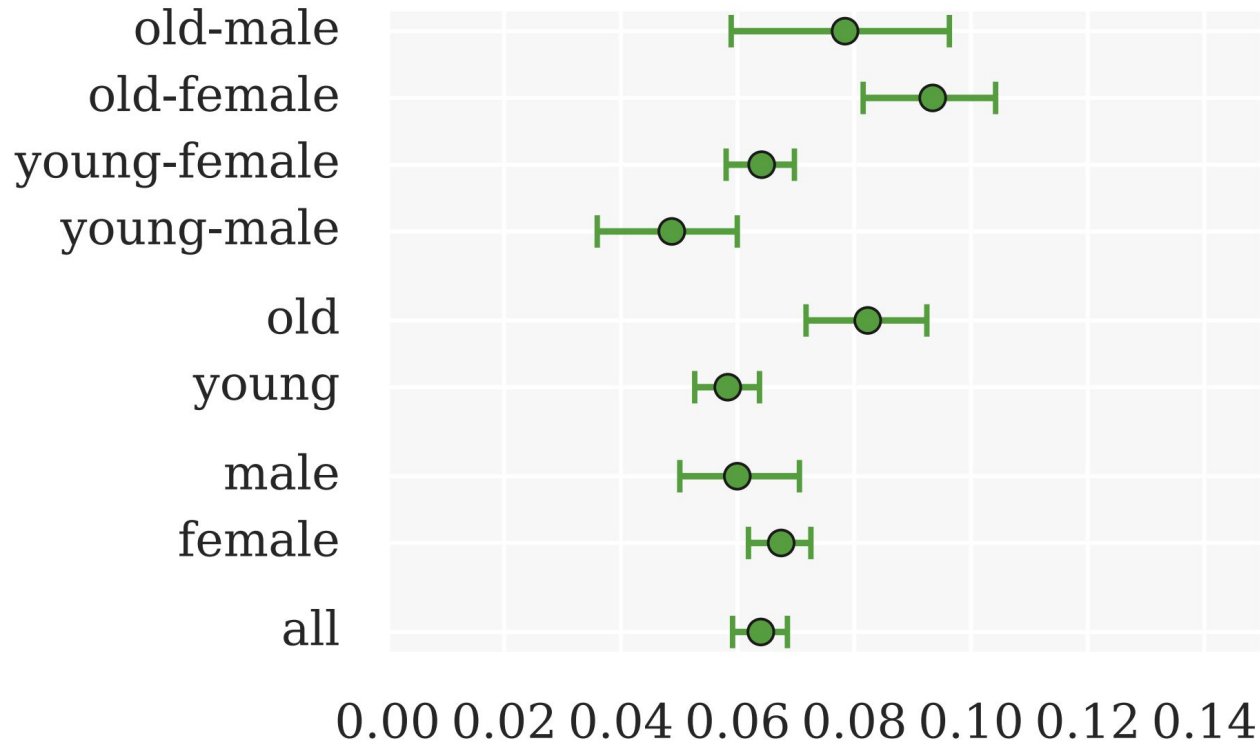
Different Types of Fairness

- Group Fairness
 - “treat different groups equally”
 - E.g., demographic parity across groups (along age, gender, race, etc.)

- Individual Fairness
 - “treat similar examples similarly”
 - E.g., counterfactual fairness (if we switch the gender, does the prediction change?)

Group Fairness

False Positive Rate @ 0.5



Individual Fairness

```
text_to_sentiment("My name is Emily")
```

```
2.2286179364745311
```

```
text_to_sentiment("My name is Heather")
```

```
1.3976291151079159
```

```
text_to_sentiment("My name is Yvette")
```

```
0.9846380213298556
```

```
text_to_sentiment("My name is Shaniqua")
```

```
-0.47048131775890656
```

<http://blog.conceptnet.io/posts/2017/how-to-make-a-racist-ai-without-really-trying/>



Measuring Algorithmic Fairness/Bias

Evaluate for Fairness & Inclusion

Disaggregated Evaluation

Create for each (subgroup, prediction) pair.

Compare across subgroups.

Evaluate for Fairness & Inclusion

Disaggregated Evaluation

Create for each (subgroup, prediction) pair.
Compare across subgroups.

Example: women, face detection
men, face detection

Evaluate for Fairness & Inclusion

Intersectional Evaluation

Create for each (subgroup1, subgroup2, prediction) pair. Compare across subgroups.

Example: black women, face detection
white men, face detection



Evaluate for Fairness & Inclusion

Female Patient Results

True Positives (TP) = 10	False Positives (FP) = 1
False Negatives (FN) = 1	True Negatives (TN) = 488

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{10}{10 + 1} = 0.909$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{10}{10 + 1} = 0.909$$

Male Patient Results

True Positives (TP) = 6	False Positives (FP) = 3
False Negatives (FN) = 5	True Negatives (TN) = 48

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{6}{6 + 3} = 0.667$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{6}{6 + 5} = 0.545$$

Evaluate for Fairness & Inclusion

Female Patient Results

True Positives (TP) = 10 False Positives (FP) = 1

False Negatives (FN) = 1 True Negatives (TN) = 488

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{10}{10 + 1} = 0.909$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{10}{10 + 1} = 0.909$$

Male Patient Results

True Positives (TP) = 6 False Positives (FP) = 3

False Negatives (FN) = 5 True Negatives (TN) = 48

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{6}{6 + 3} = 0.667$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{6}{6 + 5} = 0.545$$

**“Equality of Opportunity” fairness criterion:
Recall is equal across subgroups**

Evaluate for Fairness & Inclusion

Female Patient Results

True Positives (TP) = 10	False Positives (FP) = 1
False Negatives (FN) = 1	True Negatives (TN) = 488

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{10}{10 + 1} = 0.909$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{10}{10 + 1} = 0.909$$

Male Patient Results

True Positives (TP) = 6	False Positives (FP) = 3
False Negatives (FN) = 5	True Negatives (TN) = 48

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{6}{6 + 3} = 0.667$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{6}{6 + 5} = 0.545$$

**“Predictive Parity” fairness criterion:
Precision is equal across subgroups**

Choose your evaluation metrics in light
of acceptable tradeoffs between
False Positives and **False Negatives**

False Positives Might be Better than False Negatives

Privacy in Images

False Positive: Something that doesn't need to be blurred gets blurred.

Can be a bummer.



False Negative: Something that needs to be blurred is not blurred.

Identity theft.



False Negatives Might Be Better than False Positives

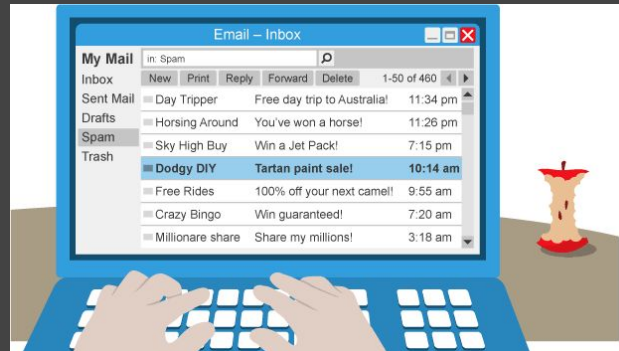
Spam Filtering

False Negative: Email that is SPAM is not caught, so you see it in your inbox.

False Positive: Email flagged as SPAM is removed from your inbox.

Usually just a bit annoying.

If it is an interview call?



Can we computationally remove
undesirable biases?

- Debiasing Meaning Representations

Methods to “de-bias” NLP models

- Gender De-Biasing

- Bolukbasi et al. **Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings.** *NIPS* (2016)
- Zhao, Jieyu, et al. **Men also like shopping: Reducing gender bias amplification using corpus-level constraints.** arXiv (2017)
- Park, et al. **Reducing gender bias in abusive language detection.** arXiv (2018)
- Zhao, Jieyu, et al. **Learning gender-neutral word embeddings.** arXiv (2018)
- Anne Hendricks, et al. **Women also snowboard: Overcoming bias in captioning models.** ECCV. (2018)

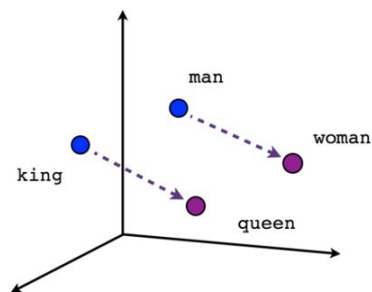
- General De-Biasing

- Beutel et al. **Data Decisions and Theoretical Implications when Adversarially Learning Fair Representations.** FATML (2017)
- Zhang, et al. **Mitigating unwanted biases with adversarial learning.** AIES, 2018
- Elazar and Goldberg. **Adversarial removal of demographic attributes from text data.** arXiv (2018)
- Hu and Strout. **Exploring Stereotypes and Biased Data with the Crowd.** arXiv (2018)

Gender Bias in Word Embeddings

$$\vec{\text{man}} - \vec{\text{woman}} \approx \vec{\text{computer programmer}} - \vec{\text{homemaker}}$$

$$\min \cos(\text{he} - \text{she}, x - y) \text{ s.t. } \|x - y\|_2 < \delta$$



Male-Female

surgeon vs. nurse

architect vs. interior designer

shopkeeper vs. housewife

superstar vs. diva

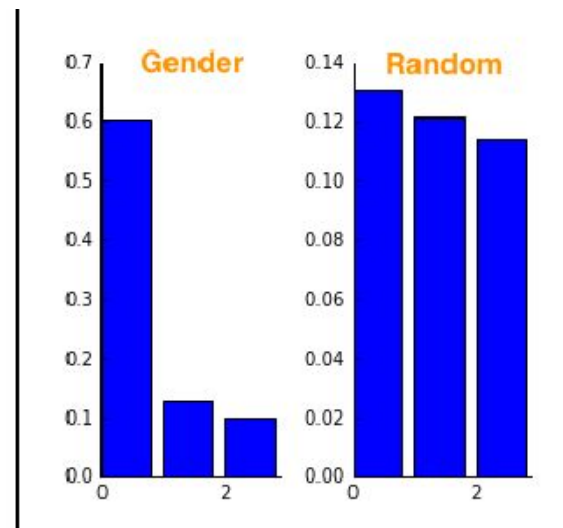
....

Towards Debiasing

1. Identify gender subspace: B

Gender Subspace

$\vec{\text{she}} - \vec{\text{he}}$
 $\vec{\text{her}} - \vec{\text{his}}$
 $\vec{\text{woman}} - \vec{\text{man}}$
 $\vec{\text{Mary}} - \vec{\text{John}}$
 $\vec{\text{herself}} - \vec{\text{himself}}$
 $\vec{\text{daughter}} - \vec{\text{son}}$
 $\vec{\text{mother}} - \vec{\text{father}}$
 $\vec{\text{gal}} - \vec{\text{gu\u00fd}}$
 $\vec{\text{girl}} - \vec{\text{boy}}$
 $\vec{\text{female}} - \vec{\text{male}}$

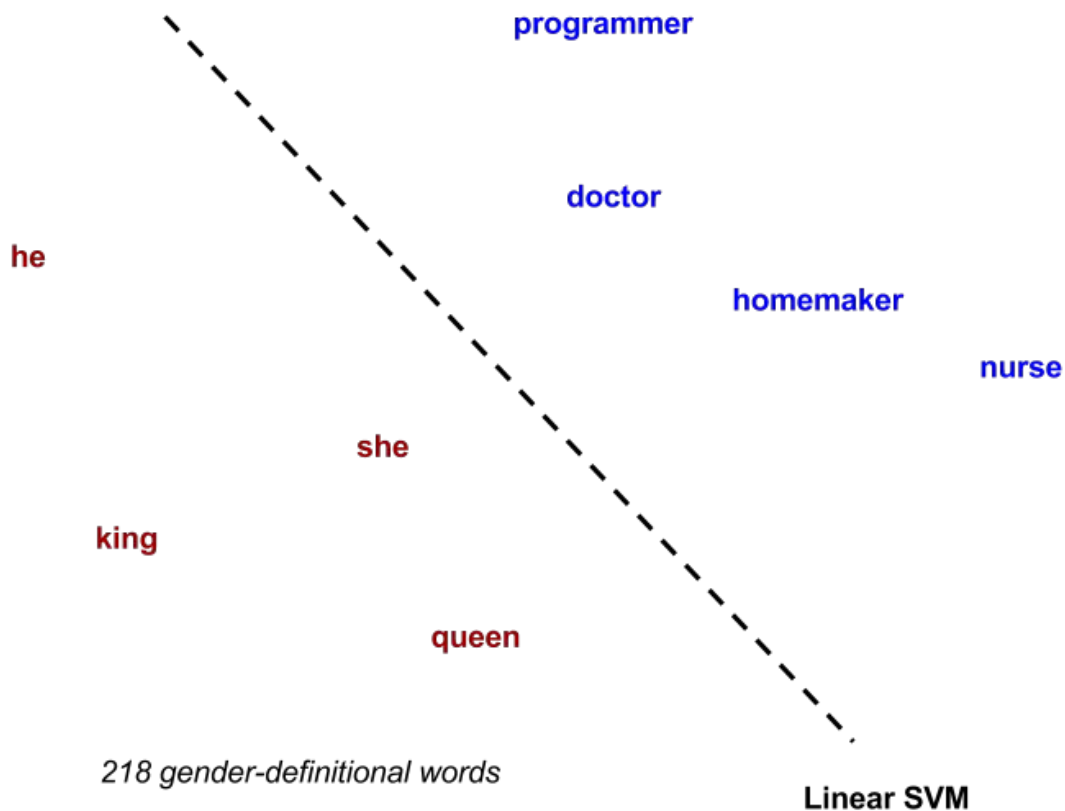


The top PC captures the gender subspace

Towards Debiasing

1. Identify gender subspace: B
2. **Identify gender-definitional (S) and gender-neutral words (N)**

Gender-definitional vs. Gender-neutral Words



Plus
Bootstrapping

Towards Gender Debiasing

1. Identify gender subspace: B
2. Identify gender-definitional (S) and gender-neutral words (N)

Towards Gender Debiasing

1. Identify gender subspace: B
2. Identify gender-definitional (S) and gender-neutral words (N)
3. Apply transform matrix (T) to the embedding matrix (W) such that
 - a. Project away the gender subspace B from the gender-neutral words N
 - b. But, ensure the transformation doesn't change the embeddings too much

$$\min_T \underbrace{\| (TW)^T (TW) - W^T W \|_F^2}_{\text{Don't modify embeddings too much}} + \lambda \underbrace{\| (TN)^T (TB) \|_F^2}_{\text{Minimize gender component}}$$

T - the desired debiasing transformation
 W - embedding matrix

B - biased space
 N - embedding matrix of gender neutral words

Can we computationally remove
undesirable biases?

- Debiasing Meaning Representations
 - Debiasing Model Predictions

Debiasing using Adversarial Learning

Bias Mitigation

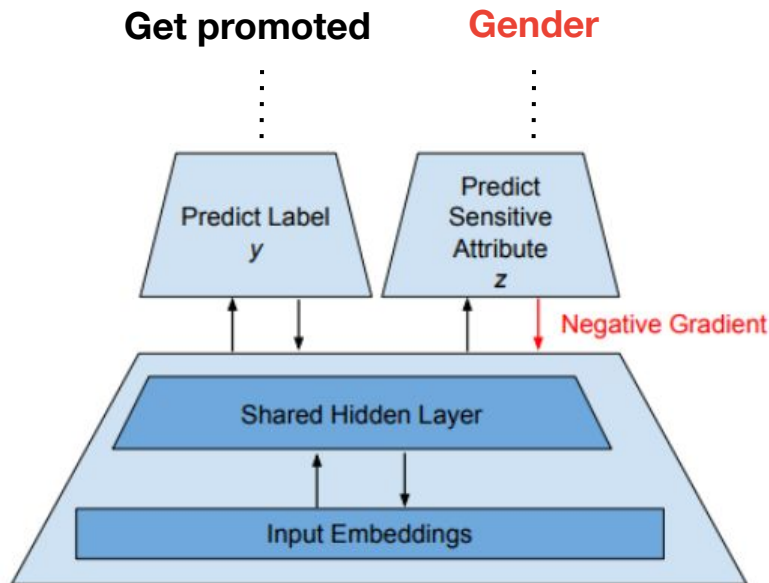
- Handling biased predictions
- Removing signal for problematic variables
 - Stereotyping
 - Sexism, Racism, *-ism

Debiasing using Adversarial Learning

Bias Mitigation

- Handling biased predictions
- Removing signal for problematic variables
 - Stereotyping
 - Sexism, Racism, *-ism

Adversarial Multi-task Learning



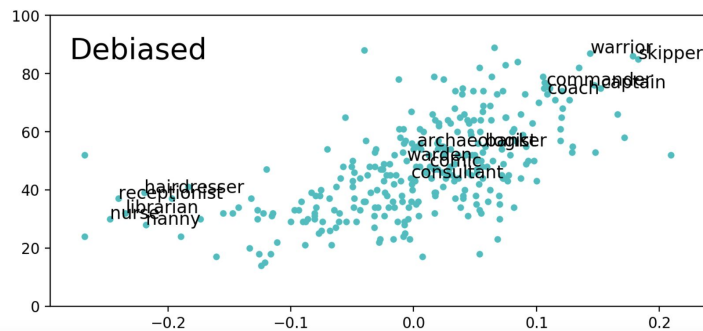
Can we computationally remove
undesirable biases?

YES!

Are we done?

Issues with relying entirely on ‘debiasing’

- Gonen, et al. **Lipstick on a Pig: Debiasing Methods Cover up Systematic Gender Biases in Word Embeddings But do not Remove Them.** NAACL (2019).



So...

What should we do?

Can we **computationally** remove
undesirable biases?

Critically **examine** cases where we
categorize humans

Towards a Critical Race Methodology in Algorithmic Fairness

Alex Hanna*

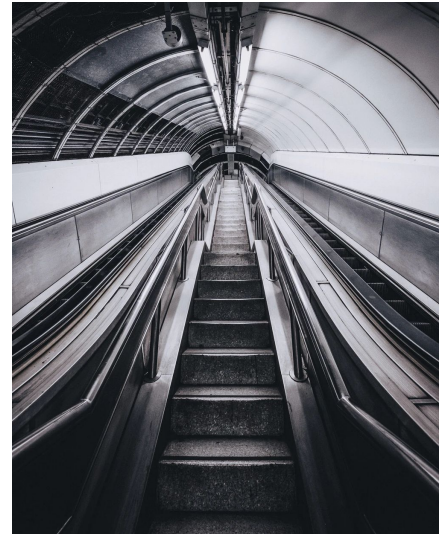
Emily Denton*

Andrew Smart

Jamila Smith-Loud

{alexhanna,dentone,andrewsmart,jsmithloud}@google.com

**Acknowledging the hierarchical,
stratified nature of racial groups**



Towards a Critical Race Methodology in Algorithmic Fairness

Alex Hanna*

Emily Denton*

Andrew Smart

Jamila Smith-Loud

{alexhanna,dentone,andrewsmart,jsmithloud}@google.com

Centering the process of
conceptualizing and
operationalizing race

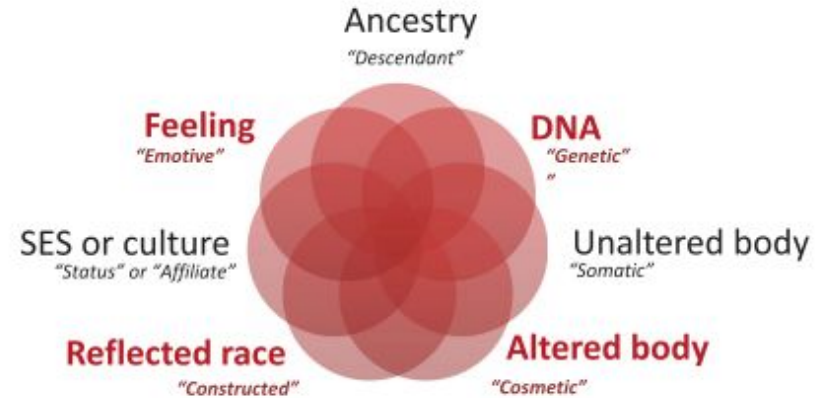


Figure 1. Core and periphery: Claimed attributes and “types” of race member.

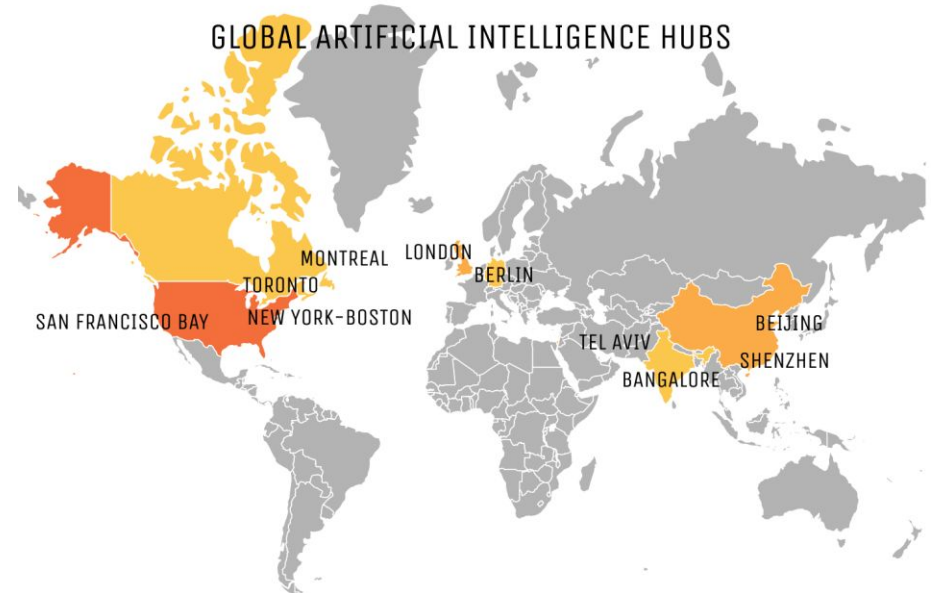
Note: New bases and types of racial membership appear in bold.

Morning. 2018. "[Kaleidoscope: contested identities and new forms of race membership.](#)" *Ethnic and Racial Studies*.

Acknowledge meta issues:

- Lack of **stakeholder perspectives**
 - Lack of **global notions** of value systems or injustices
-

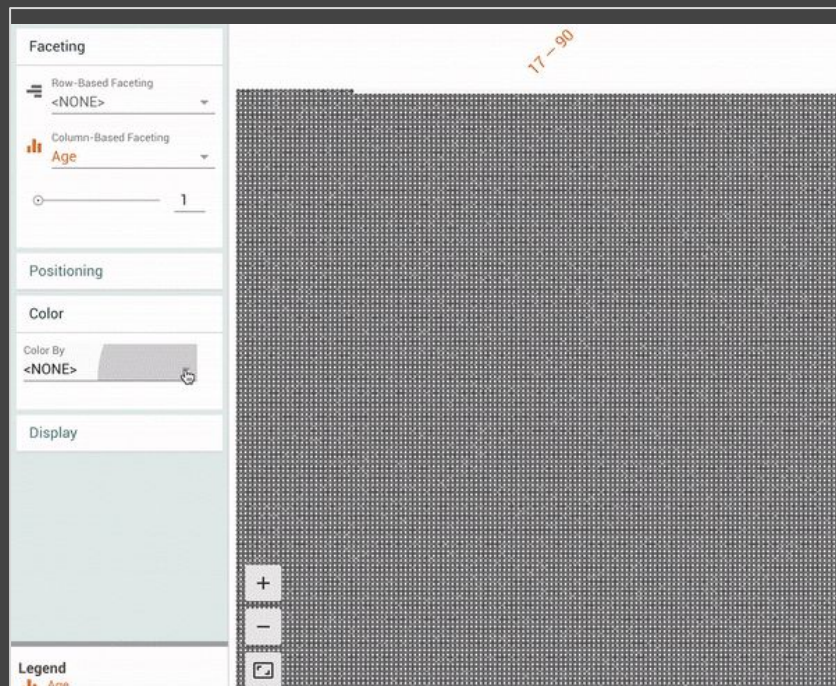
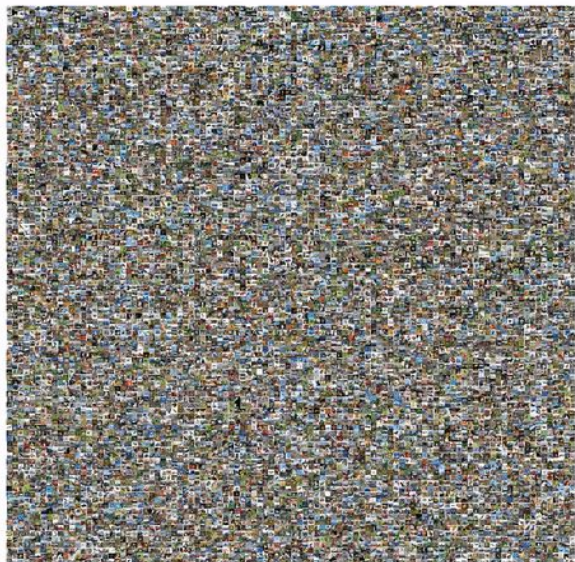
- **Who** is answering these questions?
- **What data** is used to study and answer these questions?
- **Whose value systems** inform interventions?





Data Really, Really Matters

Understand Your Data Skews



Datasheets for Datasets

**Timnit Gebru¹ Jamie Morgenstern² Briana Vecchione³ Jennifer Wortman Vaughan¹ Hanna Wallach¹
Hal Daumé III¹⁴ Kate Crawford¹⁵**

Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science

Emily M. Bender
Department of Linguistics
University of Washington
ebender@uw.edu

Batya Friedman
The Information School
University of Washington
batya@uw.edu

Datasheets for Datasets

Motivation for Dataset Creation

Why was the dataset created? (e.g., were there specific tasks in mind, or a specific gap that needed to be filled?)

What (other) tasks could the dataset be used for? Are there obvious tasks for which it should *not* be used?

Has the dataset been used for any tasks already? If so, where are the results so others can compare (e.g., links to published papers)?

Who funded the creation of the dataset? If there is an associated grant, provide the grant number.

Any other comments?

Dataset Composition

What are the instances? (that is, examples; e.g., documents, images, people, countries) Are there multiple types of instances? (e.g., movies, users, ratings; people, interactions between them; nodes, edges)

Are relationships between instances made explicit in the data (e.g., social network links, user/movie ratings, etc.)?

How many instances of each type are there?

Data Collection Process

How was the data collected? (e.g., hardware apparatus/sensor, manual human curation, software program, software interface/API; how were these constructs/measures/methods validated?)

Who was involved in the data collection process? (e.g., students, crowdworkers) How were they compensated? (e.g., how much were crowdworkers paid?)

Over what time-frame was the data collected? Does the collection time-frame match the creation time-frame?

How was the data associated with each instance acquired? Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part of speech tags; model-based guesses for age or language)? If the latter two, were they validated/verified and if so how?

Does the dataset contain all possible instances? Or is it, for instance, a sample (not necessarily random) from a larger set of instances?

If the dataset is a sample, then what is the population? What was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)? Is the sample representative of the larger set (e.g., geographic coverage)? If not, why not (e.g., to cover a more diverse range of instances)? How does this affect possible uses?

Dataset Fact Sheet

Metadata



Title COMPAS Recidivism Risk Score Data

Author Broward County Clerk's Office, Broward County Sheriff's Office, Florida

Email browardcounty@florida.usa

Description Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat.

DOI 10.5281/zenodo.1164791

Time Feb 2013 - Dec 2014

Keywords risk assessment, parole, jail, recidivism, law

Records 7214

Variables 25

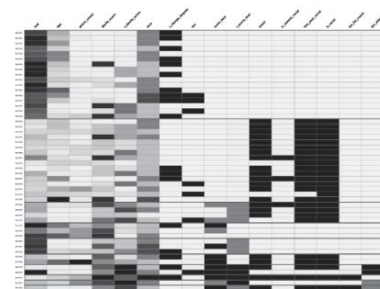
priors_count: *Ut enim ad minim veniam, quis nostrud exercitation* numerical

two_year_prior: *Lorem ipsum dolor sit amet conseq*

Probabilistic Modeling

Analysis

12





Dependency Probability **Pearson R**






Release Your Models Responsibly



Transparency for Electronics Components

Products Manufacturers Applications Services & Tools Help Order History Log In Register 

All ▾  In Stock RoHS

[All Products](#) > [Passive Components](#) > [Capacitors](#) > [Tantalum Capacitors](#) > [Tantalum Capacitors - Polymer SMD](#) > [See an Error?](#)

T520B107M006ATE040



[Enlarge](#)

Images are for reference only
See Product Specifications

[Share](#)

Mouser #:	80-T520B107M6ATE40
Mfr. #:	T520B107M006ATE040
Mfr.:	KEMET
Customer #:	<input type="text"/>
Description:	Tantalum Capacitors - Polymer SMD 6.3volts 100uF 20% ESR=40 Available in MultiSIM BLUE View Simulation and SPICE Model in K-SIM
Datasheet:	T520B107M006ATE040 Datasheet
More Information:	Learn more about KEMET T520B107M006ATE040

In Stock: 7,998

Stock:	7,998 Can Ship Immediately
On Order:	2000 View Delivery Dates
Factory Lead-Time:	21 Weeks
Enter Quantity:	Minimum: 1 Multiples: 1 <input type="text"/> Buy

Pricing (USD)

Qty.	Unit Price	Ext. Price
1	\$1.22	\$1.22
10	\$0.838	\$8.38
100	\$0.644	\$64.40

“Operating Characteristics” of a component



Miniature Aluminum Electrolytic Capacitors

XRL Series

■ FEATURES

- Low impedance characteristics
- Case sizes are smaller than conventional general-purpose capacitors, with very high performance
- Can size larger than 9mm diameter has safety vents on rubber end seal
- RoHS Compliant



■ CHARACTERISTICS

Item	Characteristics
Operating Temperature Range	-40°C ~ +85°C
Capacitance Tolerance	±20% at 120Hz, 20°C
Leakage Current	≤100V $I = 0.01CWV$ or $3\mu A$ whichever is greater after 2 minutes of applied rated DC working voltage at 20°C Where: C = rated capacitance in μF ; WV = rated DC working voltage
	>100V $CWV \leq 1000 \mu F$: $I = 0.03 CWV + 15\mu A$; $C > 1000 \mu F$: $I = 0.02 CWV + 25\mu A$; WV = rated DC working voltage in V
Dissipation Factor (Tan δ , at 20°C 120Hz)	Working voltage (WV) 6.3 10 16 25 35 50 63 100 160 250 350 450
	Tan δ 0.23 0.20 0.16 0.14 0.12 0.10 0.09 0.08 0.12 0.17 0.20 0.25 For capacitors whose capacitance exceeds 1,000 μF , the specification of tan δ is increased by 0.02 for every addition of 1,000 μF .
Surge Voltage	Working voltage (WV) 6.3 10 16 25 35 50 63 100 160 250 350 450
	Surge voltage (SV) 8 13 20 32 44 63 79 125 200 300 400 500
Low Temperature Characteristics (Imp. ratio @ 120Hz)	Working voltage (WV) 6.3 10 16 25 35 50 63 100 160 250 350 450
	Z(-25°C)/Z(+20°C) $\alpha D \leq 16$ 6 4 3 3 2 2 2 2 3 8 12 16
	$\alpha D \geq 16$ 8 6 4 4 3 3 3 3 3 8 12 16
Load Test	Z(-40°C)/Z(+20°C) $\alpha D \leq 16$ 10 8 6 6 4 3 3 3 4 10 16 20
	$\alpha D \geq 16$ 18 16 12 10 8 8 6 6 4 10 16 20
Self Life Test	When returned to +20°C after 2,000 hours application of working voltage at +85°C, the capacitor will meet the following limits: Capacitance change is $\leq \pm 20\%$ of initial value; tan δ is $< 200\%$ of specified value; leakage current is within specified value.

■ PART NUMBERING SYSTEM

1	4	0	-	X	R	L	1	6	V	1	0	0	-	R	C
Prefix	Series			Voltage Actual Value			Capacitance (μF) Actual Value			Suffix RoHS Compliant					

■ RIPPLE CURRENT AND FREQUENCY MULTIPLIERS

Capacitance (μF)	Frequency (Hz)				
	60 (50)	120	500	1K	±10K
<100	0.70	1.0	1.30	1.40	1.50
100 ~ 1000	0.75	1.0	1.20	1.30	1.35
>1000	0.80	1.0	1.10	1.12	1.15

■ RIPPLE CURRENT AND TEMPERATURE MULTIPLIERS

Temperature (°C)	<50	70	85
Multiplier	1.78	1.4	1.0

XICON PASSIVE COMPONENTS • (800) 628-0544



XC-600178 Specifications are subject to change without notice. No liability or warranty implied by this information. Environmental compliance based on producer documentation. Date Revised: 1/8/07



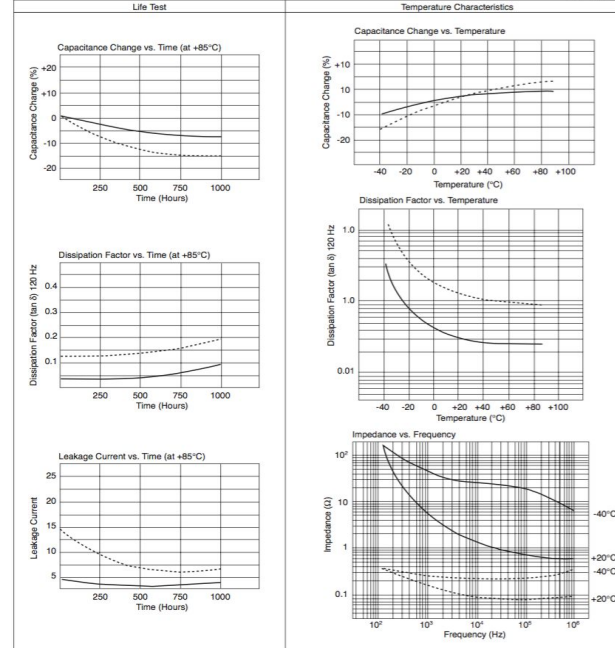
Miniature Aluminum Electrolytic Capacitors

XRL Series

■ TYPICAL PERFORMANCE CHARACTERISTICS

----- 1000 μF 16V

1 μF 50V



XICON PASSIVE COMPONENTS • (800) 628-0544



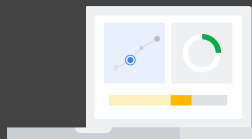
XC-600178 Specifications are subject to change without notice. No liability or warranty implied by this information. Environmental compliance based on producer documentation. Date Revised: 1/8/07

Model Cards for Model Reporting

- Currently no common practice of reporting how well a model works when it is released

Model Cards for Model Reporting

Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, Timnit Gebru
{mmitchellai,simonewu,andrewzaldivar,parkerbarnes,lucyvasserman,benhutch,espitzer,tgebru}@google.com
deborah.rajai@mail.utoronto.ca



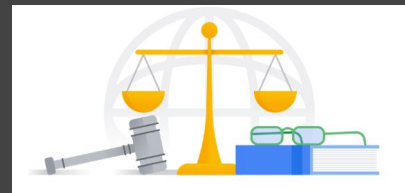
What It Does

A report that focuses on transparency in model performance to encourage responsible AI adoption and application.



How It Works

It is an easily discoverable and usable artifact presented at important steps of a user journey for a diverse set of users and public stakeholders.



Why It Matters

It keeps model developer accountable to release high quality and fair models.

Intended Use, Factors and Subgroups

Example Model Card - Toxicity in Text	
Model Details	Developed by Jigsaw in 2017 as a convolutional neural network trained to predict the likelihood that a comment will be perceived as toxic.
Intended Use	Supporting human moderation, providing feedback to comment authors, and allowing comment viewers to control their experience.
Factors	Identity terms referencing frequently attacked groups focusing on the categories of sexual orientation, gender identity and race.

Metrics and Data

Metrics	<i>Pinned AUC</i> , which measures threshold-agnostic separability of toxic and non-toxic comments for each group, within the context of a background distribution of other groups.
Evaluation Data	A synthetic test set generated using a template-based approach, where identity terms are swapped into a variety of template sentences.
Training Data	Includes comments from a variety of online forums with crowdsourced labels of whether the comment is “toxic”. “Toxic” is defined as, “a rude, disrespectful, or unreasonable comment that is likely to make you leave a discussion”.

Considerations, Recommendations

Ethical Considerations	A set of values around community, transparency, inclusivity, privacy and topic-neutrality to guide their work.
Caveats & Recommendations	Synthetic test data covers only a small set of very specific comments. While these are designed to be representative of common use cases and concerns, it is not comprehensive.

Disaggregated Intersectional Evaluation

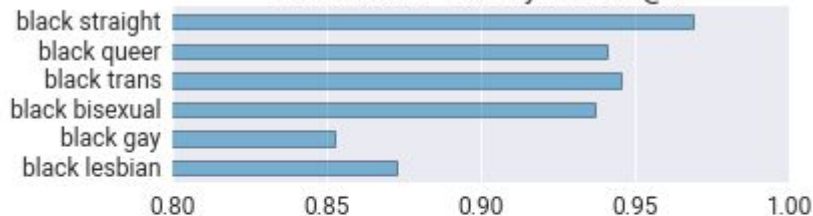
Toxicity @1

Identity groups	Subgroup AUC	BPSN AUC	BNSP AUC
lesbian	0.93	0.74	0.98
gay	0.94	0.65	0.99
queer	0.98	0.96	0.93
straight	0.99	1.00	0.87
bisexual	0.96	0.95	0.92
homosexual	0.87	0.53	0.99
heterosexual	0.96	0.94	0.92
cis	0.99	1.00	0.87
trans	0.97	0.96	0.91
nonbinary	0.99	0.99	0.90
black	0.91	0.85	0.95
white	0.91	0.88	0.94

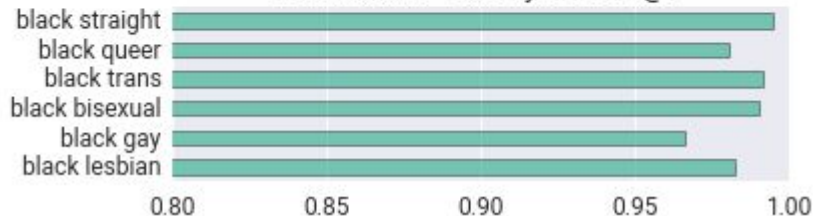
0.5 0.6 0.7 0.8 0.9 1.0



Pinned AUC Toxicity Scores @1



Pinned AUC Toxicity Scores @5



Jigsaw



The False Positive

In Summary...

- **Question why** should we build NLP model X, and who it may **harm**
- Always **be mindful** of various sorts of biases in the NLP models and the data
- Consider this an **iterative process**, than something that has a “done” state
- Explore “debiasing” techniques, but **be cautious**
- Identify **fairness interventions that matter** for your problem
- Be **transparent** about your model and its performance in different settings

Closing Note

“Fairness and justice are properties of social and legal systems”

“To treat fairness and justice as terms that have meaningful application to technology separate from a social context is therefore [...] an abstraction error”

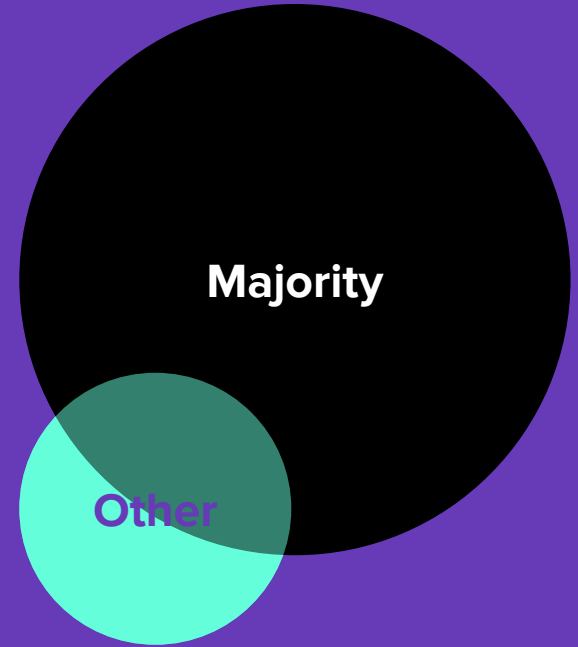


Questions?



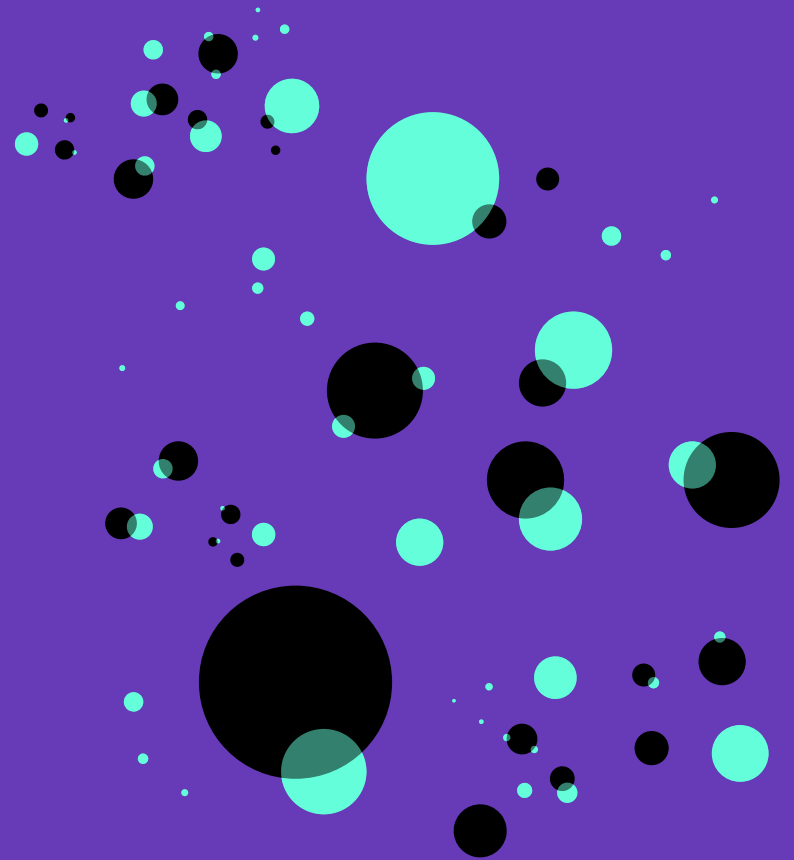
BACKUP Slides

Moving from majority
representation...



Moving from majority
representation...

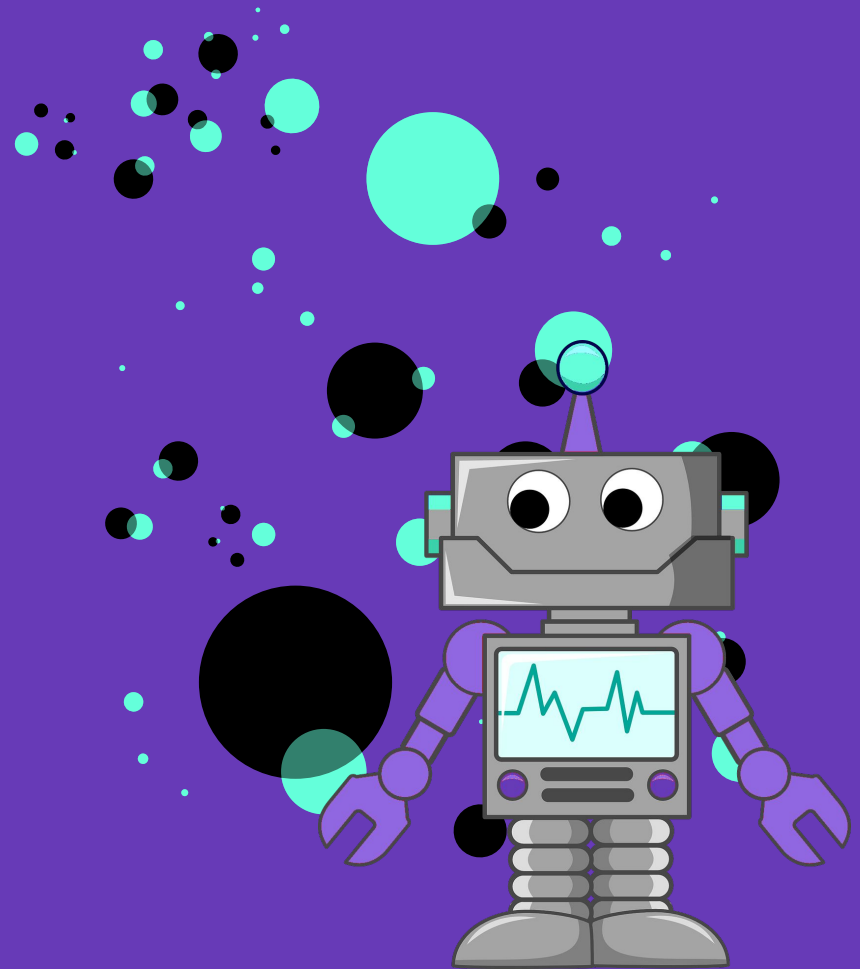
...to diverse
representation



Moving from majority
representation...

...to diverse
representation

...for ethical AI



Thanks!

margarmitchell@gmail.com

m-mitchell.com

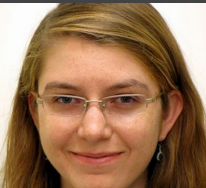
Need MOAR? ml-fairness.com



Andrew
Zaldivar



Me



Simone
Wu



Parker
Barnes



Lucy
Vasserman



Ben
Hutchinson



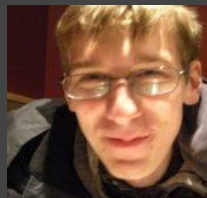
Elena
Spitzer



Deb
Raji



Timnit Gebru



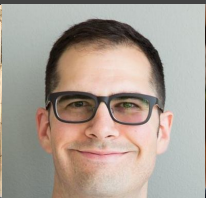
Adrian
Benton



Brian
Zhang



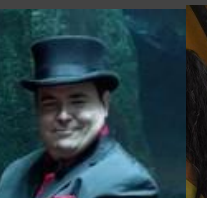
Dirk
Hovy



Josh
Lovejoy



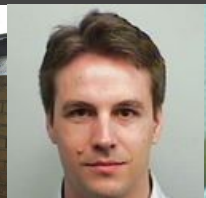
Alex
Beutel



Blake
Lemoine



Hee Jung
Ryu



Hartwig
Adam



Blaise
Agüera y
Arcas

More free, hands-on tutorials on how to build more inclusive ML

Measuring and Mitigating Unintended Bias in Text Classification

John Li
jetpack@google.com

Lucas Dixon
ldixon@google.com

Nithum Thain
nthain@google.com

Lucy Vasserman
lucyvasserman@google.com

Jeffrey Sorensen
sorenj@google.com

Mitigating Unwanted Biases with Adversarial Learning

Brian Hu Zhang
Stanford University
Stanford, CA
bhz@stanford.edu

Blake Lemoine
Google
Mountain View, CA
lemoine@google.com

Margaret Mitchell
Google
Mountain View, CA
mmitchellai@google.com

The screenshot shows a Colab notebook interface. The title bar reads "pinned_auc_demo.ipynb". The main content area is titled "Conversation AI's Pinned AUC Unintended Model Bias Demo". It includes an author list: "Author: ldixon@google.com, jetpack@google.com, sorenj@google.com, nthain@google.com, lucyvasserman@google.com". A summary section states: "This notebook demonstrates Pinned AUC as an unintended model bias metric for Conversation AI wikipedia models." A disclaimer section contains two bullet points: "This notebook contains experimental code, which may be changed without notice." and "The ideas here are some ideas relevant to fairness - they are not the whole story!". The bottom of the notebook shows a code cell with the command: `!pip install -U -q git+https://github.com/conversationai/unintended-ml-bias-analysi`

The screenshot shows a Colab notebook interface. The title bar reads "Debiasing Word Embeddings using Fair Adversarial Networks (FANs)". The main content area is titled "Debiasing Word Embeddings using Fair Adversarial Networks (FANs)". It includes authors: "Authors: lemoine@, zhangbrian@, benhutch@, guajardo@" and contributors: "Contributors: mmitchellai@, andrewzaldivar@". A summary section states: "This Colab was put together as part of the ML-fairness inspired hackathon in late August 2017 to demonstrate how to mitigate bias in word embeddings using an adversarial network." A disclaimer section contains two bullet points: "This notebook contains experimental code, which may be changed without notice." and "The ideas here are some ideas relevant to fairness - they are not the whole story!". The bottom of the notebook shows a code cell with the command: `!pip install -U -q git+https://github.com/conversationai/unintended-ml-bias-analysi`

Get Involved

- Find free machine-learning tools open to anyone at ai.google/tools
- Check out Google's ML Fairness codelab at ml-fairness.com
- Explore educational resources at ai.google/education
- Take a free, hands-on Machine Learning Crash Course at <https://developers.google.com/machine-learning/crash-course/>
- Share your feedback: acceleratewithgoogle@google.com

Build for everyone





Measuring Algorithmic Fairness/Bias

Evaluate for Fairness & Inclusion

Disaggregated Evaluation

Create for each (subgroup, prediction) pair.

Compare across subgroups.

Evaluate for Fairness & Inclusion

Disaggregated Evaluation

Create for each (subgroup, prediction) pair.

Compare across subgroups.

Example: women, face detection
men, face detection

Evaluate for Fairness & Inclusion: Confusion Matrix

Model Predictions

References

Evaluate for Fairness & Inclusion: Confusion Matrix

		Model Predictions	
		Positive	Negative

References	Positive
	Negative

Evaluate for Fairness & Inclusion: Confusion Matrix

		Model Predictions	
		Positive	Negative
References	Positive	<ul style="list-style-type: none">● Exists● Predicted True Positives	
	Negative		<ul style="list-style-type: none">● Doesn't exist● Not predicted True Negatives

Evaluate for Fairness & Inclusion: Confusion Matrix

		Model Predictions	
		Positive	Negative
References	Positive	<ul style="list-style-type: none">● Exists● Predicted True Positives	<ul style="list-style-type: none">● Exists● Not predicted False Negatives
	Negative	<ul style="list-style-type: none">● Doesn't exist● Predicted False Positives	<ul style="list-style-type: none">● Doesn't exist● Not predicted True Negatives

Evaluate for Fairness & Inclusion: Confusion Matrix

		Model Predictions		
		Positive	Negative	
References	Positive	<ul style="list-style-type: none">ExistsPredicted True Positives	<ul style="list-style-type: none">ExistsNot predicted False Negatives	Recall, False Negative Rate
	Negative	<ul style="list-style-type: none">Doesn't existPredicted False Positives	<ul style="list-style-type: none">Doesn't existNot predicted True Negatives	False Positive Rate, Specificity
		Precision, False Discovery Rate	Negative Predictive Value, False Omission Rate	LR+, LR-

Evaluate for Fairness & Inclusion

Female Patient Results

True Positives (TP) = 10	False Positives (FP) = 1
False Negatives (FN) = 1	True Negatives (TN) = 488

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{10}{10 + 1} = 0.909$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{10}{10 + 1} = 0.909$$

Male Patient Results

True Positives (TP) = 6	False Positives (FP) = 3
False Negatives (FN) = 5	True Negatives (TN) = 48

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{6}{6 + 3} = 0.667$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{6}{6 + 5} = 0.545$$

Evaluate for Fairness & Inclusion

Female Patient Results

True Positives (TP) = 10	False Positives (FP) = 1
False Negatives (FN) = 1	True Negatives (TN) = 488

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{10}{10 + 1} = 0.909$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{10}{10 + 1} = 0.909$$

Male Patient Results

True Positives (TP) = 6	False Positives (FP) = 3
False Negatives (FN) = 5	True Negatives (TN) = 48

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{6}{6 + 3} = 0.667$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{6}{6 + 5} = 0.545$$

**“Equality of Opportunity” fairness criterion:
Recall is equal across subgroups**

Evaluate for Fairness & Inclusion

Female Patient Results

True Positives (TP) = 10	False Positives (FP) = 1
False Negatives (FN) = 1	True Negatives (TN) = 488

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{10}{10 + 1} = 0.909$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{10}{10 + 1} = 0.909$$

Male Patient Results

True Positives (TP) = 6	False Positives (FP) = 3
False Negatives (FN) = 5	True Negatives (TN) = 48

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{6}{6 + 3} = 0.667$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{6}{6 + 5} = 0.545$$

**“Predictive Parity” fairness criterion:
Precision is equal across subgroups**