# NLP and Society:
# Towards Socially Responsible NLP

## Vinodkumar Prabhakaran
Research Scientist

Google

# What's in this talk...

- Motivation for Machine Learning (ML) Fairness research

- Why and how ML models may be unfair

- Fairness issues in ML-based Natural Language Processing

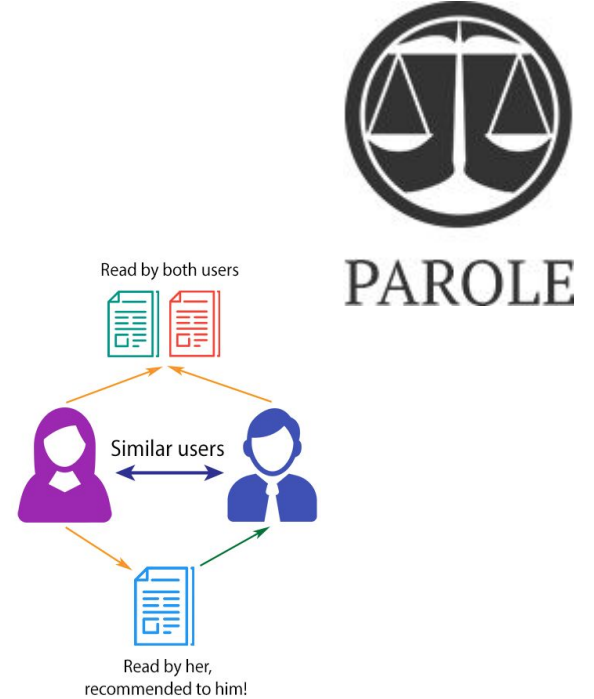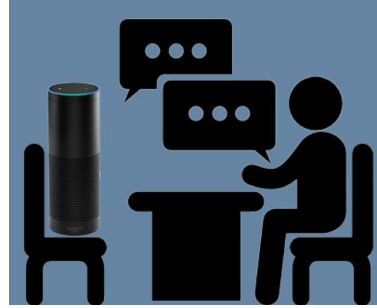- What can/should we do?

# What's **NOT** in this talk...

- Definitive answers to fairness/ethical questions

- Prescriptive solutions to fix ML/NLP (un)fairness

- Focus on research done by myself, my team, or Google.

- ...

# What's **<u>also</u>** in this talk...

- Research done in academia and various industry labs

- Research from other disciplines, including Psychology, Philosophy, and Social Sciences in general …

- Uncomfortable impacts of technology on society

# Machine Learning is Everywhere!!!

# Machine Learning is Everywhere!!!

"It's true that they can follow instructions at superhuman speed, with superhuman fidelity and over unimaginable quantities of data. **But these instructions don't come from nowhere.** Although neural networks might be said to write their own programs, they do so towards g**oals set by humans, using data collected for human purposes**. If the data is skewed, even by accident, the computers will amplify injustice."

— The Guardian

"It's true that they can follow instructions at superhuman speed, with superhuman fidelity and over unimaginable quantities of data. **But these instructions don't come from nowhere**. Although neural networks might be said to write their own programs, they do so towards g**oals set by humans, using data collected for human purposes**. If the data is skewed, even by accident, the computers will *amplify injustice*."

— The Guardian
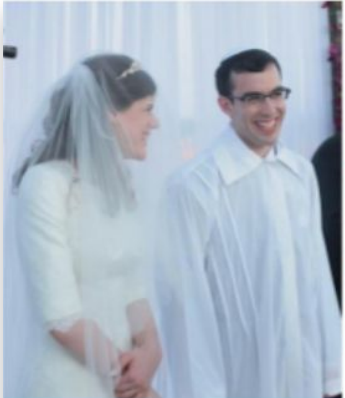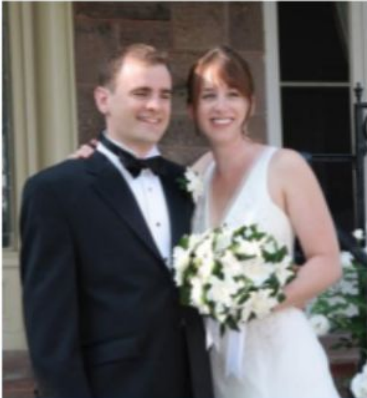
# Fairness in Machine Learning
## A Few Case Studies

# Photo captioning
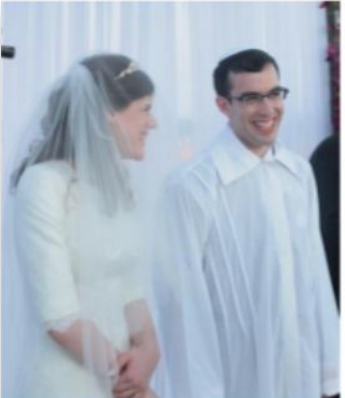


ceremony, wedding, bride, man, groom, woman, dress

bride, ceremony, wedding, dress, woman

ceremony, bride, wedding, man, groom, woman, dress
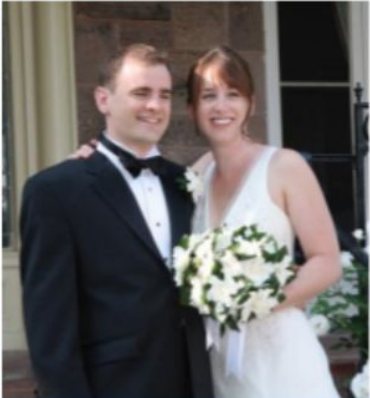
# Photo captioning



ceremony, wedding, bride, man, groom, woman, dress

bride, ceremony, wedding, dress, woman

ceremony, bride, wedding, man, groom, woman, dress

person, people

# Predicting Sexual Orientation



*Original Paper*: "Deep neural networks are more accurate than humans at detecting sexual orientation from facial images" Wang and Kosinsky, 2017. PsyArXiv

# Predicting Sexual Orientation



"Differences between lesbian or gay and straight faces in selfies relate to grooming, presentation, and lifestyle—**that is, differences in culture, not in facial structure**."

"Do Algorithms Reveal Sexual Orientation or Just Expose our Stereotypes?" Medium, Blaise Agüera y Arcas, Alexander Todorov and Margaret Mitchell
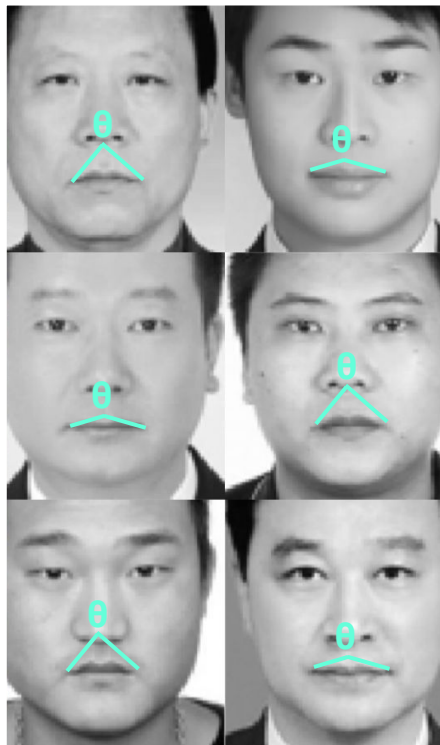
# Predicting criminality



"Automated Inference on Criminality using Face Images"
Wu and Zhang, 2016.  arXiv

# Predicting criminality

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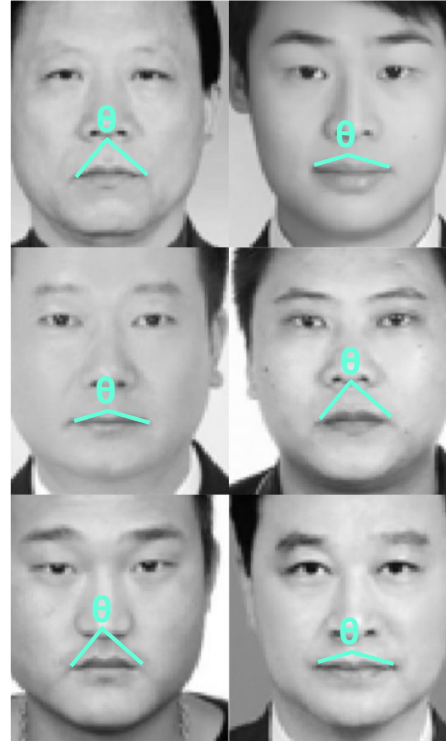


"Automated Inference on Criminality using Face Images"
Wu and Zhang, 2016.  arXiv

# Predicting criminality: physiognomy?

"[...] angle ϴ from nose tip to two mouth corners is on average 19.6% smaller for criminals than for non-criminals ..."

[Physiognomy's New Clothes](#) (Medium Blog Post) - by Blaise Agüera y Arcas, Margaret Mitchell and Alexander Todorov

*"Deep learning based on superficial features is decidedly not a tool that should be deployed to "accelerate" criminal justice; attempts to do so will instead perpetuate injustice."*



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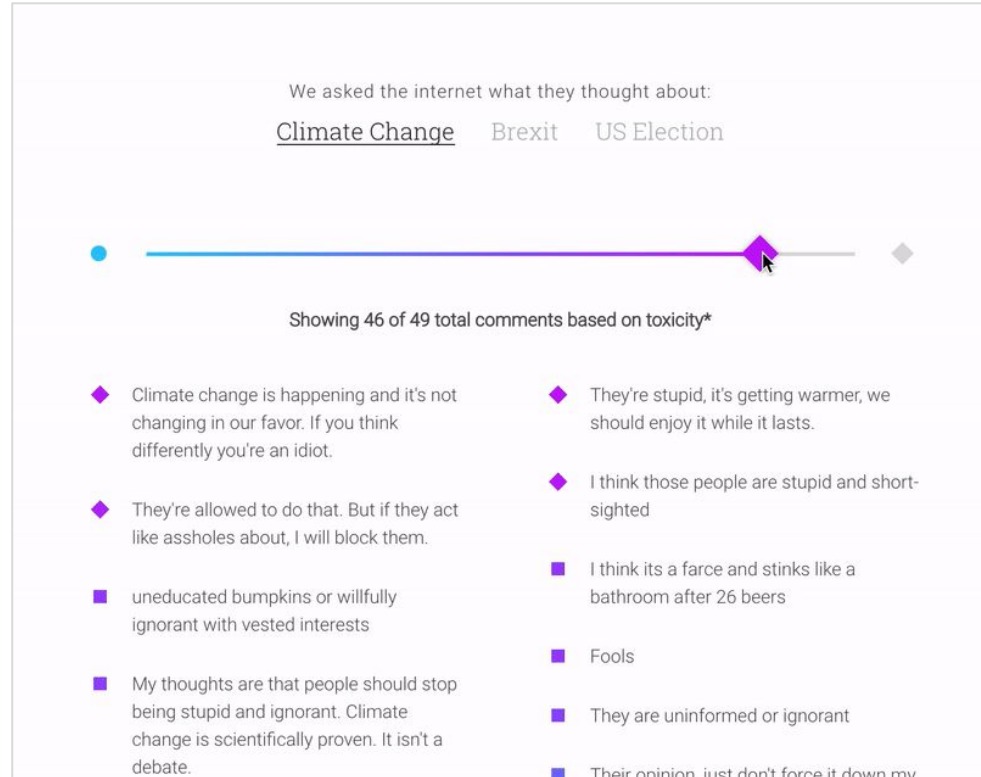Wu and Zhang, 2016.  [arXiv](#)

# Toxicity Classification



Jigsaw

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 allowed to do that. But if they act like assholes about, I will block them.

■ uneducated bumpkins or willfully ignorant with vested interests

■ My thoughts are that people should stop being stupid and ignorant. Climate change is scientifically proven. It isn't a debate.

◆ They're stupid, it's getting warmer, we should enjoy it while it lasts.

◆ I think those people are stupid and short-sighted

■ I think its a farce and stinks like a bathroom after 26 beers

■ Fools

■ They are uninformed or ignorant

■ Their opinion, just don't force it down my
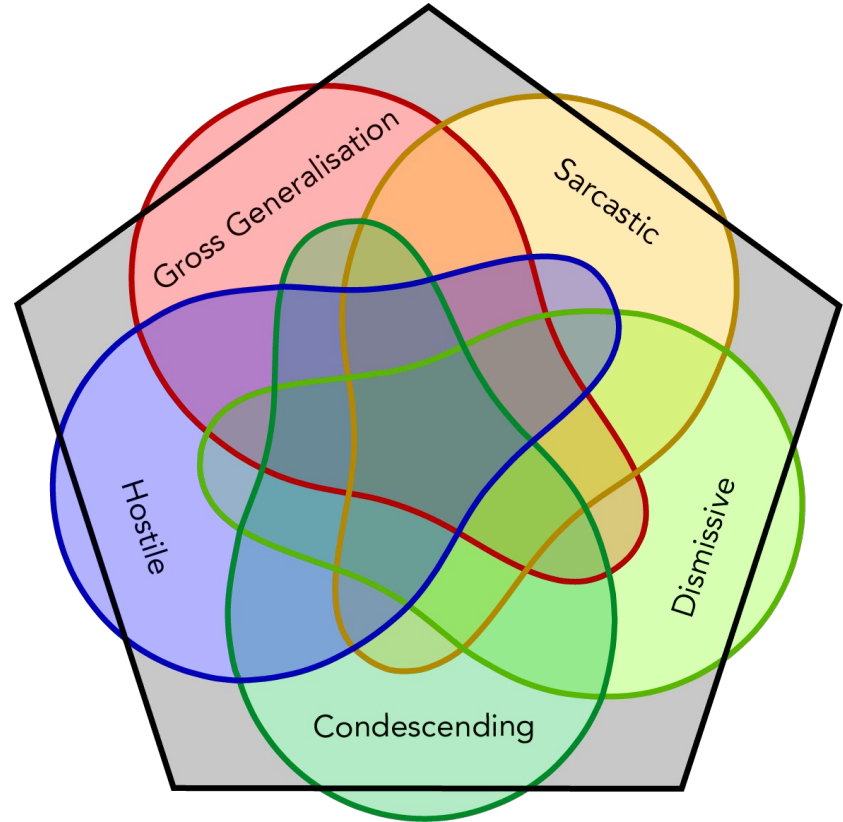
# Toxicity Classification

Toxic 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

**A naively trained model** will have **strong unintended biases** as illustrated by these examples:

| 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

**A naively trained model** will have **strong unintended biases** as illustrated by these examples:

| Comment | Toxicity Score |
|---|---|
| I hate Justin Timberlake. | 0.90 |
| I hate Rihanna. | 0.69 |

- Prabhakaran et al. (2019). "Perturbation Sensitivity Analysis to Detect Unintended Model Biases" EMNLP 2019

# Toxicity Classification

A **naively trained model** will have **strong unintended biases** as illustrated by these examples:

| Comment | Toxicity Score |
|---|---|
| I am a person. | 0.08 |
| I am a tall person. | 0.03 |

- Hutchinson et al. (2019). *Unintended Machine Learning Biases as Social Barriers for Persons with Disabilities.* SIGACCESS ASSETS AI Fairness Workshop 2019.

# Toxicity Classification

**A naively trained model** will have **strong unintended biases** as illustrated by these examples:

| 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 |

- Hutchinson et al. (2019). *Unintended Machine Learning Biases as Social Barriers for Persons with Disabilities.* SIGACCESS ASSETS AI Fairness Workshop 2019.

# Toxicity Classification

A **naively trained model** will have **strong unintended biases** as illustrated by these examples:

| 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 |

- Hutchinson et al. (2019). *Unintended Machine Learning Biases as Social Barriers for Persons with Disabilities*. SIGACCESS ASSETS AI Fairness Workshop 2019.

# Allocative Harm

*"when a system allocates or withholds a certain opportunity or resource"*

# Associative Harm

*"when systems reinforce the subordination of some groups along the lines of identity"*

# Why do these things happen?

# Machine Learning "sequence"

**Collect and annotate training data.**

**Train model.**

**Filter, rank, aggregate, or generate content.**

**People see output.**

Google

# Potential biases

**Collect and annotate training data.**

**Human Biases in Data**

**Reporting bias**

**Selection bias**

**Overgeneralization**

**Out-group homogeneity bias**

**Unconscious bias from "the world"** that we might reflect in ML when using some of the world's data

**Human Biases in Collection and Annotation**

**Confidence bias / Overconfidence effect**

**Confirmation bias**

**Experimenter's bias**

**Unconscious bias in our procedures** that we might reflect in our ML

Google

# Unconscious bias interferes



Unconscious bias gets reinforced in the training data

Collect and annotate training data.

Train model.

Filter, rank, aggregate, or generate content.

People see output

Unconscious bias affects the way we collect and classify data, design, and write code

Google

# Fairness in Natural Language Processing
## A Deeper Dive

The common misconception is that language has to do with **words** and what they mean.

**It doesn't.**

**It has to do with people and what they mean.**

Herbert H. Clark & Michael F. Schober, 1992

# Fairness in Natural Language Processing

## A Deeper Dive

- Is my data biased?

# Selection Bias: World Englishes



60M Speakers

125M Speakers

251M Speakers

90M Speakers

79M Speakers

# Selection Bias: Gender Equity



Female vs. male internet gender gaps

26% gender gap globally

4% Europe & Central Asia

21% Middle East & North Africa

4% Latin America & Caribbean

70% South Asia

4% East Asia & Pacific

34% Sub-Saharan Africa

GSMA, 2018

# Selection Bias: Gender Equity

- 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)

# A case study:
## Language Identification

# Sampling Bias in Language Identification (LID)

- Most NLP applications employ off-the-shelf LID systems as the first step

(Jurgens et al. ACL'17)

# Sampling Bias in Language Identification (LID)

- Most NLP applications employ off-the-shelf LID systems as the first step

Example Application:
- Public Health Monitoring

# 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.

# World Englishes

# World Englishes

**The Royal Family** ✓
@RoyalFamily
**Follow**

Taking place this week on the river Thames is 'Swan Upping' – the annual census of the swan population on the Thames.

**da'Rah-zingSun**
@TIME7SS
**Follow**

@kimguilfoyle prblm I hve wit ur reportng is its 2 literal, evry1 knos pple tlk diffrnt evrywhere, u kno wut she means jus like we do!

**Mooktar**
@bossmukky
**Follow**

"@Ecstatic_Mi: @bossmukky Ebi like say I wan dey sick sef wlh 'Flu' my whole body dey weak"uw gee...

**Ebenezer·**
@Physique_cian
**Follow**

@Tblazeen R u a wizard or wat gan sef : in d mornin- u tweet, afternoon - u tweet, nyt gan u dey tweet.beta get ur IT placement wiv twitter

- Language identification degrades significantly on African American Vernacular English
  (Blodgett et al. 2016)

# LID Usage Example: Public Health Monitoring

# Socioeconomic Bias in Language Identification

- 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
systemically <span style="color:red">excluded</span> from this technology

# Better Social Representation through Network-based Sampling

- Re-sampling from strategically-diverse corpora

**Topical**

**Geographic**

**Social**

**Multilingual**

classifier

— langid.py
— CLD2
— EquiLID

Estimated accuracy for English tweets (y-axis: 0.6 to 1.0)

Human Development Index of text's origin country (x-axis: 0.4 to 1.0)

(Jurgens et al. ACL'17)

# Fairness in Natural Language Processing

## A Deeper Dive

- Is my data biased?
- Is my model biased?

# Bias in NLP Models

1. Bolukbasi T... ...aligrama V., Kalai A. (2016) **Man is to ...** ...woman is to Homemaker? Debiasing Word** **En...**

2. C...liskan, ..., Brys...n, 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.*

Slide from SRNLP Tutorial at NAACL 2018

# Bias in NLP Models

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. Garg et al. **Word embeddings quantify 100 years of gender and ethnic stereotypes.** *PNAS.* (2018)
4. Zhao, Jieyu, et al. **Men also like shopping: Reducing gender bias amplification using corpus-level constraints.** *arXiv* (2017)

**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)

**2019**
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. …

# Where to look for biases?



Input Text → (Input/Embedding Layer) (Hidden Layers) (Output Layer) → Prediction

**Bias in Input Representations?**

# Input Representation: Word Embeddings



**Neural Language Model** (Bengio et al, `03)

**word2vec** (Mikolov et al, `03)

Skip-gram

CBOW

*i*-th output = $P(w_t = i \mid context)$

softmax

most computation here

tanh

$C(w_{t-n+1})$   $C(w_{t-2})$   $C(w_{t-1})$

Table look-up in $C$   Matrix $C$ shared parameters across words

index for $w_{t-n+1}$   index for $w_{t-2}$   index for $w_{t-1}$

INPUT   PROJECTION   OUTPUT

w(t)   w(t-2) w(t-1) w(t+1) w(t+2)

w(t-2) w(t-1) w(t+1) w(t+2)   SUM   w(t)

Documents

Terms  A  =  U  Σ  V$^T$

$m \times n$   $m \times r$   $r \times r$   $r \times n$

A   =   U   D   V$^T$

**Latent Semantic Analysis**
(Deerwester et al, `90, Turney & Pantel `10)

BERT (Ours)   OpenAI GPT   ELMo

**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

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

# Social Stereotypes → Word Embeddings?

# 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 recog... skinned faces. It often reveals an automatic preference for light-skin relative to da...

**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...

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

Greenwald et al. 1998

**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.

# 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, …*}
    - $Y$ = {*nurse, teacher, librarian, …*}
  - Attribute words
    - $A$ = {*man, male, …* }
    - $B$ = {*woman, female, …*}

# Word Embedding Association Test

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

Association of a word $w$ with an attribute: $\quad s(w, A, B) = \text{mean}_{a \in A} \cos(\vec{w}, \vec{a}) - \text{mean}_{b \in B} \cos(\vec{w}, \vec{b})$

Association of two sets of target words with an attribute: $\quad s(X, Y, A, B) = \sum_{x \in X} s(x, A, B) - \sum_{y \in Y} s(y, A, B)$

The effect size of bias: $\quad \dfrac{\text{mean}_{x \in X} s(x, A, B) - \text{mean}_{y \in Y} s(y, A, B)}{\text{std-dev}_{w \in X \cup Y} s(w, A, B)}$

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 | $10^{-8}$ | $25 \times 2$ | $25 \times 2$ | 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*, Tvree, 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                                                        WEAT

| 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 | $10^{-5}$ | $32 \times 2$ | $25 \times 2$ | 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

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

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

surgeon vs. nurse

architect vs. interior designer

shopkeeper vs. housewife

superstar vs. diva

....

man

woman

king

queen

Male-Female

# 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.

# Social Stereotypes → Word Embeddings?
## Yes, they do!

But aren't they just reflecting Society?

# Gender bias in Occupations

# Gender bias in Adjectives over the decades



Height of women's movements in 1960s-70s

# But aren't they just reflecting Society?

## Yup!

Oisin Deery & Katherine Bailey
Ethics in NLP workshop. NAACL '18

Shouldn't we then just leave them as is?

# Shouldn't we then just leave them as is?

1. Would that harm certain groups of people?

# 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?



Input Text

(Input/Embedding Layer)

(Hidden Layers)

(Output Layer)

**Bias in Predictions?**

Prediction

**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 identitiy 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?

1. Would that harm certain groups of people?

2. 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?



FAST COMPANY

CO.DESIGN | TECH | WORK LIFE | CREATIVITY | IMPACT | AUDIO | VIDEO

05.08.18 | THE FUTURE OF WORK

**The Potential Hidden Bias In Automated Hiring Systems**

More companies are using machine-learning software to screen candidates, but it may be unwittingly perpetuating past bias.

Bloomberg

Business

**Artificial Intelligence Is Coming for Hiring, and It Might Not Be That Bad**

Even with all of its problems, AI is a step up from the notoriously biased recruiting process.

# Examples of Harm from NLP Bias

## Compounding imbalances

**Surgeons**

females in data:
**14.6%**

females in true positives:
**11.6%**

Males:
**71% recall**

Females:
**54% recall**

Slide credit: Maria De-Arteaga

# 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")
```

0.98463802132985556

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

-0.47048131775890656

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

# 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

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

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

surgeon vs. nurse

architect vs. interior designer

shopkeeper vs. housewife

superstar vs. diva

....

Male-Female

# Towards Debiasing

1. Identify gender subspace: B

# Gender Subspace



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



programmer

doctor

he

homemaker

nurse

she

king

queen

Plus
Bootstrapping

*218 gender-definitional words*

**Linear SVM**

# 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

$$min_T||(TW)^T(TW) - W^TW||_F^2 + \lambda||(TN)^T(TB)||_F^2$$

Don't modify embeddings too much

Minimize gender component

T - the desired debiasing transformation    B - biased space
W - embedding matrix    N - embedding matrix of gender neutral words

# Can we computationally remove undesirable biases?

- **Debiasing Meaning Representations**

  - **Debiasing Model Predictions**

# Debiasing using Adversarial Learning

Beutel et al. (2017)
Zhang et al. (2018)

**Bias Mitigation**

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

# Debiasing using Adversarial Learning

Beutel et al. (2017)
Zhang et al. (2018)

**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?

# Recommendations

- Always **be mindful** of various sorts of biases in the NLP models and the data

- Explore "debiasing" techniques, **but be cautious**

- Think about the biases that matter for your problem and **test for those biases**

- Be **transparent** about the models you release to the world

# Speaking of Transparency...

# Transparency for Electronics Components



Slide by Timnit Gebru

# Transparency for Electronics Components



Slide by Timnit Gebru

# Speaking of Transparency...

- **Data Sheets for Datasets**

# Datasheets for Datasets

- Gebru et al. (2019)
  - https://arxiv.org/pdf/1803.09010.pdf

- Key questions for each stage:
  - Motivation
  - Composition
  - Collection Process
  - Preprocessing/cleaning/labeling
  - Uses
  - Distribution
  - Maintenance

- For dataset creators:
  - Encourage reflection on the process and assumptions

- For dataset consumers:
  - Provide information for making informed decisions

# Speaking of Transparency...

- **Data Sheets for Datasets**
- **Model Cards for model reporting**

# Model Card for Toxicity Model

## Model Card - Toxicity in Text

**Model Details**
- The TOXICITY classifier provided by Perspective API [32], trained to predict the likelihood that a comment will be perceived as toxic.
- Convolutional Neural Network.
- Developed by Jigsaw in 2017.

**Intended Use**
- Intended to be used for a wide range of use cases such as supporting human moderation and providing feedback to comment authors.
- Not intended for fully automated moderation.
- Not intended to make judgments about specific individuals.

**Factors**
- Identity terms referencing frequently attacked groups, focusing on sexual orientation, gender identity, and race.

**Metrics**
- Pinned AUC, as presented in [11], which measures threshold-agnostic separability of toxic and non-toxic comments for each group, within the context of a background distribution of other groups.

**Ethical Considerations**
- Following [31], the Perspective API uses a set of values to guide their work. These values are Community, Transparency, Inclusivity, Privacy, and Topic-neutrality. Because of privacy considerations, the model does not take into account user history when making judgments about toxicity.

**Training Data**
- Proprietary from Perspective API. Following details in [11] and [32], this includes comments from a online forums such as Wikipedia and New York Times, 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."
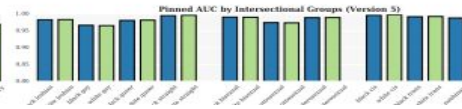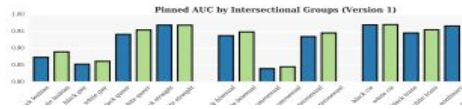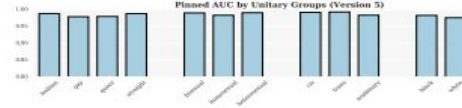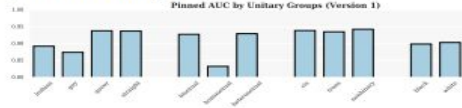
**Evaluation Data**
- A synthetic test set generated using a template-based approach, as suggested in [11], where identity terms are swapped into a variety of template sentences.
- Synthetic data is valuable here because [11] shows that real data often has disproportionate amounts of toxicity directed at specific groups. Synthetic data ensures that we evaluate on data that represents both toxic and non-toxic statements referencing a variety of groups.

**Caveats and 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.

**Quantitative Analyses**

# 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"

Selbst et al., Fairness and Abstraction in Sociotechnical Systems. FAT* 2018

# Thank You!

## Acknowledgments:

Team

Internal          External