# Bias in bios: fairness in a high-stakes machine-learning setting

## Maria De-Arteaga

Joint PhD Student, Machine Learning & Public Policy

Advisors: Artur Dubrawski & Alexandra Chouldechova

**ML** MACHINE LEARNING DEPARTMENT

**Carnegie Mellon University** **Heinz College**
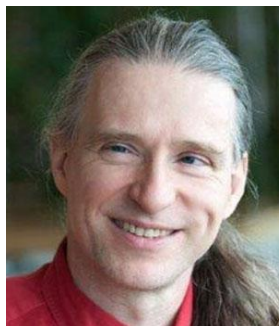
Microsoft Research

Alexey Romanov

Adam Kalai

Hanna Wallach

Jennifer Chayes

Christian borgs

Alexandra Chouldechova

Krishnaram Kenthapadi

Sahin Geyik

Max Leiserson

Nathaniel Swinger, Neil Thomas Heffernan IV

**What are the biases in our data?**

**Why do they matter?**

**What can we do about them?**

# What are the biases in my data?

**What are the biases in my word embedding? (AIES 2019)**

Nathaniel Swinger[=] (Lexington HS),  Maria De-Arteaga[=] (CMU), Neil Thomas Heffernan IV (Shrewsbury HS), Mark Leiserson (UMD), Adam Kalai (MSR)

# Why do they matter?

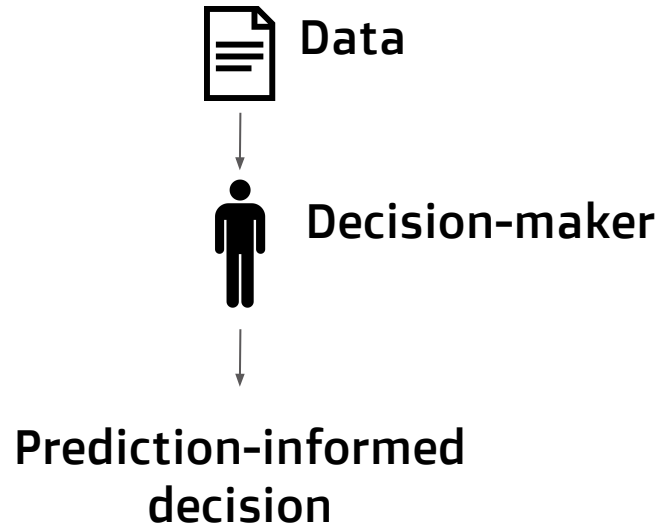**Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting (FAT\* 2019)**

Maria De-Arteaga (CMU), Alexey Romanov (UMASS), Hanna Wallach (MSR), Jennifer Chayes (MSR), Christian Borgs (MSR), Alexandra Chouldechova (CMU), Sahin Geyik (LinkedIn), Krishnaram Kenthapadi (LinkedIn), Adam Kalai (MSR)

# What can we do about them?

**What's in a Name? Reducing Bias in Bios without Access to Protected Attributes (NAACL 2019)**

Alexey Romanov (UMASS), Maria De-Arteaga (CMU), Hanna Wallach (MSR), Jennifer Chayes (MSR), Christian Borgs (MSR), Alexandra Chouldechova (CMU), Sahin Geyik (LinkedIn), Krishnaram Kenthapadi (LinkedIn), Anna Rumshisky (UMASS), Adam Kalai (MSR) Best Thematic Paper :)

# **Humans** and high-stakes predictions

Data

Decision-maker

**Prediction-informed
decision**

# **Humans** and high-stakes predictions

**Defendant's record**

**Judge**

**Bail?**

# **Humans** and high-stakes predictions

Defendant's record → Judge → Bail?

Candidate's CV → Recruiter → Interview? Hire?

# **Humans** and high-stakes predictions

# Humans, **machines** and high-stakes predictions



Data

Machine prediction

Human decision

# Humans, **machines** and high-stakes predictions

Data

Machine prediction

Human decision

**Machines are better than humans at making predictions!**
[Meehl'54, Dawes'89, Grove'00]
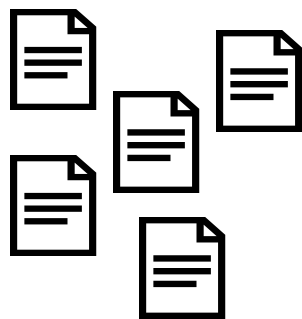
# Humans, **machines** and high-stakes predictions
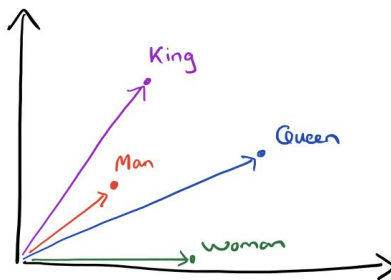
Data

Machine prediction

Human decision

But what happens when available data embeds societal biases?
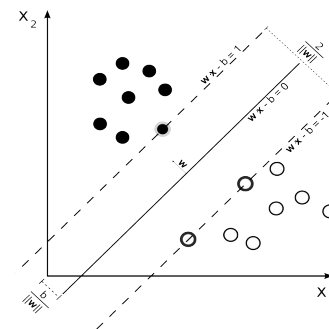
# In this talk...

What are the risks of semantic representation bias?
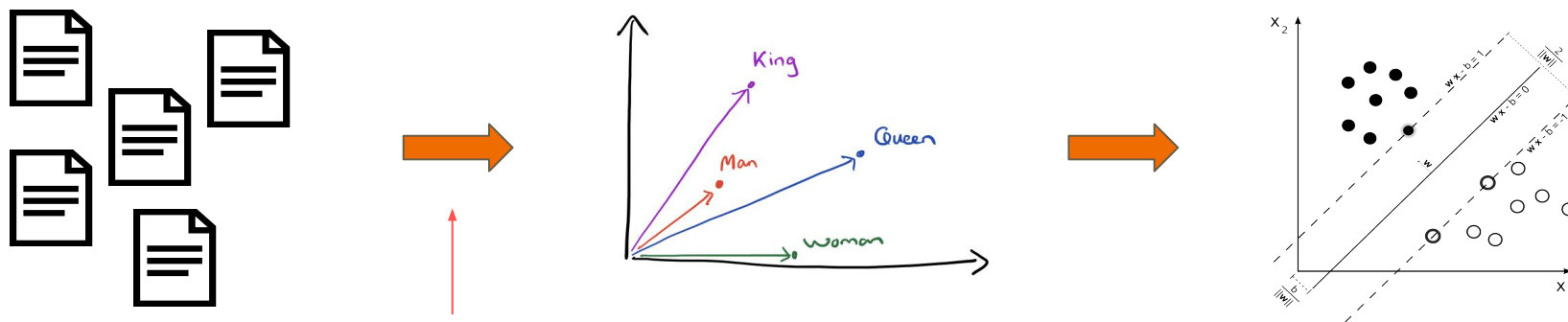


Input data



Semantic representation



Machine learning algorithm

# In this talk...

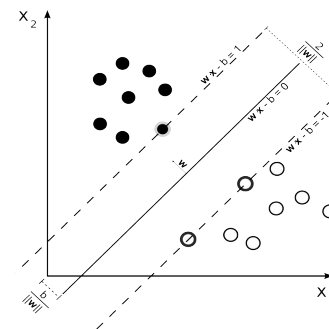## What are the risks of semantic representation bias?



**Part 1:** Representational harms

**What are the biases in my word embedding? (AIES 2019)**

Nathaniel Swinger[=] (Lexington HS),  Maria De-Arteaga[=] (CMU), Neil Thomas Heffernan IV (Shrewsbury HS), Mark Leiserson (UMD), Adam Kalai (MSR)

# In this talk…

## What are the risks of semantic representation bias?



**Part 2:** Allocative harms

**Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting (FAT* 2019)**
Maria De-Arteaga (CMU), Alexey Romanov (UMASS), Hanna Wallach (MSR), Jennifer Chayes (MSR), Christian Borgs (MSR), Alexandra Chouldechova (CMU), Sahin Geyik (LinkedIn), Krishnaram Kenthapadi (LinkedIn), Adam Kalai (MSR)

# In this talk...

### What are the risks of semantic representation bias?



**Part 3:** Mitigating allocative harms

**What's in a Name? Reducing Bias in Bios without Access to Protected Attributes (NAACL 2019)**
Alexey Romanov (UMASS), Maria De-Arteaga (CMU), Hanna Wallach (MSR), Jennifer Chayes (MSR), Christian Borgs (MSR), Alexandra Chouldechova (CMU), Sahin Geyik (LinkedIn), Krishnaram Kenthapadi (LinkedIn), Anna Rumshisky (UMASS), Adam Kalai (MSR) Best Thematic Paper :)

# Word embeddings

# Word embeddings

Man :: computer programmer

Woman ::

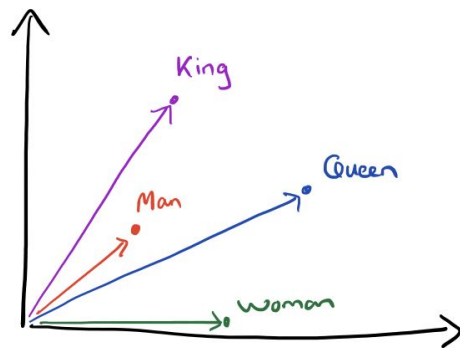| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | the | 0.056 | 0.043 | 0.051 | 0.08 | 0.006 | 0.041 | 0.032 | 0.011 | 0.057 | 0.004 | 0.083 |
| 2 | cat | 0.072 | 0.076 | 0.1 | 0.085 | 0.055 | 0.082 | 0.058 | 0.017 | 0.011 | 0.062 | 0.027 |
| 3 | dog | 0.088 | 0.099 | 0.028 | 0.059 | 0.06 | 0.059 | 0.039 | 0.09 | 0.001 | 0.031 | 0.071 |
| 4 | nurse | 0.03 | 0.018 | 0.058 | 0.074 | 0.055 | 0.028 | 0.025 | 0.054 | 0.094 | 0.052 | 0.093 |
| 5 | doctor | 0.097 | 0.093 | 0.035 | 0.057 | 0.044 | 0.052 | 0.046 | 0.055 | 0.072 | 0.055 | 0.02 |
| 6 | king | 0.013 | 0.059 | 0.024 | 0.032 | 0.038 | 0.078 | 0.052 | 0.067 | 0.05 | 0.087 | 0.033 |
| 7 | queen | 0.087 | 0.072 | 0.029 | 0.042 | 0.05 | 0.083 | 0.095 | 0.012 | 0.098 | 0.009 | 0.076 |
| 8 | bird | 0.042 | 0.044 | 0.006 | 0.003 | 0.003 | 0.082 | 0.034 | 0.024 | 0.003 | 0.05 | 0.04 |

# Word embeddings

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

Tolga Bolukbasi[1], Kai-Wei Chang[2], James Zou[2], Venkatesh Saligrama[1,2], Adam Kalai[2]

| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | the | 0.056 | 0.043 | 0.051 | 0.08 | 0.006 | 0.041 | 0.032 | 0.011 | 0.057 | 0.004 | 0.083 |
| 2 | cat | 0.072 | 0.076 | 0.1 | 0.085 | 0.055 | 0.082 | 0.058 | 0.017 | 0.011 | 0.062 | 0.027 |
| 3 | dog | 0.088 | 0.099 | 0.028 | 0.059 | 0.06 | 0.059 | 0.039 | 0.09 | 0.001 | 0.031 | 0.071 |
| 4 | nurse | 0.03 | 0.018 | 0.058 | 0.074 | 0.055 | 0.028 | 0.025 | 0.054 | 0.094 | 0.052 | 0.093 |
| 5 | doctor | 0.097 | 0.093 | 0.035 | 0.057 | 0.044 | 0.052 | 0.046 | 0.055 | 0.072 | 0.055 | 0.02 |
| 6 | king | 0.013 | 0.059 | 0.024 | 0.032 | 0.038 | 0.078 | 0.052 | 0.067 | 0.05 | 0.087 | 0.033 |
| 7 | queen | 0.087 | 0.072 | 0.029 | 0.042 | 0.05 | 0.083 | 0.095 | 0.012 | 0.098 | 0.009 | 0.076 |
| 8 | bird | 0.042 | 0.044 | 0.006 | 0.003 | 0.003 | 0.082 | 0.034 | 0.024 | 0.003 | 0.05 | 0.04 |

# Embedding geometry: proximity and parallelism

Slide created by Adam Kalai

**nurse** ('nərs) n., pl., -s **1.** A woman trained to care for the sick or infirm, especially in a hospital.

**computer programmer** (kəmˈpjuːtə ˈprəʊɡræmə) n., pl., -s **1.** A man who writes programs for the operation of computers, especially as an occupation.

BAD because **compounds** biases
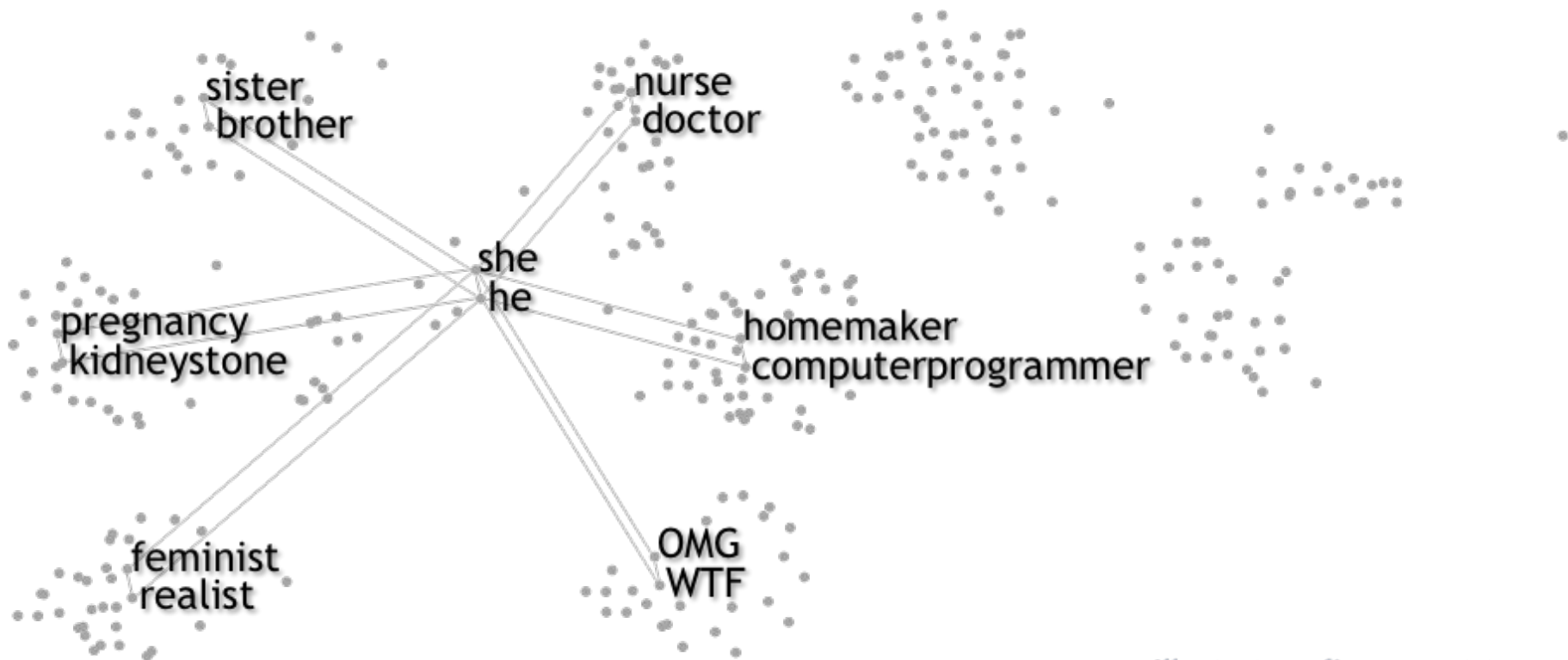
# Word embeddings

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

Tolga Bolukbasi[1], Kai-Wei Chang[2], James Zou[2], Venkatesh Saligrama[1,2], Adam Kalai[2]

|   | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | the | 0.056 | 0.000 | 0.051 | 0.08 | 0.006 | 0.041 | 0.032 | 0.011 | 0.057 | 0.004 | 0.083 |
| 2 | cat | 0.072 | 0.000 | 0.1 | 0.085 | 0.055 | 0.082 | 0.058 | 0.017 | 0.011 | 0.062 | 0.027 |
| 3 | dog | 0.088 | 0.000 | 0.028 | 0.059 | 0.06 | 0.059 | 0.039 | 0.09 | 0.001 | 0.031 | 0.071 |
| 4 | nurse | 0.03 | 0.000 | 0.058 | 0.074 | 0.055 | 0.028 | 0.025 | 0.054 | 0.094 | 0.052 | 0.093 |
| 5 | doctor | 0.097 | 0.000 | 0.035 | 0.057 | 0.044 | 0.052 | 0.046 | 0.055 | 0.072 | 0.055 | 0.02 |
| 6 | king | 0.013 | 0.000 | 0.024 | 0.032 | 0.038 | 0.078 | 0.052 | 0.067 | 0.05 | 0.087 | 0.033 |
| 7 | queen | 0.087 | 0.000 | 0.029 | 0.042 | 0.05 | 0.083 | 0.095 | 0.012 | 0.098 | 0.009 | 0.076 |
| 8 | bird | 0.042 | 0.000 | 0.006 | 0.003 | 0.003 | 0.082 | 0.034 | 0.024 | 0.003 | 0.05 | 0.04 |

Slide created by Adam Kalai

# Word embeddings

**What are the biases in my word embedding?**
**(beyond gender bias)**

# Implicit Association Test

*[Greenwald'98]*

Implicit association between categories?

# Implicit Association Test
*[Greenwald'98]*

Implicit association between categories?

# Implicit Association Test
*[Greenwald'98]*

Female

Male

**Setting 1**

Career

Family

# Implicit Association Test
*[Greenwald'98]*

Female                                              Male

Career                                              Family

Salary

# Implicit Association Test
*[Greenwald'98]*

Female                                                          Male

Career                                                          Family

Paul

# Implicit Association Test
*[Greenwald'98]*

Female                                                        Male

Career                                                        Family

Emily

# Implicit Association Test
*[Greenwald'98]*

Female                                                          Male

Career                                                          Family

Wedding

# Implicit Association Test
*[Greenwald'98]*

Female

**Setting 2**

Male

Family

Career

# Implicit Association Test
*[Greenwald'98]*

Female                                          Male

Family                                          Career




                        Salary

# Implicit Association Test
*[Greenwald'98]*

Female                                        Male

Family                                        Career

Emily

# Implicit Association Test
*[Greenwald'98]*

Female                                        Male

Family                                        Career

Wedding

# Implicit Association Test
*[Greenwald'98]*

Female                                                          Male

Family                                                          Career



John

# Implicit Association Test
*[Greenwald'98]*

Differences in average response time between **setting 1** and **setting 2**?

# Word embedding Association Test
*[Caliskan et al, 2017]*

# Word embedding Association Test
*[Caliskan et al, 2017]*



Differences in average distances between groups of words?

# Word embedding Association Test
*[Caliskan et al, 2017]*



1. Which sets $X_1$, $X_2$, $A_1$, $A_2$ should we consider?

2. How to deal with the combinatorial explosion that arises when considering intersectional groups?

# Word embedding Association Test

*[Caliskan et al, 2017]*



| Is bias X in my word embedding? *[Caliskan'17]* | What are the biases in my word embedding? [Swinger* and De-Arteaga* et al, AIES, 2019] |

**Unsupervised bias enumeration**

# Generalized Word embedding Association Test
*[Swinger* and De-Arteaga* et al, 2018]*

$$g(X_1, A_1, \ldots, X_n, A_n) \stackrel{\text{def}}{=} \sum_{i=1}^{n} (\overline{X}_i - \mu) \cdot (\overline{A}_i - \overline{\mathcal{A}})$$

$$\text{where } \mu \stackrel{\text{def}}{=} \begin{cases} \overline{\mathcal{X}} & \text{for } n = 1, \\ \sum_i \overline{X}_i / n & \text{for } n \geq 2. \end{cases}$$

# Generalized Word embedding Association Test

*[Swinger* and De-Arteaga* et al 2018]*

$$g(X_1, A_1, \ldots, X_n, A_n) \overset{\text{def}}{=} \sum_{i=1}^{n} (\overline{X}_i - \boldsymbol{\mu}) \cdot (\overline{A}_i - \overline{\mathcal{A}})$$

$$\text{where } \boldsymbol{\mu} \overset{\text{def}}{=} \begin{cases} \overline{\mathcal{X}} & \text{for } n = 1, \\ \sum_i \overline{X}_i / n & \text{for } n \geq 2. \end{cases}$$

**n=2** ⟶ **Lemma 1.** *For any* $X_1, X_2 \subseteq \mathcal{X}$ *of equal sizes* $|X_1| = |X_2|$ *and any nonempty* $A_1, A_2 \subseteq \mathcal{A}$,

$$s(X_1, A_1, X_2, A_2) = 2|X_1| \, g(X_1, A_1, X_2, A_2)$$

# Generalized Word embedding Association Test

*[Swinger* and De-Arteaga* et al 2018]*

$$g(X_1, A_1, \ldots, X_n, A_n) \overset{\text{def}}{=} \sum_{i=1}^{n} (\overline{X}_i - \boldsymbol{\mu}) \cdot (\overline{A}_i - \overline{\mathcal{A}})$$

$$\text{where } \boldsymbol{\mu} \overset{\text{def}}{=} \begin{cases} \overline{\mathcal{X}} & \text{for } n = 1, \\ \sum_i \overline{X}_i / n & \text{for } n \geq 2. \end{cases}$$

**Lemma 1.** *For any $X_1, X_2 \subseteq \mathcal{X}$ of equal sizes $|X_1| = |X_2|$ and any nonempty $A_1, A_2 \subseteq \mathcal{A}$,*

$$s(X_1, A_1, X_2, A_2) = 2|X_1| \, g(X_1, A_1, X_2, A_2)$$

n=1 ⟶ **Lemma 2.** *For any nonempty sets $X \subset \mathcal{X}$, $A \subset \mathcal{A}$, let their complements sets $X^c = \mathcal{X} \setminus X$ and $A^c = \mathcal{A} \setminus A$. Then,*

$$g(X, A) = 2g(X, A, \mathcal{X}, \mathcal{A}) = 2 \frac{|X^c|}{|\mathcal{X}|} \frac{|A^c|}{|\mathcal{A}|} g(X, A, X^c, A^c)$$

# Generalized Word embedding Association Test
*[Swinger* and De-Arteaga* et al 2018]*

$$g(X_1, A_1, \ldots, X_n, A_n) \stackrel{\text{def}}{=} \sum_{i=1}^{n} (\overline{X}_i - \mu) \cdot (\overline{A}_i - \overline{\mathcal{A}})$$

$$\text{where } \mu \stackrel{\text{def}}{=} \begin{cases} \overline{\mathcal{X}} & \text{for } n = 1, \\ \sum_i \overline{X}_i / n & \text{for } n \geq 2. \end{cases}$$

**Lemma 1.** *For any $X_1, X_2 \subseteq \mathcal{X}$ of equal sizes $|X_1| = |X_2|$ and any nonempty $A_1, A_2 \subseteq \mathcal{A}$,*

$$s(X_1, A_1, X_2, A_2) = 2|X_1|\, g(X_1, A_1, X_2, A_2)$$

**Lemma 2.** *For any nonempty sets $X \subset \mathcal{X}$, $A \subset \mathcal{A}$, let their complements sets $X^c = \mathcal{X} \setminus X$ and $A^c = \mathcal{A} \setminus A$. Then,*

$$g(X, A) = 2g(X, A, \mathcal{X}, \mathcal{A}) = 2 \frac{|X^c|}{|\mathcal{X}|} \frac{|A^c|}{|\mathcal{A}|} g(X, A, X^c, A^c)$$

n>1
(decompositi

**Lemma 3.** *For any $n > 1$ and nonempty $X_1, X_2, \ldots, X_n \subseteq \mathcal{X}$ and $A_1, A_2, \ldots, A_n \subseteq \overline{\mathcal{A}}$,*
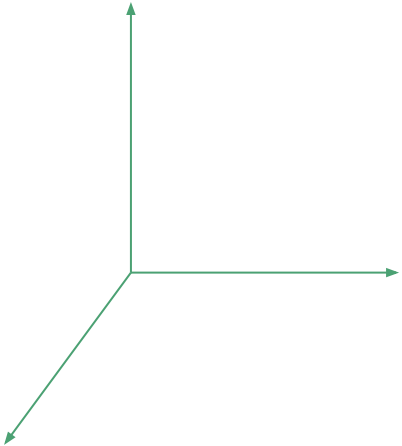
$$g(X_1, A_1, \ldots, X_n, A_n) = \sum_{i \in [n]} g(X_i, A_i) - \sum_{i,j \in [n]} \frac{g(X_i, A_j)}{n}$$

# Unsupervised Bias Enumeration (UBE) algorithm

| name | meaning | default |
|------|---------|---------|
| $WE$ | word embedding | w2v |
| $\mathcal{X}$ | set of names | SSA |
| $n$ | number of target groups | 12 |
| $m$ | number of categories | 64 |
| $M$ | number of frequent lower-case words | 30,000 |
| $t$ | number of words per WEAT | 3 |
| $\alpha$ | false discovery rate | 0.05 |

Attributes →

# Input

Step 1: Discover groups

Markisha
Latisha
Tyrique
$X_3$

Amanda
Erika
Zoe
$X_1$

Yael
Moses
Michal
$X_2$

Step 1: Discover groups

| w2v F1 | w2v F2 | w2v F3 | w2v F4 | w2v F5 | w2v F6 | w2v F7 | w2v F8 | w2v F9 | w2v F10 | w2v F11 | w2v F12 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Amanda | Janice | Marquisha | Mia | Kayla | Kamal | Daniela | Miguel | Yael | Randall | Dashaun | Keith |
| Renee | Jeanette | Latisha | Keva | Carsyn | Nailah | Lucien | Deisy | Moses | Dashiell | Jamell | Gabe |
| Lynnea | Lenna | Tyrique | Hillary | Aislynn | Kya | Marko | Violeta | Michal | Randell | Marlon | Alfred |
| Zoe | Mattie | Marygrace | Penelope | Cj | Maryam | Emelie | Emilio | Shai | Jordan | Davonta | Shane |
| Erika | Marylynn | Takiyah | Savanna | Kaylei | Rohan | Antonia | Yareli | Yehudis | Chace | Demetrius | Stan |
| +581 | +840 | +692 | +558 | +890 | +312 | +391 | +577 | +120 | +432 | +393 | +494 |

# Step 1: Discover groups

| w2v F1 | w2v F2 | w2v F3 | w2v F4 | w2v F5 | w2v F6 | w2v F7 | w2v F8 | w2v F9 | w2v F10 | w2v F11 | w2v F12 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Amanda | Janice | Marquisha | Mia | Kayla | Kamal | Daniela | Miguel | Yael | Randall | Dashaun | Keith |
| Renee | Jeanette | Latisha | Keva | Carsyn | Nailah | Lucien | Deisy | Moses | Dashiell | Jamell | Gabe |
| Lynnea | Lenna | Tyrique | Hillary | Aislynn | Kya | Marko | Violeta | Michal | Randell | Marlon | Alfred |
| Zoe | Mattie | Marygrace | Penelope | Cj | Maryam | Emelie | Emilio | Shai | Jordan | Davonta | Shane |
| Erika | Marylynn | Takiyah | Savanna | Kaylei | Rohan | Antonia | Yareli | Yehudis | Chace | Demetrius | Stan |
| +581 | +840 | +692 | +558 | +890 | +312 | +391 | +577 | +120 | +432 | +393 | +494 |
| 98% F | 98% F | 89% F | 85% F | 78% F | 65% F | 59% F | 56% F | 40% F | 27% F | 5% F | 4% F |
| 1983 | 1968 | 1978 | 1982 | 1993 | 1991 | 1985 | 1986 | 1989 | 1981 | 1984 | 1976 |
| 4% B | 8% B | 48% B | 10% B | 2% B | 7% B | 4% B | 2% B | 5% B | 10% B | 32% B | 6% B |
| 4% H | 4% H | 3% H | 9% H | 1% H | 4% H | 9% H | 70% H | 10% H | 3% H | 5% H | 3% H |
| 3% A | 3% A | 1% A | 11% A | 1% A | 32% A | 4% A | 8% A | 5% A | 4% A | 3% A | 5% A |
| 89% W | 84% W | 47% W | 69% W | 95% W | 56% W | 83% W | 21% W | 79% W | 83% W | 59% W | 86% W |

# Step 1: Discover groups

Step 2: Discover word categories
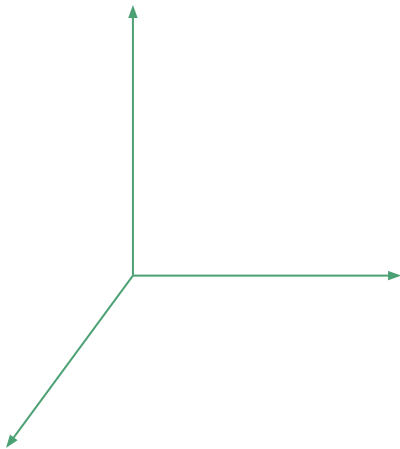
Step 2: Discover word categories

$X_3$

$X_1$

$A_j$

$X_2$

Step 3: Partition $A_j$

$$V_{ij} = \left\{ w \in \mathcal{A}_j \mid i = \arg \max_{i' \in [n]} \overline{w} \cdot \overline{X}_{i'} \right\}$$

**X₃**

*Kosher*

*Tortillas*

*Hummus*

*Tequila*

**A_j**

*Caviar*

**X₁**

**X₂**

Step 3: Partition A_j

A$_{i,j}$ contains top $t$ words s.t.

$$\max_{w \in V_{ij}} (\overline{X}_i - \mu) \cdot (\overline{w} - \overline{\mathcal{A}}_j)$$

$$V_{ij} = \left\{ w \in \mathcal{A}_j \mid i = \arg\max_{i' \in [n]} \overline{w} \cdot \overline{X}_{i'} \right\}$$

**X$_3$**

**A$_{3,j}$**

**A$_{1,j}$**

**X$_1$**

**A$_{2,j}$**

**X$_2$**

Step 3: Partition A$_j$

Is $A_{i,j}$ significantly closer to $X_i$ than it could be expected through sheer randomness?

$X_3$

$A_{3,j}$

$A_{1,j}$

$X_1$

$A_{2,j}$

$X_2$

Step 4: Establish statistical significance

Step 4: Establish statistical significance

Step 4: Establish statistical significance

Step 4: Establish statistical significance

$$\sigma_{ij} = (\overline{X}_i - \mu) \cdot (\overline{A}_{ij} - \overline{\mathcal{A}})$$

Step 4: Establish statistical significance

$$\sigma_{ij} = (\overline{X}_i - \mu) \cdot (\overline{A}_{ij} - \overline{\mathcal{A}})$$

Is $\sigma_{i,j}$ significantly large?

$\overline{X_i}$

$\overline{A_{i,j}}$

$\overline{A}$

Step 4: Establish statistical significance

**Rotational null hypothesis**

1. Rotate X: X → XUr



Step 4: Establish statistical significance

**Rotational null hypothesis**

2. Find $A_{i,j,r}$

$X_3$

$A_{3,j,r}$

$A_{1,j,r}$

$A_{2,j,r}$

$X_1$

$X_2$

Step 4: Establish statistical significance

**Rotational null hypothesis**

3. Calculate $\sigma_{i,j,r}$



Step 4: Establish statistical significance

**Rotational null hypothesis**

3.    Calculate $\sigma_{i,j,r}$



Step 4: Establish statistical significance

**Rotational null hypothesis**

3.  Calculate p-value:

$$p_{i,j} = [\, \delta(\sigma_{i,j} > \sigma_{i,j,r}) + 1 \,] \; / \; [\, R + 1 \,]$$

r = 1,2,...,10k



# Step 4: Establish statistical significance

**Rotational null hypothesis**

4.  Determine critical p-value, $\alpha$-bound guarantee on false discovery rate (*Benjamini-Hochbergh*)



Step 4: Establish statistical significance

# Disclaimer

The biases in the following slides contain offensive stereotypes.

These do not reflect our views.

| 98% F | 98% F | 89% F | 85% F | 78% F | 65% F | 59% F | 56% F | 40% F |
|---|---|---|---|---|---|---|---|---|
| 1983 | 1968 | 1978 | 1982 | 1993 | 1991 | 1985 | 1986 | 1989 |
| 4% B | 8% B | 48% B | 10% B | 2% B | 7% B | 4% B | 2% B | 5% B |
| 4% H | 4% H | 3% H | 9% H | 1% H | 4% H | 9% H | 70% H | 10% H |
| 89% W | 84% W | 47% W | 69% W | 95% W | 56% W | 83% W | 21% W | 79% W |
| 3% A | 3% A | 1% A | 11% A | 1% A | 32% A | 4% A | 8% A | 5% A |
|  | cookbook, baking, baked goods | sweet potatoes, macaroni, green beans |  |  | saffron, halal, sweets | mozzarella, foie gras, caviar | tortillas, salsa, tequila | kosher, humm[us], bagel |
| herself, hers, moms | husband, homebound, grandkids | aunt, niece, grandmother | hubby, socialite, cuddle | twin sister, girls, classmate | elder brother, dowry, refugee camp |  |  | berea[ved], immig[rant], emigr[...] |
| hostess, cheerleader, dietitian | registered nurse, homemaker, chairwoman |  | supermodel, beauty queen, stripper | helper, getter, snowboarder | shopkeeper, villager, cricketer |  | translator, interpreter, smuggler |  |
|  | log cabin, library, fairgrounds | front porch, carport, duplex | racecourse, plush, tenements | picnic tables, bleachers, concession stand | locality, mosque, slum | prefecture, chalet, sauna |  | synag[ogue], constr[...], hilltop |
|  | parish, | pastor, | goddess, |  | fatwa, | monastery, | rosary, | rabbis |

# Crowdsourcing evaluation

**Qualification**:
36 names, 3 per group
+1 per name labeled in correct group

# Crowdsourcing evaluation

**Qualification**:
36 names, 3 per group
+1 per name labeled in correct group

*If accuracy > 50%*

**Is the UBE output consistent with society's stereotypes?**
For each WEAT:
- Groups in output {X1, X2, … , Xk} and {A1, A2, …, Ak} shown
- For each name group Xi, which group Ai contains words most stereotypically associated with these names?

# Crowdsourcing evaluation

**Qualification**:

36 names, 3 per group
+1 per name labeled in correct group

*If accuracy > 50%*

**Is the UBE output consistent with society's stereotypes?**

For each WEAT:
- Groups in output {X1, X2, ... , Xk} and {A1, A2, ..., Ak} shown
- For each name group Xi, which group Ai contains words most stereotypically associated with these names?

*If most commonly chosen group matches UBE pairing*

**Is it offensive? Rate:**

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |

*Politically correct, inoffensive, or just random*

*Politically incorrect, possibly very offensive*

# Crowdsourcing evaluation

| Emb. | # significant | % accurate | % offensive |
|---|---|---|---|
| w2v | 235 | 72% | 35% |
| fast | 160 | 80% | 38% |
| glove | 442 | 48% | 24% |

# Disclaimer

The biases in the following slides contain offensive stereotypes.

These do not reflect our views or the views of crowd workers.

# Crowdsourcing evaluation

| Emb. | # significant | % accurate | % offensive |
|------|---------------|------------|-------------|
| w2v | 235 | 72% | 35% |
| fast | 160 | 80% | 38% |
| glove | 442 | 48% | 24% |

| Word2Vec trained on Google news | | | fastText trained on the Web | | | GloVe trained on the Web | | |
|---|---|---|---|---|---|---|---|---|
| **Miguel** | **Dashaun** | **Kamal** | **Marquell** | **Ahmed** | **Alejandra** | **Amina** | **Alejandra** | **Kylee** |
| **Deisy** | **Jamell** | **Nailah** | **Antwan** | **Shanti** | **Maricella** | **Yair** | **Epifanio** | **Shaye** |
| **Violeta** | **Marlon** | **Kya** | **Dakari** | **Mariyah** | **Ona** | **Rani** | **Monalisa** | **Tayla** |
| **Emilio** | **Davonta** | **Maryam** | **Pernell** | **Siddharth** | **Fabiola** | **Danial** | **Eulalia** | **Latasha** |
| **Yareli** | **Demetrius** | **Rohan** | **Jarred** | **Yasmin** | **Sulema** | **Safa** | **Alicea** | **Jessi** |
| illegal immigrant | aggravated robbery | subcontinent | n***** | jihad | s****** | turban | cartel | pornstar |
| drug trafficking | aggravated assault | tribesmen | f***** | militants | maid | saree | undocumented | hottie |
| deported | felonious assault | miscreants | dreads | caliphate | busty | hijab | culpable | nubile |

*These associations do not reflect our views or those of the crowd workers.

# Why does this matter?

- Representational harms

- Harmful bias encoded in semantic representation used for learning

- **Removing names is not enough to get rid of bias!**

  - Words in category clusters may be used as proxy for gender/race/etc

*Hostess*

*Cab driver*

*volleyball*

*cornerback*

# In this talk...

## What are the risks of semantic representation bias?



## Part 2: Allocative harms

**Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting (FAT* 2019)**
Maria De-Arteaga (CMU), Alexey Romanov (UMASS), Hanna Wallach (MSR), Jennifer Chayes (MSR), Christian Borgs (MSR), Alexandra Chouldechova (CMU), Sahin Geyik (LinkedIn), Krishnaram Kenthapadi (LinkedIn), Adam Kalai (MSR)

# An artificially intelligent headhunter?

**Forbes**

Billionaires    Innovation    Leadership    Money    Consu

Forbes CommunityVoice    Connecting expert communities to the Forbes audience.    What is

5,220 views | Jul 12, 2018, 07:00am

## Welcome To The Age Of Recruiting Automation

**n p r**

SCIENCE

Now Algorithms Are Deciding Whom To Hire, Based On Voice

4:27

**CNBC** MENU

Get ready, this year your next job interview may be with an A.I. robot

77

# An artificially intelligent headhunter?



**Forbes**

Billion...

Forbes Commu...

6,220 views

**n p r**

Wel...

Recr...

SCIENCE

Now Algorithms Are Dec...
Hire, Based On Voice

4:27

**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.

**CNBC**

Get ready, this year your next job interview may be with an A.I. robot
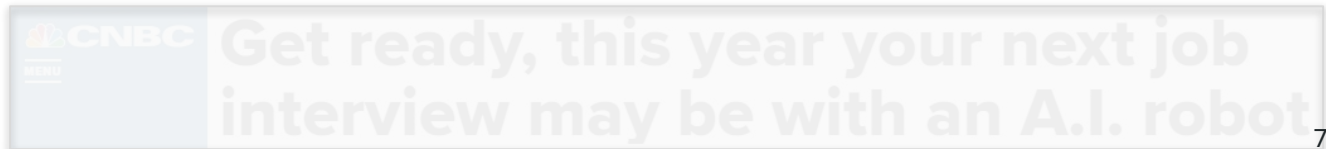
78

# 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.

79

Can we **quantify the risks** of incorporating **ML** in **hiring and recruiting pipelines**?

Can we **quantify the risks** of incorporating **ML** in **hiring and recruiting pipelines**?

Can we **characterize** the **effects**?

Can we **quantify the risks** of incorporating **ML** in **hiring and recruiting pipelines**?

Can we **characterize** the **effects**?

**Our findings**:
- Gender accuracy gap in large-scale study
- "Scrubbing" gender indicators ≠ gender blindness
- Compounding imbalances

# Computer Programmer 🔍

# Computer Programmer

Slide created by Adam Kalai

# Computer Programmer

# Computer Programmer



**SOFTWARE ENGINEER** JANE_DOE.ORG

**JANE DOE**

**OBJECTIVE**

Writing solid software for meaningful applications that have a positive impact on the world.

**EXPERIENCE**

**DEVELOPER • MICROSOFT • 2007-2013**
Wrote software for cloud platform involving distributed computing, databases, and logging.

**BLACK FEMALE**

**LEADERSHIP**

**SOFTBALL TEAM CAPTAIN • SPELMAN C**003
Led team to division champi ponsible for coordinating

Java, Python, C++, SQL,

# Computer Programmer

**Adam Kalai**

Principal Researcher

Contact Info
- (857) 453-6323
- Email

Microsoft Research
Office 12038
Cambridge, MA 02142

About

I have been...
fun problem...
accessibility...
be less bias...

Previously...
fortunate to...
followed by...

**Sahin Cem Geyik**

Computer Science Department
Rensselaer Polytechnic Institute
TROY, NY, 12180
email: sahincem2

**Krishnaram Kenthapadi**

Krishnaram Kenthapadi is part of the AI team at LinkedIn, where he leads the fairness, transparency, explainability (AETHER) Committee. He shaped the technical roadmap and led the privacy/modeling efforts for LinkedIn Salary intersection of members, recruiters, and career opportunities. Previously, he was a Researcher at Microsoft Research Science from Stanford University in 2006, under the supervision of Professor Rajeev Motwani. Before joining Stan

Krishnaram's expertise is in the areas of fairness/transparency/explainability/privacy in AI/ML systems, algorithms 17+ years of experience (including 12+ years in industry after his PhD), working on challenging problems in these fairness/privacy, and improved business metrics for existing products via technology transfers. He has collaborated his fields of interest. He serves regularly on the program committees of KDD, WWW, WSDM, and related conference best case studies paper award, SODA best student paper award, and WWW best paper award nomination. He has t

**I have successfully completed and**
**Started working at Turn Inc. as an Applied Scientist.**

**Alexey Romanov**

A Ph.D. Student at UMass

## Hello

I am currently a second year Ph.D. student at UMass Lowell in the Text-Machine Lab working with Anna Rumshisky. My research interests at this moment are particularly focused on applying Deep Learning methods in Natural Language Processing.

# Jennifer Chayes

About  Projects  Publications  Videos

Jennifer Tour Chayes is Technical Fellow and Managing Director of Microsoft Research New England in Cambridge, Massachusetts, which she co-founded in 2008, and Microsoft Research New York City, which she co-founded in 2012, and Microsoft Research Montreal since 2017. These three laboratories are widely renowned interdisciplinary centers, bringing together computer scientists, mathematicians, physicists, social scientists, and biologists, and helping to lay the foundations of data science. Prior to founding these labs, Chayes was Research Area Manager for Mathematics, Theoretical Computer Science, and Cryptography at Microsoft Research Redmond. Chayes joined Microsoft Research in 1997, when she co-founded the Theory Group. Her research areas include phase transitions in discrete mathematics and computer science, structural and dynamical properties of large networks, mechanism design, and graph algorithms. She is the co-author of about scientific papers and the co-inventor of about 30 patents.

# Christian Borgs

Deputy Managing Director,
Microsoft Research New England

Contact Info
- Website

Research areas
Mathematics

About  Projects  Publications  Videos

Christian Borgs is deputy managing director and co-fo... ...search New England in Cambridge, Massachusetts.

**Hanna Wallach**

Principal Researcher

Contact Info
- Email
- Website
- Twitter

...Researcher at Microsoft Research New York City and an Adjunct Professor in the College of In... ...ences at the University of Massachusetts Amherst. She is also a member of UMass's Comp... ...ence Institute. Hanna develops machine learning methods for analyzing the structure, content, and dynamics of... ...esses. Her work is inherently interdisciplinary: she collaborates with political scientists, sociologists, and journalists to understand how organizations work by analyzing publicly available interaction data, such as email networks, document collections, press releases, meeting transcripts, and news articles. To complement this agenda, she also studies issues of fairness, accountability, and transparency as they relate to machine learning. Hanna's research has had broad impact in machine learning, natural language processing, and computational social science. In 2014, she was named

Maria De-Arteaga    About  CT  Publications

## About

I am a fifth year PhD student in the joint Machine Learning and Public Policy program at Carnegie Mellon University's Machine Learning Department and Heinz College. I am co-advised by Prof. Artur Dubrawski and Prof. Alexandra Chouldechova, and I am part of the Auton Lab.

Currently, my main focus is algorithmic fairness, studying how to measure and prevent bias and discrimination that may arise when using machine learning for decision support. I am passionate about understanding the roadblocks that prevent the effective use of machine learning to advance global development, and conducting machine learning research to tackle those challenges.

**Alexandra Chouldechova**

Assistant Professor of Statistics and Public Policy

Heinz College, Carnegie Mellon University
Office: Hamburg Hall 2224
Email: achould(at)cmu.edu
Phone: 412-268-4414

### Education

Ph.D. in Statistics, Stanford University, 2014
B.Sc. in Mathematical Statistics, University of Toronto, 2005-2009

### Research

My research focuses on problems related to fairness in predictive modeling. I work on better understanding how to assess black-box predictors for potentially unanticipated biases that could lead to discriminatory practices. Questions that I am actively investigating include:

Under what conditions can disparate impact arise?

How can we quantitatively characterize fairness?

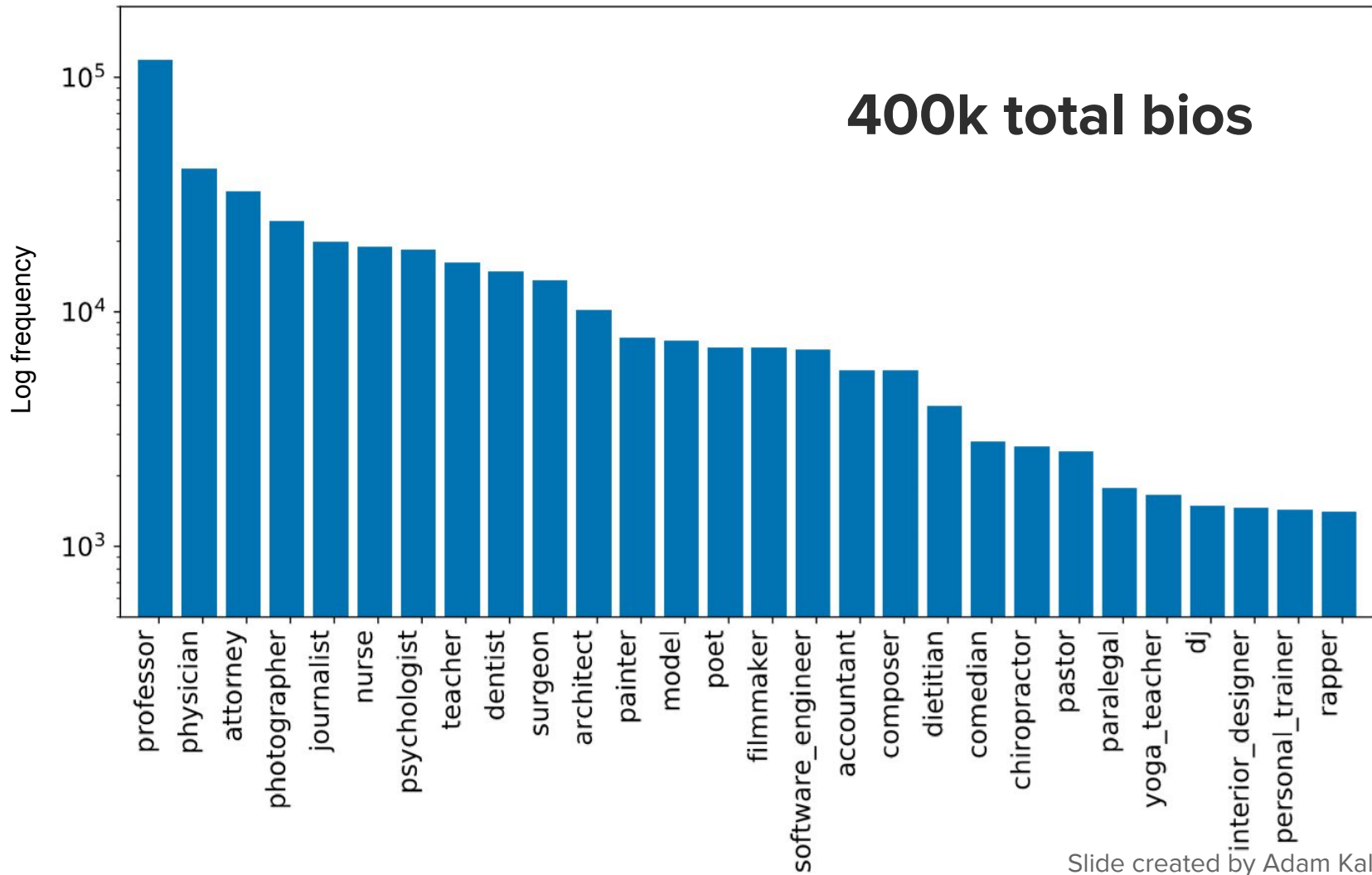How can we use such characterizations to develop improved systems that are less likely to result in disparate impact?

88

# Bias in bios: Biographies dataset

- 400,000 third-person web bios from Common Crawl.

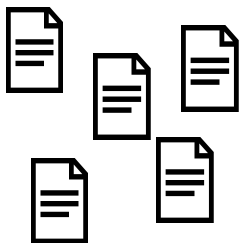*"Xxx Xxx is a(n) (xxx) [title]...he/she..."* *title* ∈ {common BLS SOC titles}

*Alexandra Chouldechova is an Assistant* **Professor** *of Statistics and Public Policy at Carnegie Mellon University's Heinz College of Informations Systems and Public Policy. She received her B.Sc. from the University of Toronto in 2009, and in 2014 she completed her Ph.D. in Statistics at Stanford University. While at Stanford, she also worked at Google and Symantec on developing statistical assessment methods for information retrieval systems.*
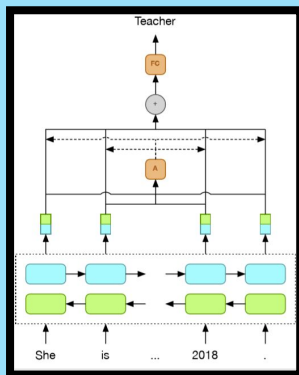
- Classification problem: 28 title-from-bio-text

89

**400k total bios**

90

# Learning pipeline

Input data:
Biographies

**Semantic representations:**

1. Bag-of-words
2. Word embedding
3. Deep neural network (GRU) with attention

Objective:
Predict Y = *Occupation*

**Gender sensitivity**: How do predictions change if explicit gender indicators are swapped?

[Bertrand, Mulliainathan'04]

# Biases in bios

**Enter the bio**

She is a fifth year PhD student in the joint Machine Learning and Public Policy program at Carnegie Mellon University's Machine Learning Department and Heinz College. She is co-advised by Prof. Artur Dubrawski and Prof. Alexandra Chouldechova, and she is part of the Auton Lab.

Currently, her main focus is algorithmic fairness, studying how to measure and prevent bias and discrimination that may arise when using machine learning for decision support. She is passionate about developing novel machine learning algorithms that are

**PREDICT TITLE**  **SHE**  **HE**

she is a fifth year phd student in the joint machine learning and public policy program at carnegie mellon university <unk> s machine learning department and heinz college . she is co-advised by prof. artur <unk> and prof. alexandra chouldechova , and she is part of the auton lab . currently , her main focus is algorithmic fairness , studying how to measure and prevent bias and discrimination that may arise when using machine learning for decision support . she is passionate about developing novel machine learning algorithms that are motivated by existing policy problems , and understanding how machine learning can better help us overcome important societal challenges . prior to graduate school she received her b.sc . in mathematics from universidad nacional de colombia and worked as a journalist for one of colombia <unk> s main news magazine , semana . she is the recipient of a microsoft

## teacher

93

# Biases in bios

*She → he*

**Enter the bio**

He is a fifth year PhD student in the joint Machine Learning and Public Policy program at Carnegie Mellon University's Machine Learning Department and Heinz College. He is co-advised by Prof. Artur Dubrawski and Prof. Alexandra Chouldechova, and he is part of the Auton Lab.

Currently, his main focus is algorithmic fairness, studying how to measure and prevent bias and discrimination that may arise when using machine learning for decision support. He is passionate about developing novel machine learning algorithms that are motivated

**PREDICT TITLE**    **SHE**    **HE**

he is a fifth year phd student in the joint machine learning and public policy program at carnegie mellon university <unk> s machine learning department and heinz college . he is co-advised by prof. artur <unk> and prof. alexandra chouldechova , and he is part of the auton lab . currently , his main focus is algorithmic fairness , studying how to measure and prevent bias and discrimination that may arise when using machine learning for decision support . he is passionate about developing novel machine learning algorithms that are motivated by existing policy problems , and understanding how machine learning can better help us overcome important societal challenges . prior to graduate school he received his b.sc . in mathematics from universidad nacional de colombia and worked as a journalist for one of colombia <unk> s main news magazine , semana . he is the recipient of a microsoft

software_engineer

94

# Biases in bios

**Enter the bio**

He is a fifth year PhD student in the joint Machine Learning and Public Policy program at Carnegie Mellon University's Machine Learning Department and Heinz College. He is co-advised by Prof. Artur Dubrawski and Prof. Alexandra Chouldechova, and he is part of the Auton Lab.

Curr

usin

**P**

he is

learni

lab . c

using

existi

gradu

colombia <unk> s main news magazine , semana . he is the recipient of a microsoft

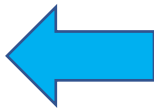| $y^1$ | $y^2$ | $\Pi_{\text{female},(y^1,y^2)}$ |
|---|---|---|
| model | rapper | 14.7% |
| teacher | pastor | 8.5% |
| professor | software engineer | 6.5% |
| professor | surgeon | 4.8% |
| physician | surgeon | 3.8% |

| $y^1$ | $y^2$ | $\Pi_{\text{male},(y^1,y^2)}$ |
|---|---|---|
| attorney | paralegal | 7.1% |
| architect | interior designer | 4.7% |
| professor | dietitian | 4.3% |
| photographer | interior designer | 3.5% |
| teacher | yoga teacher | 3.3% |

software_engineer

Beyond explicit gender indicators: **the gender accuracy gap**

**Compounding** gender imbalance

More accurate on F

More accurate on M

TPR GENDER GAP — ACCURACY ON FEMALES – ACCURACY ON MALES

% FEMALE

Compounding gender imbalance — scatter plot of ACCURACY ON FEMALES − ACCURACY (TPR GENDER GAP) vs % FEMALE, with professions labeled. Blue arrows indicate "More accurate on F" (top) and "More accurate on M" (bottom).
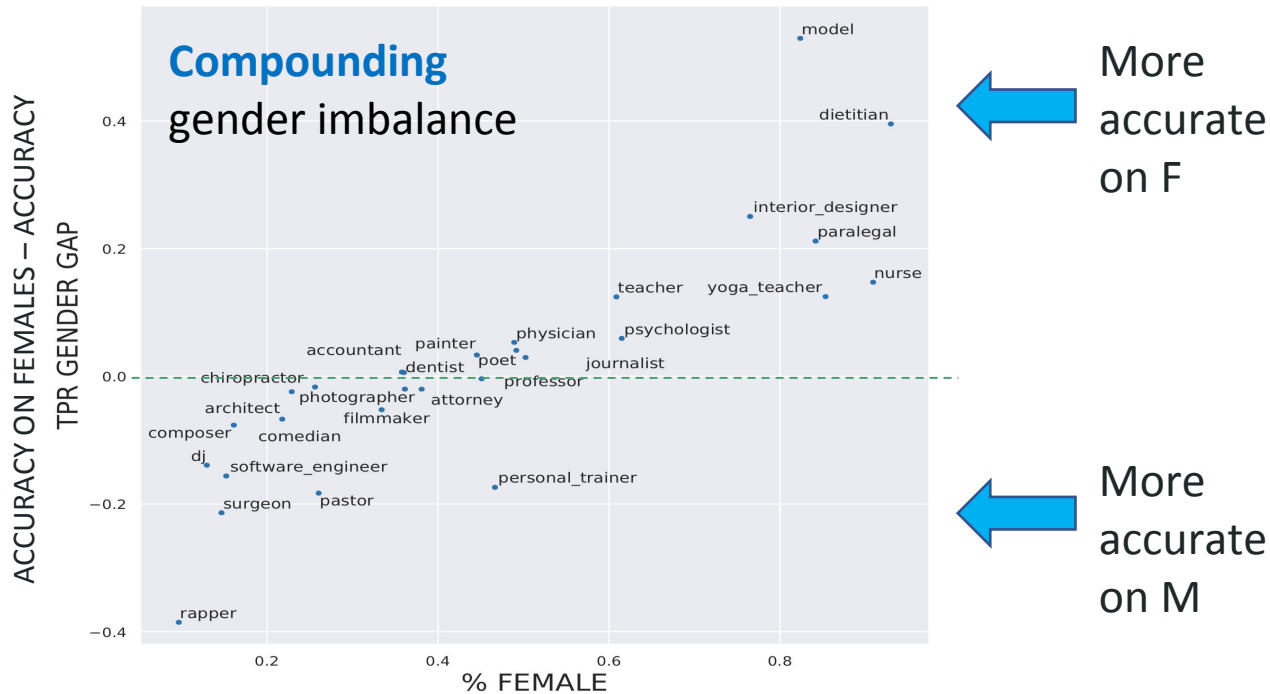
# Compounding imbalance

If female fraction p < 0.5 and gender gap < 0 for title, then female fraction in true positives < p (similarly for males)
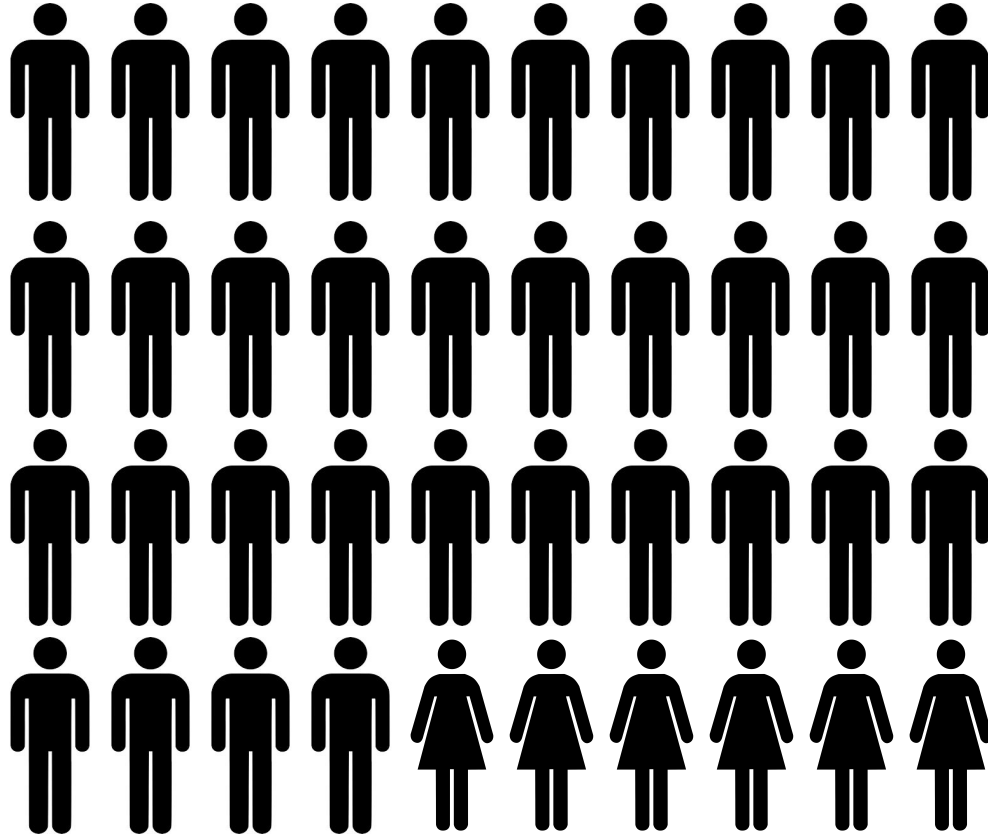
Compounding gender imbalance scatter plot. X-axis: % FEMALE. Y-axis: ACCURACY ON FEMALES − ACCURACY / TPR GENDER GAP. Labeled points include: model, dietitian, interior_designer, paralegal, nurse, teacher, yoga_teacher, physician, psychologist, accountant, painter, poet, journalist, dentist, chiropractor, professor, photographer, attorney, architect, filmmaker, composer, comedian, dj, software_engineer, personal_trainer, surgeon, pastor, rapper. Arrows: "More accurate on F" (top) and "More accurate on M" (bottom).

# Compounding injustice [Hellman'18]

If initial imbalance constitutes injustice: Model's prediction is informed by, and compounds, previous injustice

# Compounding imbalances



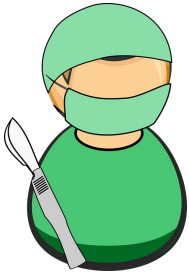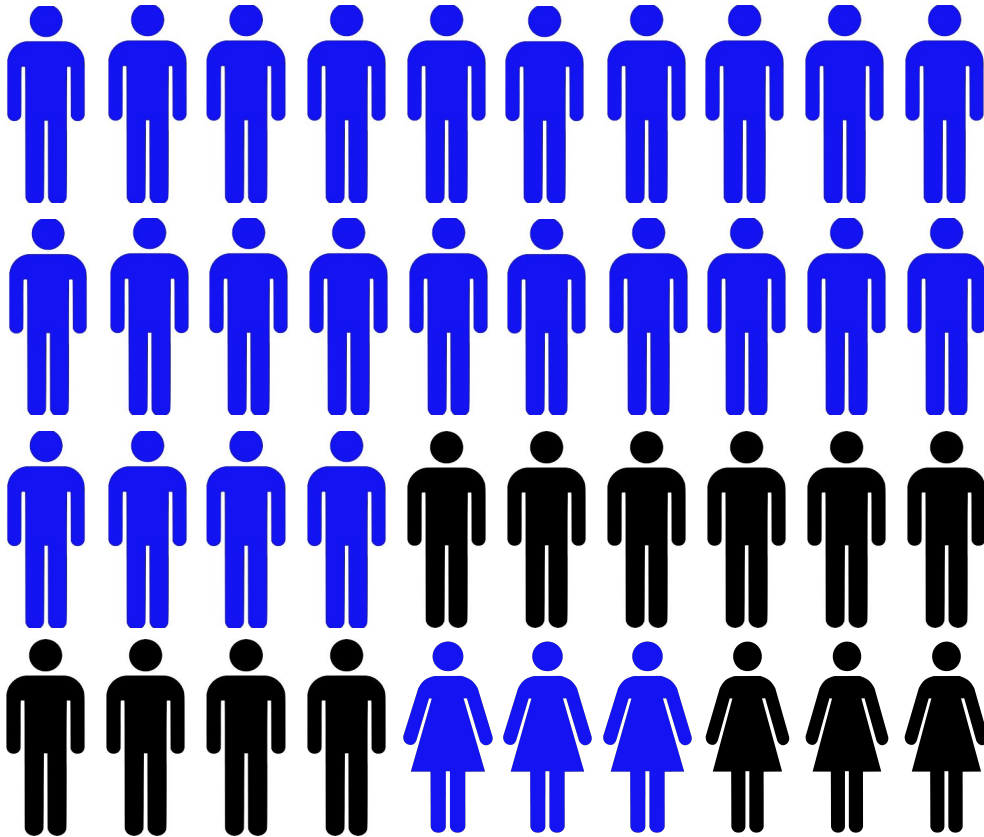**Surgeons**

females in data:
**14.6%**
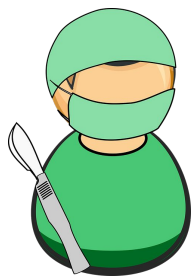
# Compounding imbalances

**Surgeons**

females in data:
**14.6%**



Males:
**71%** recall

Females:
**54%** recall

# Compounding imbalances

**Surgeons**

females in data:
**14.6%**

females in true positives:
**11.6%**

Males:
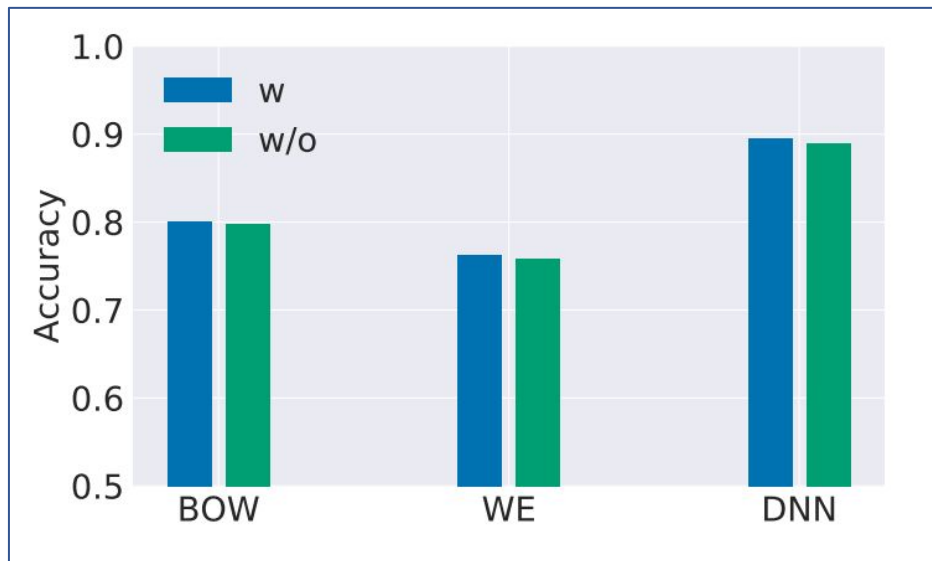**71%** recall

Females:
**54%** recall

103

# "Scrub" explicit gender indicators?

≈ same accuracy
with/without
explicit gender indicators

no scrub
scrub



104

# Compounding imbalances



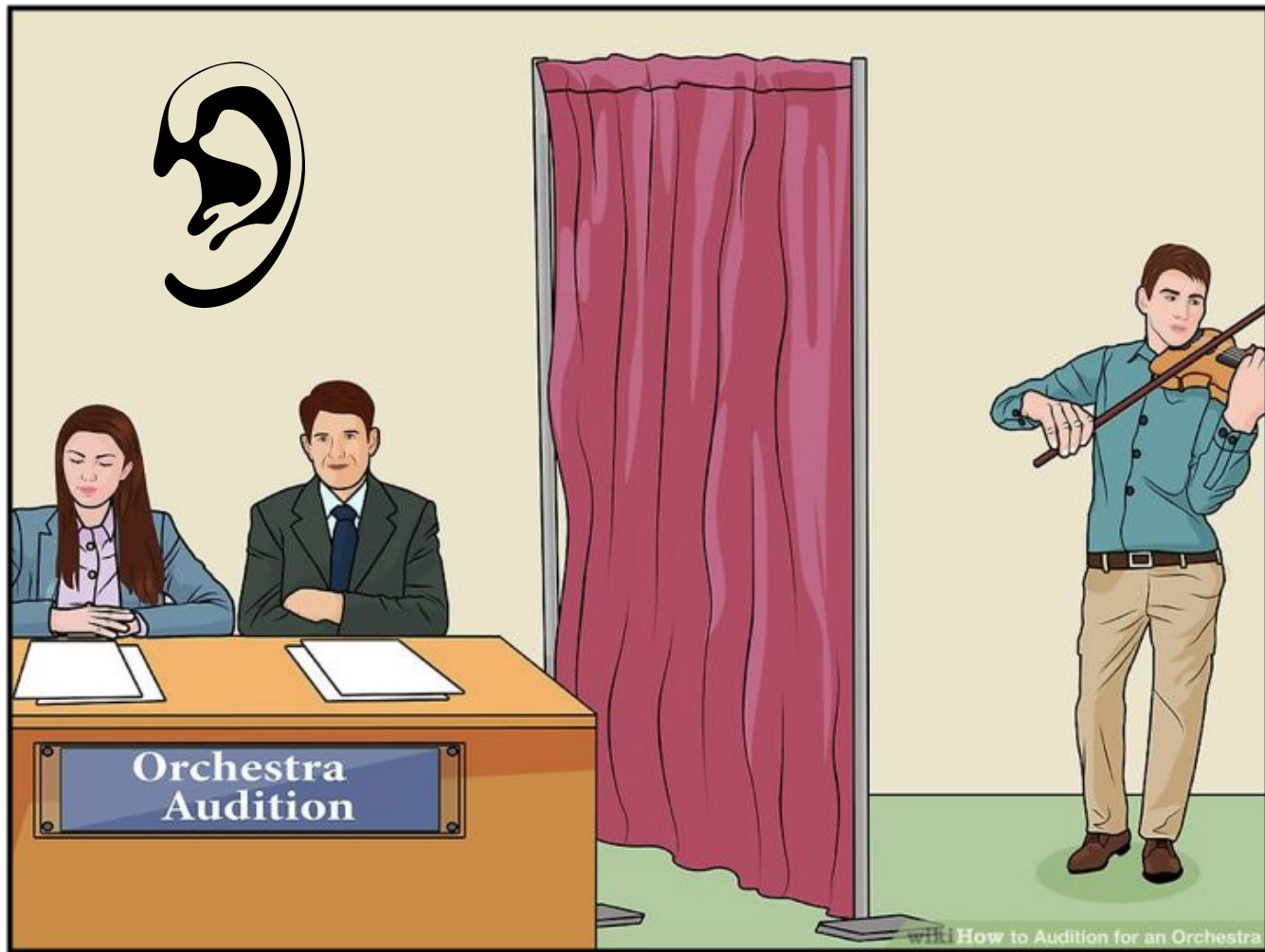Representation = BOW    Representation = WE    Representation = DNN

GAG = (ACC F) - (ACC M)

% FEMALE

—— no scrub
—— scrub

Accuracy

w
w/o

BOW    WE    DNN

105

Orchestra Audition

wikiHow to Audition for an Orchestra

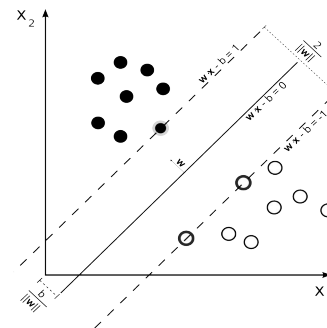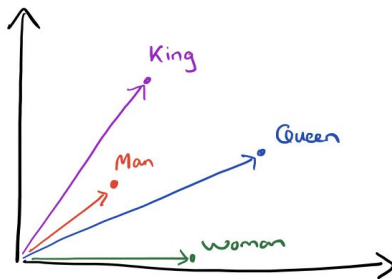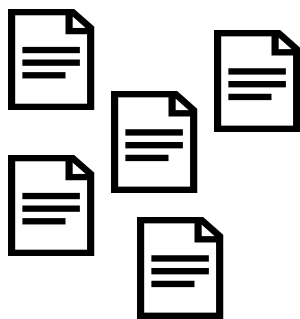Slide created by Adam Kalai

wikiHow to Audition for an Orchestra

# Can we mitigate this problem?

- Additional challenges:
  - Sensitive attributes may be unavailable, or it may be illegal to use them
  - Need to consider several attributes and their intersection
    - ➢ Race, gender, ethnicity, . . .

# In this talk…

What are the risks of semantic representation bias?



**Part 3:** Mitigating allocative harms

**What's in a Name? Reducing Bias in Bios without Access to Protected Attributes (NAACL 2019)**
Alexey Romanov (UMASS), Maria De-Arteaga (CMU), Hanna Wallach (MSR), Jennifer Chayes (MSR), Christian Borgs (MSR), Alexandra Chouldechova (CMU), Sahin Geyik (LinkedIn), Krishnaram Kenthapadi (LinkedIn), Anna Rumshisky (UMASS), Adam Kalai (MSR) Best Thematic Paper :)

Names encode societal biases, and...

"What's in a name? That which we call a rose

By any other name would smell as sweet."
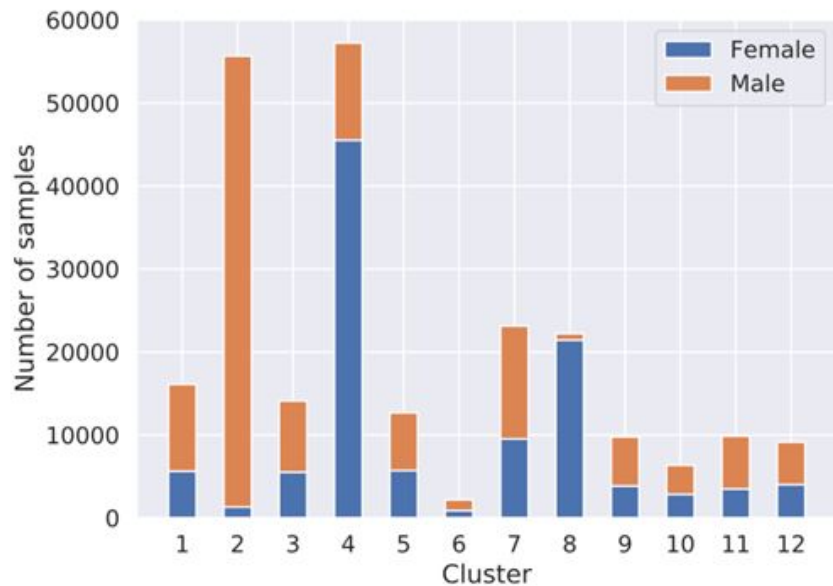
*William Shakespeare, Romeo and Juliet*

# Main idea

- Leverage biases presented in word embeddings
    - Use embeddings of names as "universal proxies"
    - No need to define protected groups
- Embeddings are used only in the loss calculation
    - No need for names or protected attributes during deployment
    - Gains extend to individuals who are poorly proxied

Credit: Alexey Romanov

# Names are indeed "universal proxies"



(a) Gender membership.

(b) Race membership.

Slide created by Alexey Romanov

# Algorithms: regularize accuracy gaps

- Training data $x_1, y_1, \dots, x_n, y_n \in X \times \{1, 2, \dots, T\}$
- Model parameters $\theta$, regularization parameter $R \geq 0$
- $\mathcal{L}(\theta)$ = standard misclassification loss, e.g., $-\frac{1}{n}\sum_i \log p_\theta(y_i|x_i)$
- Minimize $\mathcal{L}(\theta) + R \cdot \mathcal{L}_{\text{CluCL}}(\theta)$
- Cluster constrained loss: cluster names into $K$ groups
- $\mathcal{L}_{k,t}(\theta)$ = group $k$ loss for title $t$
- $\mathcal{L}_{\text{CluCL}}(\theta) = \dfrac{\sum_{j,k,t}\left(\mathcal{L}_{j,t}(\theta) - \mathcal{L}_{k,t}(\theta)\right)^2}{TK(K-1)}$
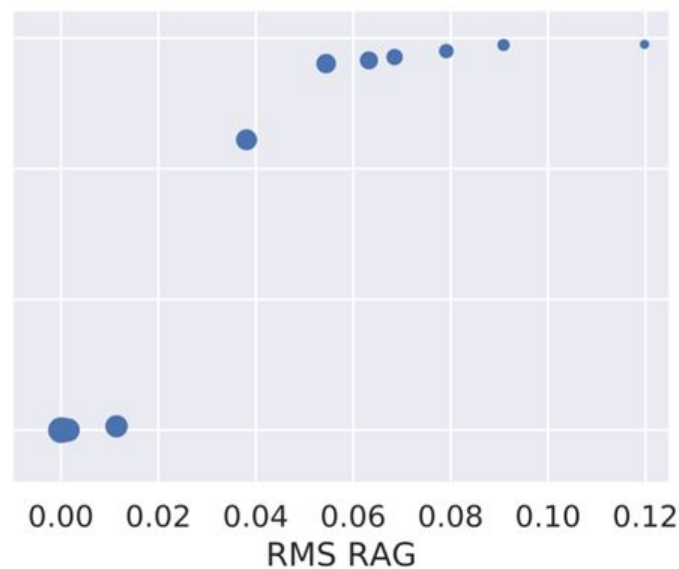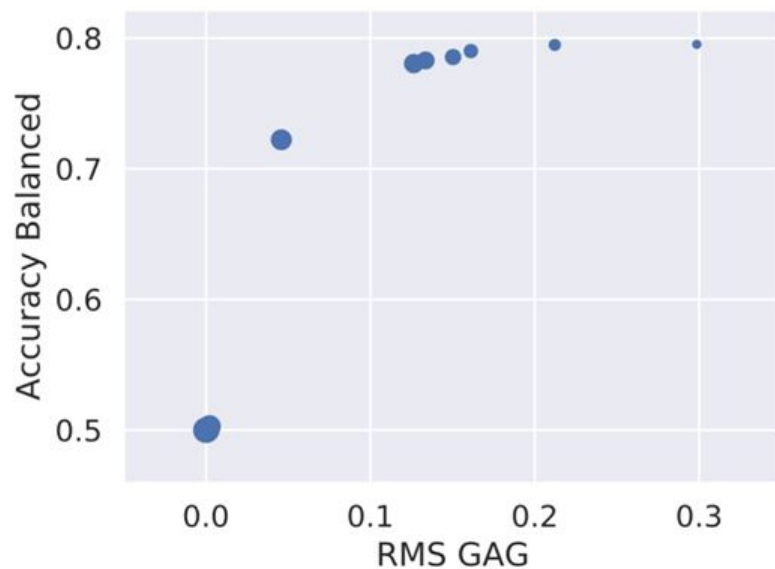
# Algorithms: regularize accuracy gaps

- Training data $x_1, y_1, \ldots, x_n, y_n \in X \times \{1, 2, \ldots, T\}$
- Model parameters $\theta$, regularization parameter $R \geq 0$
- $\mathcal{L}(\theta)$ = standard misclassification loss, e.g., $-\frac{1}{n}\sum_i \log p_\theta(y_i|x_i)$
- Minimize $\mathcal{L}(\theta) + R \cdot \mathcal{L}_{\text{CoCL}}(\theta)$

- Or balance covariance constrained loss $\mathcal{L}_{\text{CoCL}}(\theta) =$

$$\frac{1}{T}\sum_t \left\| \mathrm{E}_{i:y_i=t}\left[(v_{\text{name}_i} - \bar{v}_t) \cdot (p(t|x_i) - \bar{p}_t)\right] \right\|$$

- Intuition: minimize correlation between errors and name vectors

# Accuracy/fairness tradeoff on UCI Adult dataset



$$\text{Gap}_{r,c} = \text{TPR}_{r,c} - \text{TPR}_{\sim r,c}$$

$$\text{Gap}_r^{\text{RMS}} = \sqrt{\frac{1}{|C|} \sum_{c \in C} \text{Gap}_{r,c}^2}$$

# Bios dataset

| | | | Root Mean Square Gender Accuracy Gap | Root Mean Square Race Accuracy Gap | | |
|---|---|---|---|---|---|---|
| **Model** | **R** | **Accuracy Balanced** | **RMS GAG** | **RMS RAG** | **Max GAG** | **Max RAG** |
| Regular | 0 | **0.788** | 0.173 | 0.051 | 0.511 | 0.121 |
| CluCL | 1 | 0.784 | 0.168 | 0.048 | 0.494 | 0.120 |
| CluCL | 2 | 0.781 | **0.165** | **0.047** | **0.486** | 0.114 |
| CoCL | 1 | 0.785 | 0.168 | 0.048 | 0.507 | **0.109** |
| CoCL | 2 | 0.779 | 0.169 | 0.048 | 0.512 | 0.116 |

# UCI Adult dataset

| Model | R | Accuracy Balanced | RMS GAG | RMS RAG | Max GAG | Max RAG |
|---|---|---|---|---|---|---|
| Regular | 0 | **0.795** | 0.299 | 0.120 | 0.303 | 0.148 |
| CluCL | 1 | 0.788 | 0.278 | 0.121 | 0.297 | 0.145 |
| CluCL | 2 | 0.793 | 0.259 | 0.085 | 0.282 | 0.114 |
| CoCL | 1 | 0.794 | 0.215 | 0.091 | 0.251 | 0.119 |
| CoCL | 2 | 0.790 | **0.163** | **0.080** | **0.201** | **0.109** |

Root Mean Square Gender Accuracy Gap → RMS GAG

Root Mean Square Race Accuracy Gap → RMS RAG

# Several prominent biases are reduced



(a) The *Adult* dataset

(b) The *Bios* dataset, occupation "surgeon"

# Summary

- Unsupervised bias enumeration algorithm for word embeddings

    - Problematic societal biases encoded in widely used embeddings

- Link between accuracy gap and compounding injustices

- Large-scale dataset of online bios for occupation classification*

    - Gender imbalance compounded, even if explicit indicators "scrubbed"

- Bias in word embeddings can be leveraged to mitigate bias without access to protected attributes

*Code to reproduce dataset publicly available: **aka.ms/biasbios**

# Open problems

**Open problems in word embeddings**

- Debias contextual word embeddings (e.g., ELMo, BERT)
- Simultaneously reduce multiple biases in word embeddings
- Are biases in other languages different? ← **we now have some results for Spanish!**

**Open problems in occupation classification**

- Generate explainable classifications
- Better understand causes of differences
- Better understand fairness needs, e.g., affirmative action
- Fair candidate ranking

**Thanks!**

mdeartea@andrew.cmu.edu