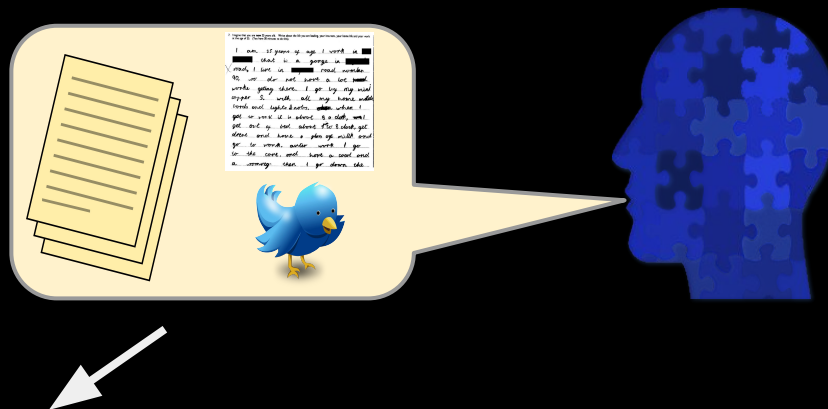# Human-Centered NLP and Ethics in NLP

CSE 354
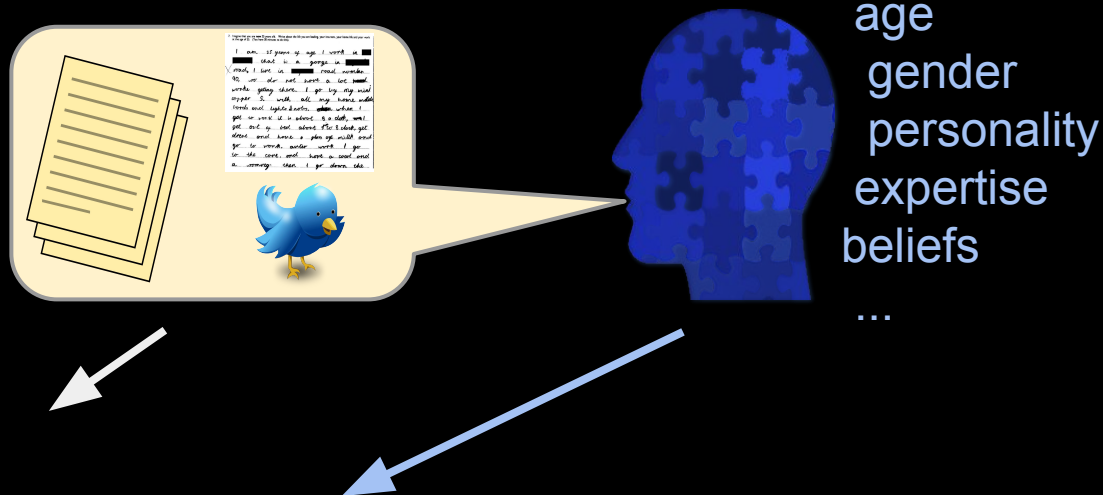
# The "Task" of human-centered NLP

Most NLP Tasks. E.g.
- POS Tagging
- Document Classification
- Sentiment Analysis
- Stance Detection
- Mental Health Risk Assessment
- …
  (language modeling, QA, …

# The "Task" of human-centered NLP



age
gender
personality
expertise
beliefs
...

Most NLP Tasks. E.g.
- POS Tagging
- Document Classification
- Sentiment Analysis
- Stance Detection
- Mental Health Risk Assessment
- …
  (language modeling, QA, …

# The "Task" of human-centered NLP



age
gender
personality
expertise
beliefs
…

Most NLP Tasks. E.g.

- POS Tagging
- Document Classification
- Sentiment Analysis
- Stance Detection
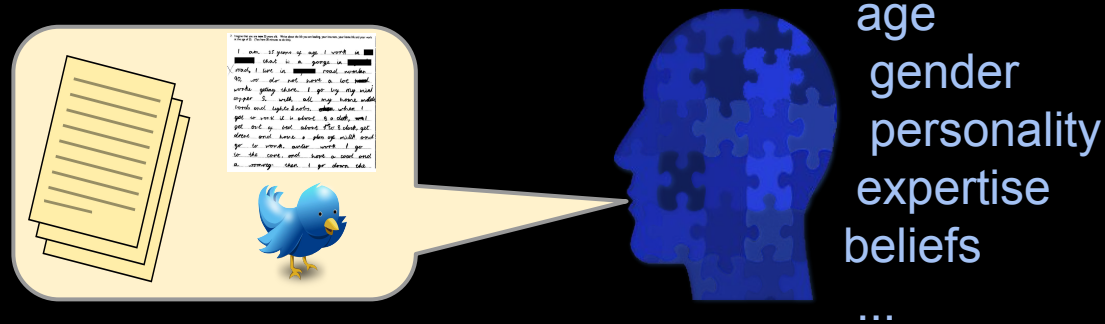- Mental Health Risk Assessment
- …
  (language modeling, QA, …

How to include extra-linguistics?

- Additive Inclusion
- Adaptive Extralinguistics
  - Adapting Embeddings
  - Adapting Models
- Correcting for bias

Natural Language Processing
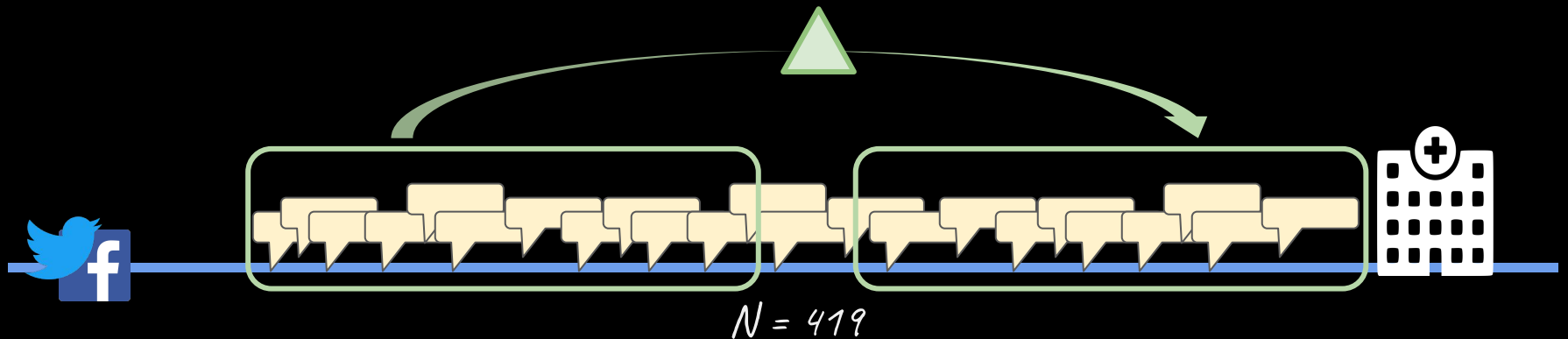
Psychological & Health Sciences

correlation strength

relative frequency

Schwartz, H. A., Eichstaedt, ... & Ungar. (2013). Personality, gender, and age in the language of social media: The open-vocabulary approach. *PloS one, 8(9)*.

Natural Language Processing

Psychological & Health Sciences

correlation strength

relative frequency

Schwartz, H. A., Eichstaedt, ... & Ungar. (2013). Personality, gender, and age in the language of social media: The open-vocabulary approach. *PloS one, 8(9)*.

Guntuku, S. C., Schwartz, H. A., Kashyap, A., Gaulton, J. S., Stokes, D. C., Asch, D. A., ... & Merchant, R. M. (2020). Variability in Language used on Social Media prior to Hospital Visits. *Nature - Scientific Reports*, *10(1),* 1-9.

Guntuku, S. C., Schwartz, H. A., Kashyap, A., Gaulton, J. S., Stokes, D. C., Asch, D. A., ... & Merchant, R. M. (2020). Variability in Language used on Social Media prior to Hospital Visits. *Nature - Scientific Reports*, *10(1),* 1-9.

Guntuku, S. C., Schwartz, H. A., Kashyap, A., Gaulton, J. S., Stokes, D. C., Asch, D. A., ... & Merchant, R. M. (2020). Variability in Language used on Social Media prior to Hospital Visits. *Nature - Scientific Reports, 10(1),* 1-9.

Guntuku, S. C., Schwartz, H. A., Kashyap, A., Gaulton, J. S., Stokes, D. C., Asch, D. A., ... & Merchant, R. M. (2020). Variability in Language used on Social Media prior to Hospital Visits. *Nature - Scientific Reports, 10(1),* 1-9.

# Problem

Natural language is written by

# Problem

Natural language is written by **people.**

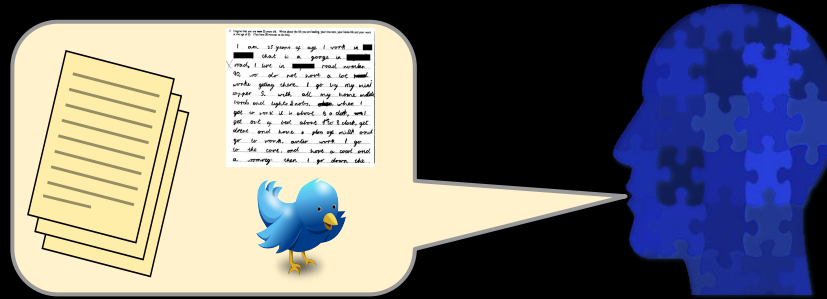# Problem

Natural language is written by **people.**

# Problem

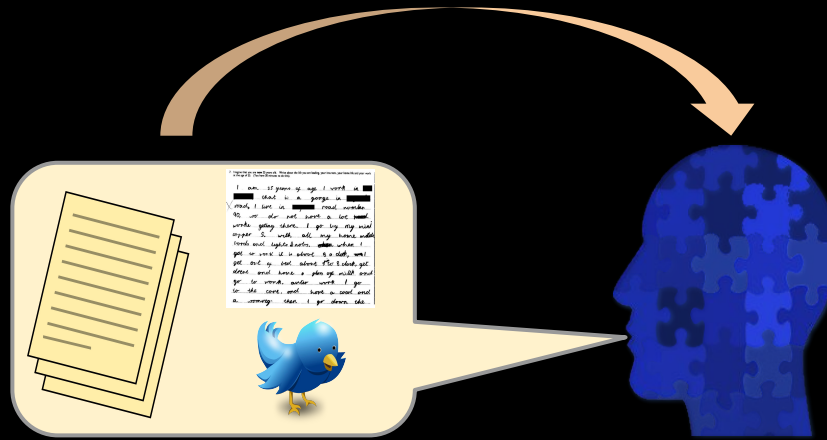Natural language is written by **people.**
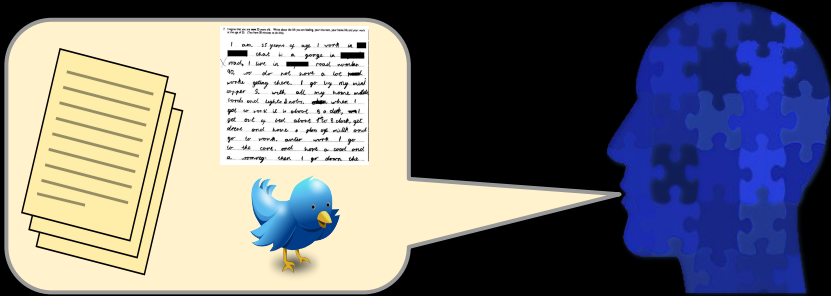
# Natural language is generated by *people*.



People have different beliefs, backgrounds, styles, vocabularies, preferences, knowledge, personalities, …

# Natural language is generated by *people*.



People have different beliefs, backgrounds, styles, vocabularies, preferences, knowledge, personalities, …,

and our language reflects these differences.

# Natural language is generated by people.
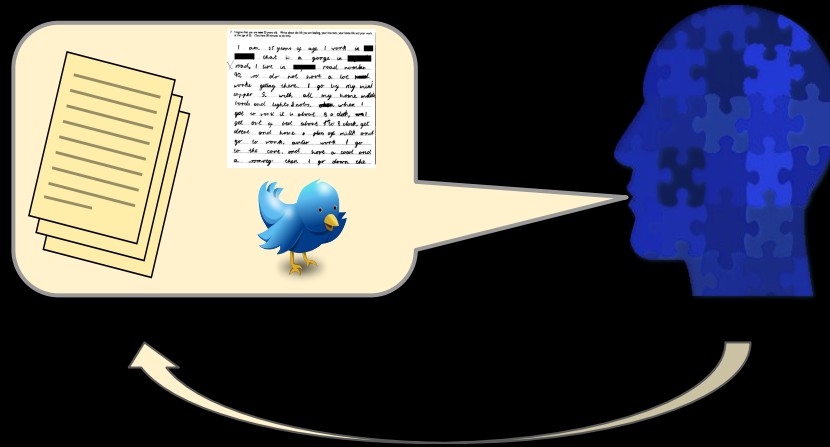


People have different beliefs, backgrounds, styles, vocabularies, preferences, knowledge, personalities, …,

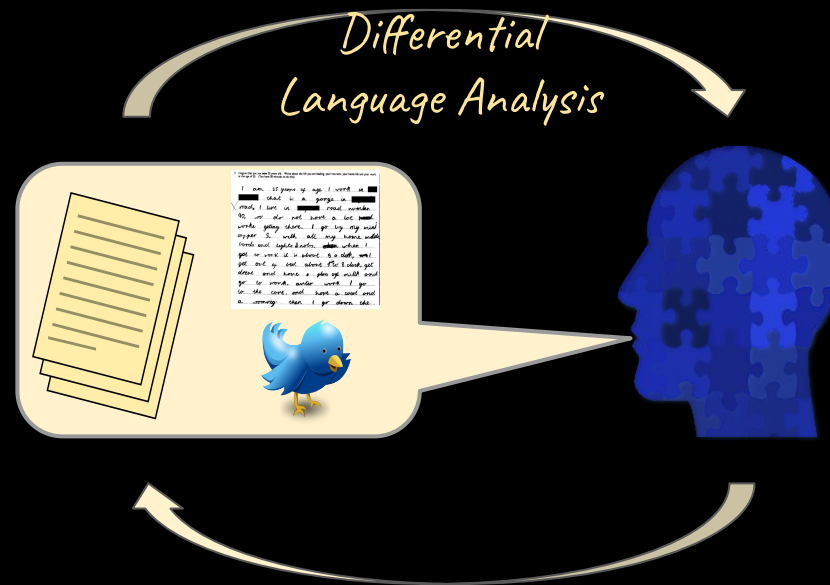and our language reflects these differences.

# Human Centered NLP:

# Human Centered NLP:

1. Model language as a human process

**Human Centered NLP:**
1. Model language as a human process
2. Use language to better understand humans.

# Differential Language Analysis

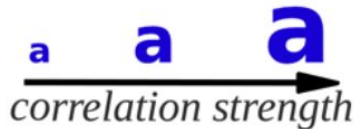Input:

    Linguistic features

    Human or community attribute

Output:

    Features distinguishing attribute

Goal: Data-driven insights about an attribute

# E.g. Words distinguishing communities with increases in real estate prices.



secret san improve texas post prices super web international starbucks california create companies downtown company pro credit tbh stoked media access tips cheap internet technology style bomb results tour experience na guide sales price cali tax source industry content nn per hellamarketing ou law followback blog search deal

a a **a**
*correlation strength*

*relative frequency*

# Differential Language Analysis

Input:
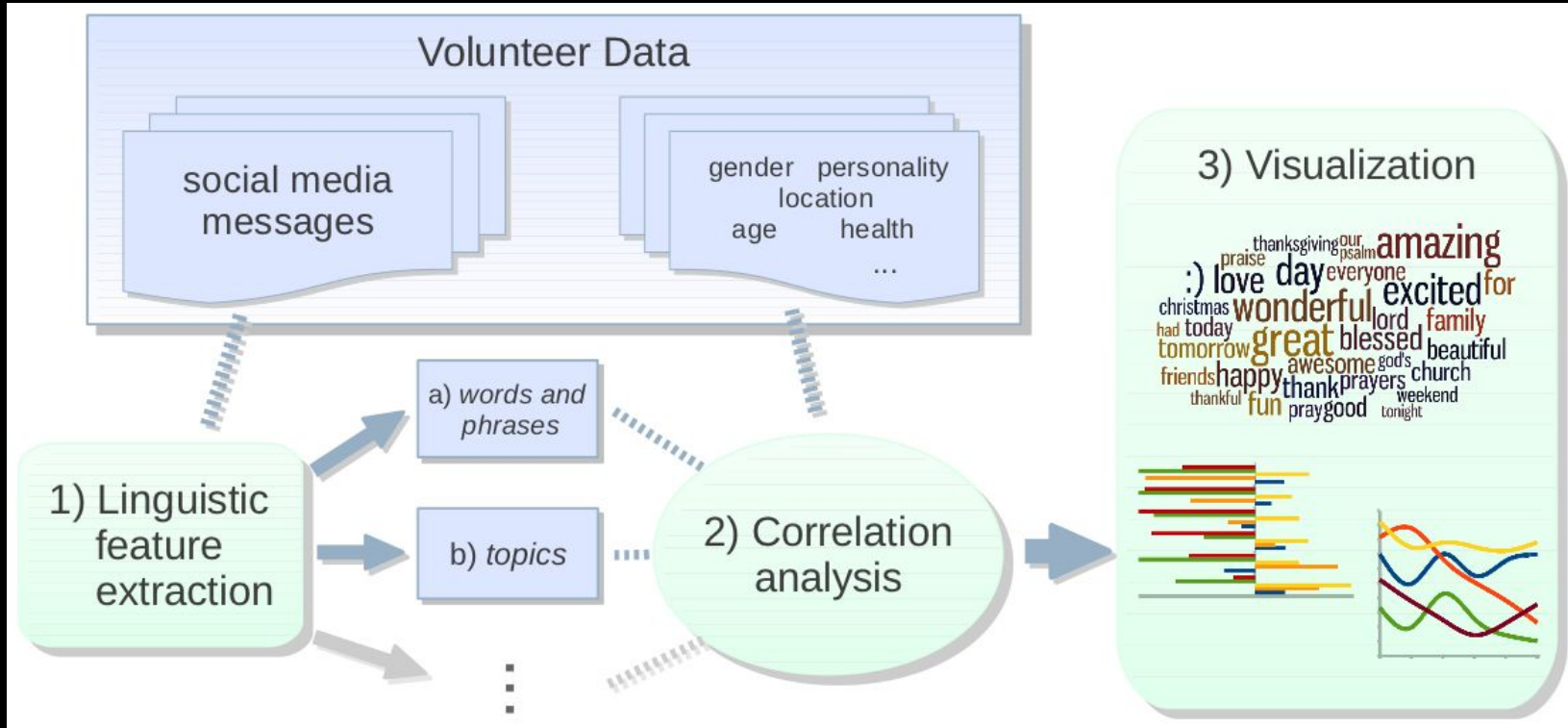
    Linguistic features

    Human or community attribute

Output:

    Features distinguishing attribute

Goal: Data-driven insights about an attribute

# Differential Language Analysis

# Differential Language Analysis

Methods of Correlation Analysis:

- Pearson Product-Moment Correlation
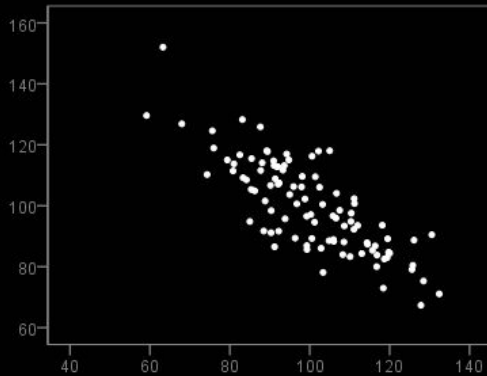    Limitation: Doesn't handle controls

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation
    - Limitation: Doesn't handle controls



r = -0.8

r = 0.5     © 2017 www.s|

r = 0.1

# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation
  Limitation: Doesn't handle controls

- Standardized Multivariate Linear Regression
  Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_m X_{m1} + \epsilon_i$$
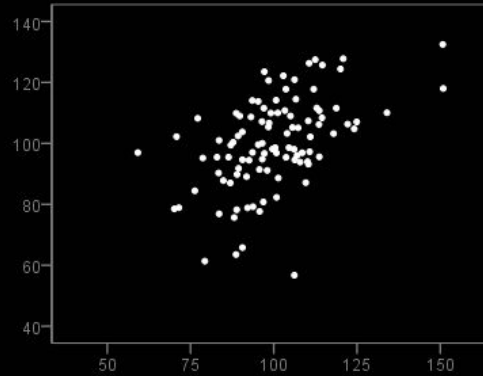
# Differential Language Analysis

Methods of Correlation Analysis:

- Pearson Product-Moment Correlation
  Limitation: Doesn't handle controls

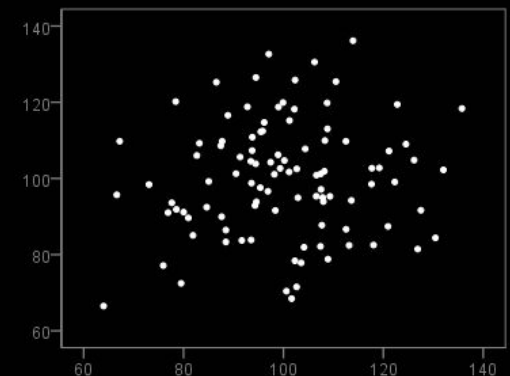$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

- **Standardized** Multivariate Linear Regression
  Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_m X_{m1} + \epsilon_i$$

  Adjust all variables to have "mean center" and "unit variance":

# Differential Language Analysis

Methods of Correlation Analysis:

- Pearson Product-Moment Correlation
    Limitation: Doesn't handle controls

$$r_{xy} = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^{n} (y_i - \bar{y})^2}}$$

- **Standardized** Multivariate Linear Regression
Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_m X_{m1} + \epsilon_i$$

Adjust all variables to have "mean center" and "unit variance":

$$z = \frac{x - \mu}{\sigma}$$

$\mu =$ Mean
$\sigma =$ Standard Deviation

# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation
     Limitation: Doesn't handle controls

- Standardized Multivariate Linear Regression
  Fit the model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_m X_{m1} + \epsilon_i$$

  Option 1: Gradient Descent:

  $J = \sum (y - \hat{y})^2$ -- "Sum of Squares" Error

# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation
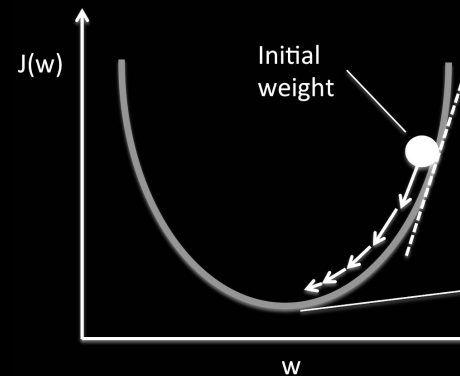  Limitation: Doesn't handle controls

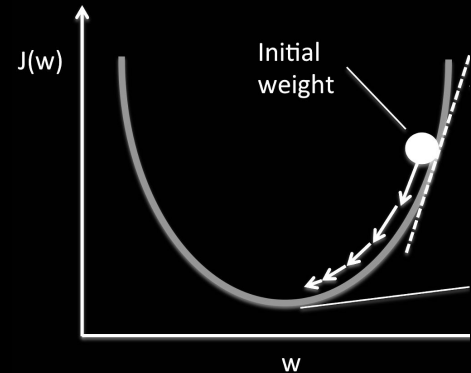- Standardized Multivariate Linear Regression
  Fit the model:

  $$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_m X_{m1} + \epsilon_i$$

  Option 1: Gradient Descent:

  $$J = \sum (y - \hat{y})^2 \quad \text{-- "Sum of Squares" Error}$$

  Option 2: Matrix model:

  $$Y = X\beta + \epsilon$$



J(w), Initial weight, w

# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation
     Limitation: Doesn't handle controls

- Standardized Multivariate Linear Regression
   Fit the model:
   Option 1: Gradient Descent:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \ldots + \beta_m X_{m1} + \epsilon_i$$

$J = \sum (y - \hat{y})^2$  -- "Sum of Squares" Error

   Option 2: Matrix model:       $Y = X\beta + \epsilon$
   Matrix Computation Solution:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

J(w)

Initial
weight

w

# Differential Language Analysis

Methods of Correlation Analysis:

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$

- Pearson Product-Moment Correlation
  Limitation: Doesn't handle controls

- Standardized Multivariate Linear Regression
  Fit the model:
  Option 1: Gradient Descent:

$$Y_i = \beta_0 + \boxed{\beta_1}X_{i1} + \beta_2 X_{i2} + \ldots + \beta_m X_{m1} + \epsilon_i$$
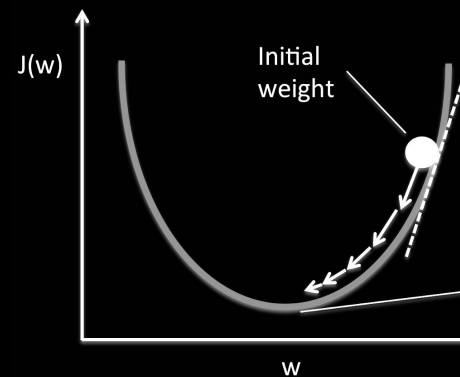
$$J = \sum (y - \hat{y})^2 \quad \text{-- "Sum of Squares" Error}$$

  Option 2: Matrix model:
  
  $$Y = X\beta + \epsilon$$
  
  Matrix Computation Solution:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$



J(w)

Initial weight

w

# Differential Language Analysis

Methods of "Correlation" Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio

$$\frac{\dfrac{countA("horrible")}{NA}}{1-\dfrac{countA("horrible")}{NA}}$$

$$\frac{\dfrac{countB("horrible")}{NB}}{1-\dfrac{countB("horrible")}{NB}}$$

# Differential Language Analysis

Methods of "Correlation" Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio

$$\frac{\dfrac{countA(\text{"horrible"})}{NA}}{1-\dfrac{countA(\text{"horrible"})}{NA}} \Bigg/ \frac{\dfrac{countB(\text{"horrible"})}{NB}}{1-\dfrac{countB(\text{"horrible"})}{NB}} \propto log\left(\frac{\dfrac{countA(\text{"horrible"})}{NA}}{1-\dfrac{countA(\text{"horrible"})}{NA}}\right) - log\left(\frac{\dfrac{countB(\text{"horrible"})}{NB}}{1-\dfrac{countB(\text{"horrible"})}{NB}}\right)$$

# Differential Language Analysis

Methods of "Correlation" Analysis for binary outcomes:

- Logistic Regression over Standardized variables

- Odds Ratio

$$\frac{\frac{countA(\text{"horrible"})}{NA}}{1-\frac{countA(\text{"horrible"})}{NA}}{\frac{countB(\text{"horrible"})}{NB}}{1-\frac{countB(\text{"horrible"})}{NB}} \quad \propto \quad log\left(\frac{\frac{countA(\text{"horrible"})}{NA}}{1-\frac{countA(\text{"horrible"})}{NA}}\right) - log\left(\frac{\frac{countB(\text{"horrible"})}{NB}}{1-\frac{countB(\text{"horrible"})}{NB}}\right)$$

$$= \quad log\left(\frac{countA(\text{"horrible"})}{NA-countA(\text{"horrible"})}\right) - log\left(\frac{countB(\text{"horrible"})}{NB-countB(\text{"horrible"})}\right)$$

# Differential Language Analysis

$$log\left(\frac{countA(\text{"horrible"})}{NA - countA(\text{"horrible"})}\right) - log\left(\frac{countB(\text{"horrible"})}{NB - countB(\text{"horrible"})}\right)$$

- Odds Ratio using Informative Dirichlet Prior

$$\delta_w^{(i-j)} = \log\left(\frac{f_w^i + \alpha_w}{n^i + \alpha_0 - (f_w^i + \alpha_w)}\right) - \log\left(\frac{f_w^j + \alpha_w}{n^j + \alpha_0 - (f_w^j + \alpha_w)}\right) \quad (20.9)$$

(Monroe et al., 2010; Jurafsky, 2017)

# Differential Language Analysis

$$log\left(\frac{countA(\text{"horrible"})}{NA-countA(\text{"horrible"})}\right) - log\left(\frac{countB(\text{"horrible"})}{NB-countB(\text{"horrible"})}\right)$$

- Odds Ratio using **Informative Dirichlet Prior**

$$\delta_w^{(i-j)} = \log\left(\frac{f_w^i + \alpha_w}{n^i + \alpha_0 - (f_w^i + \alpha_w)}\right) - \log\left(\frac{f_w^j + \alpha_w}{n^j + \alpha_0 - (f_w^j + \alpha_w)}\right) \quad (20.9)$$

(where $n^i$ is the size of corpus $i$, $n^j$ is the size of corpus $j$, $f_w^i$ is the count of word $w$ in corpus $i$, $f_w^j$ is the count of word $w$ in corpus $j$, $\alpha_0$ is the size of the background corpus, and $\alpha_w$ is the count of word $w$ in the background corpus.)

(Monroe et al., 2010; Jurafsky, 2017)

# Differential Language Analysis

$$log \left( \frac{countA("horrible")}{NA - countA("horrible")} \right) - log \left( \frac{countB("horrible")}{NB - countB("horrible")} \right)$$

- Odds Ratio using **<u>Informative Dirichlet Prior</u>**

$$\delta_w^{(i-j)} = log \left( \frac{f_w^i + \alpha_w}{n^i + \alpha_0 - (f_w^i + \alpha_w)} \right) - log \left( \frac{f_w^j + \alpha_w}{n^j + \alpha_0 - (f_w^j + \alpha_w)} \right) \quad (20.9)$$

(where $n^i$ is the size of corpus $i$, $n^j$ is the s...  ...pus $j$, $f_w^i$ is the count of word $w$ in corpus $i$, $f_w^j$ is the count of word $w$ in ...  ...is the size of the background corpus, and $\alpha_w$ is the count of word $w$ in ...  ...d corpus.)

(M...

Bayesian term for "smoothing": accounts for uncertainty as a function of event frequency (i.e. words observed less) by integrating "prior" beliefs mathematically.

# Differential Language Analysis

$$log \left( \frac{countA("horrible")}{NA-countA("horrible")} \right) - log \left( \frac{countB("horrible")}{NB-countB("horrible")} \right)$$

- Odds Ratio using **Informative Dirichlet Prior**

$$\delta_w^{(i-j)} = \log \left( \frac{f_w^i + \alpha_w}{n^i + \alpha_0 - (f_w^i + \alpha_w)} \right) - \log \left( \frac{f_w^j + \alpha_w}{n^j + \alpha_0 - (f_w^j + \alpha_w)} \right) \quad (20.9)$$

(where $n^i$ is the size of corpus $i$, $n^j$ is the s... pus $j$, $f_w^i$ is the count of word $w$ in corpus $i$, $f_w^j$ is the count of word $w$ in ... is the size of the background corpus, and $\alpha_w$ is the count of word $w$ in ... corpus.)

Bayesian term for "smoothing": accounts for uncertainty as a function of event frequency (i.e. words observed less) by integrating "prior" beliefs mathematically.
"Informative": the prior is based on past evidence. Here, the total frequency of the word.

(M...

# Differential Language Analysis

$$log \left( \frac{countA("horrible")}{NA - countA("horrible")} \right) - log \left( \frac{countB("horrible")}{NB - countB("horrible")} \right)$$

- Odds Ratio using Informative Dirichlet Prior

$$\delta_w^{(i-j)} = \log \left( \frac{f_w^i + \alpha_w}{n^i + \alpha_0 - (f_w^i + \alpha_w)} \right) - \log \left( \frac{f_w^j + \alpha_w}{n^j + \alpha_0 - (f_w^j + \alpha_w)} \right) \quad (20.9)$$

(where $n^i$ is the size of corpus $i$, $n^j$ is the size of corpus $j$, $f_w^i$ is the count of word $w$ in corpus $i$, $f_w^j$ is the count of word $w$ in corpus $j$, $\alpha_0$ is the size of the background corpus, and $\alpha_w$ is the count of word $w$ in the background corpus.)

Final score is standardized (z-scored): $\dfrac{\hat{\delta}_w^{(i-j)}}{\sqrt{\sigma^2 \left( \hat{\delta}_w^{(i-j)} \right)}}$ , where $\sigma^2 \left( \hat{\delta}_w^{(i-j)} \right) \approx \dfrac{1}{f_w^i + \alpha_w} + \dfrac{1}{f_w^j + \alpha_w}$

(Monroe et al., 2010; Jurafsky, 2017)

Natural language is generated by *people.*

# Natural language is generated by *people.*



"*The common misconception is that language has got to do with words and what they mean. It does not. It has to do with people and what they mean.*"

Shannon, 1948

Mosteller & Wallace 1963

Clark & Schober, 1992

Mairesse, Walker, et al., 2007

Hovy & Soogaard, 2015

# Natural language is generated by *people.*

## Yet, our models:



model

class
probs

# Natural language is generated by *people*.

## Yet, our models:



e.g. Document Classification: Stance

e.g., pro gun control?
yes, no

**model**

?

class probs

Mohammad, S., Kiritchenko, S., Sobhani, P., Zhu, X., & Cherry, C. (2016).
**Semeval-2016 task 6: Detecting stance in tweets.** In *Proceedings of the 10th International Workshop on Semantic Evaluation*.

# Natural language is generated by *people*.



e.g.  Document Classification: Stance

e.g., pro gun
control?
yes,  no

**model**

?

**class
probs**

Mohammad, S., Kiritchenko, S., Sobhani, P., Zhu, X., & Cherry, C. (2016).
**Semeval-2016 task 6: Detecting stance in tweets.** In *Proceedings of the
10th International Workshop on Semantic Evaluation*.

# Natural language is generated by *people.*



- personality
- demographics
- emotional states
- political ideology
- ...
- linguistic style
  (Pennebaker, 2007)
- latent user traits
  (Kulkarni et al., 2018)

model **?**

e.g.  Document Classification: Stance

e.g., pro gun
control?
yes，no

class
probs

Mohammad, S., Kiritchenko, S., Sobhani, P., Zhu, X., & Cherry, C. (2016).
**Semeval-2016 task 6: Detecting stance in tweets.** In *Proceedings of the
10th International Workshop on Semantic Evaluation.*

# Natural language is generated by *people*.

**What this means for NLP:**

1. Our data are inherently multi-level.

2. Often, there are "already-available" human attributes.

3. Our data and models are (human) biased.

# Natural language is generated by *people*.

**What this means for NLP:**

1. Our data are inherently multi-level.

2. Often, there are "already-available" human attributes.

3. Our data and models are (human) biased.

# Natural language is generated by *people.*

**What this means for NLP:**

1. Our data are inherently multi-level.

2. **Often, there are "already-available" human attributes.**

3. Our data and models are (human) biased.

# Approaches to Human Factor Inclusion

1.  Adaptive: Allow meaning if language to change depending on human context. (also called "compositional")
    (e.g. "sick" said from a young individual versus old individual)

2.  Additive: Include direct effect of human factor on outcome.
    (e.g. age and distinguishing PTSD from Depression)

3.  Bias Correction: Optimize so as not to pick up on unwanted relationships.

    (e.g. image captioner label pictures of men in kitchen as women)

# Approaches to Human Factor Inclusion

1. [What are human "factors"?] g if language to change depending
   o called "compositional")
   (ndividual versus old individual)

2. Additive: Include direct effect of human factor on outcome.
   (e.g. age and distinguishing PTSD from Depression)

3. Bias Correction: Optimize so as not to pick up on
   unwanted relationships.
   (e.g. image captioner label pictures of men in kitchen as women)

# Adaptation Approach: Domain Adaptation

Features for:  source          target

$$\Phi^s(x) = \langle x, x, 0 \rangle, \quad \Phi^t(x) = \langle x, 0, x \rangle$$

## Frustratingly Easy Domain Adaptation

**Hal Daumé III**
School of Computing
University of Utah
Salt Lake City, Utah 84112
me@hal3.name

### Abstract

We describe an approach to domain adaptation that is appropriate exactly in the case

supervised case. The fully supervised case models the following scenario. We have access to a large, annotated corpus of data from

# Adaptation Approach: Domain Adaptation

Features for: source      target

$$\Phi^s(x) = \langle x, x, 0 \rangle, \quad \Phi^t(x) = \langle x, 0, x \rangle$$

```
newX = []
for all x in source_x:
   newX.append(x + x + [0]*len(x))
for all x in target_x
   newX.append(x + [0]*len(x), x)

newY = source_y + target_y

model = model.train(newX,newY)
```

**Frustratingly Easy Domain Adaptation**

**Hal Daumé III**
School of Computing
University of Utah
Salt Lake City, Utah 84112
me@hal3.name

### Abstract

We describe an approach to domain adaptation that is appropriate exactly in the case

supervised case. The fully supervised case models the following scenario. We have access to a large, annotated corpus of data from

# Human Factors

--- Any attribute, represented as a continuous or discrete variable, of the humans generating the natural language.

E.g.
- Gender
- Age
- Personality
- Ethnicity
- Socio-economic status

# Adaptation Approach: Factor Adaptation

## Human Centered NLP with User-Factor Adaptation

Veronica E. Lynn, Youngseo Son, Vivek Kulkarni
Niranjan Balasubramanian and H. Andrew Schwartz
Stony Brook University
Stony Brook, NY
{velynn, yson, vvkulkarni, niranjan, has}@cs.stonybrook.edu

### Abstract

We pose the general task of *user-factor adaptation* — adapting supervised learning models to real-valued user factors inferred from a background of their lan...

and Costa Jr., 1989; Ruscio and Ruscio, 2000; Widiger and Samuel, 2005).

Here, we ask how one can adapt NLP models to real-valued human *factors* – continuous valued attributes that capture fine-grained differences be-

## Residualized Factor Adaptation
## for Community Social Media Prediction Tasks

Mohammadzaman Zamani,[1] H. Andrew Schwartz,[1] Veronica E. Lynn,[1]
Salvatore Giorgi,[2] and Niranjan Balasubramanian[1]
[1] Computer Science Department, Stony Brook University
[2] Department of Psychology, University of Pennsylvania
mzamani@cs.stonybrook.edu

### Abstract

Predictive models over social media language ... promise in capturing community ...

linked to socio-demographic factors (age, gender, race, education, income levels) with many social scientific studies supporting their predictive ...

# Adaptation



typically requires putting people into discrete bins

"*most latent variables of interest to psychiatrists and personality and clinical psychologists are dimensional [continuous]*"
(Haslam et al., 2012)

Type A

Type B

*"most latent variables of interest to psychiatrists and personality and clinical psychologists are dimensional [continuous]"*
(Haslam et al., 2012)

"*most latent variables of interest to psychiatrists and personality and clinical psychologists are dimensional [continuous]*"
(Haslam et al., 2012)

Less *Factor* A

More *Factor* A

# Our Method: Continuous Adaptation



(Lynn et al., 2017)

# Our Method: Continuous Adaptation



(Lynn et al., 2017)

# Our Method: Continuous Adaptation



(Lynn et al., 2017)

# User Factor Adaptation: Handling multiple factors

Replicate features for each factor:

A compositional function $c$ combines $d$ user factor scores $f_{u,d}$ with original feature values $\mathbf{x}$:

$$\Phi(\mathbf{x}, u) = \langle \mathbf{x}, c(f_{u,1}, \mathbf{x}), c(f_{u,2}, \mathbf{x}), \cdots, c(f_{u,d}, \mathbf{x}) \rangle$$

(Lynn et al., 2017)

# User Factor Adaptation: Handling multiple factors

Replicate features for each factor:

A compositional function $c$ combines $d$ user factor scores $f_{u,d}$ with original feature values $\mathbf{x}$:

$$\Phi(\mathbf{x}, u) = \langle \mathbf{x}, c(f_{u,1}, \mathbf{x}), c(f_{u,2}, \mathbf{x}), \cdots, c(f_{u,d}, \mathbf{x}) \rangle$$

| User | Factor Classes | Augmented Instance $\Phi(\mathbf{x}, u)$ |
|---|---|---|
| User 1 | $F_1$ | $\langle \mathbf{x}, \mathbf{x}, \mathbf{0}, \mathbf{0}, \cdots, 0 \rangle$ |
| User 2 | $F_2$ | $\langle \mathbf{x}, \mathbf{0}, \mathbf{x}, \mathbf{0}, \cdots, 0 \rangle$ |
| User 3 | $F_1, F_3$ | $\langle \mathbf{x}, \mathbf{x}, \mathbf{0}, \mathbf{x}, \cdots, 0 \rangle$ |
| User 4 | $F_k$ | $\langle \mathbf{x}, \mathbf{0}, \mathbf{0}, \cdots, 0, \mathrm{x} \rangle$ |

Table 1: Discrete Factor Adaptation: Augmentations of an original instance vector $\mathbf{x}$ under different factor class mappings. With $k$ domains the augmented feature vector is of length $n(k+1)$.

(Lynn et al., 2017)

# User Factor Adaptation: Handling multiple factors

Replicate features for each factor:

A compositional function $c$ combines $d$ user factor scores $f_{u,d}$ with original feature values $\mathbf{x}$:

$$\Phi(\mathbf{x}, u) = \langle \mathbf{x}, c(f_{u,1}, \mathbf{x}), c(f_{u,2}, \mathbf{x}), \cdots, c(f_{u,d}, \mathbf{x}) \rangle$$



| User | Factor Classes | Augmented Instance $\Phi(\mathbf{x}, u)$ |
|------|------|------|
| User 1 | $F_1$ | $\langle \mathbf{x}, \mathbf{x}, \mathbf{0}, \mathbf{0}, \cdots, 0 \rangle$ |
| User 2 | $F_2$ | $\langle \mathbf{x}, \mathbf{0}, \mathbf{x}, \mathbf{0}, \cdots, 0 \rangle$ |
| User 3 | $F_1, F_3$ | $\langle \mathbf{x}, \mathbf{x}, \mathbf{0}, \mathbf{x}, \cdots, 0 \rangle$ |
| User 4 | $F_k$ | $\langle \mathbf{x}, \mathbf{0}, \mathbf{0}, \cdots, 0, \mathbf{x} \rangle$ |

Table 1: Discrete Factor Adaptation: Augmentations of an original instance vector $\mathbf{x}$ under different factor class mappings. With $k$ domains the augmented feature vector is of length $n(k+1)$.

(Lynn et al., 2017)

# Main Results

Adaptation improves over unadapted baselines (Lynn et al., 2017)

| Task | Metric | No Adaptation | Gender | Personality | Latent (User Embed) |
|---|---|---|---|---|---|
| Stance | F1 | 64.9 | **65.1 (+0.2)** | **66.3 (+1.4)** | **67.9 (+3.0)** |
| Sarcasm | F1 | 73.9 | **75.1 (+1.2)** | **75.6 (+1.7)** | **77.3 (+3.4)** |
| Sentiment | Acc. | 60.6 | **61.0 (+0.4)** | **61.2 (+0.6)** | **60.7 (+0.1)** |
| PP-Attach | Acc. | 71.0 | 70.7 (-0.3) | 70.2 (-0.8) | 70.8 (-0.2) |
| POS | Acc. | 91.7 | **91.9 (+0.2)** | 91.2 (-0.5) | 90.9 (-0.8) |

# Example: How Adaptation Helps

Women
more adjectives→sarcasm

Men
more adjectives→no sarcasm



**more "male"**                                    **more "female"**

# Problem

User factors are not always available.

# Solution: User Factor Inference

## past tweets



Niranjan @b_niranjan · Sep 2
There must be a word for trending #hashtags that you know you will regret if you click. Is there?

Niranjan @b_niranjan · Aug 31
Passwords spiral: Forget password for the acnt you use twice a year. Ask for reset. Can't use previous. Create a new one to forget later.

Niranjan @b_niranjan · Jul 31
Thrilled to hear @acl2017's diversity efforts as the first thing in the conference.

♡ 1

→ **inferred factors**

**Known**
Age       (Sap et al. 2014)
Gender (Sap et al. 2014)
Personality (Park et al. 2015)

**Latent**
User Embeddings
   (Kulkarni et al. 2017)
*Word2Vec*
*TF-IDF*

# Background Size

Using more background tweets to infer factors produces larger gains

# Full User Factors Adaptation Pipeline: with latent factors from training

| doc-id | user-id | document | | |
|--------|---------|----------|---|---|
| 1 | 42 | text… | → | **emb**dngs |
| 2 | 42 | text… | → | **emb**dngs |
| 3 | 16 | text… | → | **emb**dngs |
| 4 | 42 | text… | → | **emb**dngs |
| 5 | 12 | text… | → | **emb**dngs |
| 6 | 16 | text… | → | **emb**dngs |
| … | … | … | | |

d = 128

total documents

# Full User Factors Adaptation Pipeline: with latent factors from training

| doc-id | user-id | document |
|--------|---------|----------|
| 1 | 42 | text… |
| 2 | 42 | text… |
| 3 | 16 | text… |
| 4 | 42 | text… |
| 5 | 12 | text… |
| 6 | 16 | text… |
| … | … | … |

total documents

d = 128

**emb**dngs
**emb**dngs
**emb**dngs
**emb**dngs
**emb**dngs
**emb**dngs

avg
avg
avg

users x avg_embeddings

d = 128

N users

**Step 1: Create User Factors**

PCA

user x factors

| 42 | f1, f2, f3 |
|----|------------|
| 16 | f1, f2, f3 |
| 12 | f1, f2, f3 |
| … | … |

d = 3
(or other lower dimension)

**Full User Factors Adaptation Pipeline:** with latent factors from training

Step 2: Create User-adapted Features

# Full User Factors Adaptation Pipeline: with latent factors from training



| doc-id | user-id | document |
|--------|---------|----------|
| 1 | 42 | text… |
| 2 | 42 | text… |
| 3 | 16 | text… |
| 4 | 42 | text… |
| 5 | 12 | text… |
| 6 | 16 | text… |
| … | … | … |

d = 128

**emb**dngs → avg

users x avg_embeddings

d = 128

**Step 3: Train Model**

| doc-id | rating |
|--------|--------|
| 1 | 2 |
| 2 | 3 |
| 3 | 4 |
| 4 | 1 |

**PCA**

Then feed these as features into your document level classifier or regressor.

total documents

| doc-id | user-adapted embedding |
|--------|------------------------|
| 1 | emb x f1; emb x f2; emb x f3 |
| 2 | emb x f1; emb x f2; emb x f3 |
| 3 | emb x f1; emb x f2; emb x f3 |
| 4 | emb x f1; emb x f2; emb x f3 |
| … | … |

user x factors

d = 3
(or other lower dimension)

| | |
|----|------------|
| 42 | f1, f2, f3 |
| 16 | f1, f2, f3 |
| 12 | f1, f2, f3 |
| … | … |

# Full User Factors Adaptation Pipeline: with latent factors from training

| doc-id | user-id | document |
|--------|---------|----------|
| 1 | 42 | text… |
| 2 | 42 | text… |

**This was training data; now assume test**

| 5 | 12 | text… |
| 6 | 16 | text… |
| ... | ... | ... |

d = 128

**emb**dngs → avg
**emb**dngs → avg
**emb**dngs → avg
**emb**dngs
**emb**dngs
**emb**dngs

**users x avg_embeddings**

d = 128

N users

**What about when predicting on new documents?**

total documents

| doc-id |
|--------|
| 1 |
| 2 |
| 3 |
| 4 |
| ... |

**user-adapted embeddings**

emb x f1; emb x f2; emb x f3
emb x f1; emb x f2; emb x f3
emb x f1; emb x f2; emb x f3
emb x f1; emb x f2; emb x f3
...

**PCA**

**user x factors**

| 42 | f1, f2, f3 |
| 16 | f1, f2, f3 |
| 12 | f1, f2, f3 |
| ... | ... |

d = 3
(or other lower dimension)

# Full User Factors Adaptation Pipeline: with latent factors from training

# Full User Factors Adaptation Pipeline: with latent factors from training
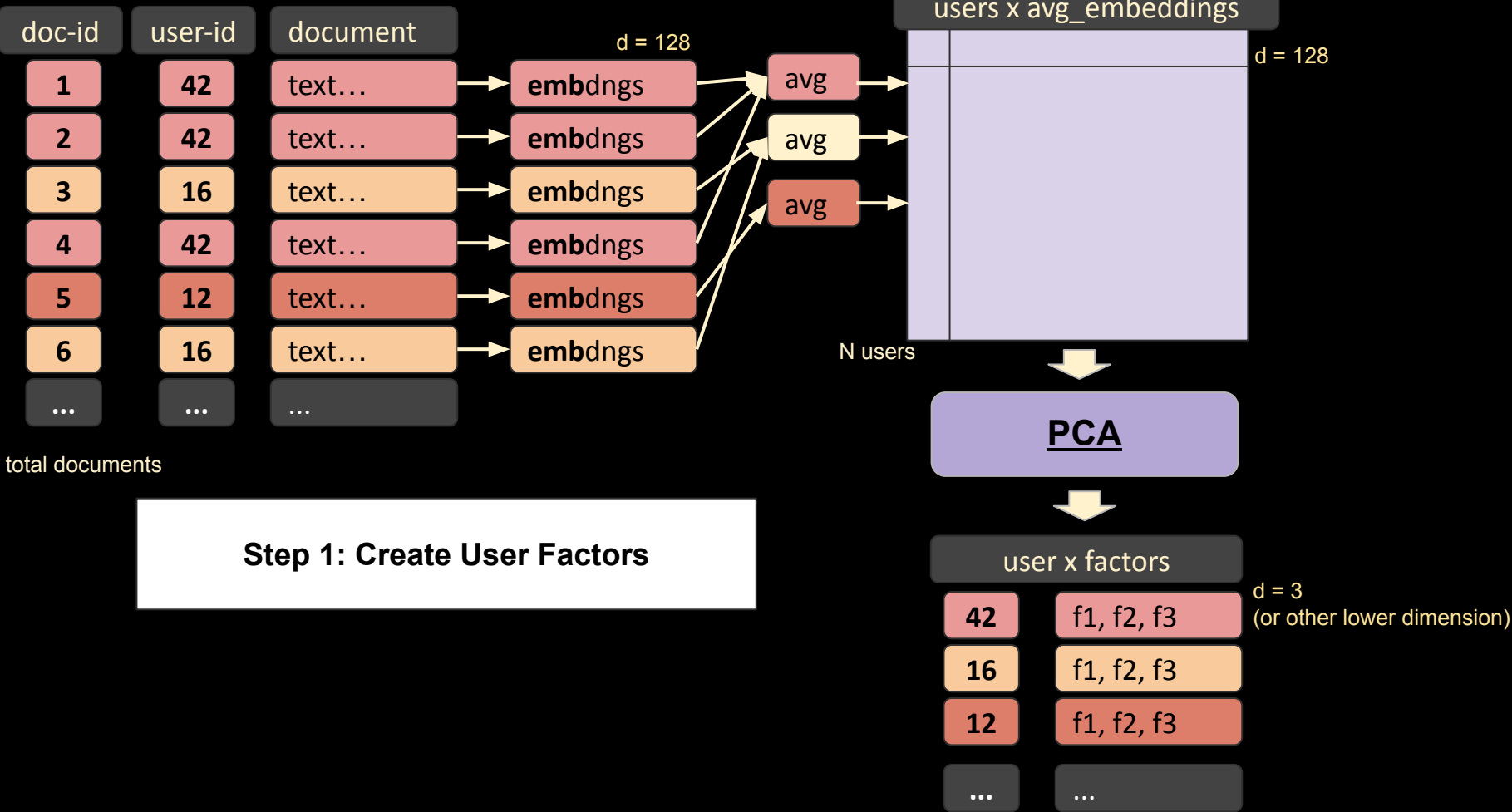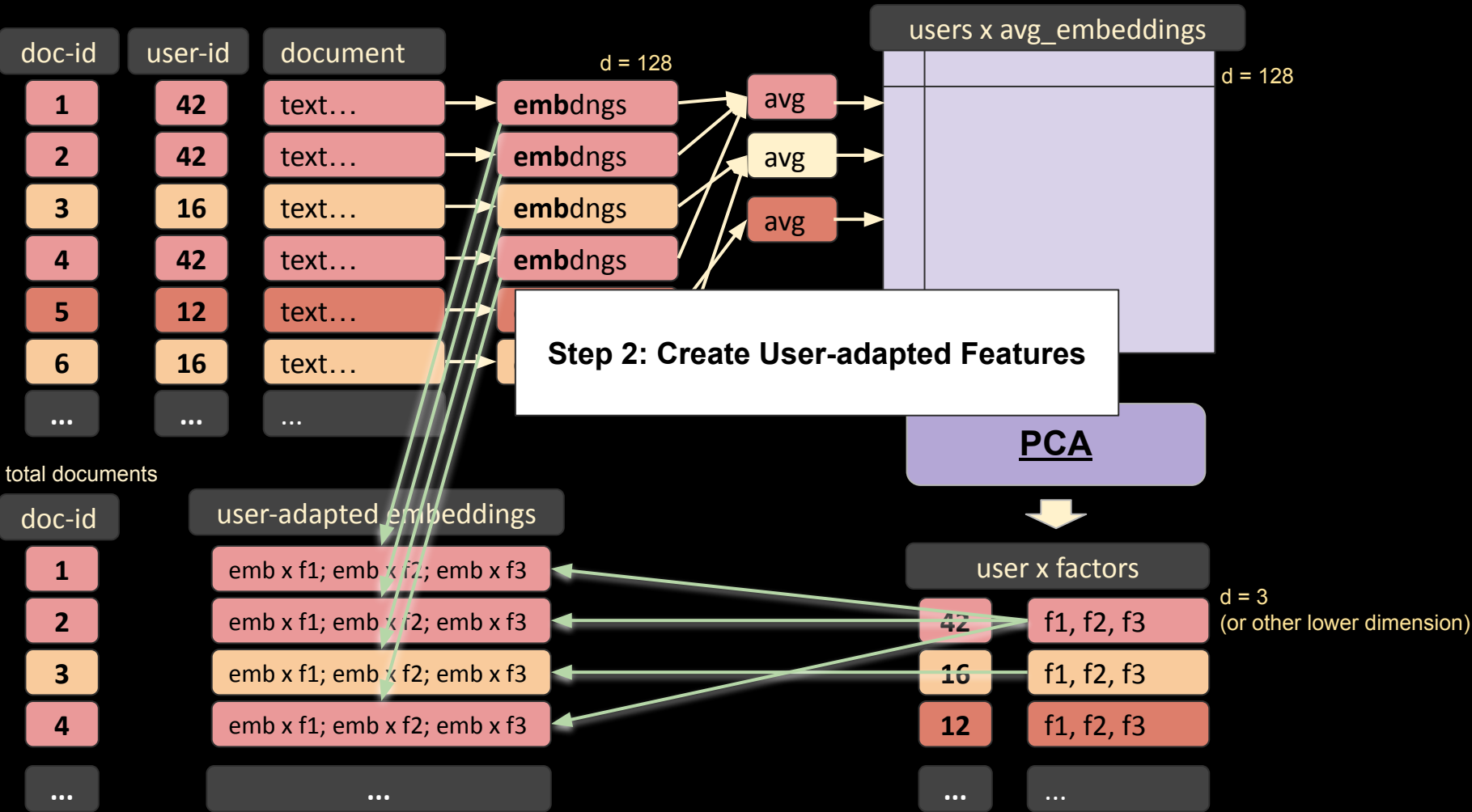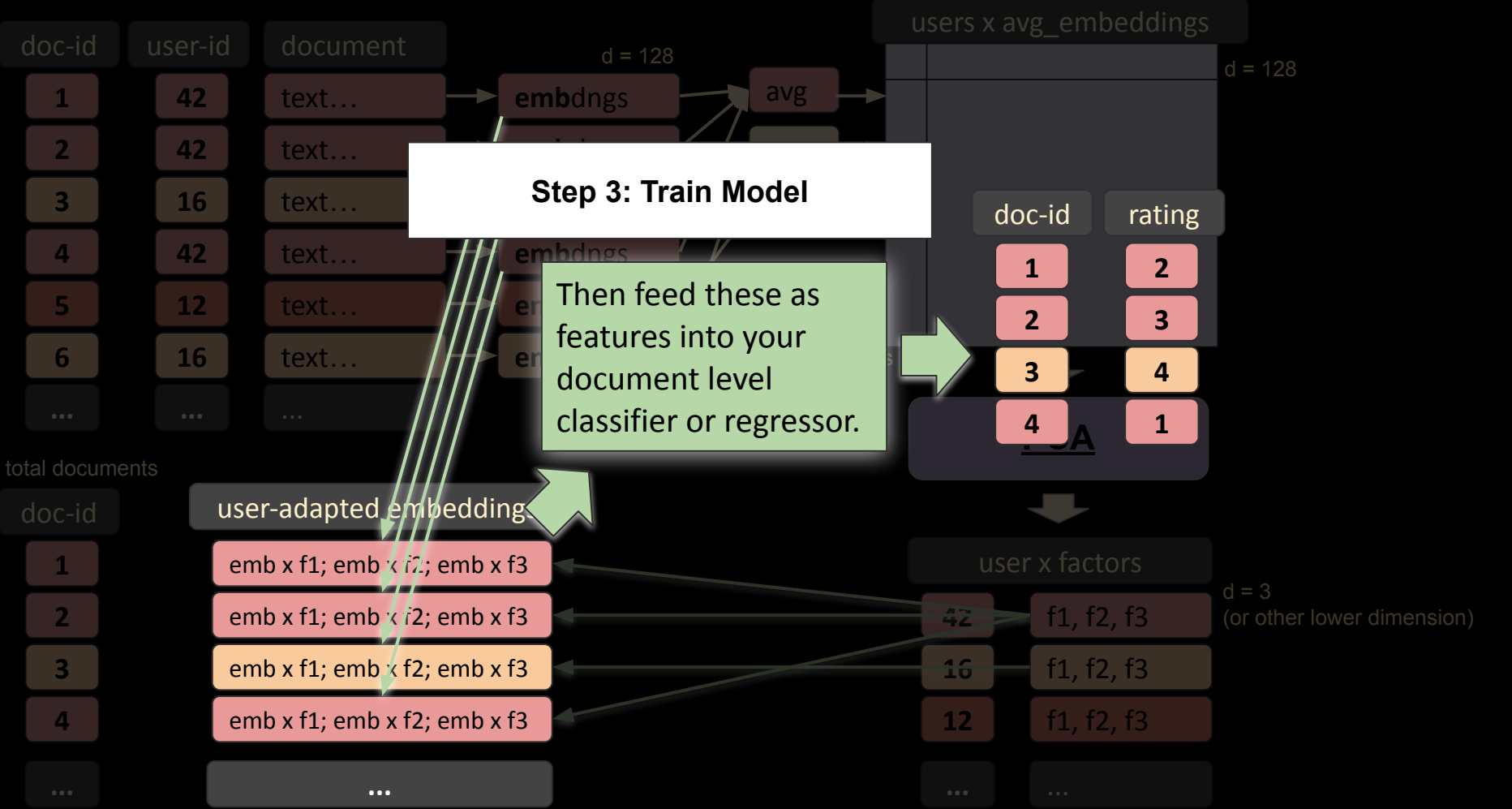
# Full User Factors Adaptation Pipeline: with latent factors from training



doc-id | user-id | document

d = 128

| 1 | 42 | text… | **emb**dngs |
| 2 | 42 | text… | **emb**dngs |
| | | | **emb**dngs |
| | | | **emb**dngs |
| 5 | 12 | text… | **emb**dngs |
| 6 | 16 | text… | **emb**dngs |
| … | … | … |

avg

avg

avg

users x avg_embeddings

N users

**This was training data; now assume test**

**What about when predicting on new documents?**
(easy as A, B, C)

**PCA**

user x factors

Transformation Matrix (V)

total documents

doc-id | user-adapted embeddings

| 1 | emb x f1; emb x f2; emb x f3 |
| 2 | emb x f1; emb x f2; emb x f3 |
| 3 | emb x f1; emb x f2; emb x f3 |
| 4 | emb x f1; emb x f2; emb x f3 |
| … | … |

| 42 | f1, f2, f3 |
| 16 | f1, f2, f3 |
| 12 | f1, f2, f3 |
| … | … |

**A.** Save the transformation (V) from PCA during training

**B.** Apply V to *user x avg_embeddings* matrix during test/trial.

**C.** Adapt document features by user factors just like in training.

(another lower dimension)

# Approaches to Human Factor Inclusion

1. Adaptive: Allow meaning if language to change depending on human context. (also called "compositional")
   (e.g. "sick" said from a young individual versus old individual)

2. Additive: Include direct effect of human factor on outcome.
   (e.g. age and distinguishing PTSD from Depression)

3. Bias Correction: Optimize so as not to pick up on unwanted relationships.
   (e.g. image captioner label pictures of men in kitchen as women)

# Approaches to Human Factor Inclusion

1. Adaptive: Allow meaning if language to change depending on human context. (also called "compositional")
   (e.g. "sick" said from a young individual versus old individual)

2. Additive: Include direct effect of human factor on outcome.
   (e.g. age and distinguishing PTSD from Depression)

3. Bias Correction: Optimize so as not to pick up on unwanted relationships.
   (e.g. image captioner label pictures of men in kitchen as women)

# Ethics in NLP

1. Adaptive: Allow meaning if language to change depending on human context. (also called "compositional")
   (e.g. "sick" said from a young individual versus old individual)

2. Additive: Include direct effect of human factor on outcome.
   (e.g. age and distinguishing PTSD from Depression)

3. Bias Correction: Optimize so as not to pick up on unwanted relationships.
   (e.g. image captioner label pictures of men in kitchen as women)

# Ethics in NLP

Bias

Privacy

Ethical Research

# Ethics in NLP - Bias

Consequences of Sociodemographic Bias in NLP Models:

- Outcome Disparity:  Predicted distribution given A,
                          are dissimilar from ideal distribution given A

- Error Disparity: Predicts less accurate for authors of given demographics.

Shah, D., Schwartz, H. A., Hovy, D. (2020). Predictive Biases in Natural Language Processing Models: A Conceptual Framework and Overview. *In ACL-2020: Proceedings of the Association for Computational Linguistics.*

# Two Examples

model accuracy

## The WSJ Effect

Jørgensen/Hovy/Søgaard, 2015
Hovy & Søgaard, 2015

distance from "standard" WSJ author demographics

# Two Examples

model
accuracy

The W



| COOKING | |
|---------|---------|
| **ROLE** | **VALUE** |
| AGENT | WOMAN |
| FOOD | FRUIT |
| HEAT | ∅ |
| TOOL | KNIFE |
| PLACE | KITCHEN |

| COOKING | |
|---------|---------|
| **ROLE** | **VALUE** |
| AGENT | WOMAN |
| FOOD | MEAT |
| HEAT | STOVE |
| TOOL | SPATULA |
| PLACE | OUTSIDE |

| COOKING | |
|---------|---------|
| **ROLE** | **VALUE** |
| AGENT | WOMAN |
| FOOD | ∅ |
| HEAT | STOVE |
| TOOL | SPATULA |
| PLACE | KITCHEN |

Zhao, Jieyu, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. "Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints." In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 2017.

distance from "standard" WSJ author demographics

# Two Examples

model accuracy

The W...



**"Outcome Disparity"**

**"Error Disparity"**

Zhao, Jieyu, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. "Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints." In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 2017.

distance from "standard" WSJ author demographics

Our data and models are (human) biased.



"Outcome Disparity"

Person-level

■ attribute = 1

■ attribute = 2

"Error Disparity"

# Our data and models are (human) biased.



"Outcome Disparity"

Person-level
- attribute = 1
- attribute = 2

"Error Disparity"

# Conceptual Framework:

# Conceptual Framework:



**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**Embedding Corpus**

features
$\theta_{embedding}$
(Pre-trained Side)

**Source Population**

features
$X_{source}$
(Model Side)

*fit*

outcomes
$Y_{source}$

**Target Population**

features
$X_{target}$
(Application Side)

*predict*

**biased outcomes**
$\hat{Y}_{target}$

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Outcome Disparity



average predicted outcome

Predicted
$Q(\hat{Y}_t|A)$

Ideal
$P(Y_t|A)$

**human attribute**

■ value1
■ value2

## outcome disparity
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**Target Population**

features
$X_{target}$
(Application Side)

*predict*

**biased outcomes**
$\hat{Y}_{target}$

## error disparity
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Outcome Disparity



cancer

Predicted $Q(\hat{Y}_t|A)$

Predicted $Q(\hat{Y}_t|A)$

Ideal $P(Y_t|A)$

human attribute
value1
value2

human attribute
woman
man

outcome disparity
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$Q(Y_t|A) \neq P(Y_t|A)$

**Target Population**

features
$X_{target}$
(Application Side)

*predict*

**biased outcomes**
$\hat{Y}_{target}$

error disparity
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$

# Error Disparity

# Error Disparity



error

WSJ Effect

Jørgensen et al. (WNUT 2015)
Hovy & Søggard (ACL 2015)

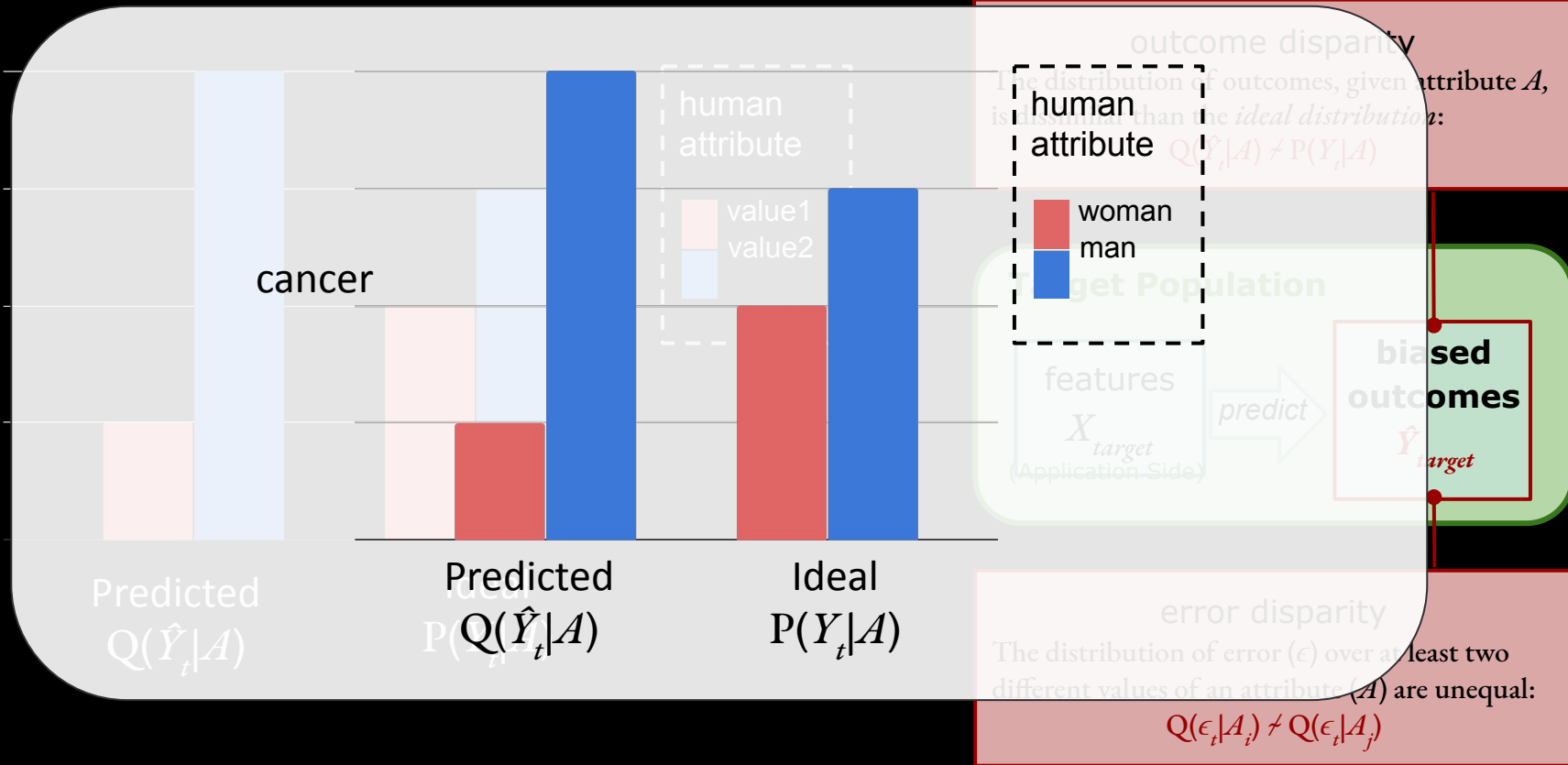**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

human attribute

value1
value2

**Target Population**

features
$X_{target}$
(Application Side)

*predict*

**biased outcomes**
$\hat{Y}_{target}$

Predicted
$Q(\hat{Y}|A)$

Ideal
$P(Y|A)$

Correlates with demographics

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

Distance from "Standard"

# Disparities





outcome disparity

The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
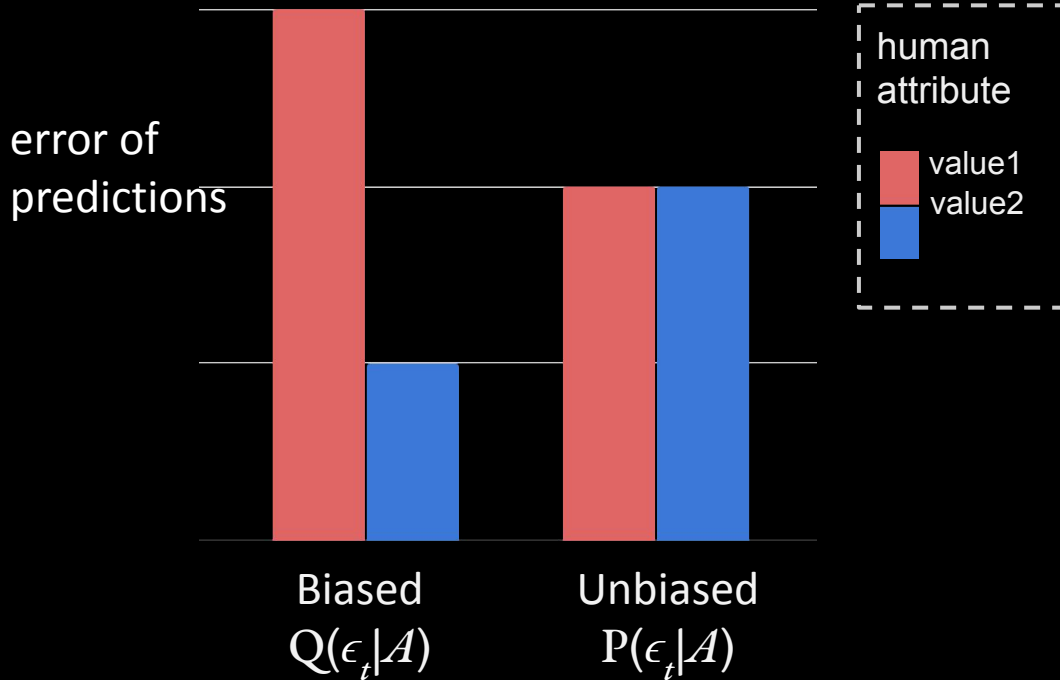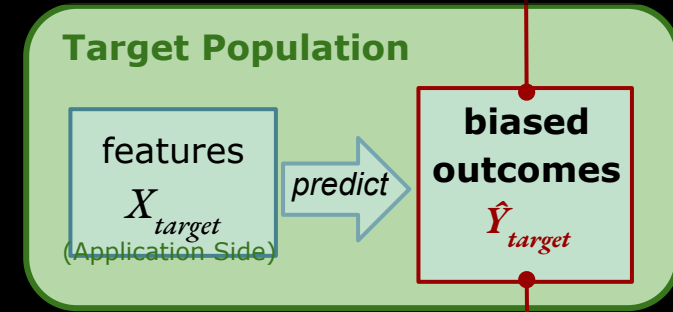$$Q(\hat{Y}_t | A) \neq P(Y_t | A)$$

**Target Population**

features
$X_{target}$
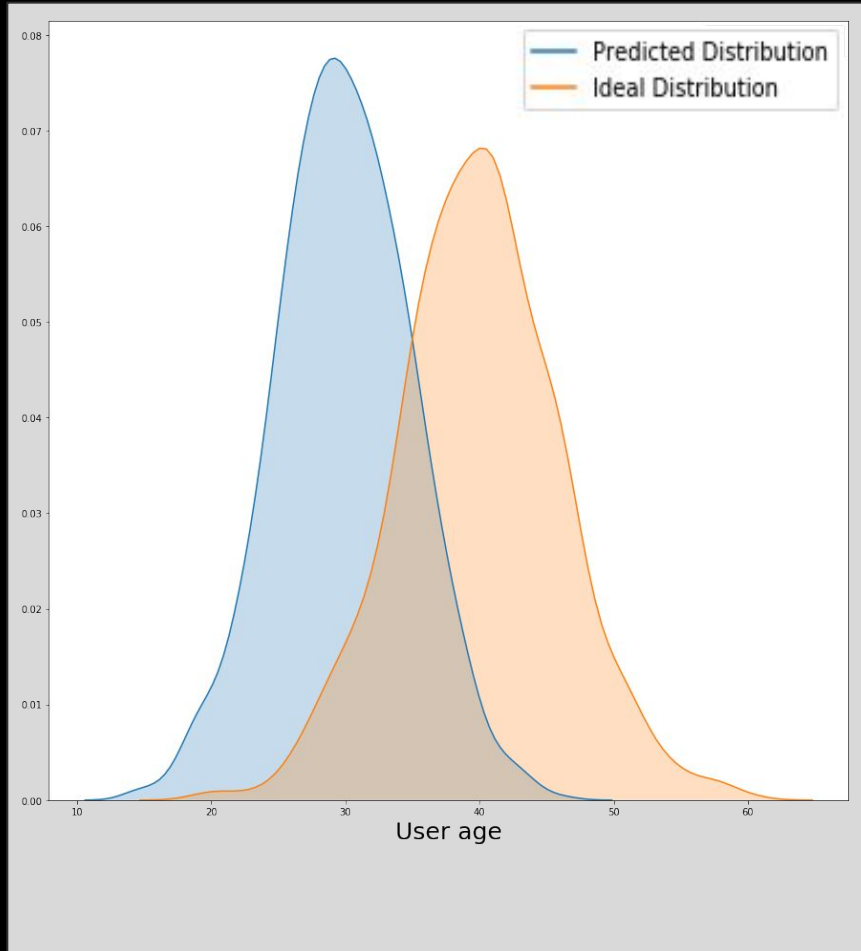(Application Side)

*predict*

**biased outcomes**
$\hat{Y}_{target}$

error disparity

The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t | A_i) \neq Q(\epsilon_t | A_j)$$

# Origins of Bias



**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**Embedding Corpus**

**features**
$$\theta_{embedding}$$
(Pre-trained Side)

**Source Population**

**features**
$$X_{source}$$
(Model Side)

*fit*

**outcomes**
$$Y_{source}$$

**Target Population**

**features**
$$X_{target}$$
(Application Side)

*predict*

**biased outcomes**
$$\hat{Y}_{target}$$

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Selection Bias



**Embedding Corpus**

features

$\theta_{embedding}$

(Pre-trained Side)

**Source Population**

**features**

$X_{source}$

(Model Side)

*fit*

**outcomes**

$Y_{source}$

**Target Population**

**features**

$X_{target}$

(Application Side)

*predict*

biased outcomes

$\hat{Y}_{target}$

outcome disparity

The distribution of outcomes, given attribute *A,* is dissimilar than the *ideal distribution*:

$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

selection bias

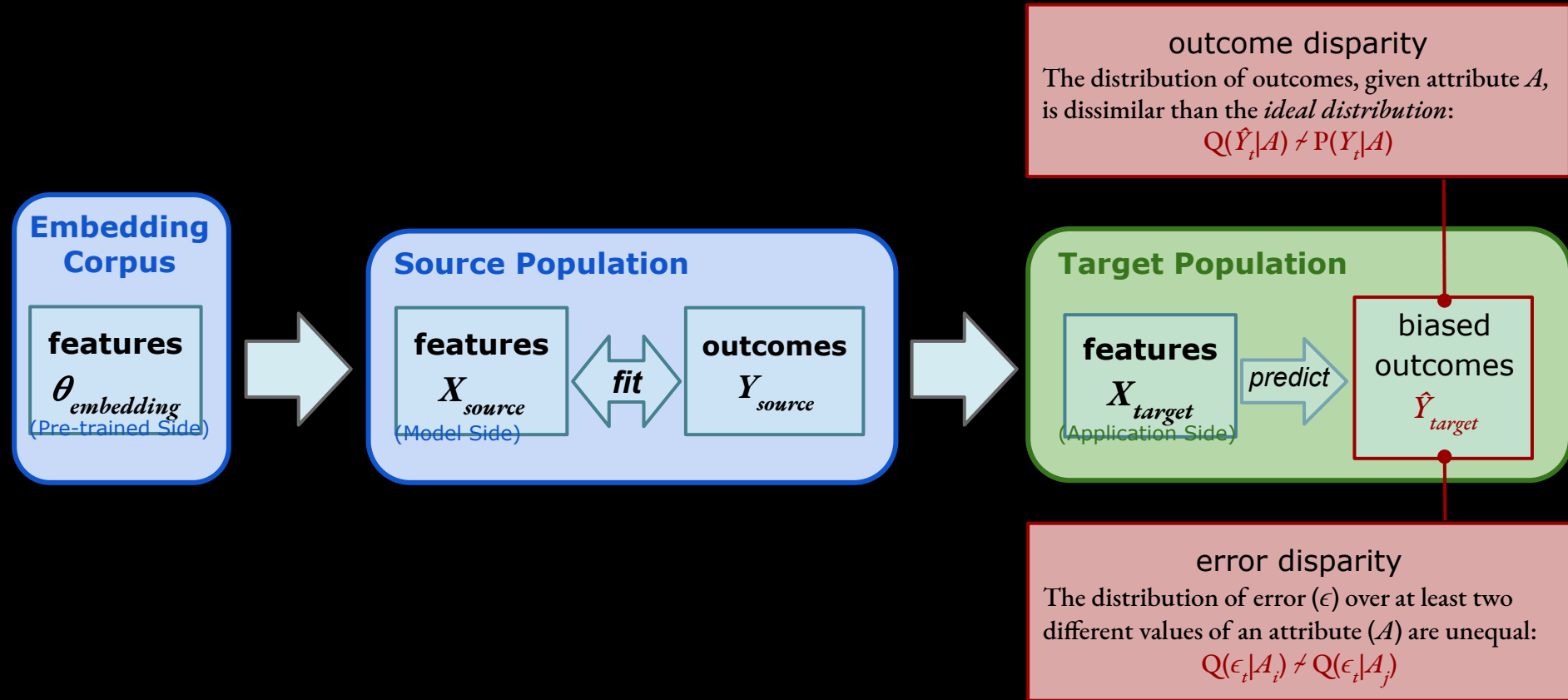The sample of observations themselves are not representative of the application population.
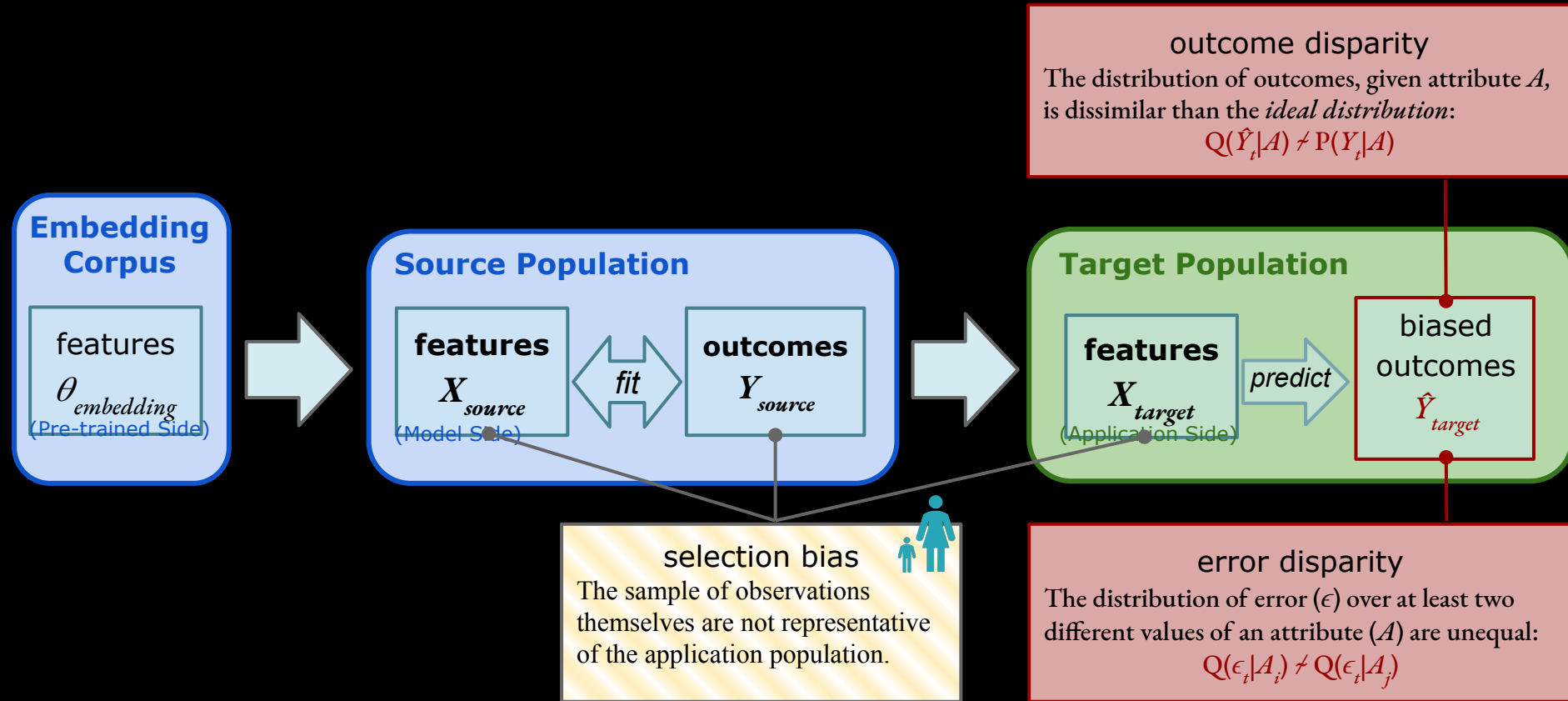
error disparity

The distribution of error ($\epsilon$) over at least two different values of an attribute (*A*) are unequal:

$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

WSJ Effect

Jørgensen et al. (WNUT 2015)
Hovy & Søggard (ACL 2015)

Selection Bias

error

**Embedding Corpus**

features

$\theta_{embedding}$

(Pre-Trained Side)

**Source Population**

**features**

$X_{source}$

fit

**outcomes**

$Y_{source}$

Correlates with demographics

**Target Population**

**features**

$X_{target}$

(Application Side)

predict

biased outcomes

$\hat{Y}_{target}$

Distance from "Standard"

outcome disparity
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**selection bias**
The sample of observations themselves are not representative of the application population.

error disparity
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
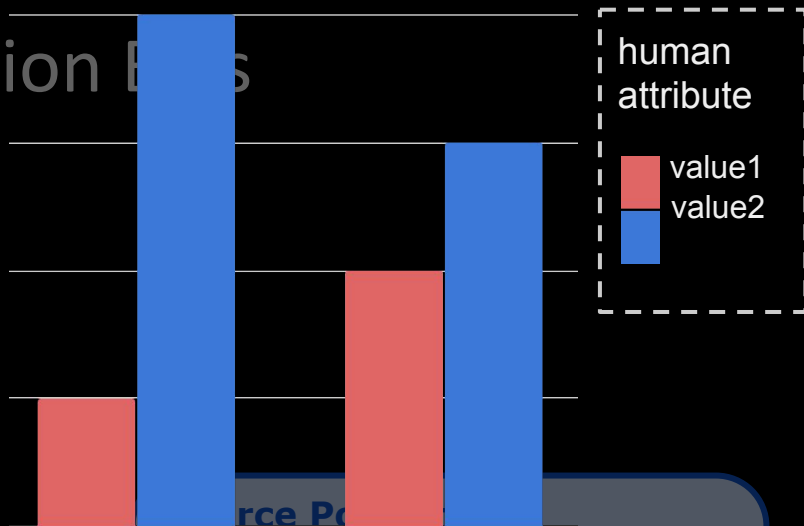$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Selection Bias



**Embedding Corpus**

features
$\theta_{embedding}$
(Pre-trained Side)

**Source Population**

features
$X_{source}$
(Model Side)

*fit*

**outcomes**
$Y_{source}$

**Target Population**

features
$X_{target}$
(Application Side)

*predict*

biased outcomes
$\hat{Y}_{target}$

**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**selection bias**
The sample of observations themselves are not representative of the application population.

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Label Bias



**label bias**
Biased annotations, interaction, or latent bias from past classifications.

**outcome disparity**
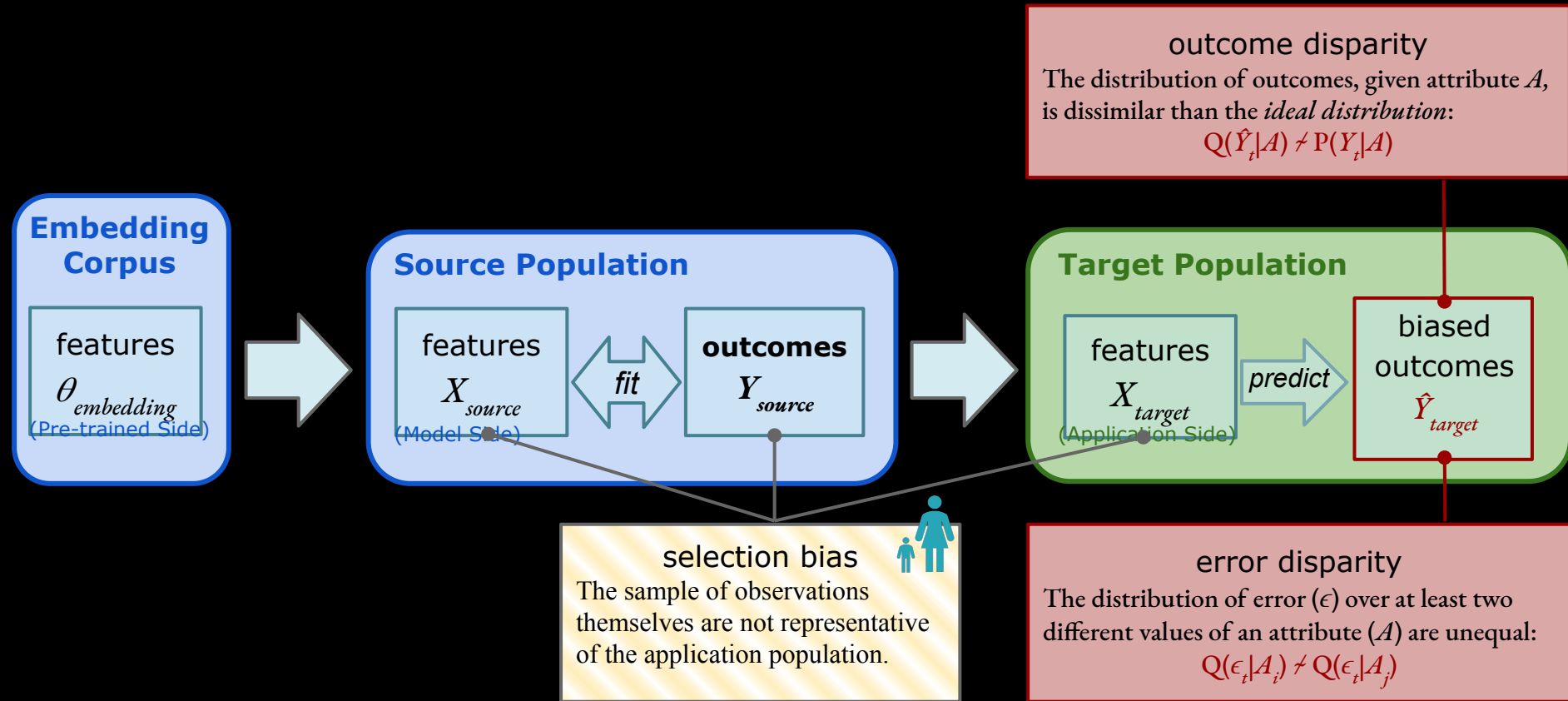The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**Embedding Corpus**

features

$\theta_{embedding}$ proportion
(Pre-trained Side) of sample

**Source Population**

features
source
(Side)

outcomes

$Y_{source}$

*fit*

human attribute

value1
value2

**Target Population**

features

$X_{target}$
(Application Side)

*predict*

biased outcomes

$\hat{Y}_{target}$

ction bias
of observations
e not representative
of the application population.

Source
$Q(Y_S|A_S)$

Ideal
$P(Y_S|A_S)$

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
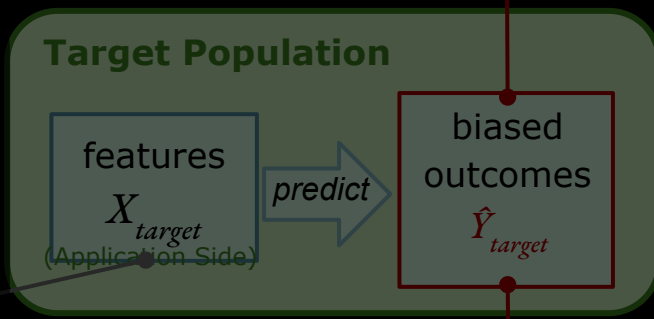$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Label Bias - Example: Label word with drawing
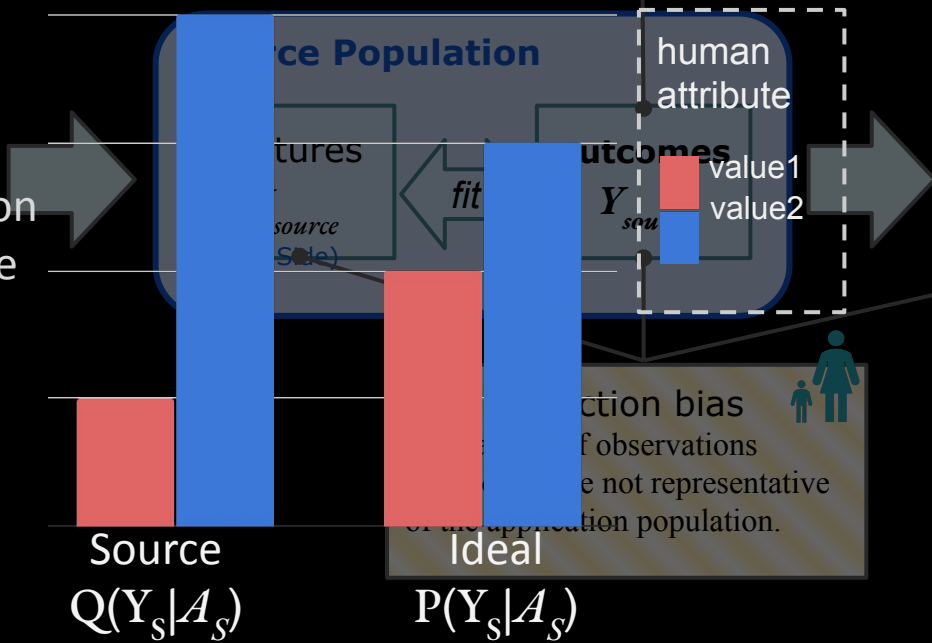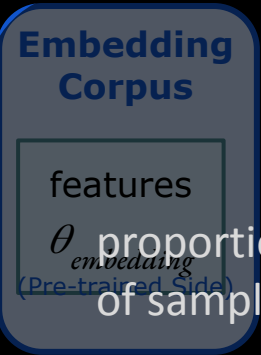


South Africa    Russia

Korea    Brazil

**label bias**
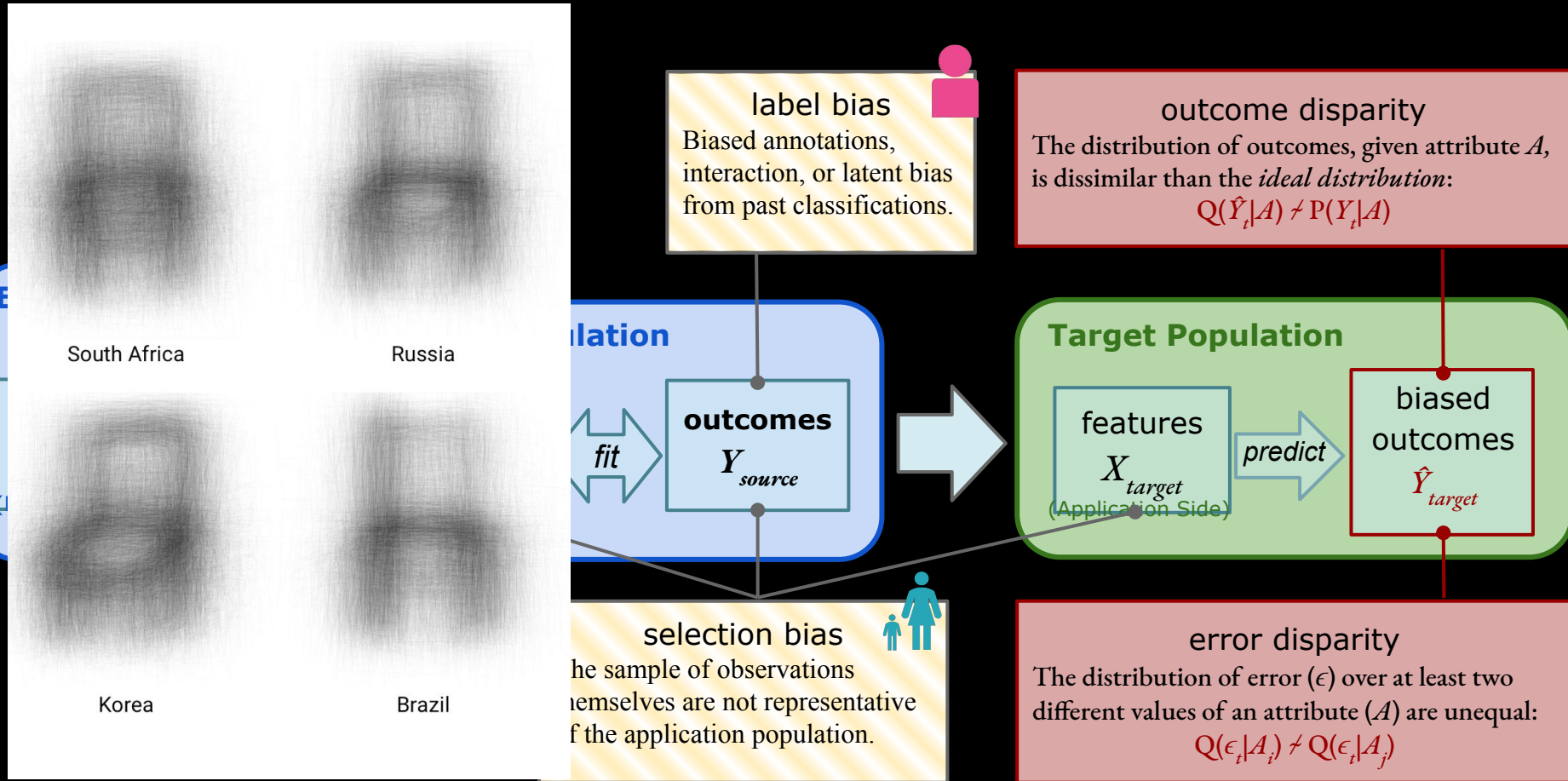Biased annotations, interaction, or latent bias from past classifications.

**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
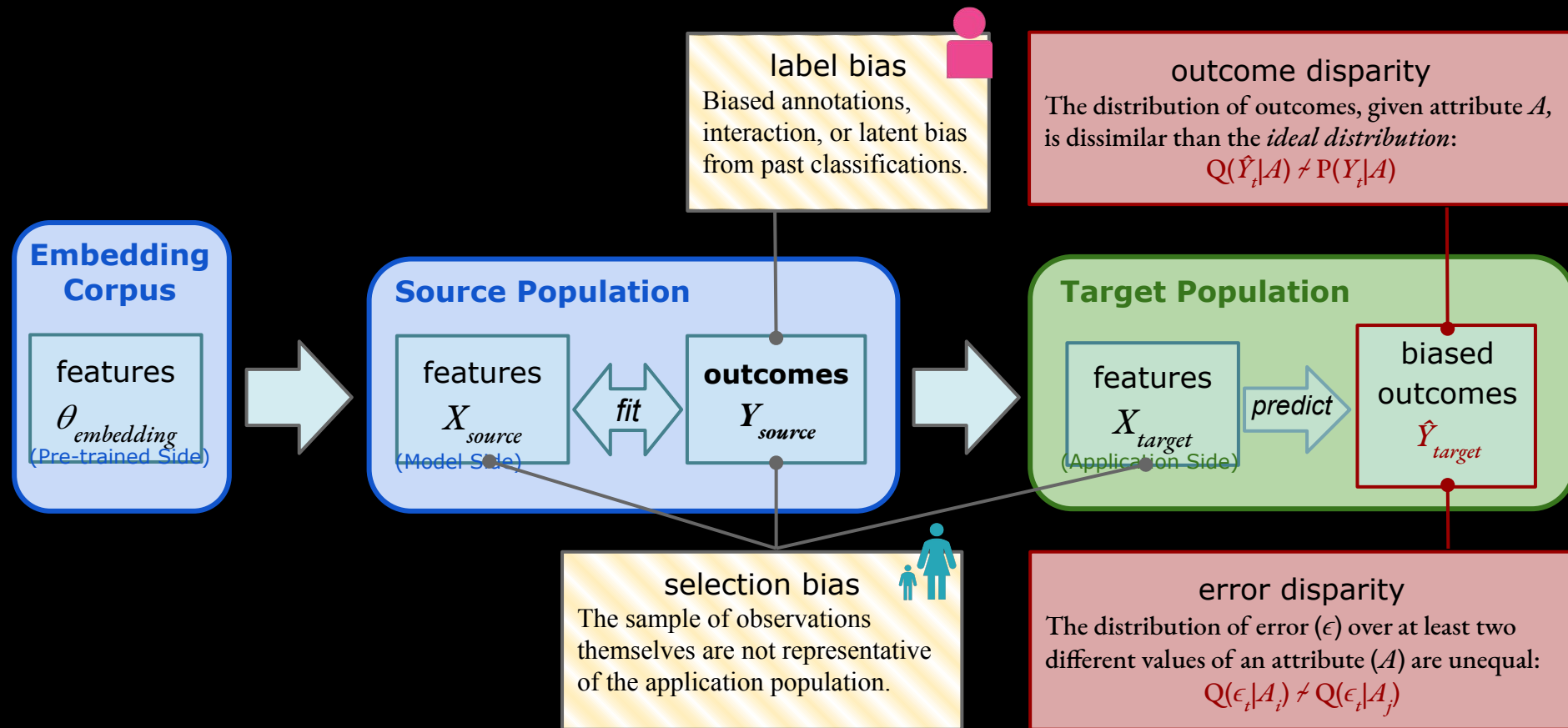$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**Population**

$\xrightarrow{fit}$ **outcomes** $Y_{source}$

**Target Population**

features $X_{target}$ (Application Side) $\xrightarrow{predict}$ biased outcomes $\hat{Y}_{target}$

**selection bias**
The sample of observations themselves are not representative of the application population.

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

Devin Coldeway. 2017. TechCrunch: Google releases millions of bad drawings for you (and your AI) to paw through
https://techcrunch.com/2017/08/25/google-releases-millions-of-bad-drawings-for-you-and-your-ai-to-paw-through/

# Label Bias



**Embedding Corpus**

features

$\theta_{embedding}$

(Pre-trained Side)

**Source Population**

features

$X_{source}$

(Model Side)

*fit*

**outcomes**

$Y_{source}$

**Target Population**

features

$X_{target}$

(Application Side)

*predict*

biased outcomes

$\hat{Y}_{target}$

**label bias**
Biased annotations, interaction, or latent bias from past classifications.

**selection bias**
The sample of observations themselves are not representative of the application population.

**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Overamplification



**over-amplification**
The model discriminates on a given human attribute beyond its source base-rate.

**label bias**
Biased annotations, interaction, or latent bias from past classifications.
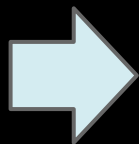
**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**Embedding Corpus**

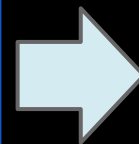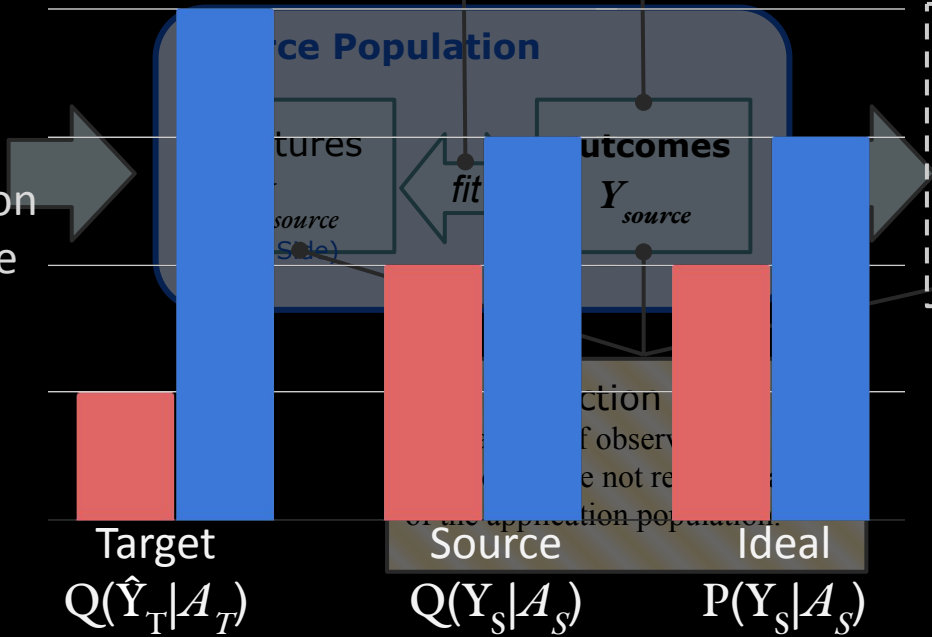features
$\theta_{embedding}$
(Pre-trained Side)

**Source Population**

features
$X_{source}$
(Model Side)

*fit*

**outcomes**
$Y_{source}$

**Target Population**

features
$X_{target}$
(Application Side)

*predict*

biased outcomes
$\hat{Y}_{target}$

**selection bias**
The sample of observations themselves are not representative of the application population.

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Overamplification



**over-amplification**
The model discriminates on a given human attribute beyond its source base-rate.

**label bias**
Biased annotations, interaction, or latent bias from past classifications.

**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**Embedding Corpus**
features
$\theta_{embedding}$
(Pre-trained Side)
proportion of sample

ce Population
features
*source*
(...de)

*fit*

utcomes
$Y_{source}$

human attribute

get Population
features
value1
value2
*target*
(Application Side)

*predict*

biased outcomes
$\hat{Y}_{target}$

ction
f observ
e not r
of the application population.

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

Target
$Q(\hat{Y}_T|A_T)$

Source
$Q(Y_S|A_S)$

Ideal
$P(Y_S|A_S)$

# Overamplification

**over-amplification**
The model discriminates on a given human attribute beyond its source base-rate.

**label bias**
Biased annotations, interaction, or latent bias from past classifications.

**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**Embedding Corpus**

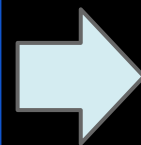features
$\theta_{embedding}$
(Pre-trained Side)

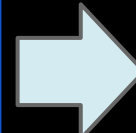**Source Population**

features
$X_{source}$
(Model Side)

*fit*

**outcomes**
$Y_{source}$

**Target Population**

features
$X_{target}$
(Application Side)

*predict*

biased outcomes
$\hat{Y}_{target}$

**selection bias**
The sample of observations themselves are not representative of the application population.

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Semantic Bias

**over-amplification**
The model discriminates on a given human attribute beyond its source base-rate.

**label bias**
Biased annotations, interaction, or latent bias from past classifications.

**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$
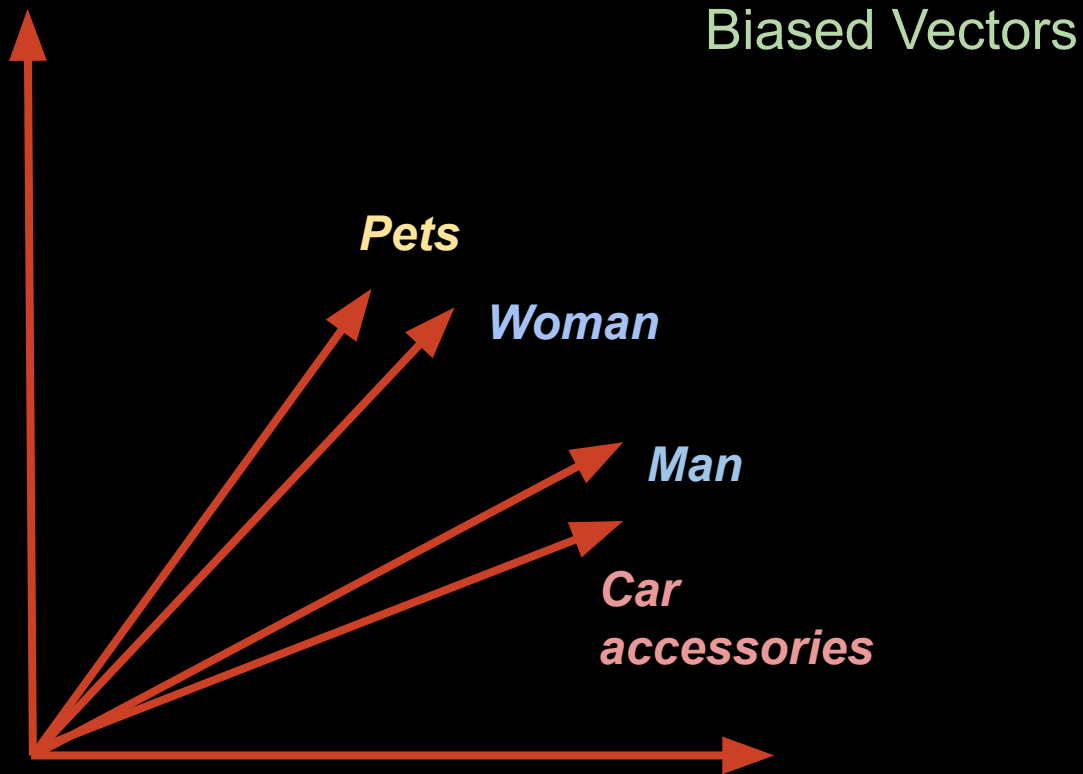
**Embedding Corpus**

**features**
$\theta_{embedding}$
(Pre-trained Side)

**Source Population**

features
$X_{source}$
(Model Side)

*fit*

outcomes
$Y_{source}$

**Target Population**

features
$X_{target}$
(Application Side)

*predict*

biased outcomes
$\hat{Y}_{target}$

**semantic bias**
Non-ideal associations between attributed lexeme (e.g. gendered pronouns) and non-attributed lexeme (e.g. occupation).

**selection bias**
The sample of observations themselves are not representative of the application population.

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

E.g. Coreference resolution:
  connecting entities to references (i.e. pronouns).

"*The doctor told Mary that she had run some blood tests.*"

**semantic bias**
Non-ideal associations between attributed lexeme (e.g. gendered pronouns) and non-attributed lexeme (e.g. occupation).

**selection bias**
The sample of observations themselves are not representative of the application population.

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t | A_i) \neq Q(\epsilon_t | A_j)$$

Shah, D., Schwartz, H. A., Hovy, D. (2020). Predictive Biases in Natural Language Processing Models: A Conceptual Framework and Overview. *In ACL-2020: Proceedings of the Association for Computational Linguistics.*

# Predictive Bias Framework for NLP

origin

consequence

**over-amplification**
The model discriminates on a given human attribute beyond its source base-rate.

**label bias**
Biased annotations, interaction, or latent bias from past classifications.

**outcome disparity**
The distribution of outcomes, given attribute $A$, is dissimilar than the *ideal distribution*:
$$Q(\hat{Y}_t|A) \neq P(Y_t|A)$$

**Embedding Corpus**

features
$\theta_{embedding}$
(Pre-trained Side)

**Source Population**

features
$X_{source}$
(Model Side)

*fit*

outcomes
$Y_{source}$

**Target Population**

features
$X_{target}$
(Application Side)

*predict*

biased outcomes
$\hat{Y}_{target}$

**semantic bias**
Non-ideal associations between attributed lexeme (e.g. gendered pronouns) and non-attributed lexeme (e.g. occupation).

**selection bias**
The sample of observations themselves are not representative of the application population.

**error disparity**
The distribution of error ($\epsilon$) over at least two different values of an attribute ($A$) are unequal:
$$Q(\epsilon_t|A_i) \neq Q(\epsilon_t|A_j)$$

# Summary of Countermeasures

| Source | Origin | Countermeasures |
|---|---|---|
| annotation | **Label Bias** | Post-stratification, Re-train annotators |
| data selection | **Selection Bias** | Stratified sampling, Post-stratification or Re-weighing techniques |
| NLP models | **Overamplification** | Synthetically match distributions, add outcome disparity to cost function |
| embeddings | **Semantic Bias** | Use above techniques and re-train embeddings |

# Bias - Takeaways

Bias, as outcome and error **disparities**, can result from many **origins**:
- the **embedding** model
- the feature **sample**
- the **fitting** process
- the **outcome** sample

Our understanding is evolving:
 This is an active area of work, both theoretically and technically!

# Ethics in NLP

Bias


Privacy


Ethical Research

# Ethics in NLP

## Privacy

- Risk Categories:
  - Revealing unintended private information
  - Targeted persuasion

# Ethics in NLP

## Privacy

- Risk Categories:
  - Revealing unintended private information
  - Targeted persuasion
- Mitigation strategies:

# Ethics in NLP

## Privacy

- Risk Categories:
  - Revealing unintended private information
  - Targeted persuasion
- Mitigation strategies:
  - Informed consent -- let participants know and opportunity to opt-in/-out
  - Do not share / secure storage
  - *Federated learning* -- obfuscate to the point of preserving privacy
  - Transparency in information targeting

    "You are being shown this ad because ..."

# Ethics in NLP

Bias

Privacy

Ethical Research

# Ethics in NLP Research

## ACM Code of Ethics; General Ethical Principles:

- Contribute to society and to human well-being, acknowledging that all people are stakeholders in computing.

- Avoid harm.

- Be honest and trustworthy.

- Be fair and take action not to discriminate.

- Respect the work required to produce new ideas, inventions, creative works, and computing artifacts.

- Respect privacy.

- Honor confidentiality.

https://www.acm.org/code-of-ethics

# Ethics in NLP

Human Subjects Research

Observational versus Interventional

# Ethics in NLP

Human Subjects Research

## Observational versus Interventional

(The Belmount Report, 1979)

 (i) Distinction of research from practice.
(ii) Risk-Benefit criteria
(iii) Appropriate selection of human subjects for participation in research
(iv) Informed consent in various research settings.