



# Applied Natural Language Processing

Info 256

Lecture 5: Bias in word embeddings (Sept. 11, 2023)

David Bamman, UC Berkeley

# Analogical inference

- Mikolov et al. 2013 show that vector representations have some potential for analogical reasoning through vector arithmetic.

apple - apples  $\approx$  car - cars

king - man + woman  $\approx$  queen

## SHARE

## REPORT



0



13

# Semantics derived automatically from language corpora contain human-like biases

Aylin Caliskan<sup>1,\*</sup>, Joanna J. Bryson<sup>1,2,\*</sup>, Arvind Narayanan<sup>1,\*</sup>

[+ See all authors and affiliations](#)

*Science* 14 Apr 2017:  
Vol. 356, Issue 6334, pp. 183-186  
DOI: 10.1126/science.aal4230



Peer Reviewed  
← see details

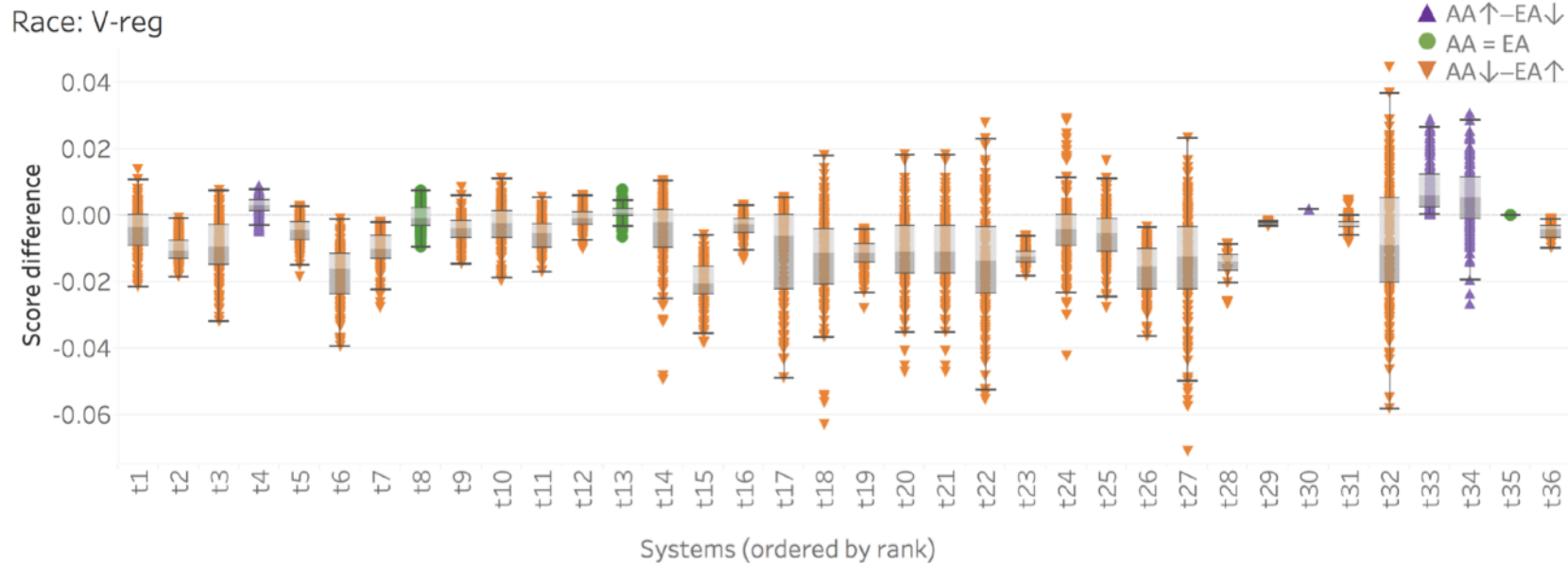
[Article](#)[Figures & Data](#)[Info & Metrics](#)[eLetters](#)[PDF](#)

# Bias

- Allocational harms: automated systems allocate resources unfairly to different groups (access to housing, credit, parole).
- Representational harms: automated systems represent one group less favorably than another (including demeaning them or erasing their existence).

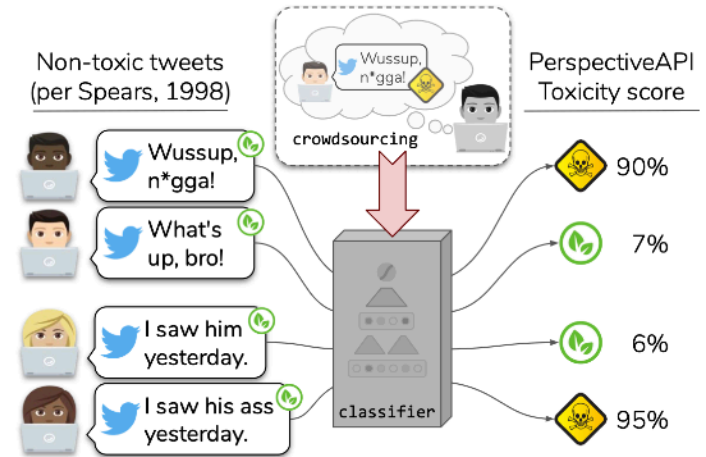
# Representations

- **Pleasant:** caress, freedom, health, love, peace, cheer, friend, heaven, loyal, pleasure, diamond, gentle, honest, lucky, rainbow, diploma, gift, honor, miracle, sunrise, family, happy, laughter, paradise, vacation.
- **Unpleasant:** abuse, crash, filth, murder, sickness, accident, death, grief, poison, stink, assault, disaster, hatred, pollute, tragedy, bomb, divorce, jail, poverty, ugly, cancer, evil, kill, rotten, vomit.
- Embeddings for African-American first names are closer to “unpleasant” words than European-American names (Caliskan et al. 2017)



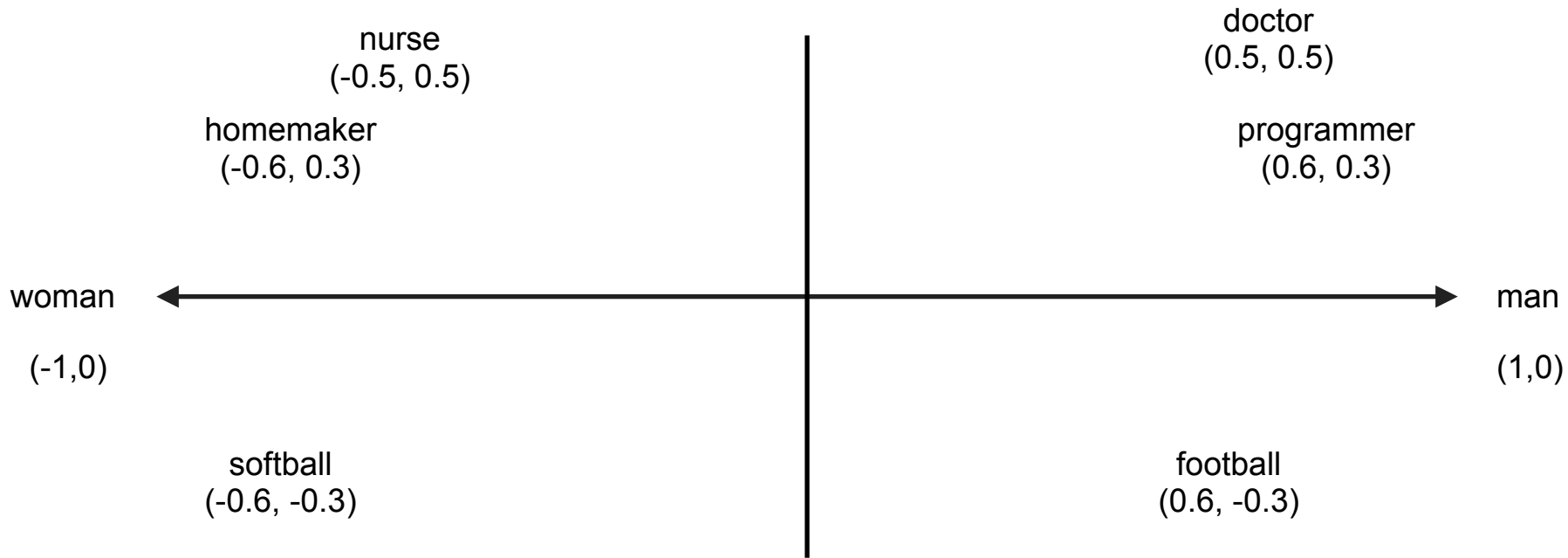
- Sentiment analysis over sentences containing African-American first names are more negative than identical sentences with European-American names.

- Toxicity detection systems score text with African-American English as more offensive
- Implicit negative perception of AAE → more AAE tweets are removed → users change language practices



Blodgett et al. (2020); Sap et al. (2019), "The risk of racial bias in hate speech detection"

# Bias





# Notation

$$x^{\top} b$$

$$x \cdot b$$

$$\text{dot}(x, b)$$

$$\sum_i x_i b_i$$

$$x = [3, 1, 2]$$

$$y = [0, 5, 2]$$

$$x \cdot y =$$

$$3 \times 0 + 1 \times 5 + 2 \times 2$$

$$= 9$$

# Cosine Similarity

$$\cos(x, y) = \frac{\sum_{i=1}^F x_i y_i}{\sqrt{\sum_{i=1}^F x_i^2} \sqrt{\sum_{i=1}^F y_i^2}}$$

# Cosine similarity

$$\cos(x, y) = \frac{\text{dot}(x, y)}{\sqrt{\text{dot}(x, x)} \times \sqrt{\text{dot}(y, y)}}$$

# Cosine similarity

$$\cos(x, y) = \frac{\text{dot}(x, y)}{\sqrt{\text{dot}(x, x)} \times \sqrt{\text{dot}(y, y)}}$$

This part can be done ahead of time by normalizing all vectors:

$$v = \frac{v}{\sqrt{\text{dot}(v, v)}}$$

If all vectors have been normalized in this way, cosine similarity is just the dot product:

$$\cos(x, y) = \text{dot}(x, y)$$

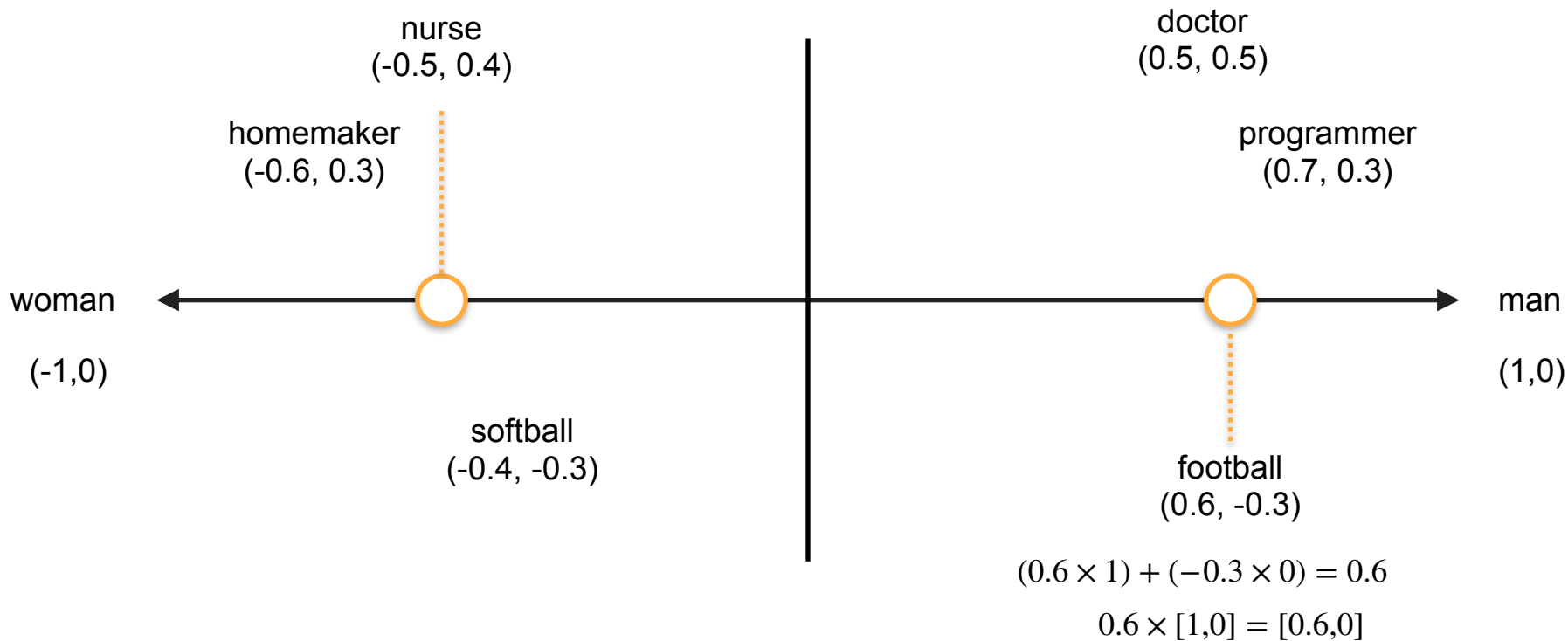
# Orthogonal projection

Assume all the vectors have  
been normalized to unit length

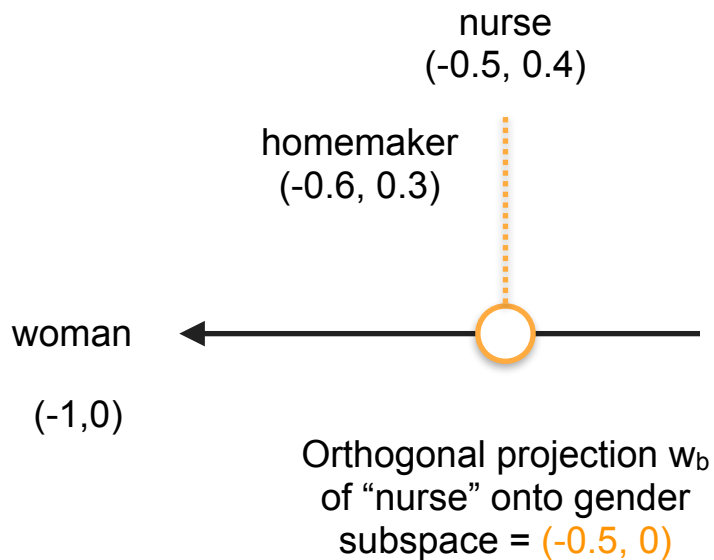
$$v = \frac{v}{\sqrt{\text{dot}(v, v)}}$$

$$x_b = \text{dot}(x, b) b$$

$$x_b = (x^\top b) b$$



# Orthogonal decomposition

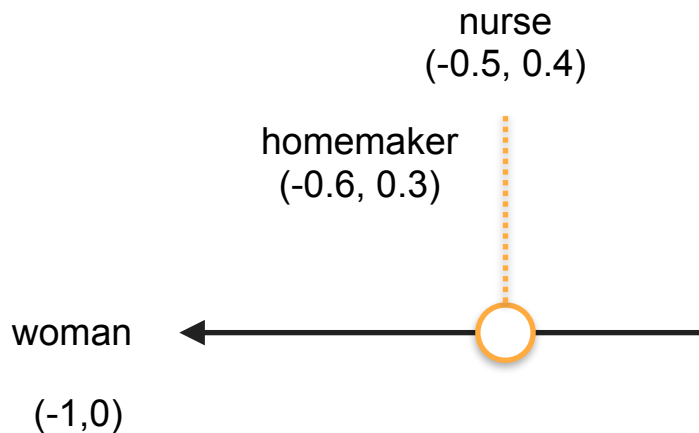


$$w = w_b + w_{b^\perp}$$

$$\begin{bmatrix} -0.5 \\ 0.4 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0.4 \end{bmatrix}$$

gender part    everything else

# Debiasing



$$\begin{aligned}w_d &= w_b^\perp \\ &= w - w_b\end{aligned}$$

$$\begin{bmatrix} -0.5 \\ 0.4 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0.4 \end{bmatrix}$$

gender part    everything else



# Bias

- The last slides illustrate this with a simple 2D subspace (where gender is effectively a 1D line).
- But the same principle (and procedure applies to any dimensionality (e.g., word embeddings of 100 dimensions)).

projection onto gender  
subspace

$$x_b = (x^\top b) b$$

debiasing by subtracting  
gender projection

$$x_d = x - x_b$$

# What's the gender subspace?

- Caliskan et al. 2018 construct this by first creating **defining sets** of gendered terms, e.g.
  - $D_1 = \{\text{man, woman}\}$
  - $D_2 = \{\text{he, she}\}$
- Performing SVD over a covariance matrix within over all terms in the defining sets (mean-normalized)
- And defining a gender subspace to be the first row of the resulting SVD.

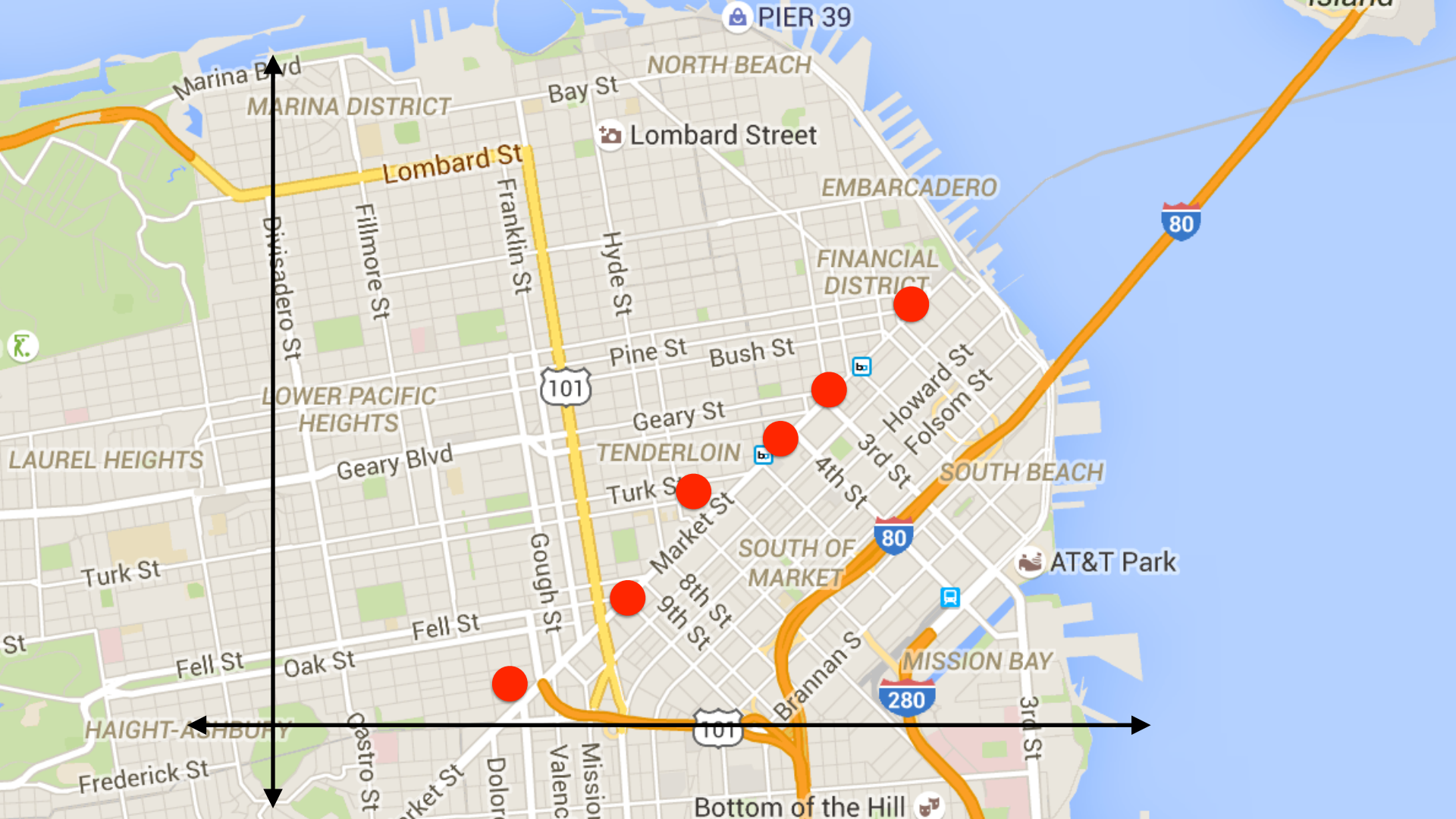
# Gender subspace

Vargas and Cotterell (2020) show that this is equivalent to PCA over the following matrix →

If each embedding is 100 dimensions, this matrix is [4 x 100] in size.

The gender subspace is then the first principle component (a 100-dimensional vector in this scenario).

man-mean(man, woman)
woman-mean(man, woman)
he-mean(he, she)
she-mean(he, she)



Marina Blvd

MARINA DISTRICT

NORTH BEACH

Lombard Street

EMBARCADERO

FINANCIAL DISTRICT

LOWER PACIFIC HEIGHTS

LAUREL HEIGHTS

TENDERLOIN

SOUTH BEACH

SOUTH OF MARKET

MISSION BAY

HAIGHT-ASHBURY

Bottom of the Hill

Divisadero St

Fillmore St

Franklin St

Hyde St

Pine St

Bush St

Geary St

Howard St

Folsom St

Geary Blvd

Turk St

Turk St

Market St

8th St

9th St

Fell St

Oak St

Fell St

Gough St

8th St

9th St

Brannan St

3rd St

Frederick St

Castro St

Market St

Dolor

Valenc

Mission

PIER 39

80

101

80

280

101

3rd St

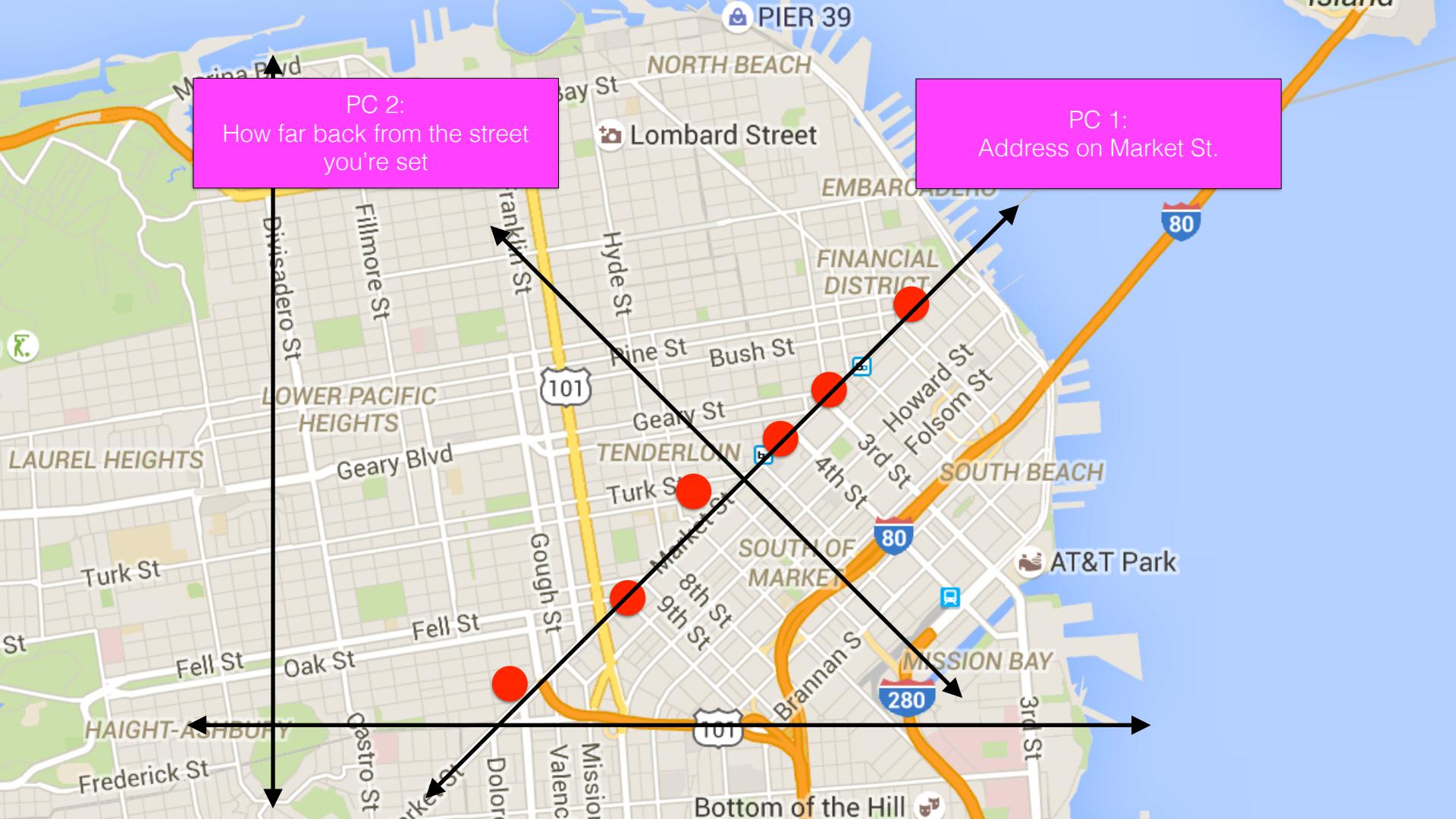
AT&T Park

# Principal Component Analysis

- Method for transforming a set of original (possible correlated) observations into new (**uncorrelated**) values.

PC 2:  
How far back from the street  
you're set

PC 1:  
Address on Market St.



- Original values: latitude and longitude (very strong correlation for these data points)
- Transformed values: street address and distance from street (no correlation)

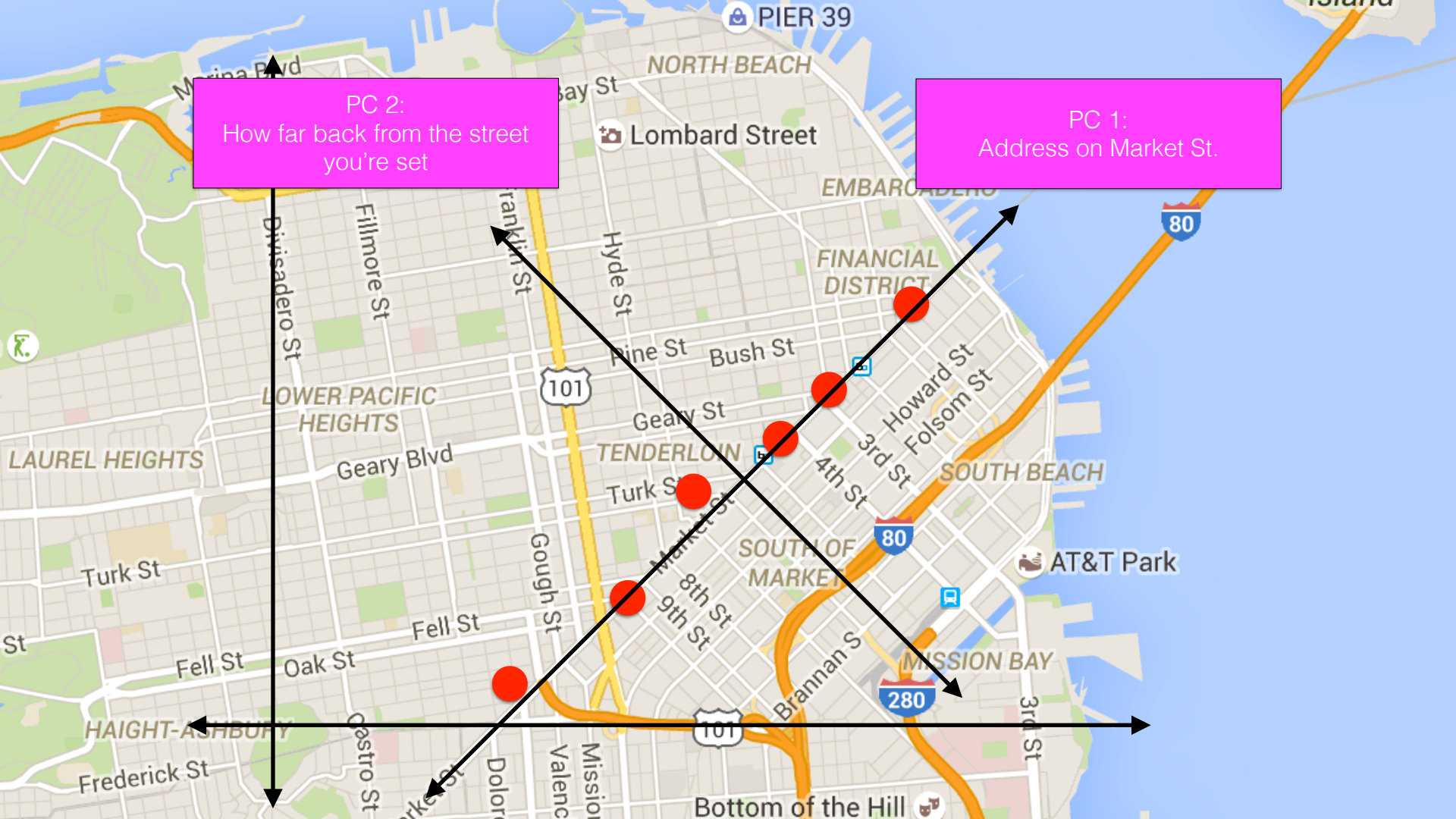
# Main idea

- Each principal component (1 ... F) is the axis that exhibits them most **variance** in the data and is uncorrelated (**orthogonal**) with earlier PCs
- The first PC explains the most variance; the second PC explains the most remaining variance, etc.



PC 2:  
How far back from the street  
you're set

PC 1:  
Address on Market St.



# Gender subspace

Vargas and Cotterell (2020) show that this is equivalent to PCA over the following matrix →

If each embedding is 100 dimensions, this matrix is [4 x 100] in size.

The gender subspace is then the first principle component (a 100-dimensional vector in this scenario).

man-mean(man, woman)
woman-mean(man, woman)
he-mean(he, she)
she-mean(he, she)

The **mean** of “man” and “woman” captures information that is common to both terms/embeddings (e.g. being people, animate, etc.). The difference is what’s left over to be explained.



# SemAxis

- Define a set of terms that comprise the endpoints of an axis of interest and average them up to form axis endpoint vectors.

$$S^- = \{v_1^-, \dots, v_n^-\}$$

{man, he, mr.}

$$V^- = \frac{1}{n} \sum_1^N v_i^-$$

$$S^+ = \{v_1^+, \dots, v_m^+\}$$

{woman, she, miss, mrs.}

$$V^+ = \frac{1}{M} \sum_1^M v_i^+$$

# SemAxis

- The axis vector is then the difference between the two endpoint vectors

{man, he, mr.}

$$V^- = \frac{1}{n} \sum_1^N v_i^-$$

{woman, she, miss, mrs.}

$$V^+ = \frac{1}{M} \sum_1^M v_i^+$$

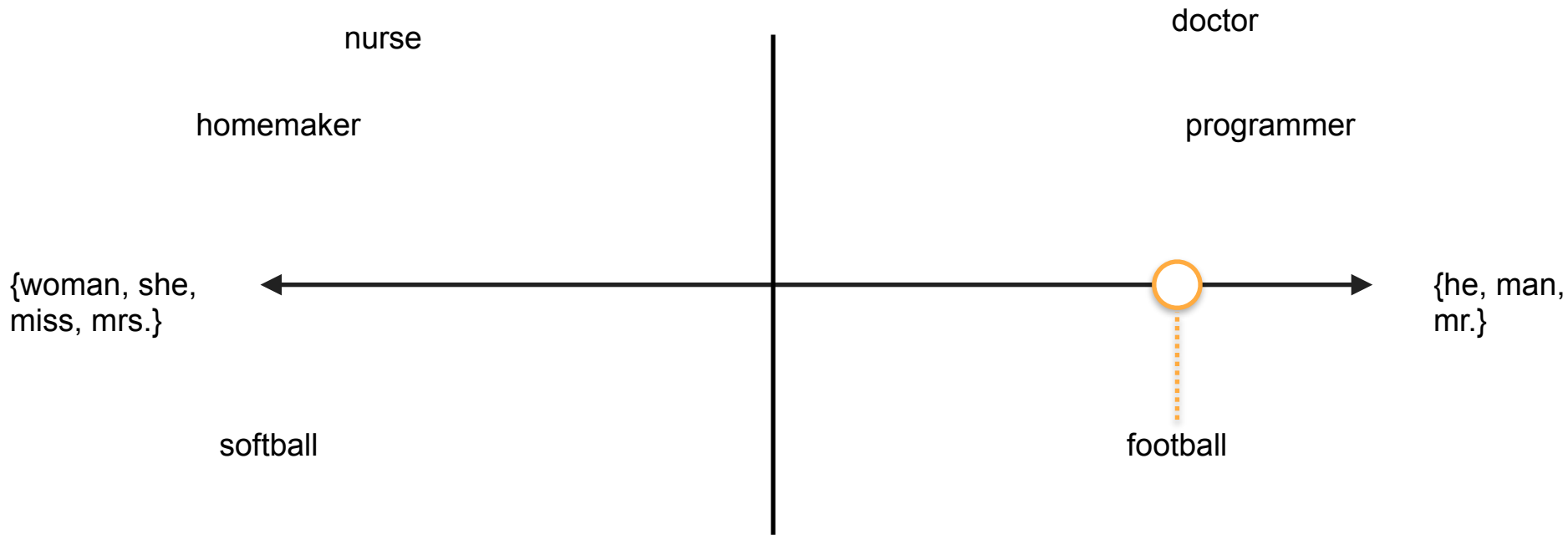
$$V_{\text{axis}} = V^+ - V^-$$

# SemAxis

- For any vector, we can find its position along this axis by taking the cosine similarity with it (or dot product if all the vectors are normalized to unit length)

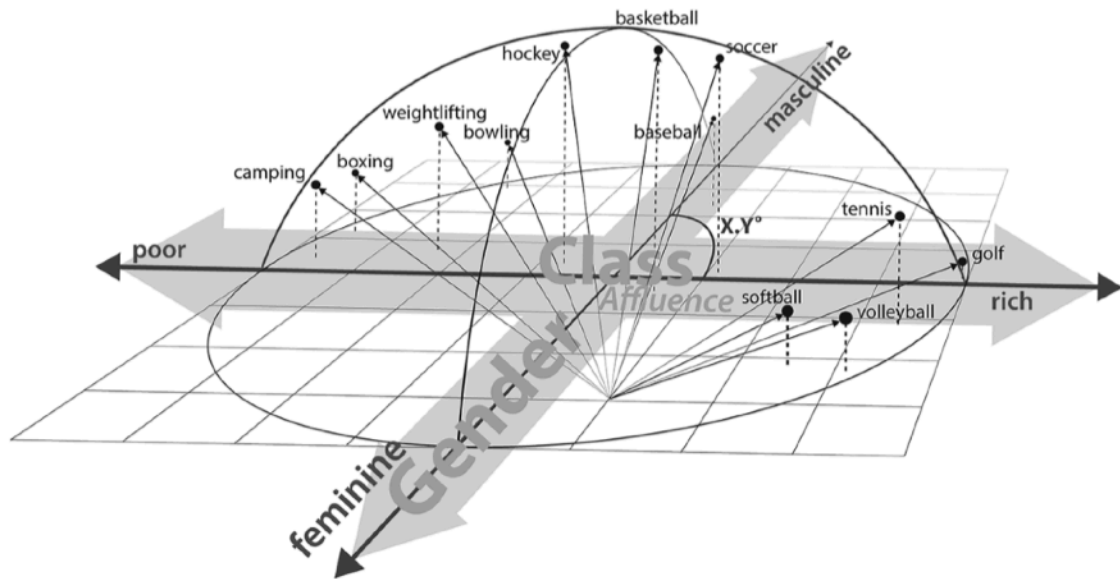
$$\text{Semaxis score} = \cos(\text{football}, V_{\text{axis}})$$

$$x_b = (x^\top b) b$$

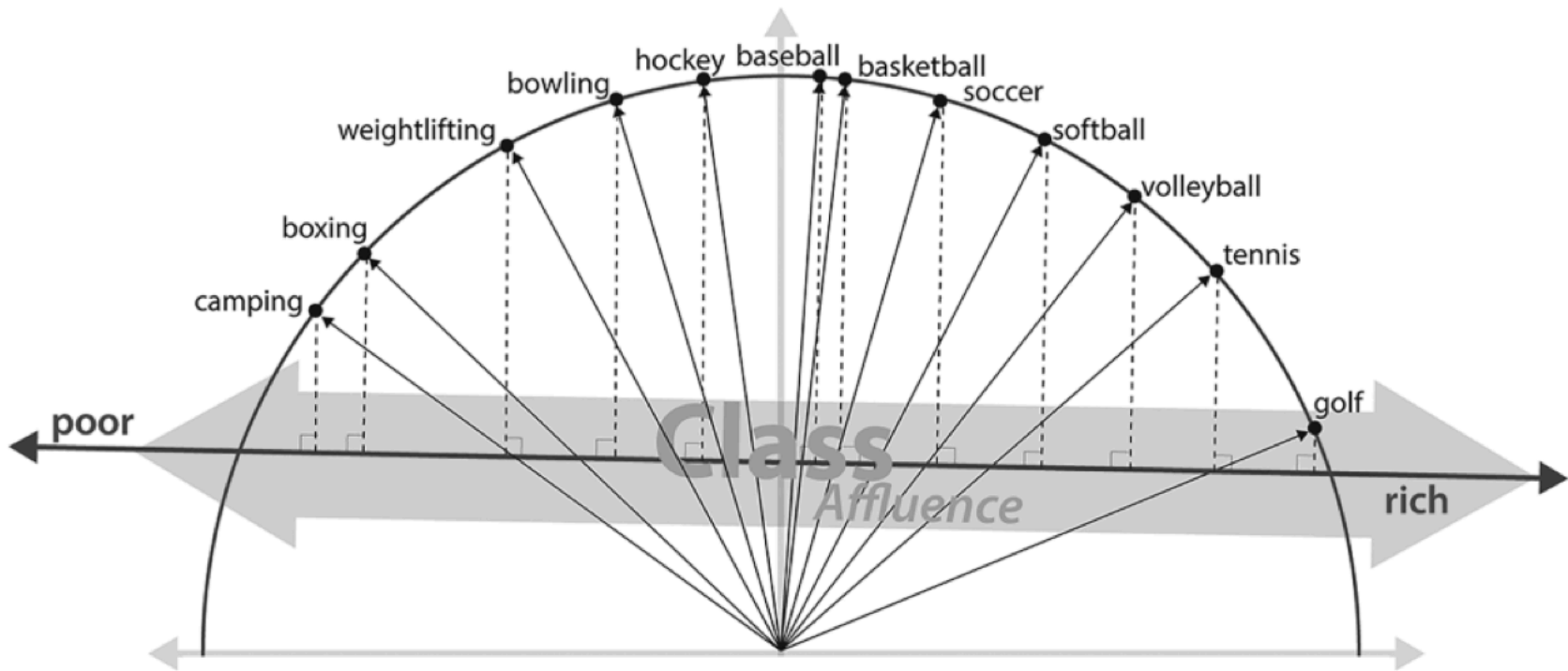


# Interrogating “bias”

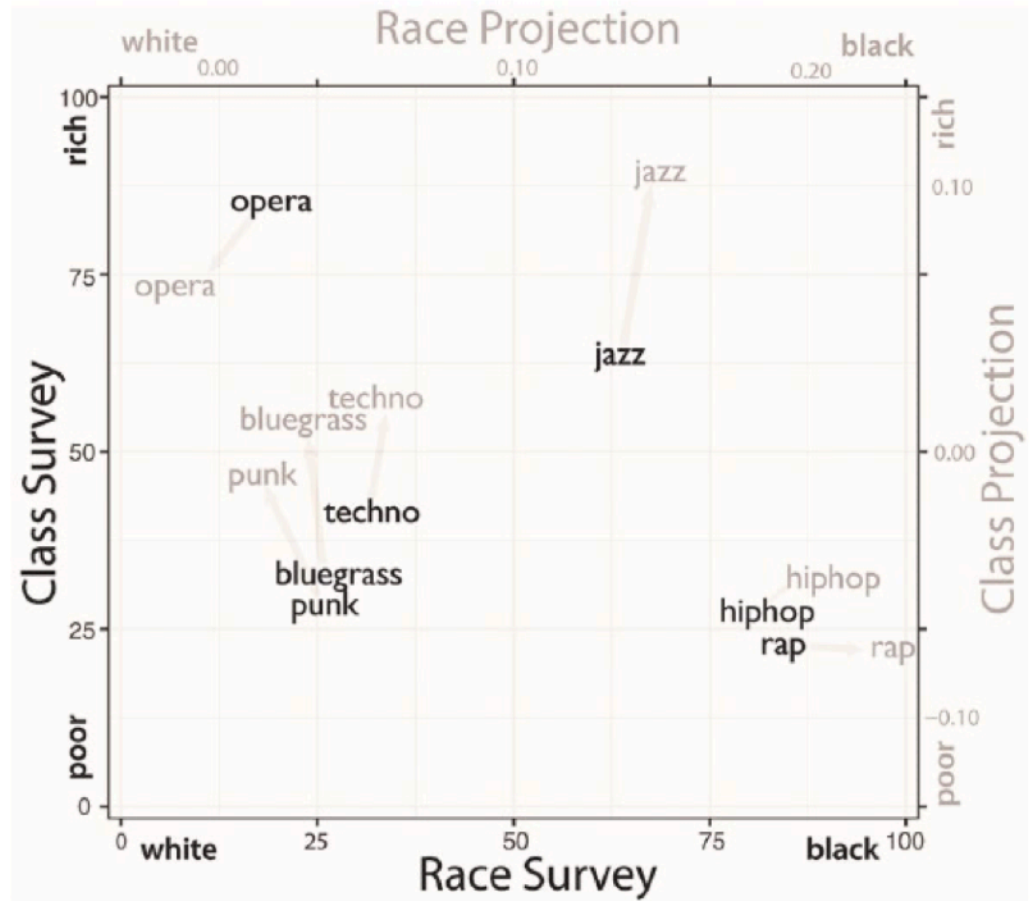
- Kozlowski et al. (2019), “The Geometry of Culture: Analyzing the Meanings of Class through Word Embeddings,” *American Sociological Review*.
- An et al. 2018, “SemAxis: A Lightweight Framework to Characterize Domain-Specific Word Semantics Beyond Sentiment”







Kozlowski et al. (2019); conceptual diagram (not real data)



Kozlowski et al. (2019)

**Table D2.** Word Pairs Used to Reconstruct 20 Semantic Differential Dimensions from Jenkins and Colleagues (1958) for Historical Survey Validation

<b>soft-hard</b> supple-tough delicate-dense pliable-rigid fluffy-firm mushy-solid softer-harder softest-hardest	<b>foolish-wise</b> dumb-smart irrational-rational stupid-thoughtful unwise-sensible silly-reasonable ridiculous-enlightened unintelligent-intelligent	<b>unimportant-important</b> inconsequential- consequential secondary-principal irrelevant-major trivial-crucial negligible-critical insignificant-significant unnecessary-essential peripheral-central	<b>fast-slow</b> quick-lagging rapid-unhurried speedy-sluggish swift-gradual quickly-slowly swiftly-gradually faster-slower fastest-slowest
<b>unusual-usual</b> different-customary abnormal-normal irregular-regular odd-standard atypical-typical unexpected-expected unconventional- conventional	<b>excitable-calm</b> volatile-tranquil nervous-still tempestuous-serene fiery-peaceful emotional-restful jumpy-sedate unsettled-settled	<b>strong-weak</b> powerful-powerless muscular-frail brawny-feeble strapping-puny sturdy-fragile robust-flimsy vigorous-languid	<b>colorful-colorless</b> brilliant-uncolored bright-pale radiant-drab vivid-pallid vibrant-lackluster colored-bleached
<b>rounded-angular</b> circular-cornered round-pointed dull-sharp smooth-jagged spherical-edged	<b>passive-active</b> immobile-mobile lethargic-energetic frail-vital subdued-vigorous static-dynamic subdued-lively	<b>true-false</b> true-untrue verifiable-erroneous veracious-fallacious accurate-inaccurate faithful-fraudulent correct-incorrect	<b>ugly-beautiful</b> unattractive-attractive unsightly-pretty hideous-handsome grotesque-gorgeous repulsive-cute

# Activity

- `SemAxis_TODO`: Implement the `SemAxis` method to define a conceptual axis using word embeddings and situate any word along that axis.
- Brainstorm other axes