



# WHAT CAN WORD VECTORS TELL US ABOUT SOCIETY?

Sandeep Soni

09/10/2024

# STORY SO FAR

- Tokenization: raw text to tokens
- Distributional hypothesis: Words used in similar contexts have similar meanings
- Many different ways to transform words into vectors to encode word semantics

# FROM LAST CLASS

	$d_1$	$d_2$	$d_3$	$\dots$	$d_m$
$w_1$	2				1
$w_2$		3			
$w_3$	1				
$\dots$					
$w_n$			5		

Term-Document Matrix

	$w_1$	$w_2$	$w_3$	$\dots$	$w_m$
$w_1$	2				1
$w_2$		3			
$w_3$	1				
$\dots$					
$w_n$			5		

Cooccurrence Matrix

	$w_1$	$w_2$	$w_3$	$\dots$	$w_n$
$w_1$	0.1				1.06
$w_2$		0.14			
$w_3$	-0.2				
$\dots$					
$w_n$			0.75		

Word Embedding Matrix

Words as distribution of counts over documents

Size =  $|V| \times |D|$

Sparse

Words as distribution of counts over cooccurrences

Size =  $|V| \times |V|$

Sparse

Words as learned vectors from cooccurrence data

Size =  $|V| \times k$

Dense

# FROM LAST CLASS

- Instead of counting, word2vec or GloVe are methods to learn vectors by predicting which words will co-occur together

x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

## Vocab

<i>service</i>
<i>wine</i>
<i>cold</i>
<i>great</i>

*wine*

## Vocab

x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

<i>service</i>
<i>wine</i>
<i>cold</i>
<i>great</i>

Vocab

X

0

*wine*

1

0

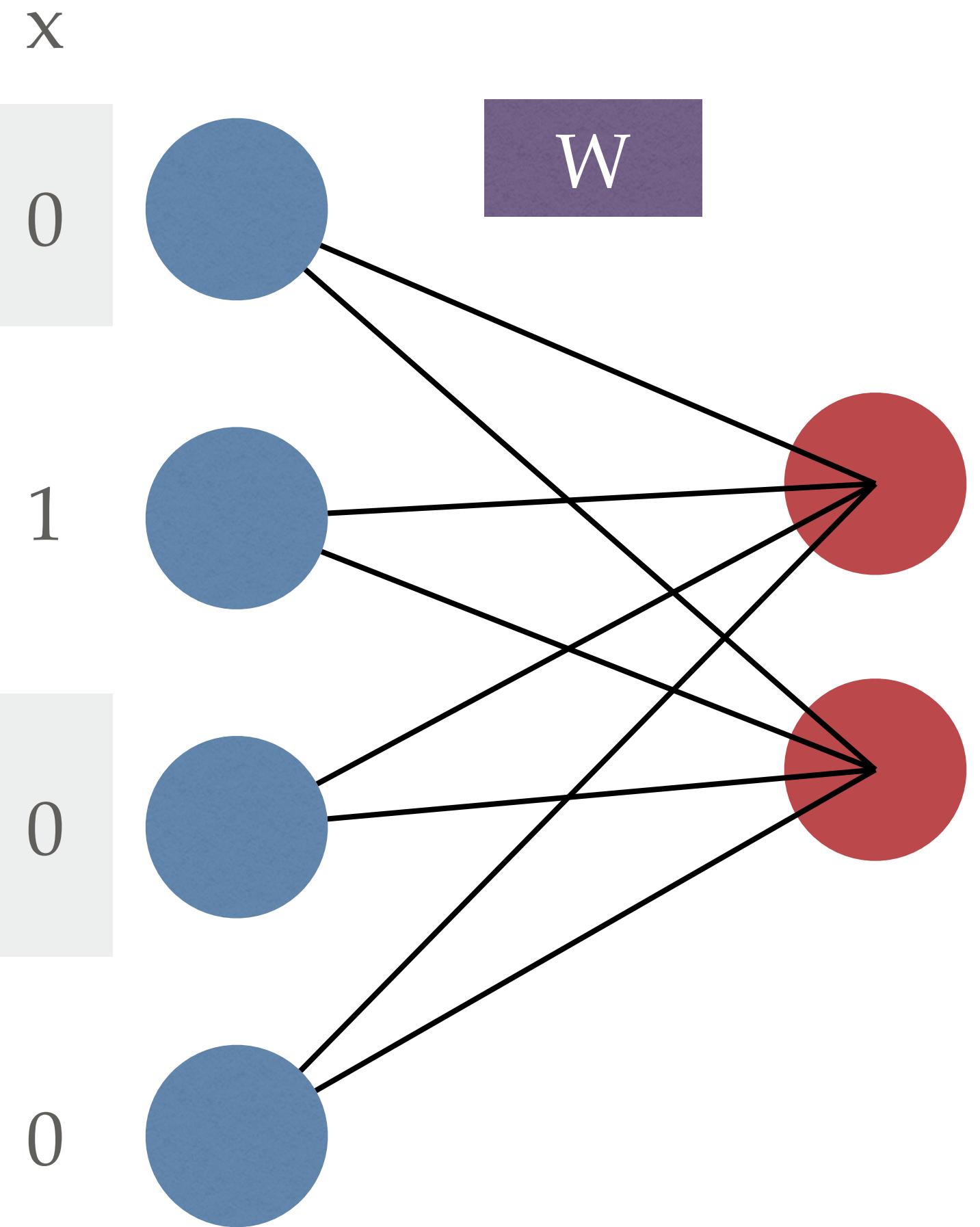
0

x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

<i>service</i>
<i>wine</i>
<i>cold</i>
<i>great</i>

Vocab

*wine*



x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

service	
wine	
cold	
great	

Vocab

wine

x

0

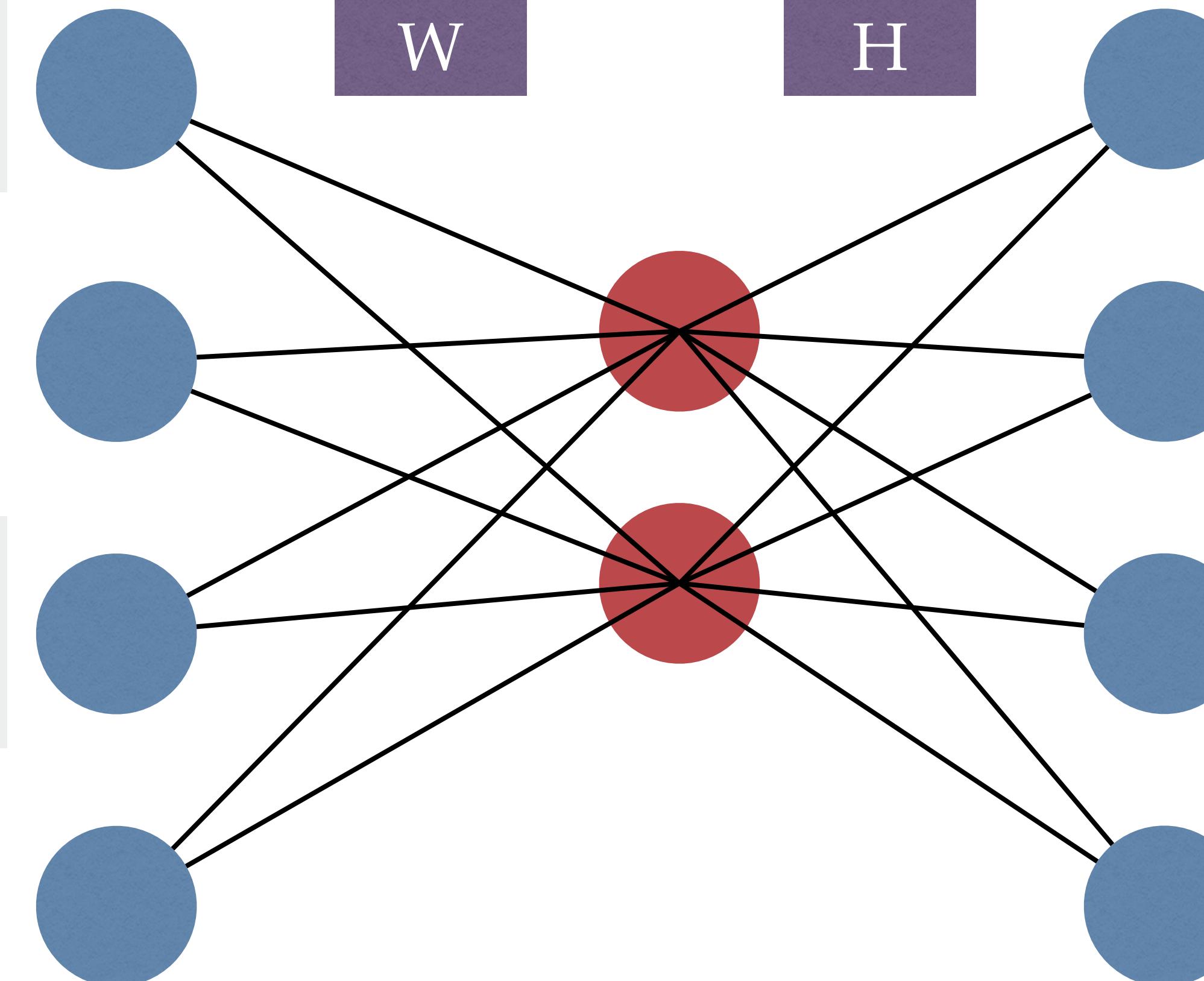
1

0

0

W

H

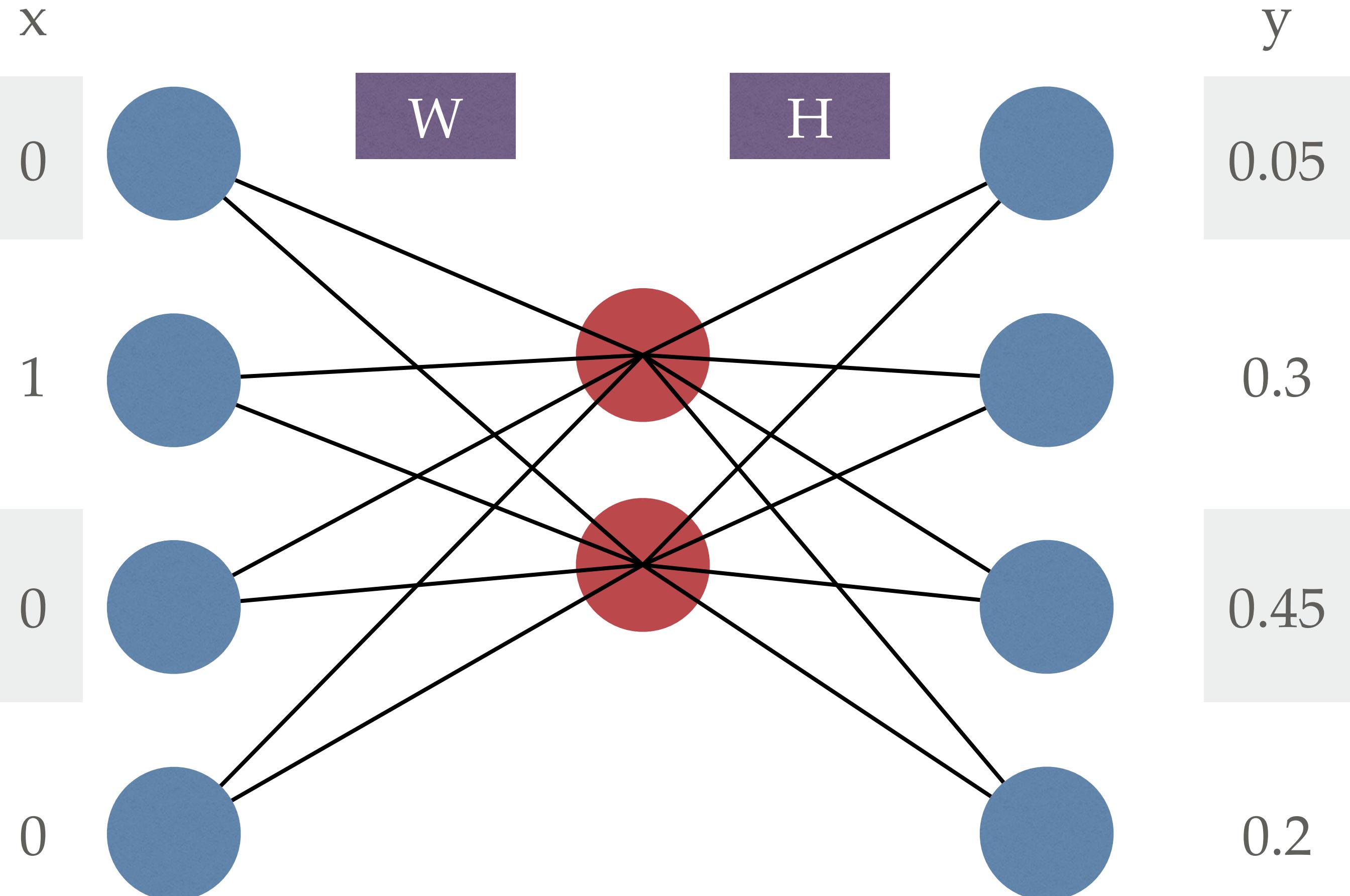


x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

<i>service</i>	
<i>wine</i>	
<i>cold</i>	
<i>great</i>	

Vocab

*wine*

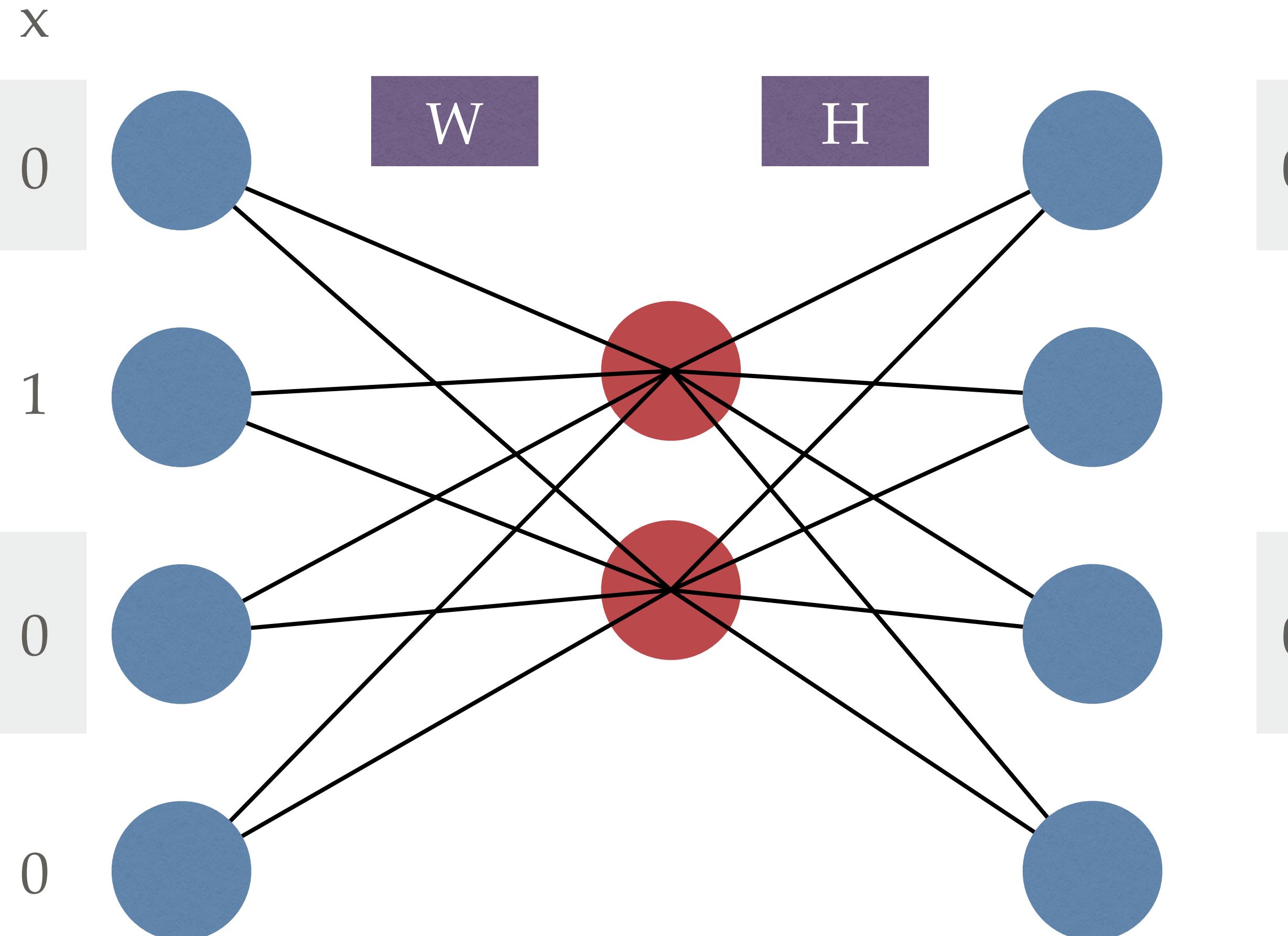


x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

<i>service</i>	
<i>wine</i>	
<i>cold</i>	
<i>great</i>	

Vocab

*wine*

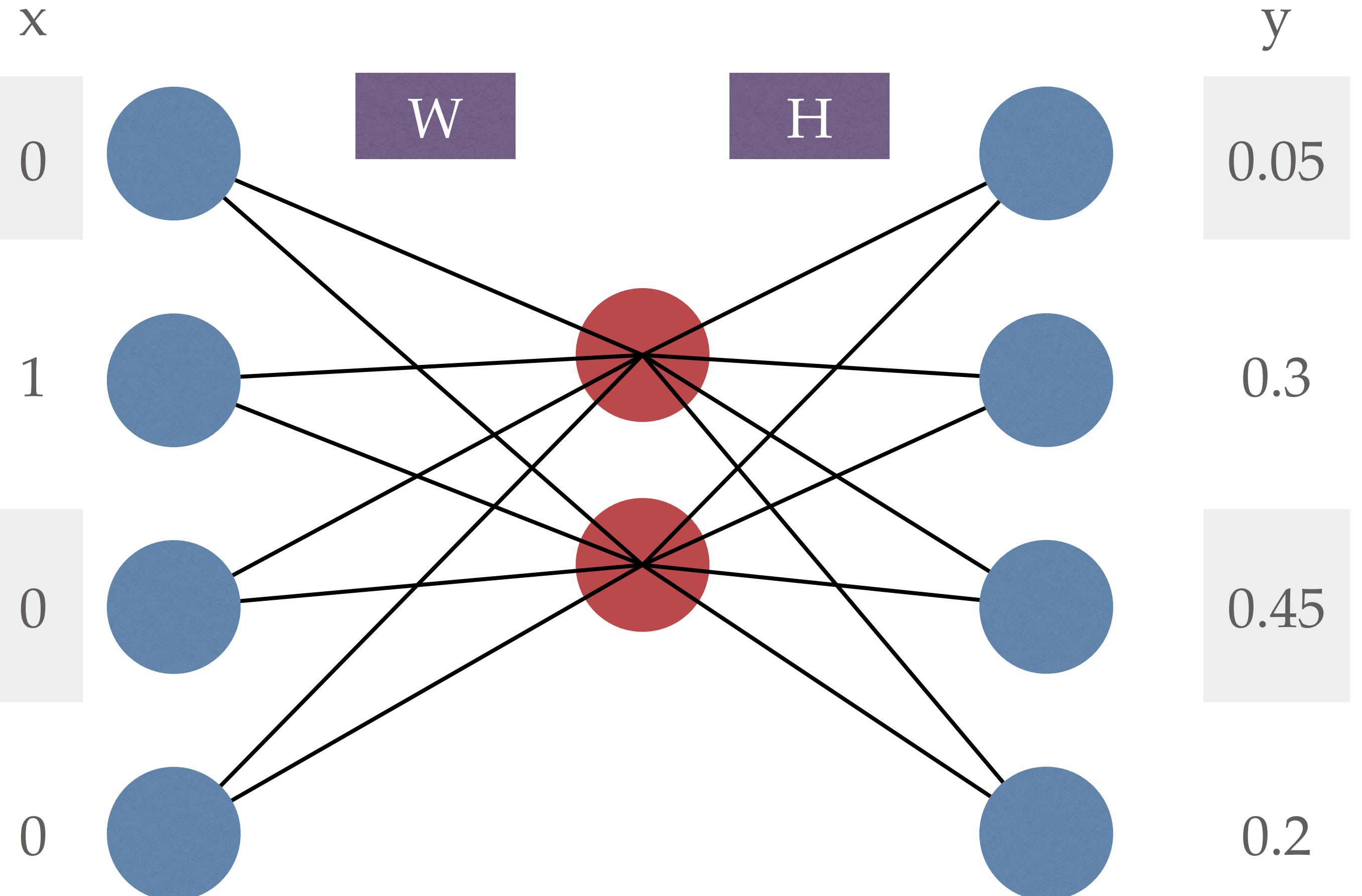


x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

<i>service</i>
<i>wine</i>
<i>cold</i>
<i>great</i>

Vocab

*wine*



W

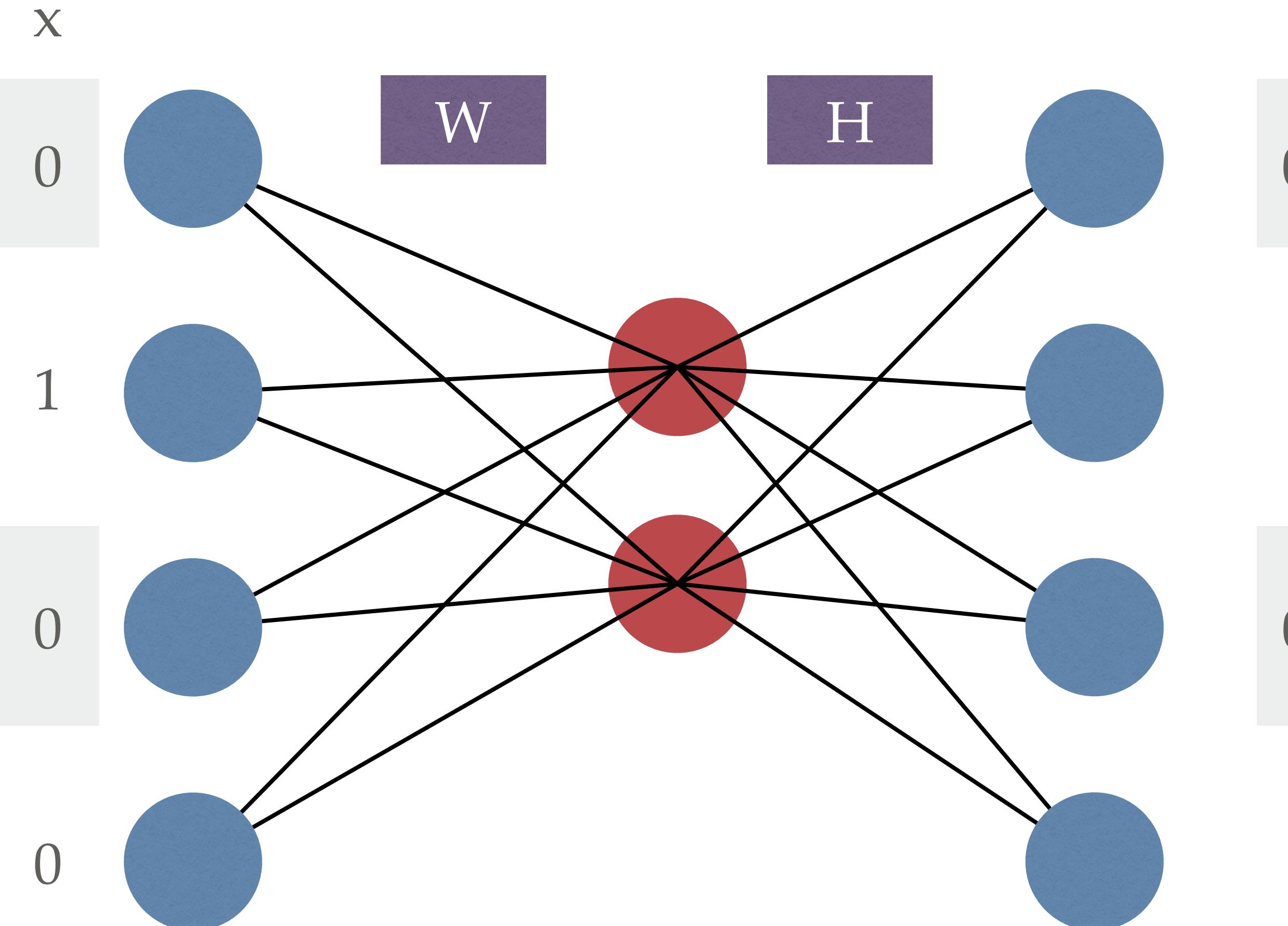
-0.3	1.2	0.5	-0.6
0.2	0.9	0.1	-0.4

x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

	x
service	0
wine	1
cold	0
great	0

Vocab

wine



x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

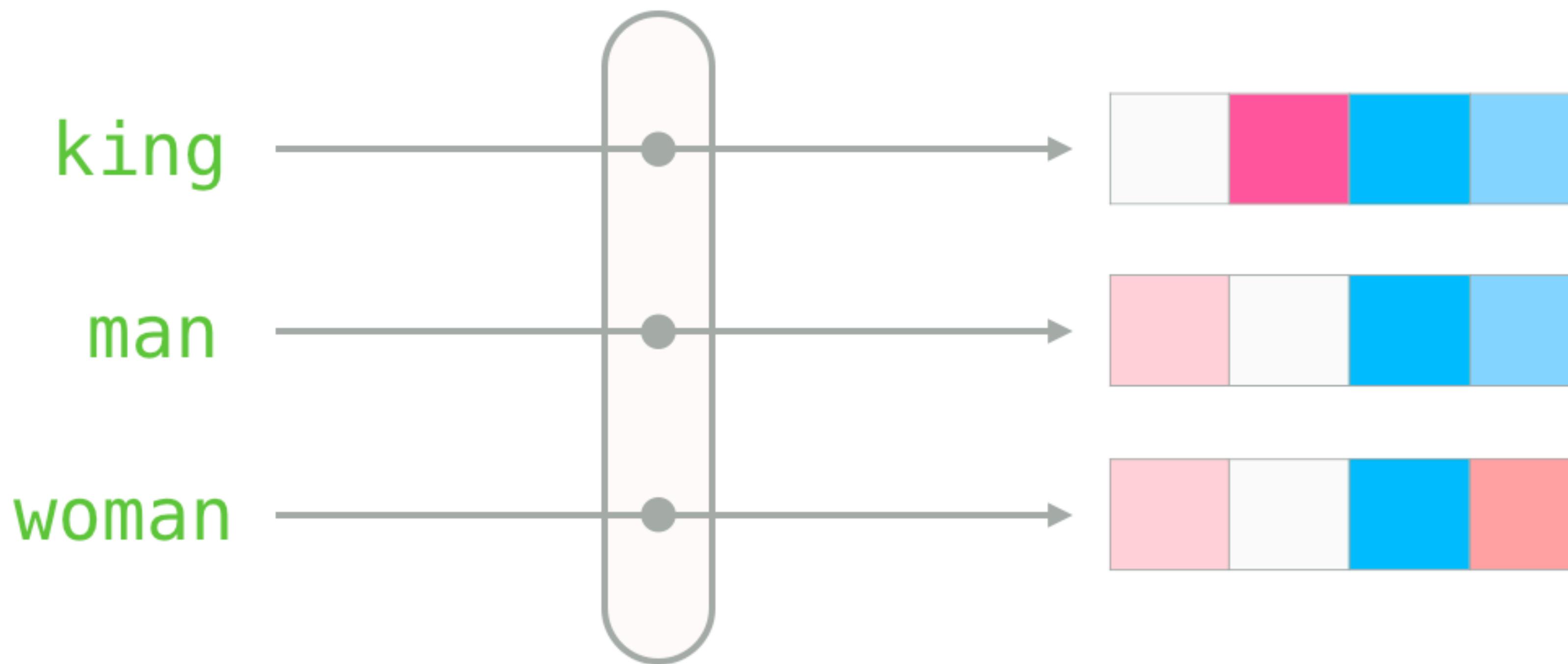
$W$

Word embeddings as columns			
-0.3	1.2	0.5	-0.6
0.2	0.9	0.1	-0.4

$H$

Context embeddings as rows	
0.1	-0.4
0.4	-0.5
0.3	-0.1
0.2	0.1

# Word2vec



<http://jalammar.github.io/illustrated-word2vec/>

# THINGS TO REMEMBER

- Embedding size
- Window size
- Corpus size
- Initialization

# WORD2VEC: FEW MORE CONSIDERATIONS

- Pretrain then finetune
- Word2Vec variations
  - CBoW  $\rightarrow P(focus\_word \mid context\_word)$
  - Skipgram  $\rightarrow P(context\_word \mid focus\_word)$
  - Hierarchical softmax or negative sampling to calculate probabilities

# IN CLASS

- Word2Vec demo

# QUESTION FOR THE DAY

“What do word embeddings encode and what can we do with them?”

# WORD EMBEDDINGS

- What is the geometry of word embeddings?
- What is their use as predictors?
- Can they be used to explain something about the world?

# GEOMETRY

# ANALOGICAL REASONING

man:woman::king:?

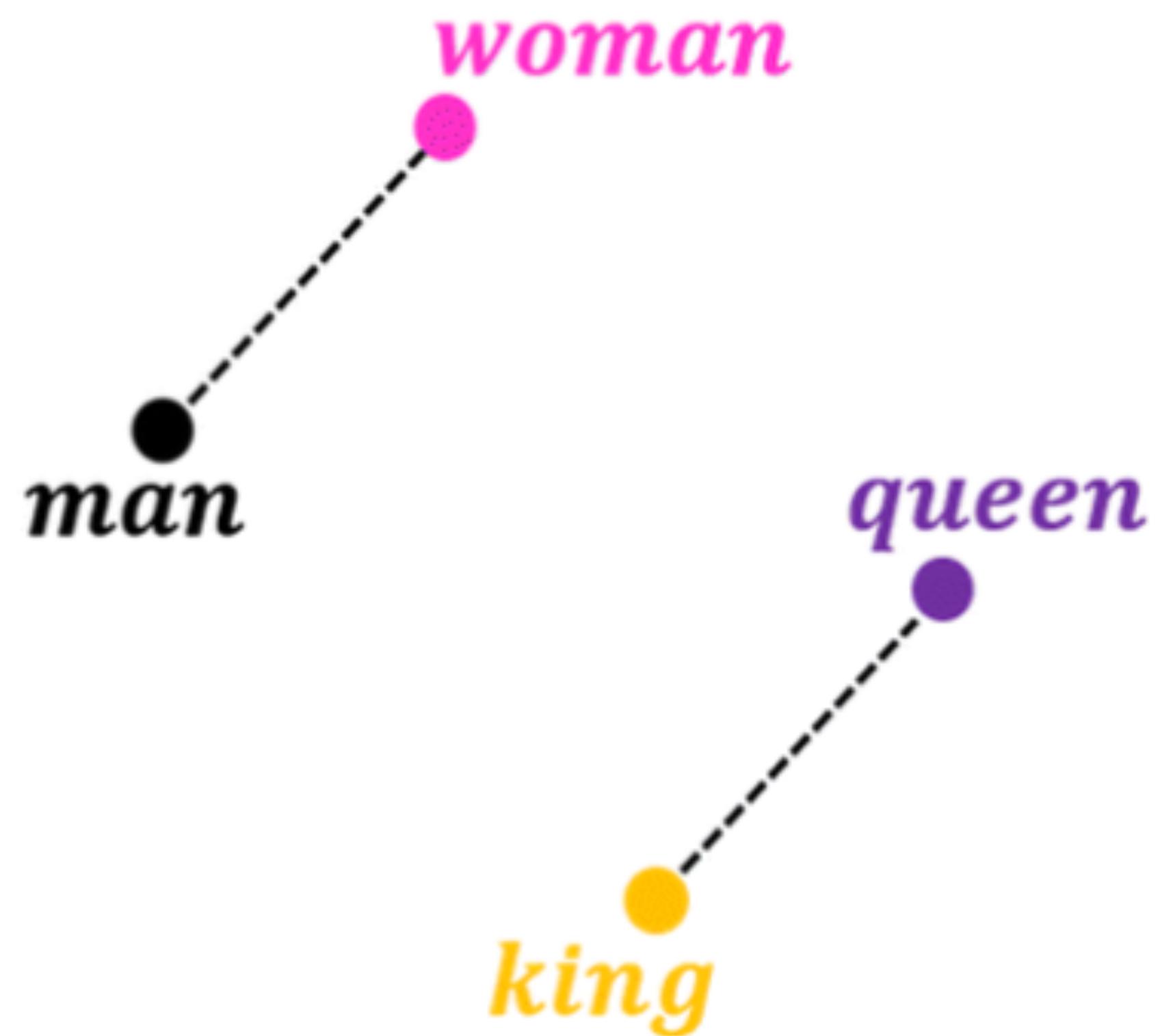
# GEOMETRY

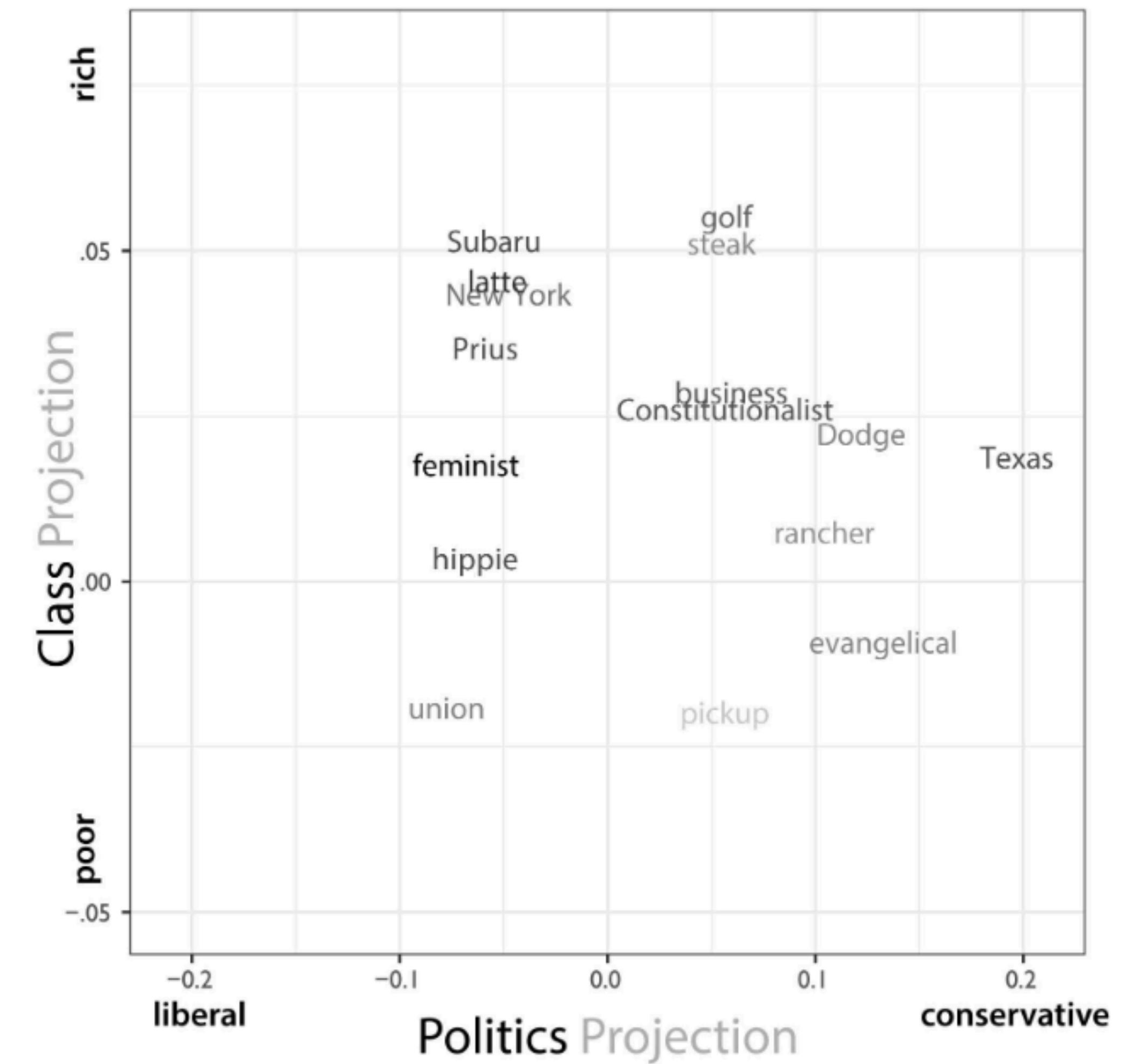
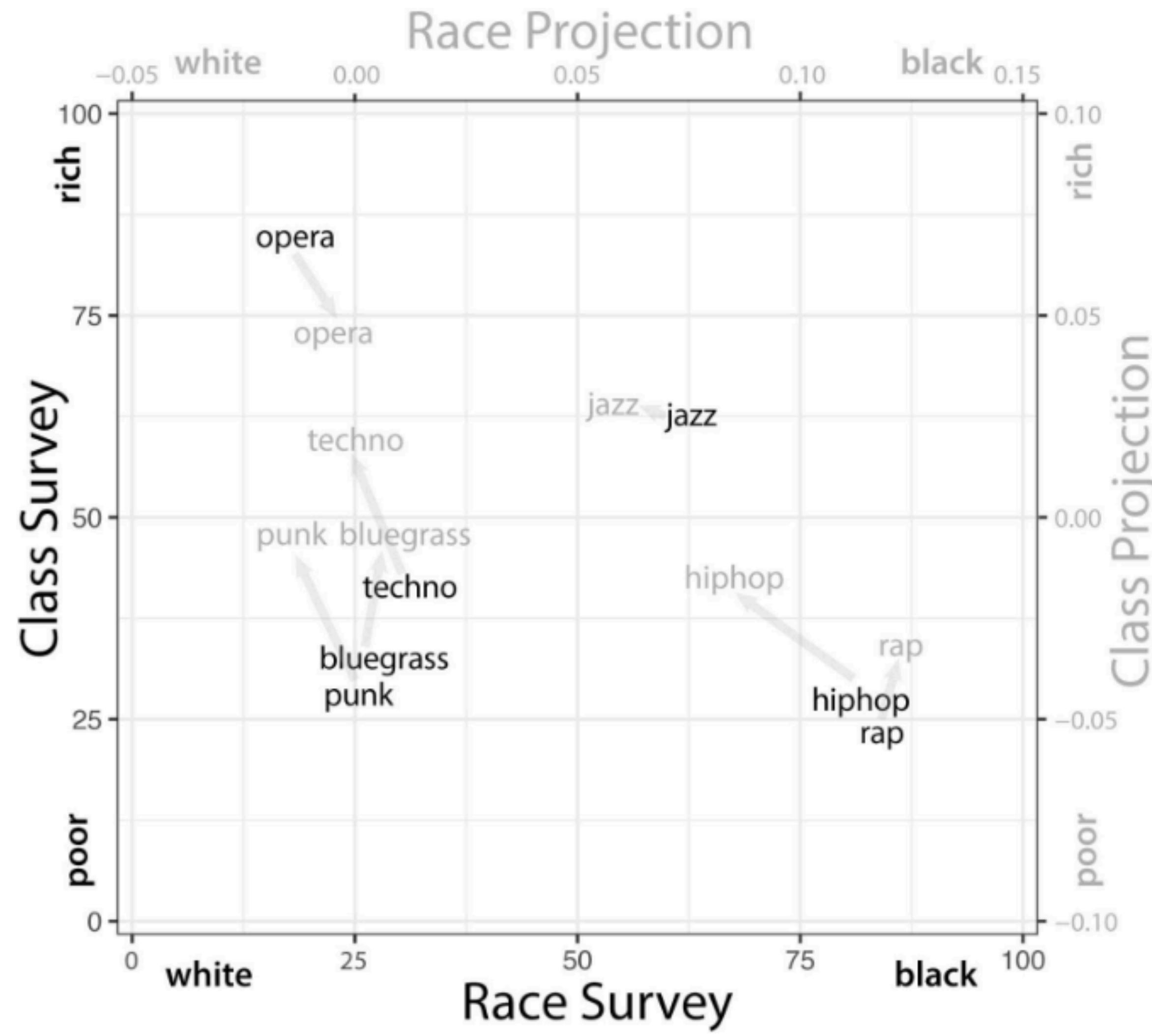
man:woman::king:queen

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{king}} - \overrightarrow{\text{queen}}$$

$$\overrightarrow{\text{king}} - \overrightarrow{\text{man}} + \overrightarrow{\text{woman}} \approx \overrightarrow{\text{queen}}$$

$$\overrightarrow{\text{king}} - \overrightarrow{\text{man}} + \overrightarrow{\text{woman}} = ?$$





# PREDICTION

# PREDICTION

ferromagnetic – NiFe + IrMn  $\approx$  antiferromagnetic

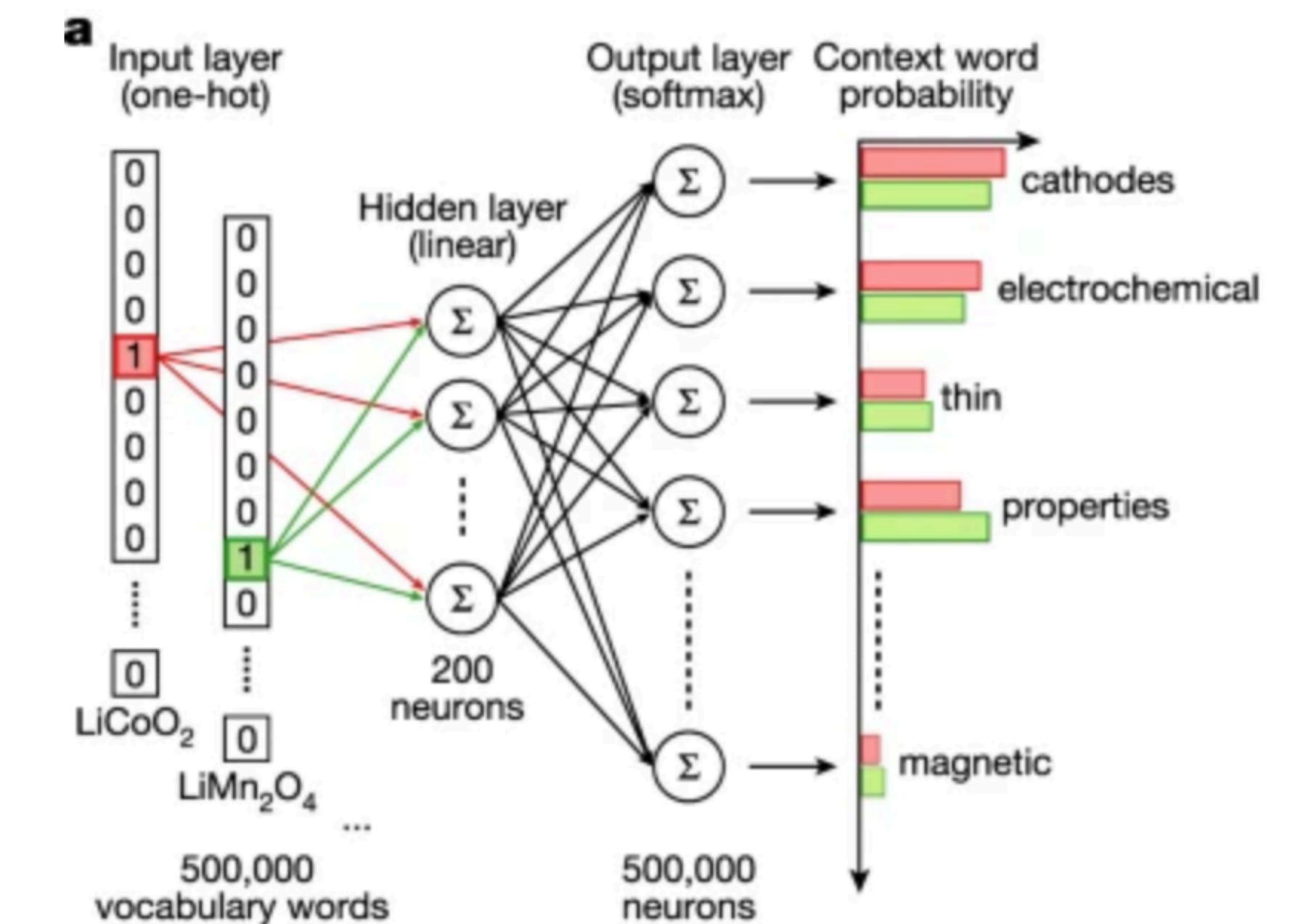
Oxides

Zr – ZrO<sub>2</sub>  $\approx$  Cr – Cr<sub>2</sub>O<sub>3</sub>  $\approx$  Ni – NiO

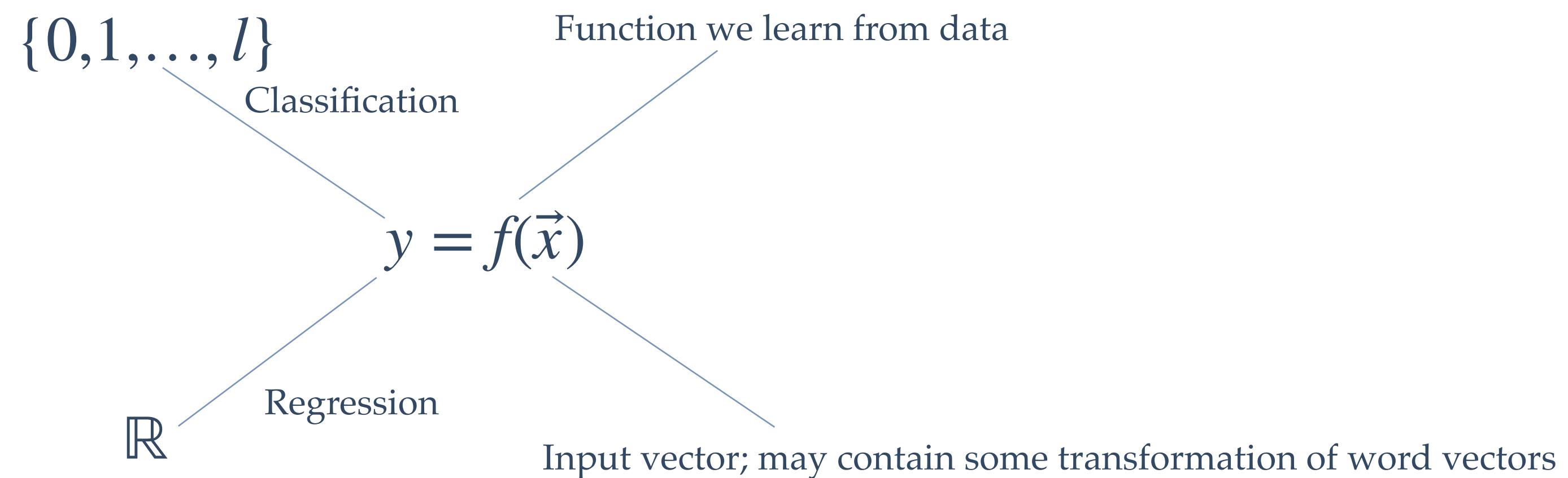
Structure

Zr – HCP  $\approx$  Cr – BCC  $\approx$  Ni – FCC

Embeddings can be used to construct knowledge bases  
that can lead to new discoveries



# PREDICTION



# PREDICTION

## Are Word Embedding-based Features Useful for Sarcasm Detection?

**Aditya Joshi**<sup>1,2,3</sup>   **Vaibhav Tripathi**<sup>1</sup>   **Kevin Patel**<sup>1</sup>

**Pushpak Bhattacharyya**<sup>1</sup>   **Mark Carman**<sup>2</sup>

<sup>1</sup>Indian Institute of Technology Bombay, India

<sup>2</sup>Monash University, Australia

<sup>3</sup>IITB-Monash Research Academy, India

{adityaj, kevin.patel, pb}@cse.iitb.ac.in, mark.carman@monash.edu

	<b>Word2Vec</b>	<b>LSA</b>	<b>GloVe</b>	<b>Dep. Wt.</b>
+S	0.835	0.86	0.918	<b>0.978</b>
+WS	<b>1.411</b>	0.255	0.192	1.372
+S+WS	<b>1.182</b>	0.24	0.845	0.795

**Table 4:** Average gain in F-Scores obtained by using intersection of the four word embeddings, for three word embedding feature-types, augmented to four prior works; Dep. Wt. indicates vectors learned from dependency-based weights

<b>Word Embedding</b>	<b>Average F-score Gain</b>
LSA	0.452
Glove	0.651
Dependency	1.048
Word2Vec	1.143

**Table 5:** Average gain in F-scores for the four types of word embeddings; These values are computed for a subset of these embeddings consisting of words common to all four

# **BIAS**



man:woman::king:?

man:woman::waiter:?

man:woman::doctor:?

man:woman::king:?

queen

okay

man:woman::waiter:?

waitress

okay

man:woman::doctor:?

nurse

huh

man:woman::king:?

queen

okay

man:woman::waiter:?

waitress

okay

man:woman::doctor:?

nurse

huh

**RESEARCH**

---

**REPORT**

**COGNITIVE SCIENCE**

# 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>

Machine learning is a means to derive artificial intelligence by discovering patterns in existing data. Here, we show that applying machine learning to ordinary human language results in human-like semantic biases. We replicated a spectrum of known biases, as measured by the Implicit Association Test, using a widely used, purely statistical machine-learning model trained on a standard corpus of text from the World Wide Web. Our results indicate that text corpora contain recoverable and accurate imprints of our historic biases, whether morally neutral as toward insects or flowers, problematic as toward race or gender, or even simply veridical, reflecting the status quo distribution of gender with respect to careers or first names. Our methods hold promise for identifying and addressing sources of bias in culture, including technology.

# BIAS

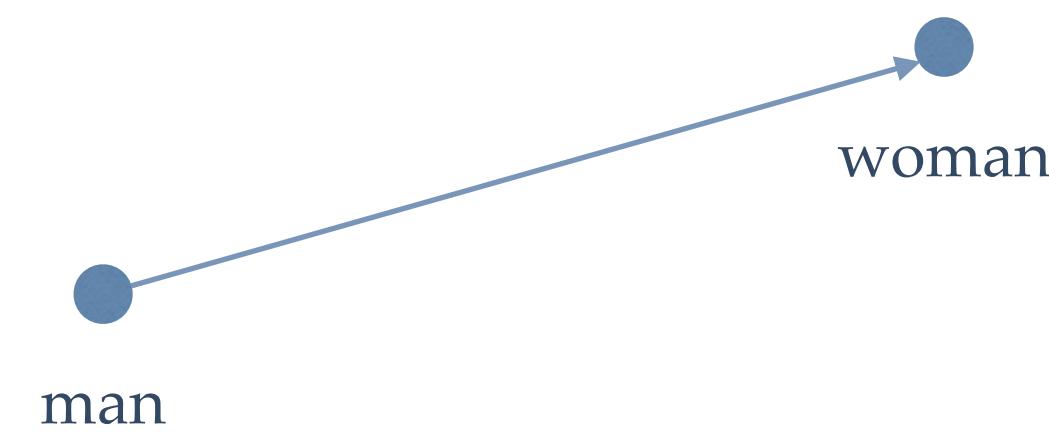
- Allocation harms: Unfair allocation of resources
- Representational harms: Unfair misrepresentation

# CALCULATING BIAS

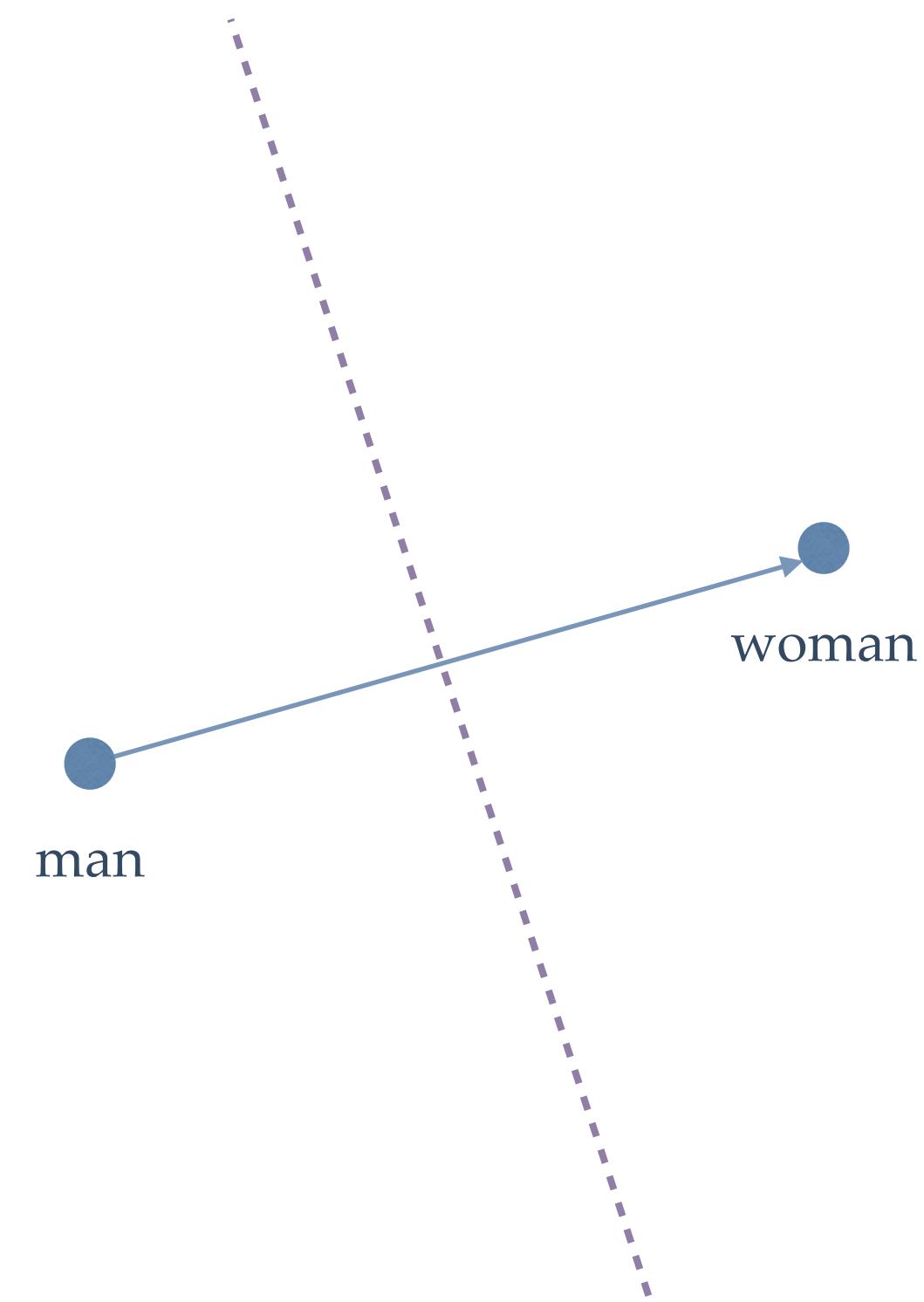
# CALCULATING BIAS



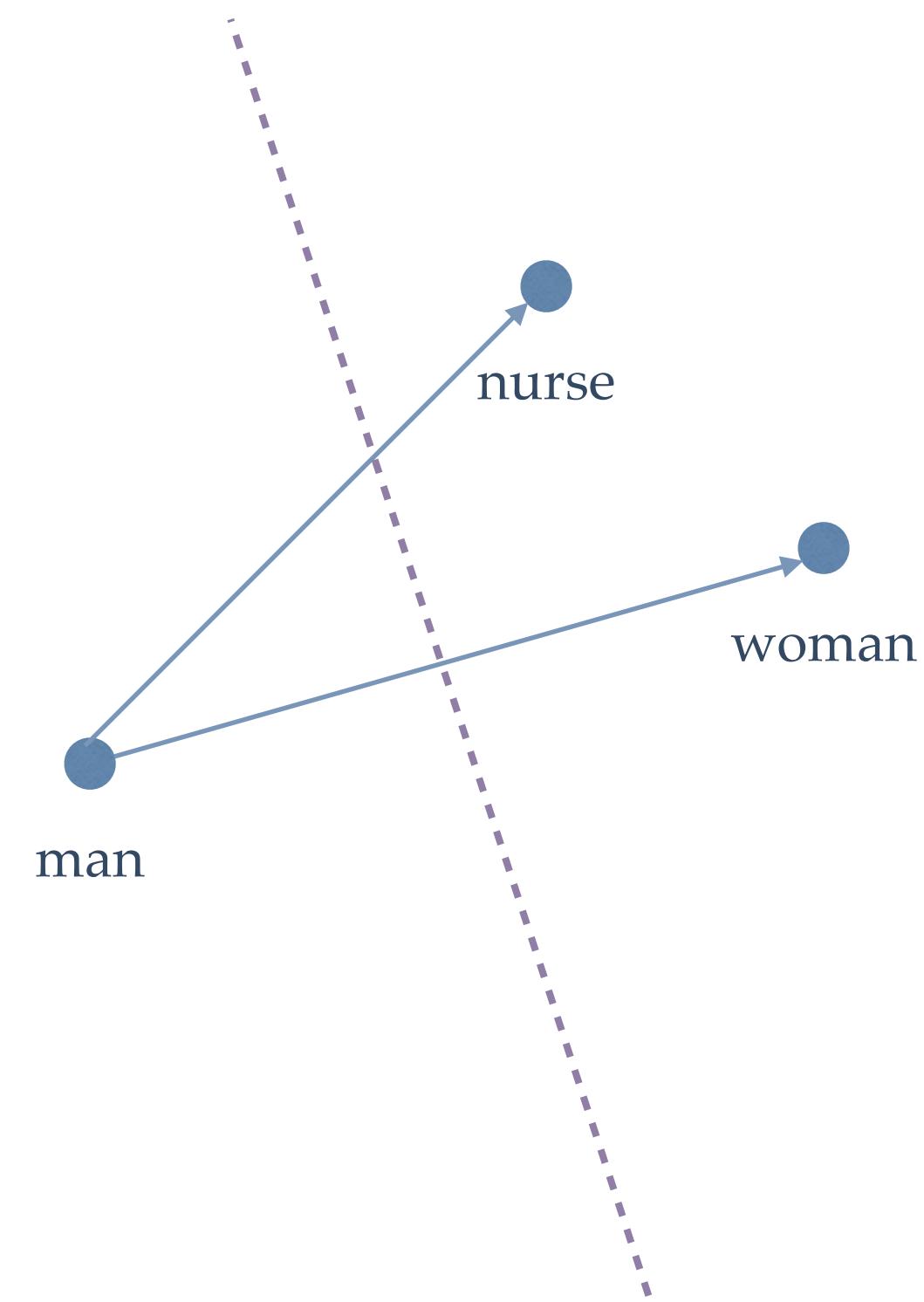
# CALCULATING BIAS



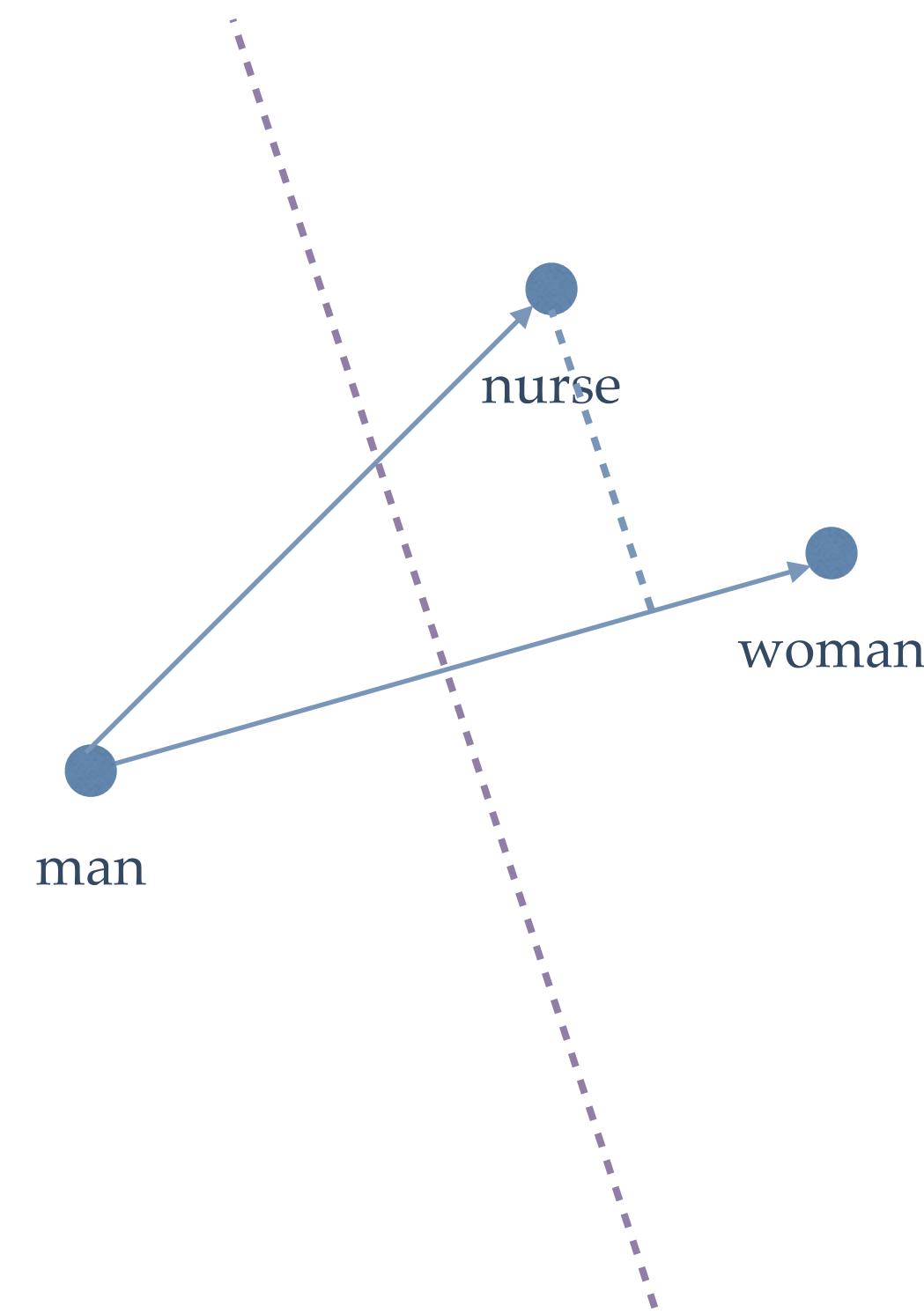
# CALCULATING BIAS



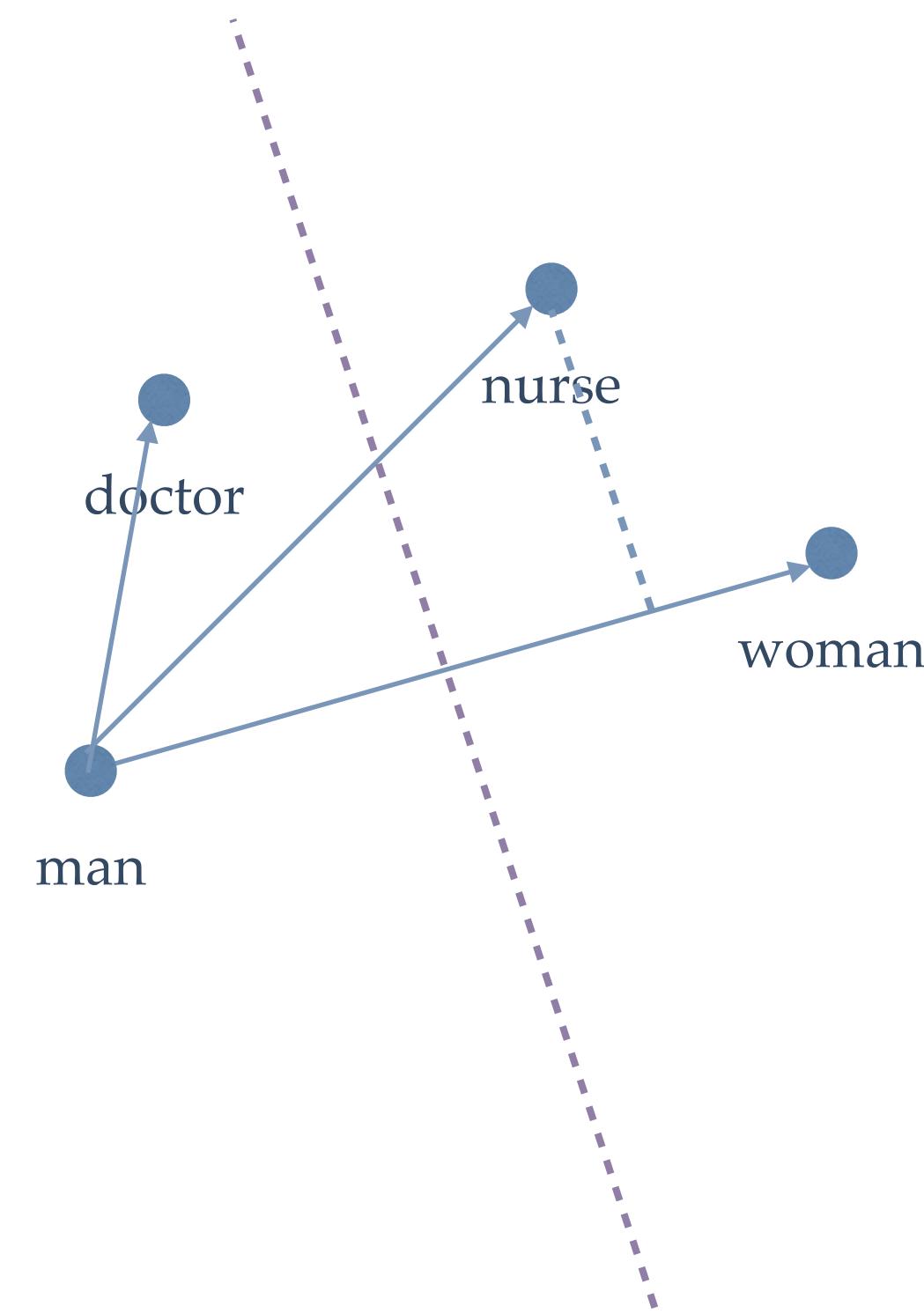
# CALCULATING BIAS



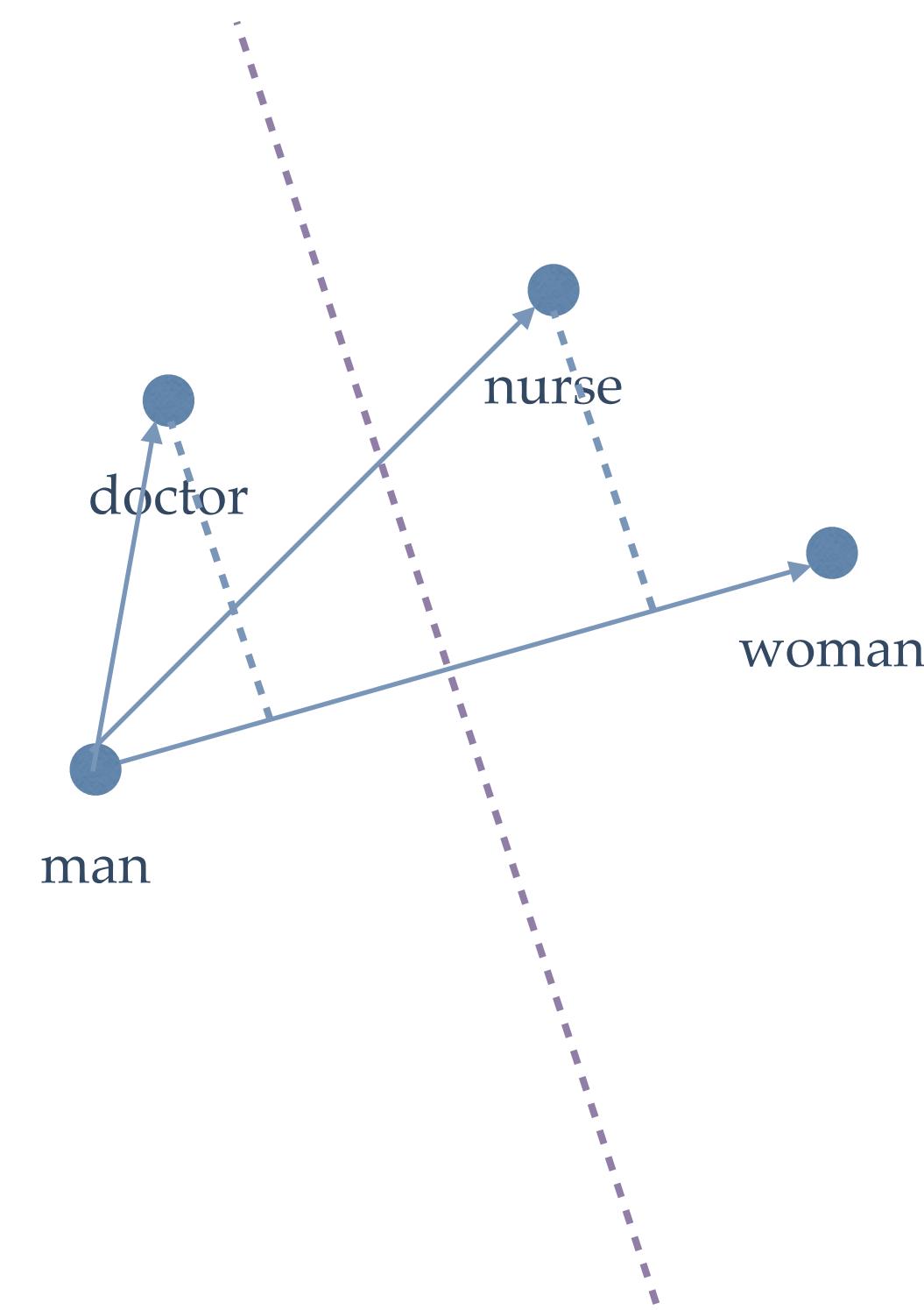
# CALCULATING BIAS



# CALCULATING BIAS



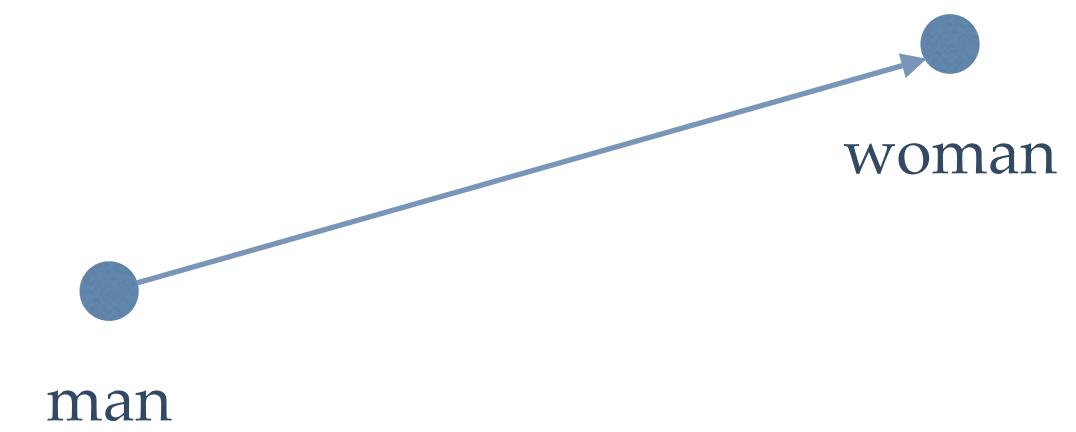
# CALCULATING BIAS



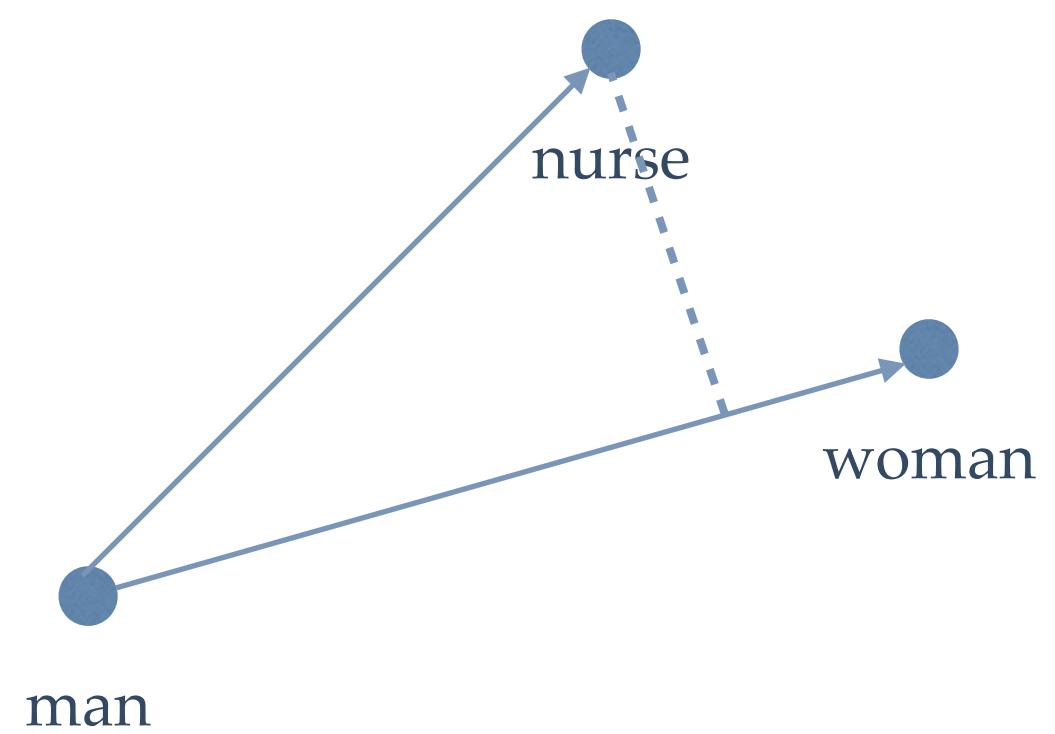
# CALCULATING BIAS

# CALCULATING BIAS

$$\vec{b} = \overrightarrow{\text{man}} - \overrightarrow{\text{woman}}$$



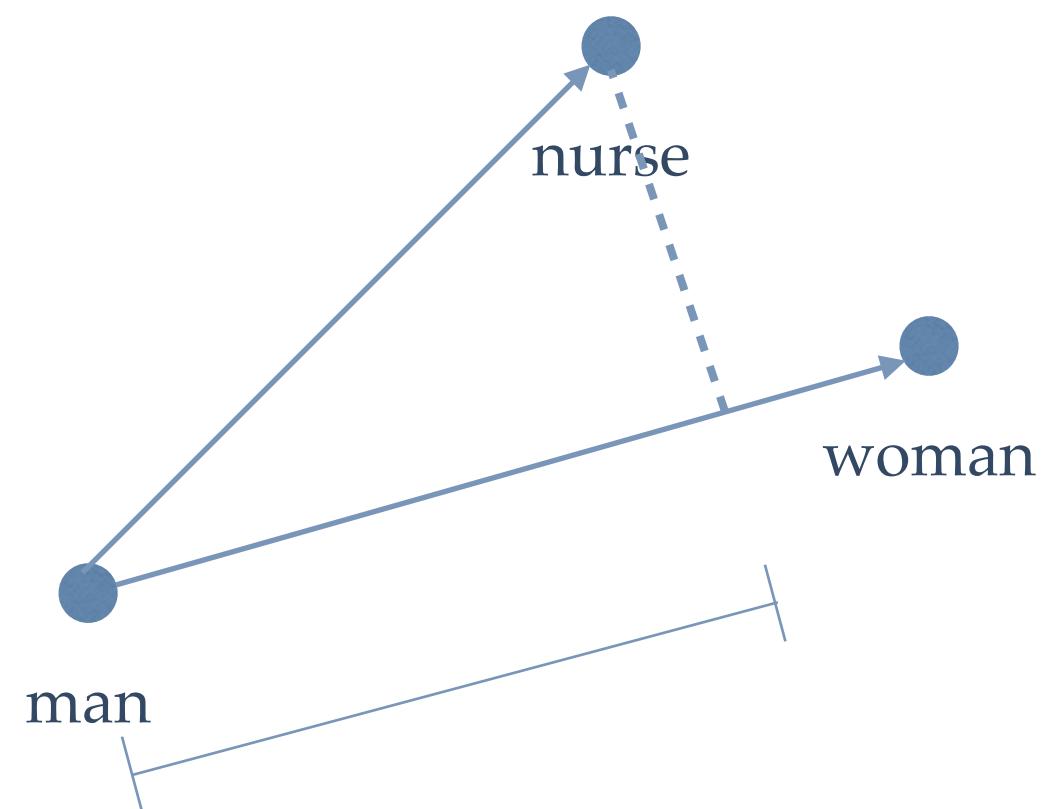
# CALCULATING BIAS



$$\vec{b} = \overrightarrow{\text{man}} - \overrightarrow{\text{woman}}$$

$$\vec{v} = \overrightarrow{\text{nurse}}$$

# CALCULATING BIAS

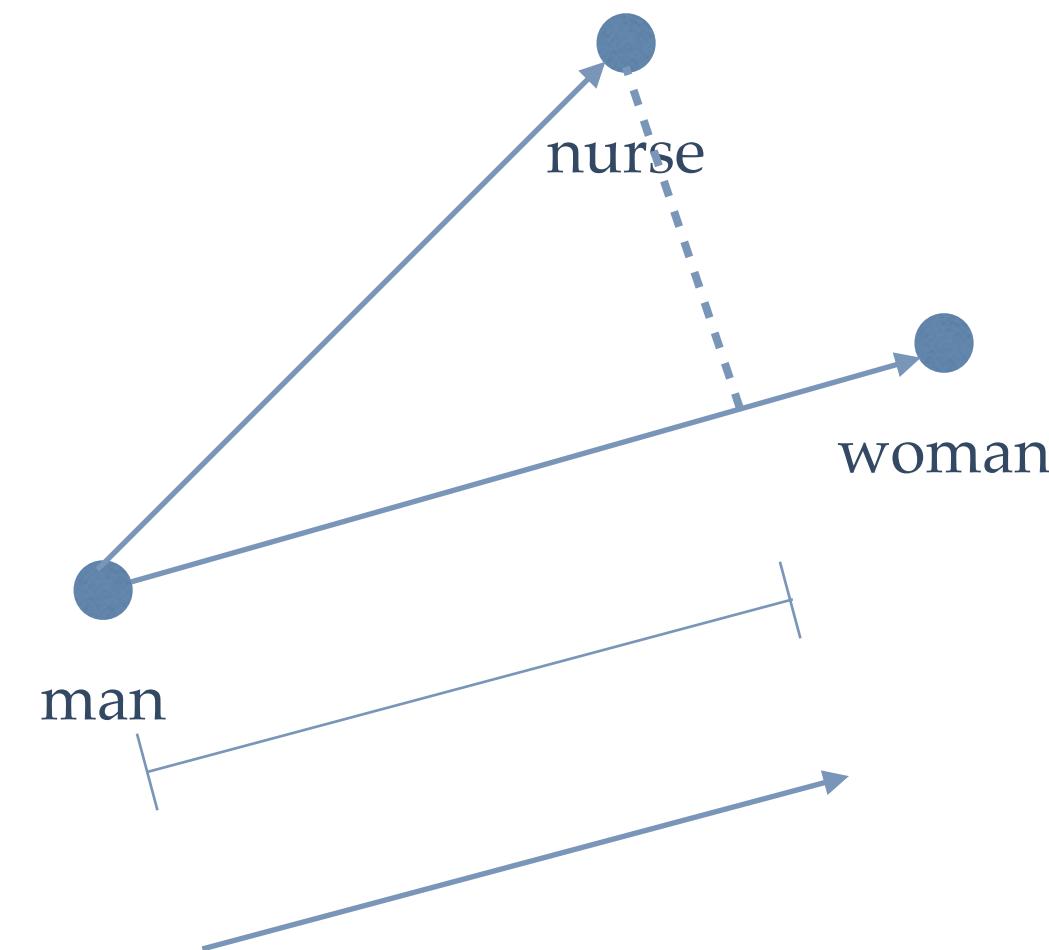


$$\vec{b} = \overrightarrow{\text{man}} - \overrightarrow{\text{woman}}$$

$$\vec{v} = \overrightarrow{\text{nurse}}$$

$$\text{Proj. length} = \vec{b} \cdot \vec{v}$$

# CALCULATING BIAS



$$\vec{b} = \overrightarrow{\text{man}} - \overrightarrow{\text{woman}}$$

$$\vec{v} = \overrightarrow{\text{nurse}}$$

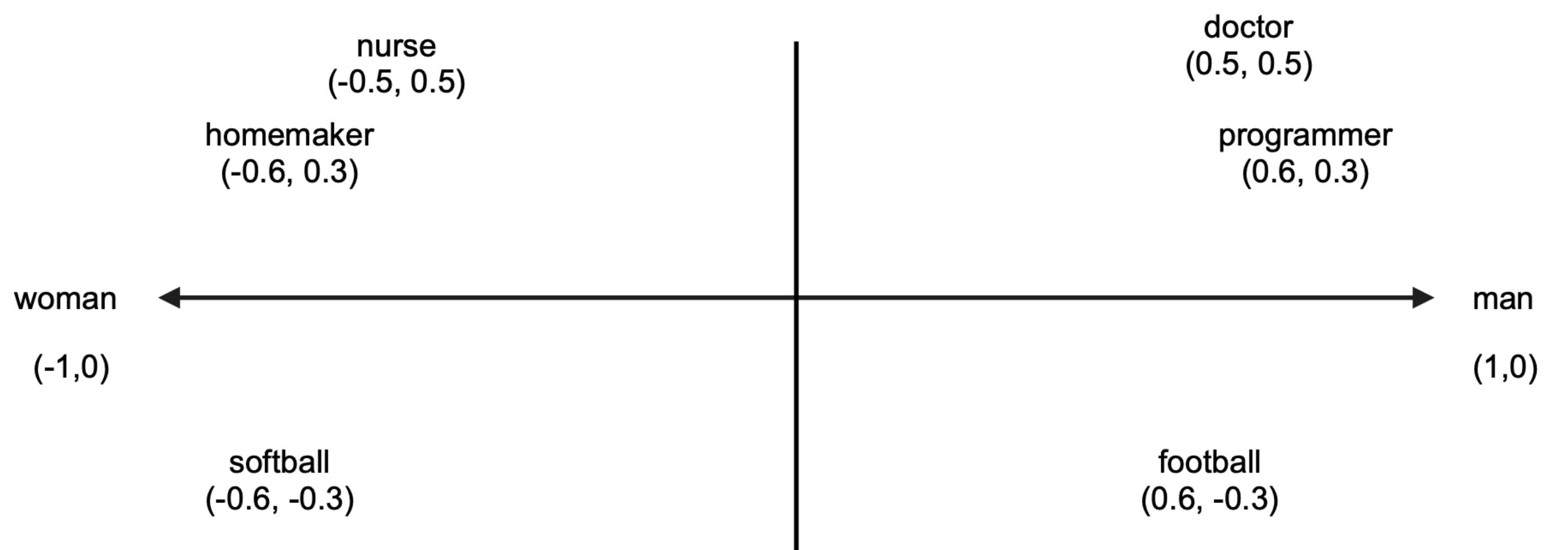
$$\text{Proj. length} = \vec{b} \cdot \vec{v}$$

$$\text{Proj. vector} = (\vec{b} \cdot \vec{v}) \vec{b}$$

# PROJECTION OF EMBEDDINGS ON A SUBSPACE

$$\vec{b} = \overrightarrow{\text{man}} - \overrightarrow{\text{woman}}$$

$$\vec{p} = (\vec{v} \cdot \vec{b})\vec{b}$$



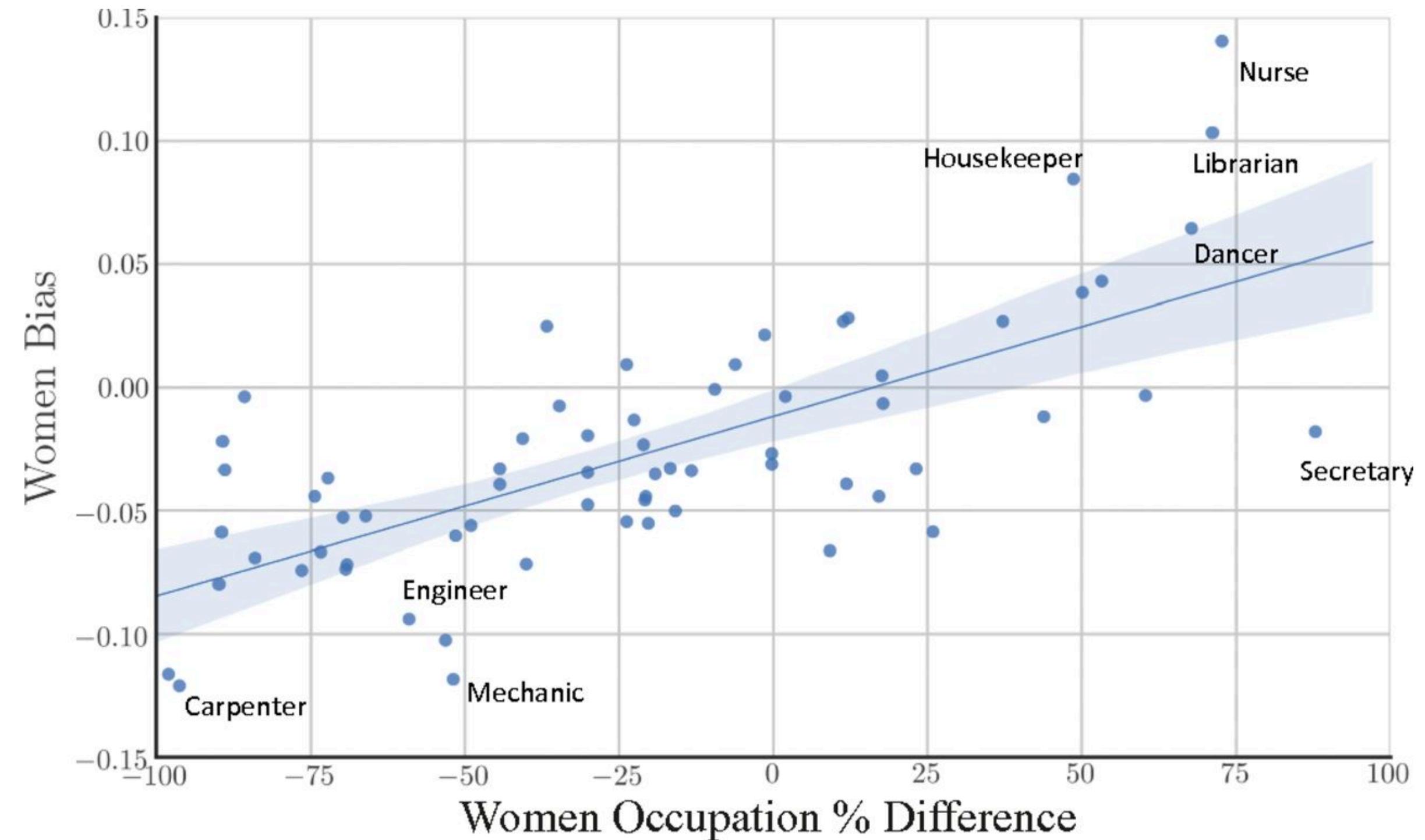


Image from Garg et. al. (2018)

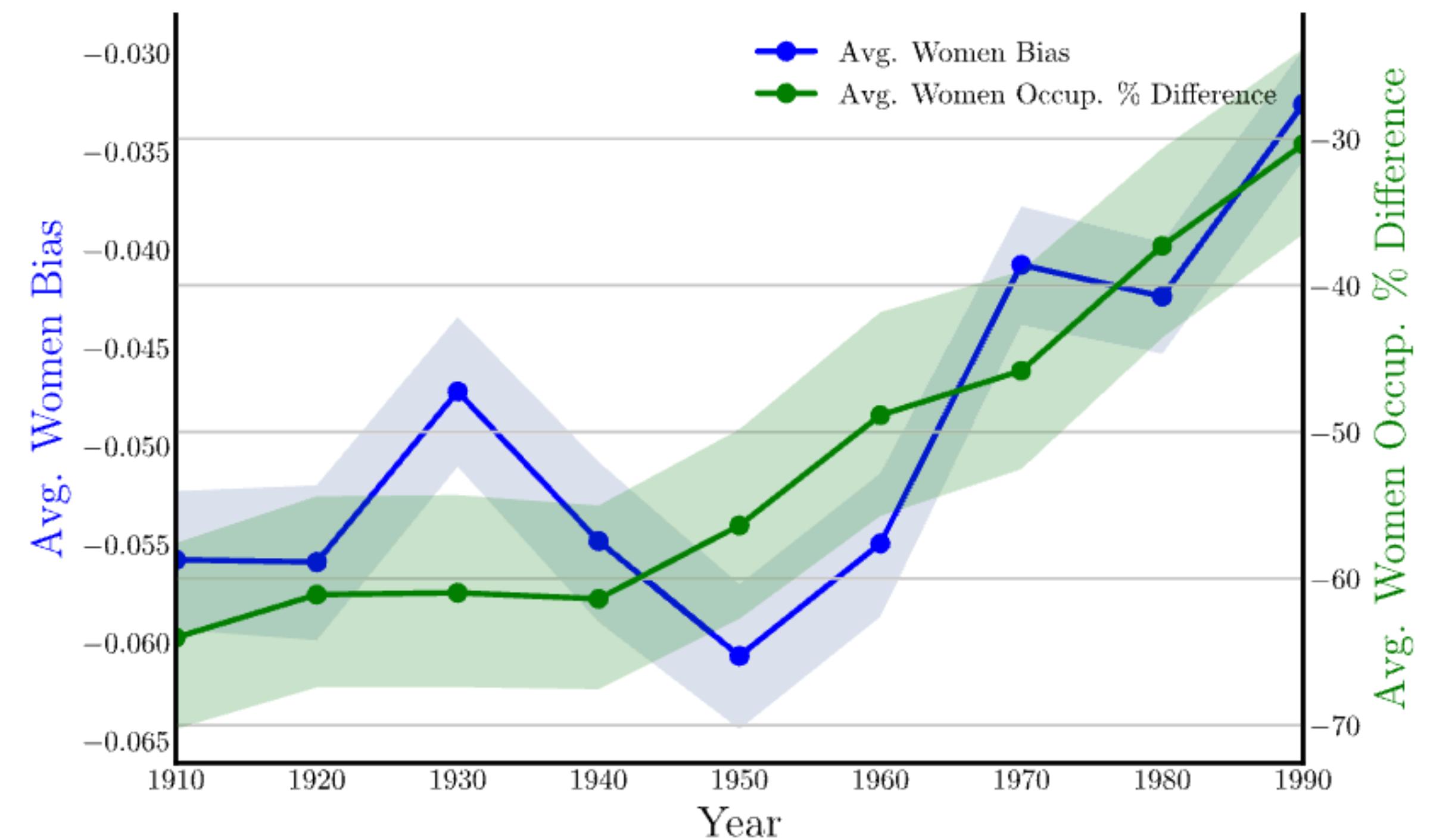


Image from Garg et. al. (2018)

# CHANGE AND VARIATION

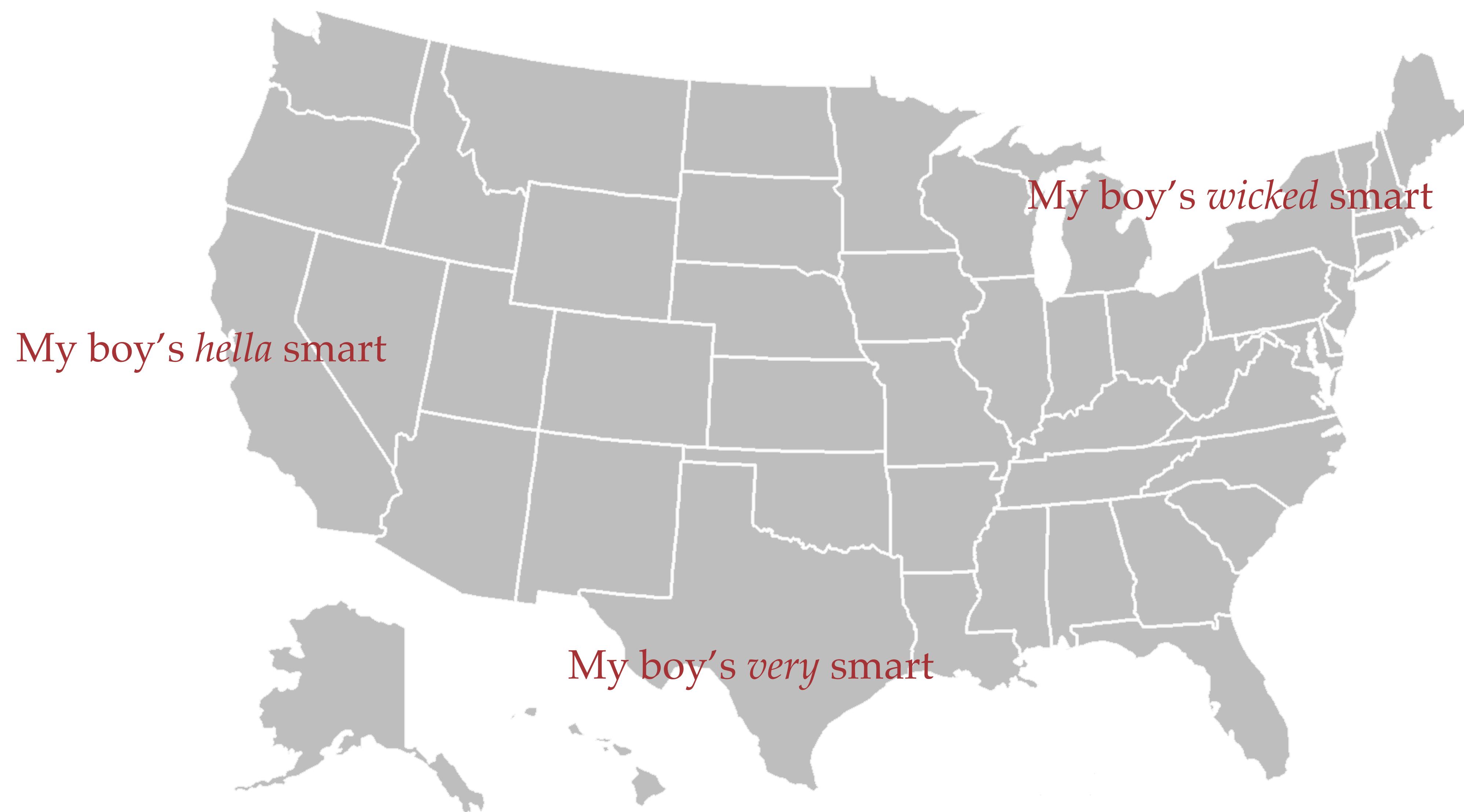
# CHANGE AND VARIATION

- Since language is situational, one can learn embeddings that depend on time, geography or other social contexts

*My boy's hella smart*

*My boy's wicked smart*

*My boy's very smart*



- David Bamman, Chris Dyer, and Noah A. Smith. 2014. Distributed Representations of Geographically Situated Language. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 828–834, Baltimore, Maryland. Association for Computational Linguistics.

David Bamman, Chris Dyer, and Noah A. Smith. 2014. [Distributed Representations of Geographically Situated Language](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 828–834, Baltimore, Maryland. Association for Computational Linguistics.



CALIFORNIA



TEXAS

x	y	State
...	...	...
wine	cold	AZ
wine	sweet	GA
hella	smart	CA
very	smart	TX
...	...	

David Bamman, Chris Dyer, and Noah A. Smith. 2014. [Distributed Representations of Geographically Situated Language](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 828–834, Baltimore, Maryland. Association for Computational Linguistics.



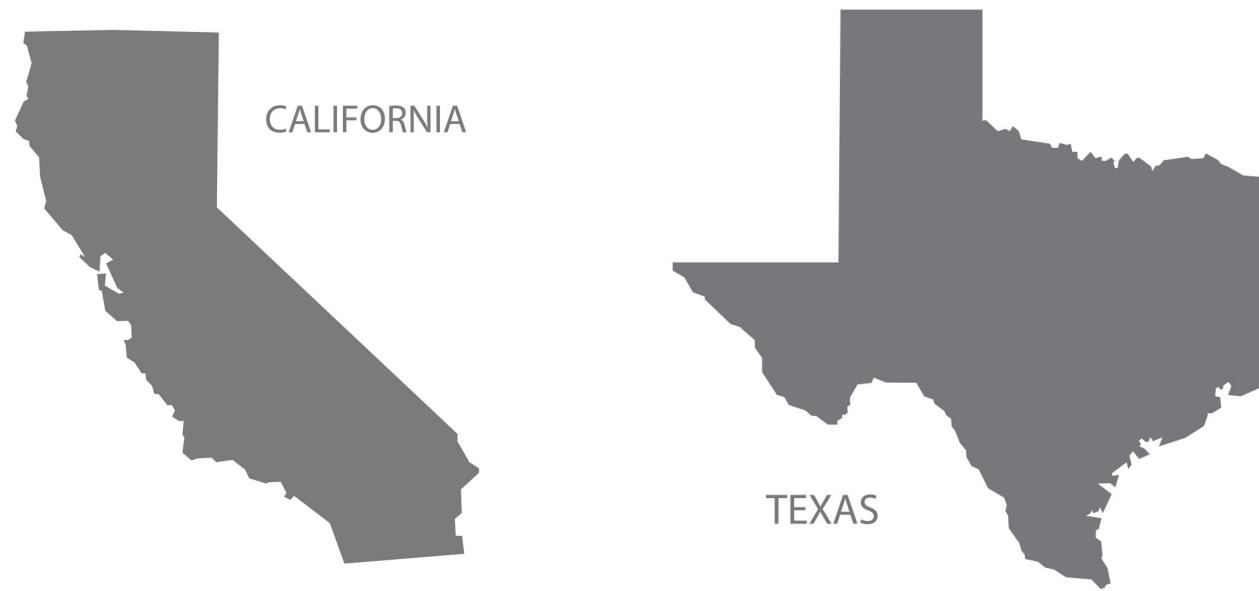
CALIFORNIA



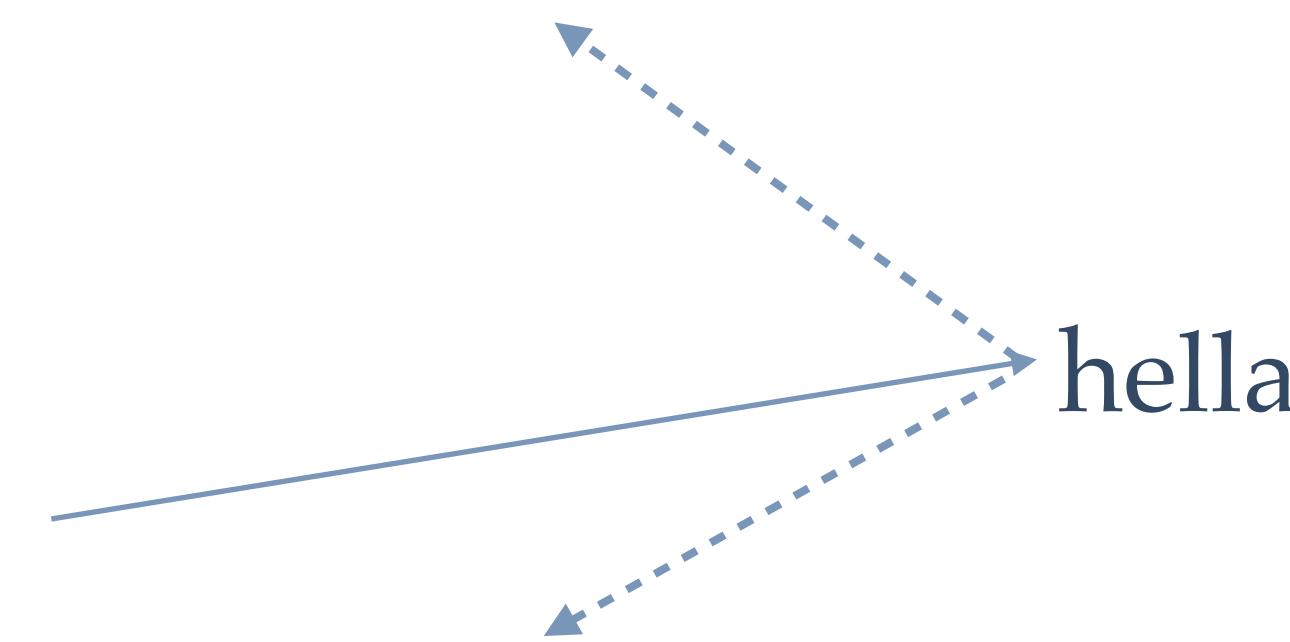
TEXAS

x	y	State
...	...	...
wine	cold	AZ
wine	sweet	GA
hella	smart	CA
very	smart	TX
...	...	

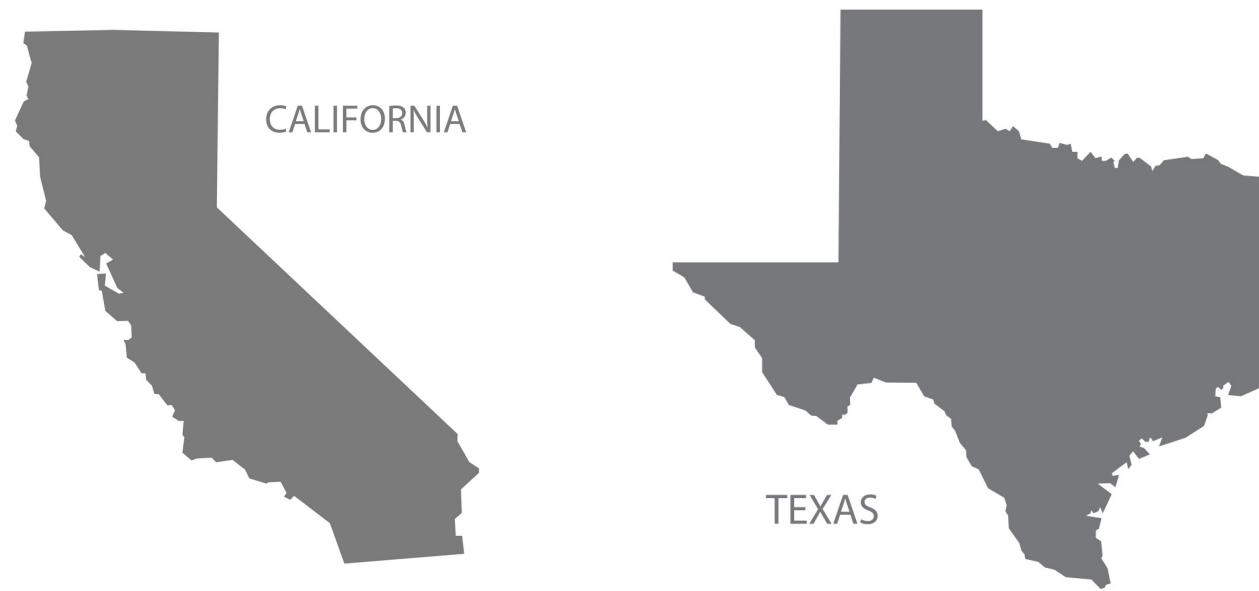
A blue arrow originates from the word "hella" and points towards the outline map of California, indicating a geographical association.



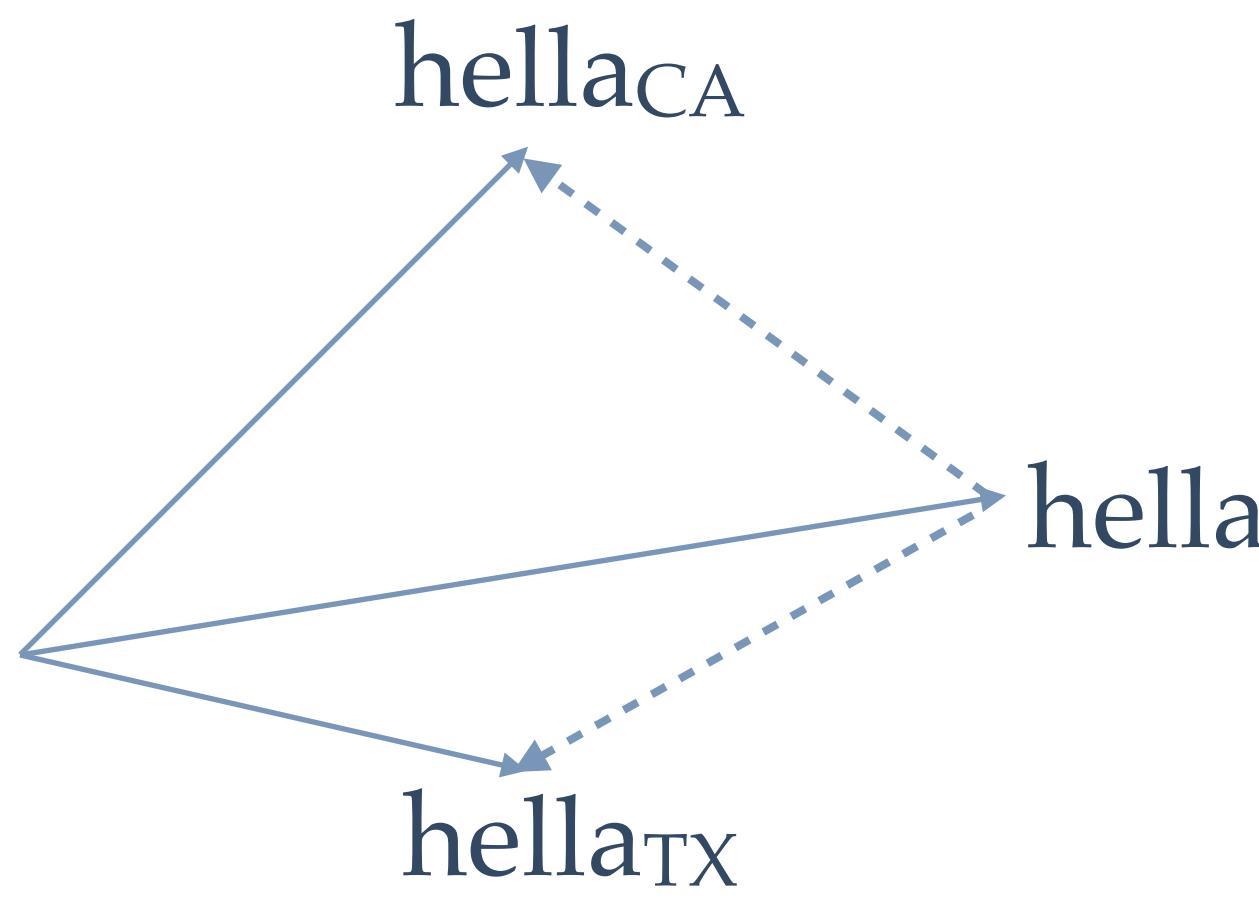
x	y	State
...	...	...
wine	cold	AZ
wine	sweet	GA
hella	smart	CA
very	smart	TX
...	...	



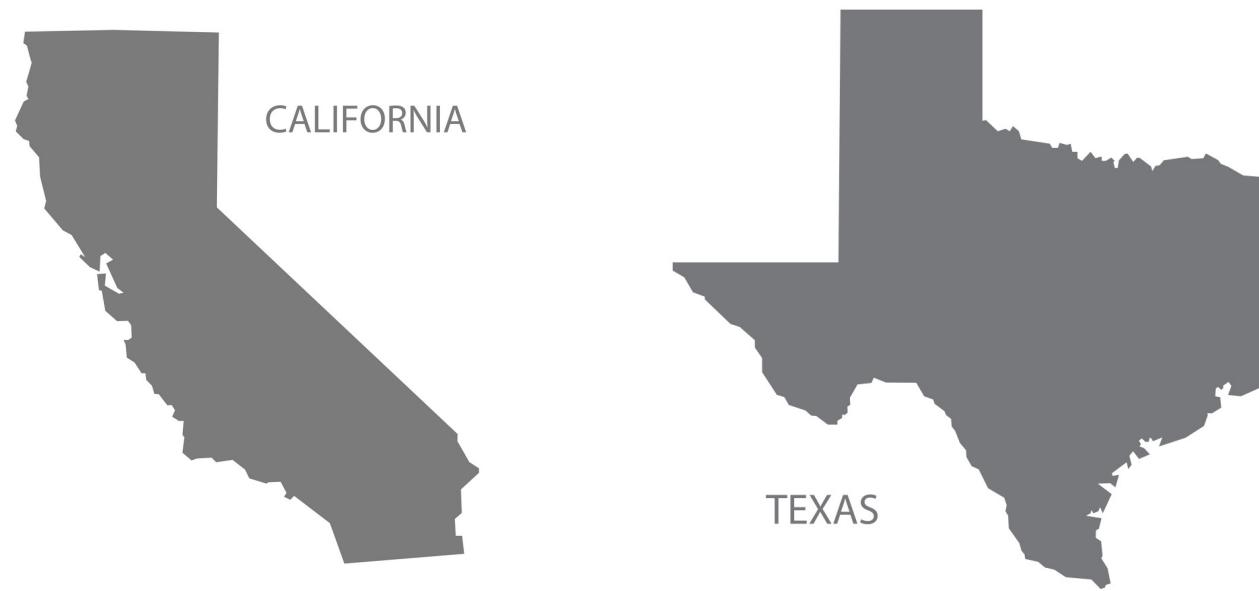
David Bamman, Chris Dyer, and Noah A. Smith. 2014. [Distributed Representations of Geographically Situated Language](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 828–834, Baltimore, Maryland. Association for Computational Linguistics.



x	y	State
...	...	...
wine	cold	AZ
wine	sweet	GA
hella	smart	CA
very	smart	TX
...	...	



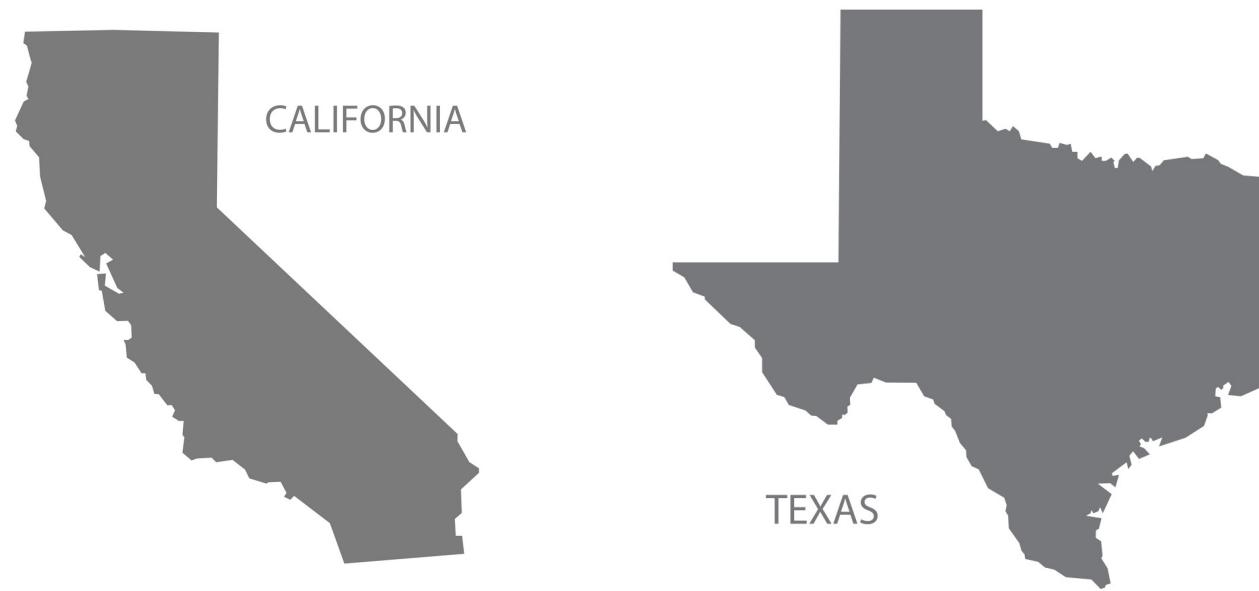
David Bamman, Chris Dyer, and Noah A. Smith. 2014. Distributed Representations of Geographically Situated Language. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 828–834, Baltimore, Maryland. Association for Computational Linguistics.



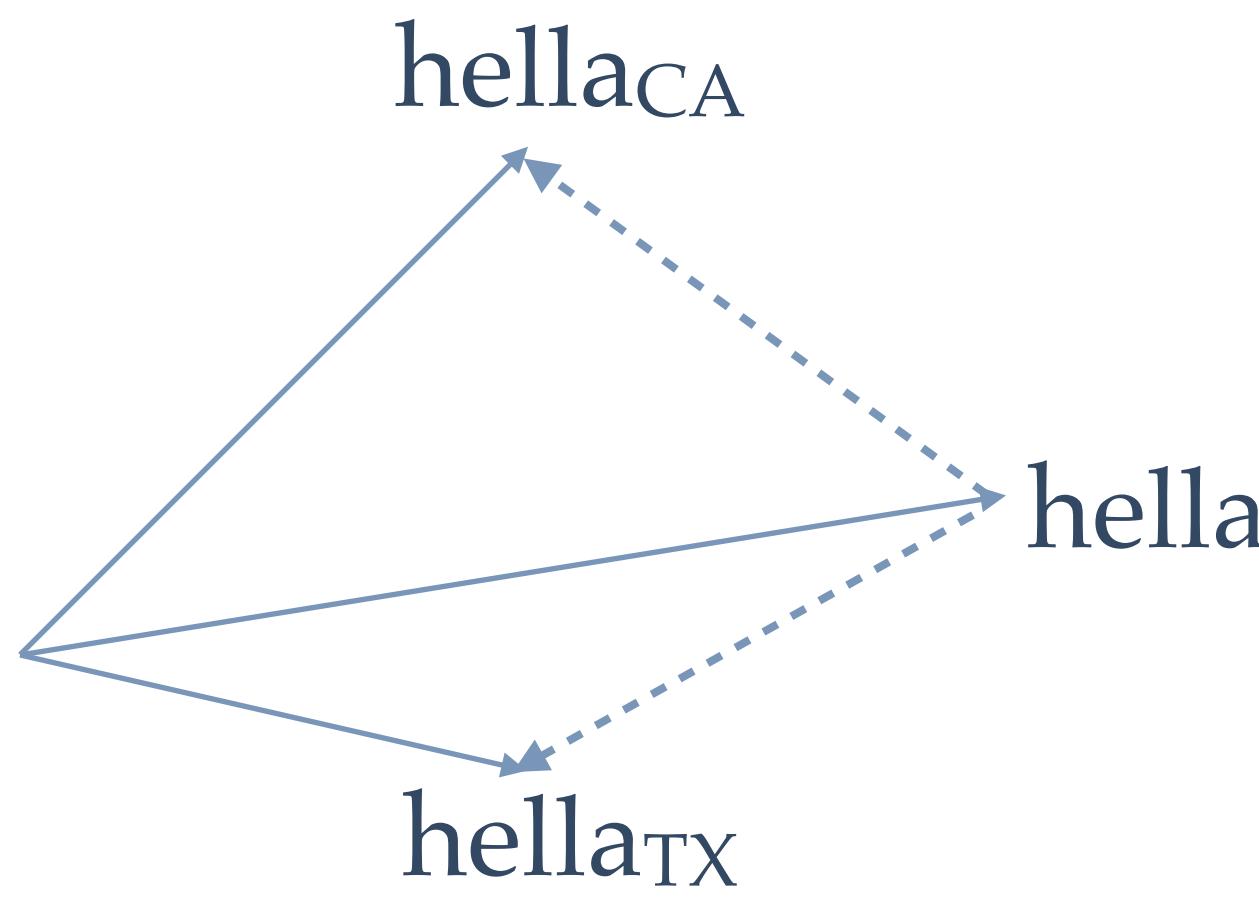
x	y	State
...	...	...
wine	cold	AZ
wine	sweet	GA
hella	smart	CA
very	smart	TX
...	...	



David Bamman, Chris Dyer, and Noah A. Smith. 2014. Distributed Representations of Geographically Situated Language. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 828–834, Baltimore, Maryland. Association for Computational Linguistics.



x	y	State
...	...	...
wine	cold	AZ
wine	sweet	GA
hella	smart	CA
very	smart	TX
...	...	



$$\overrightarrow{\text{hella}_{CA}} \stackrel{?}{\approx} \overrightarrow{\text{hella}_{TX}}$$

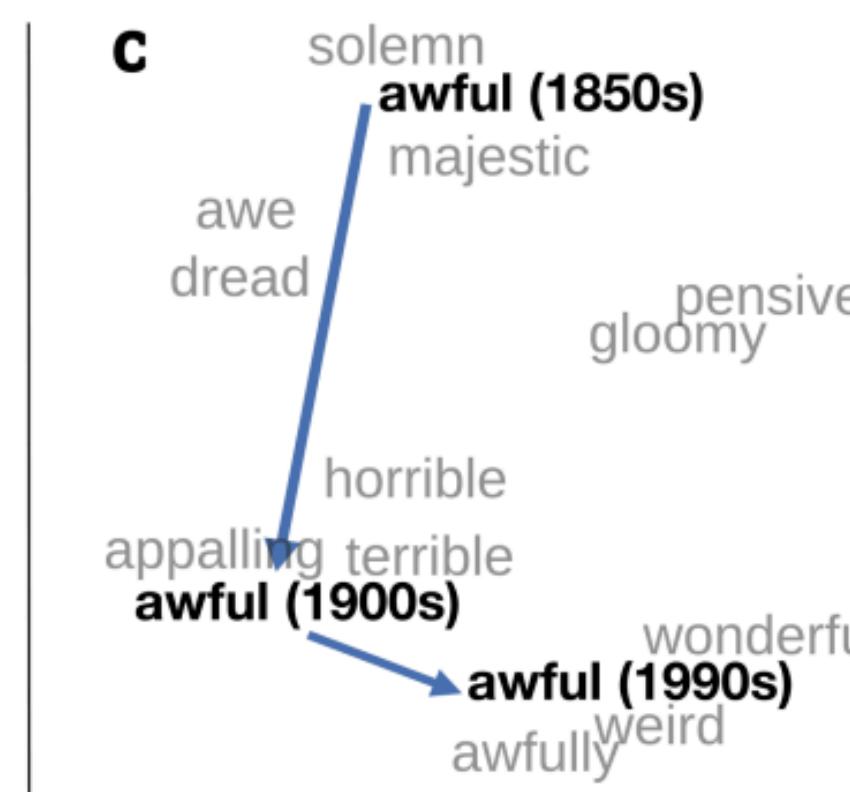
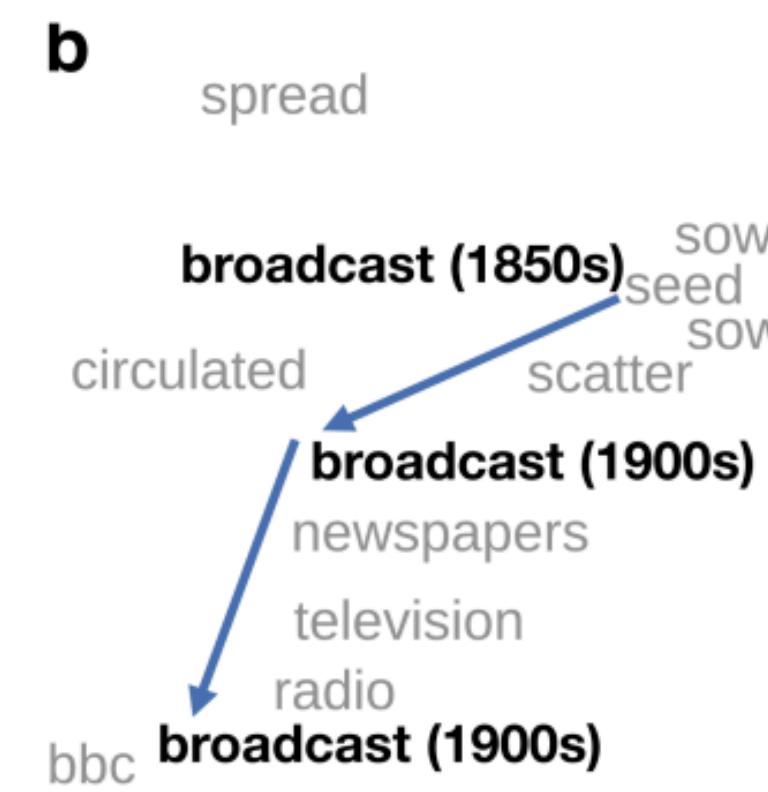
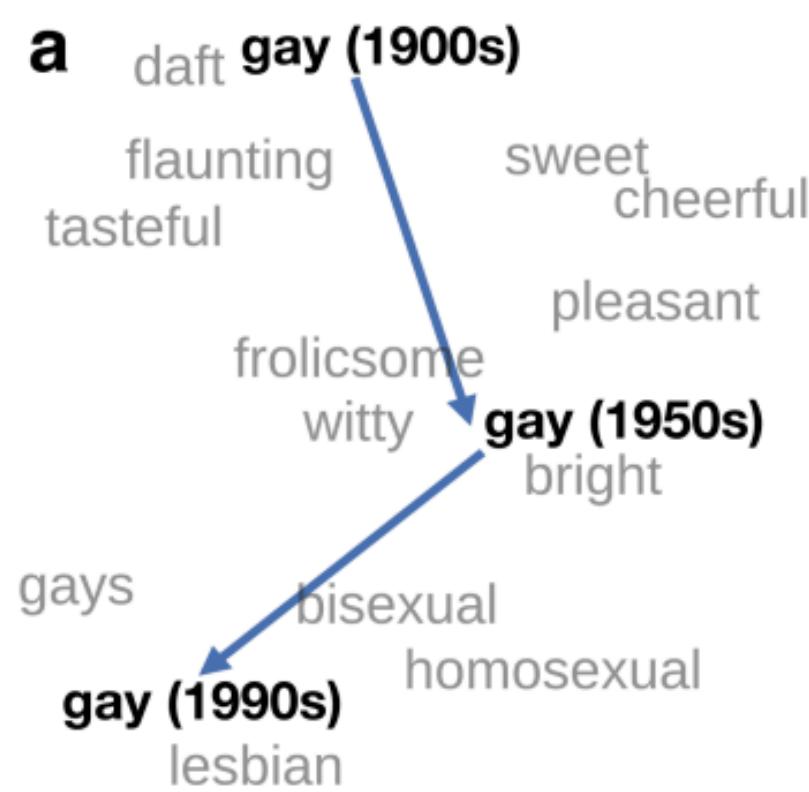
$$\overrightarrow{\text{hella}_{CA}} \stackrel{?}{\approx} \overrightarrow{\text{very}_{TX}}$$

- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. [Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1489–1501, Berlin, Germany. Association for Computational Linguistics.

x	y	Year
...	...	...
wine	cold	1990
wine	sweet	1990
gay	time	1990
gay	person	1900
...	...	

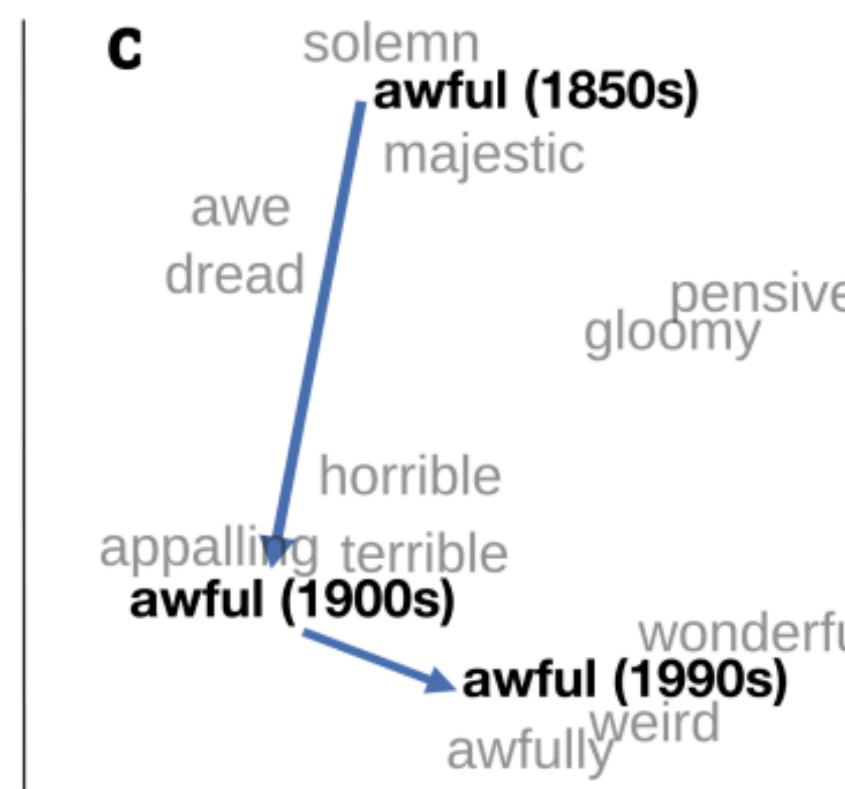
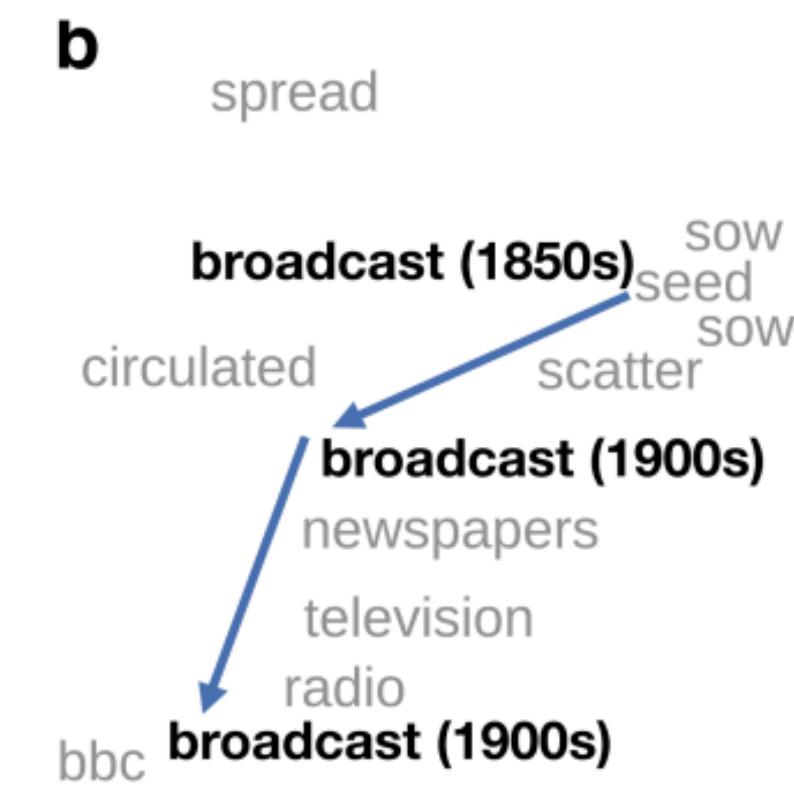
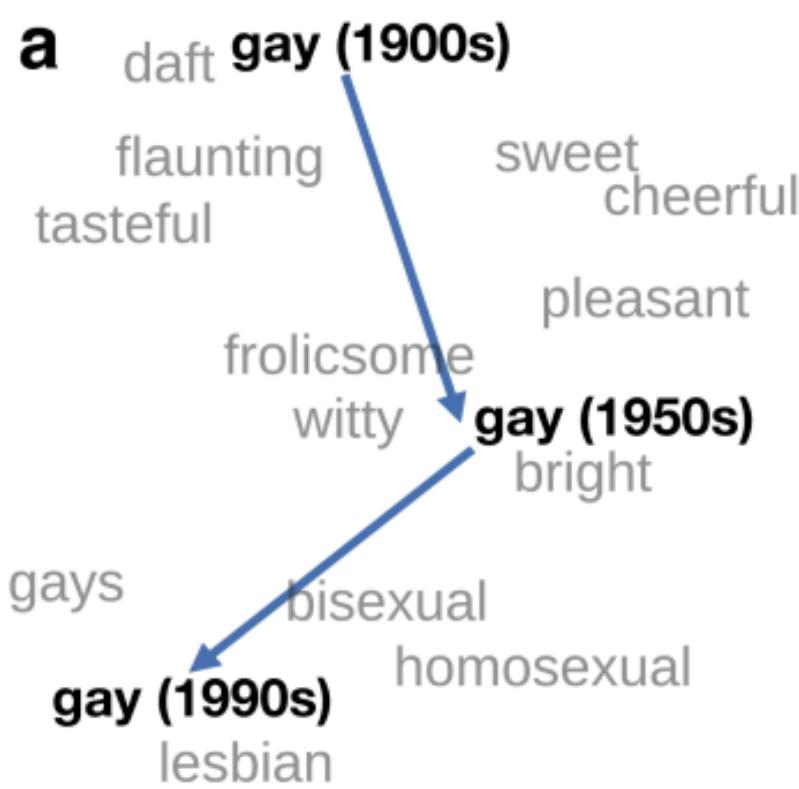
- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. [Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1489–1501, Berlin, Germany. Association for Computational Linguistics.

x	y	Year
...	...	...
wine	cold	1990
wine	sweet	1990
gay	time	1990
gay	person	1900
...	...	



- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1489–1501, Berlin, Germany. Association for Computational Linguistics.

x	y	Year
...	...	...
wine	cold	1990
wine	sweet	1990
gay	time	1990
gay	person	1900
...	...	



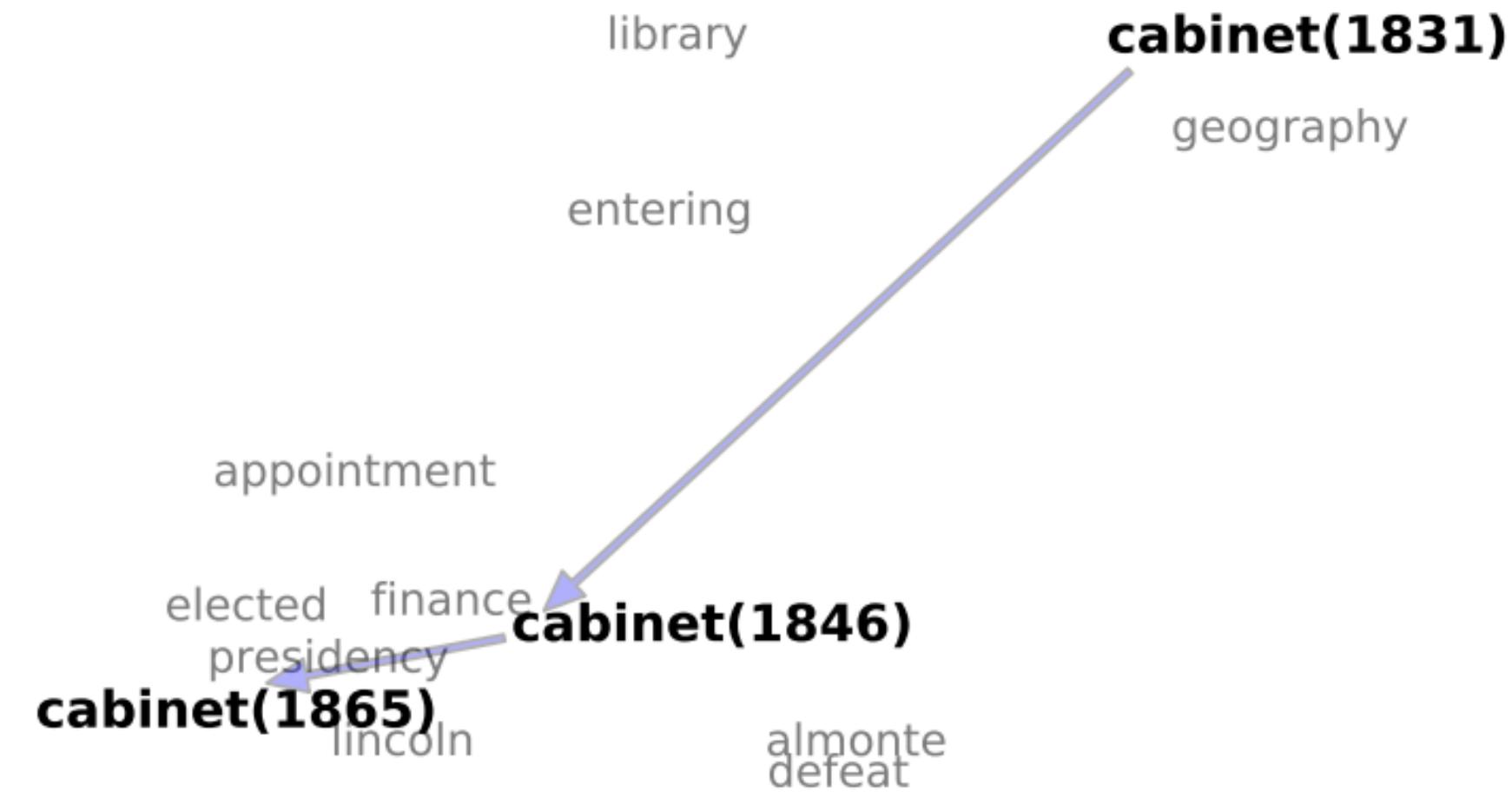
$\overrightarrow{\text{gay}}_{1900} \approx ? \overrightarrow{\text{gay}}_{1990}$

- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1489–1501, Berlin, Germany. Association for Computational Linguistics.

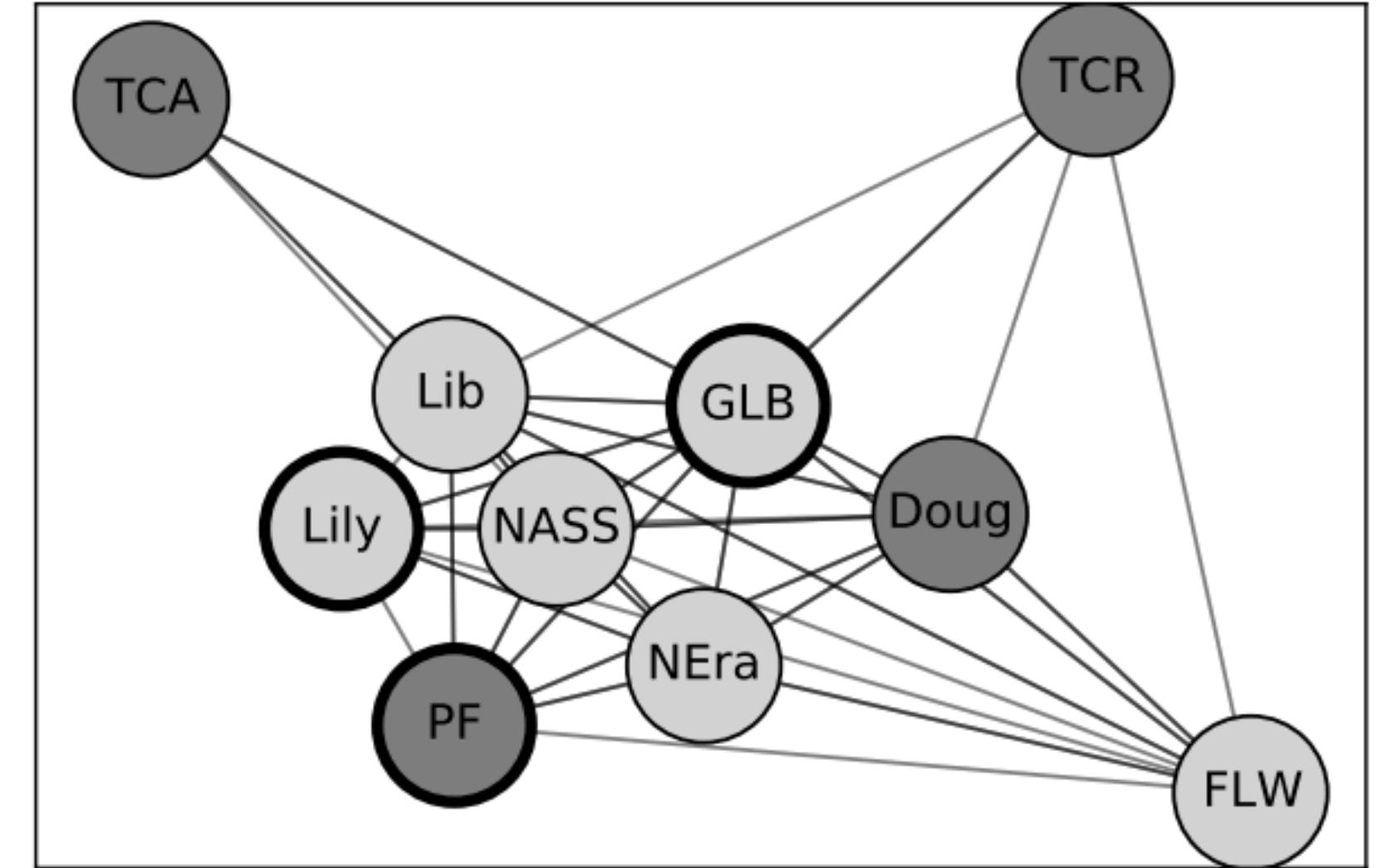
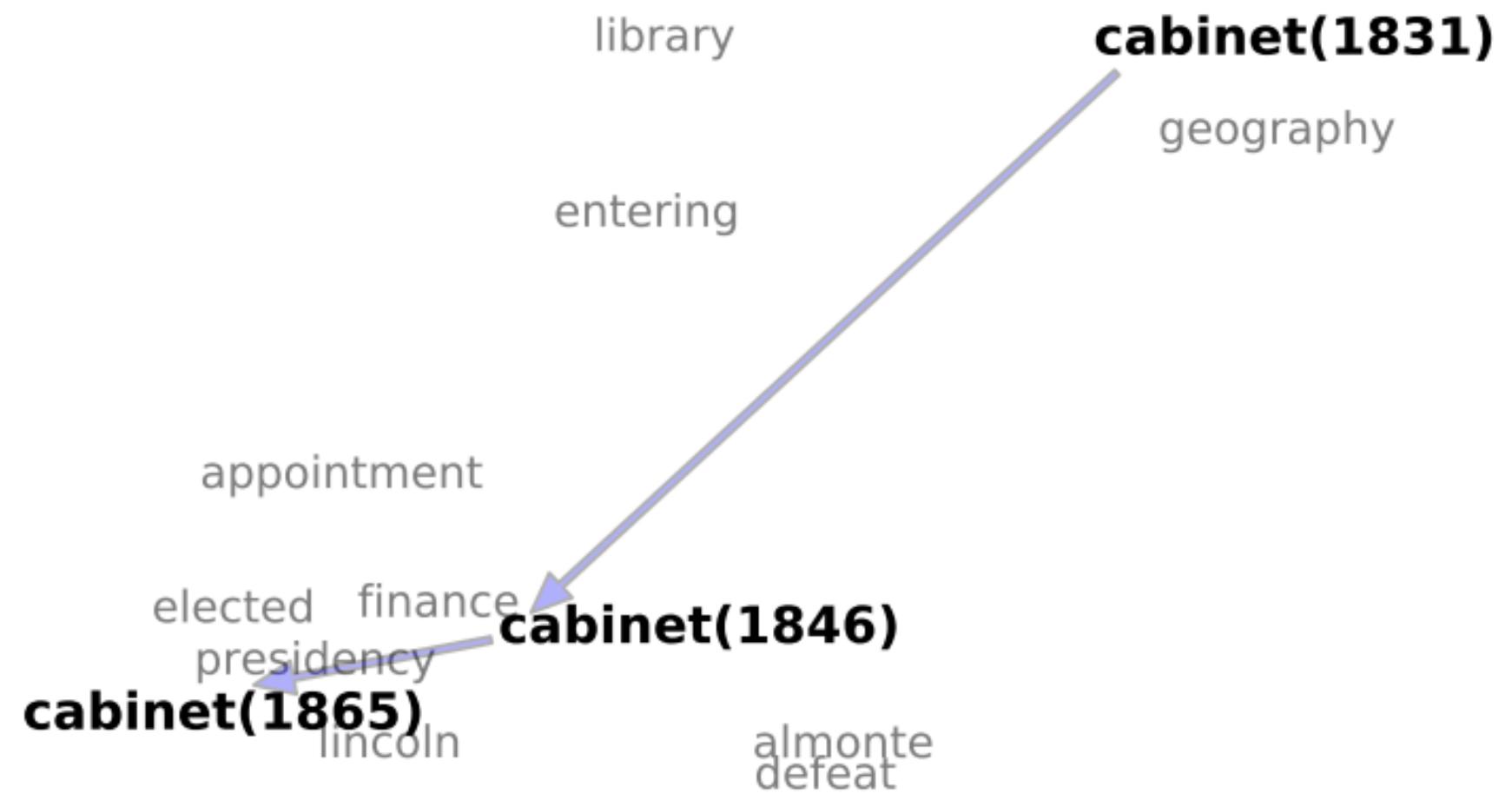




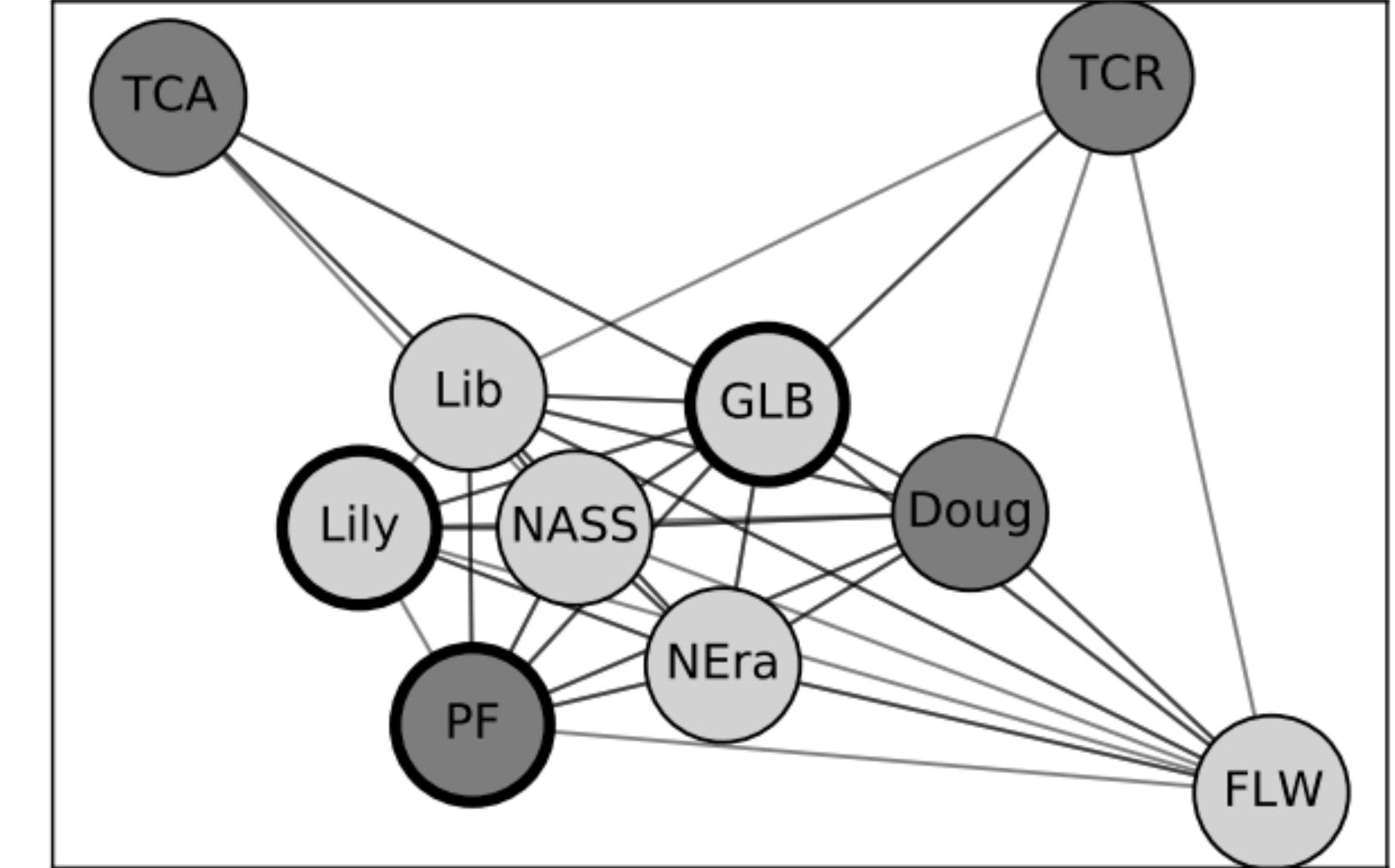
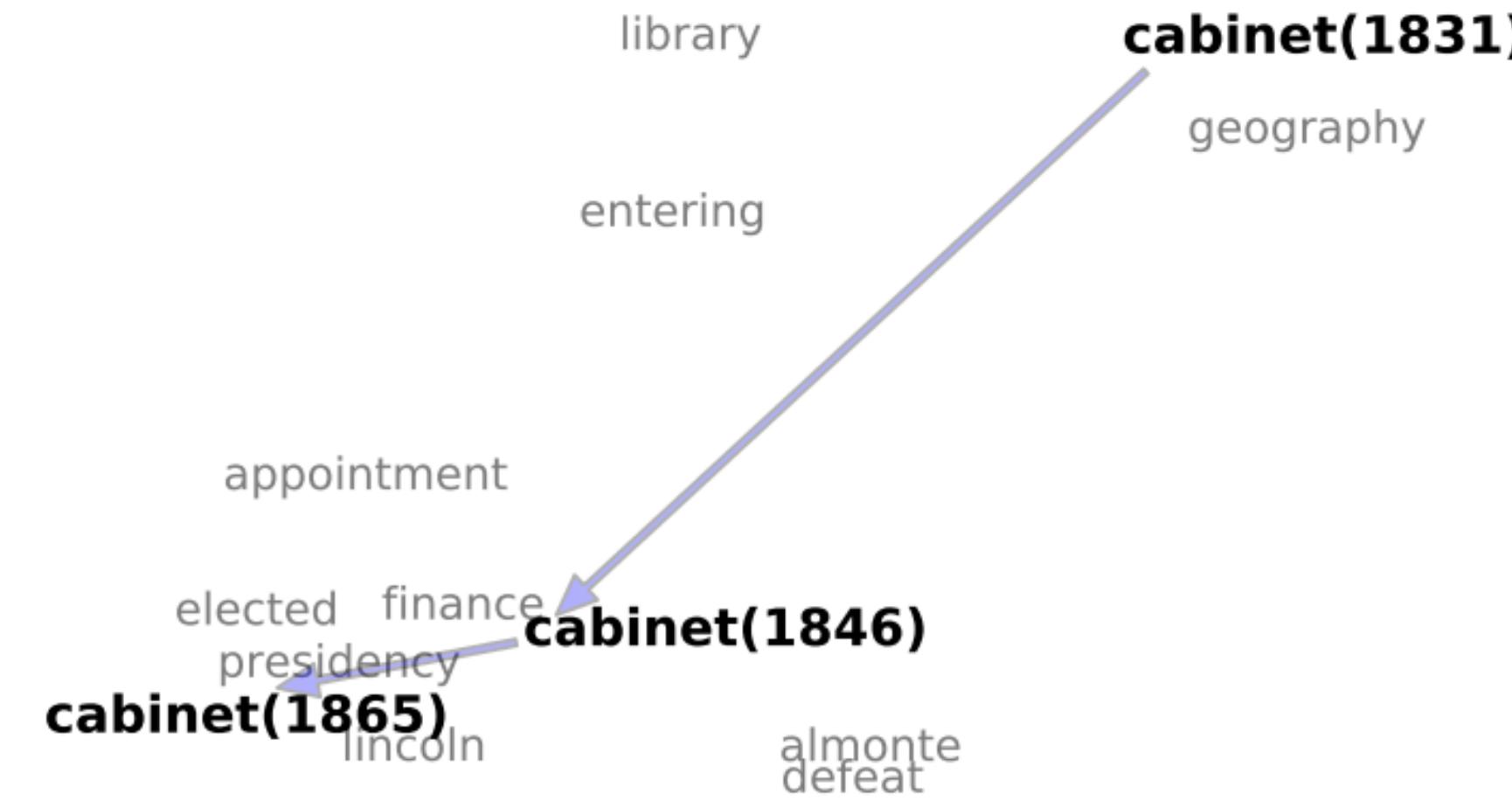
Soni, Sandeep, Lauren Klein, and Jacob Eisenstein. "Correcting whitespace errors in digitized historical texts." *Proceedings of the 3rd Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*. 2019.



Soni, Sandeep, Lauren Klein, and Jacob Eisenstein. "Correcting whitespace errors in digitized historical texts." *Proceedings of the 3rd Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*. 2019.



Soni, Sandeep, Lauren Klein, and Jacob Eisenstein. "Correcting whitespace errors in digitized historical texts." *Proceedings of the 3rd Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*. 2019.



$$\text{Lead}(\text{PF}, \text{Lib}, \text{cabinet}) = \frac{\xrightarrow{\hspace{1cm}} \text{cabinet}_{1831, \text{PF}} \cdot \text{cabinet}_{1865, \text{Lib}}}{\xrightarrow{\hspace{1cm}} \text{cabinet}_{1831, \text{Lib}} \cdot \text{cabinet}_{1865, \text{Lib}}}$$

Soni, Sandeep, Lauren Klein, and Jacob Eisenstein. "Correcting whitespace errors in digitized historical texts." *Proceedings of the 3rd Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*. 2019.

Soni, Sandeep, Lauren F. Klein, and Jacob Eisenstein. "Abolitionist Networks: Modeling Language Change in Nineteenth-Century Activist Newspapers." *Journal of Cultural Analytics* 6.1 (2021).

# IN CLASS

- Word2Vec demo