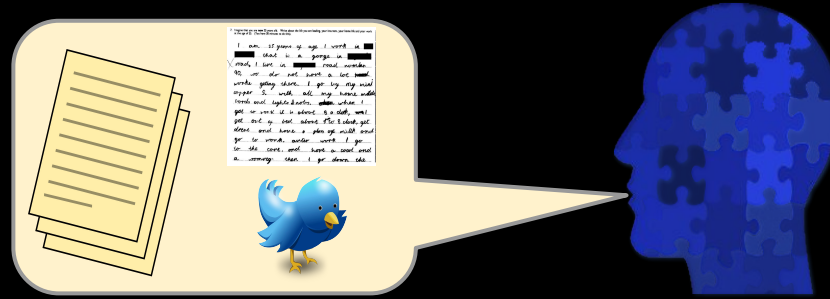


Vector Semantics and Embeddings

CSE392 - Spring 2019
Special Topic in CS

Tasks



- Vectors which represent words or sequences how?

- Dimensionality Reduction
- Recurrent Neural Network and Sequence Models

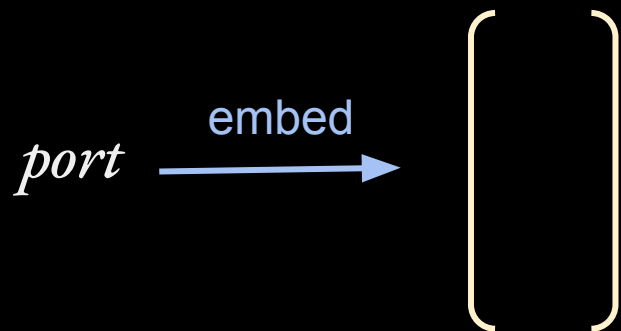
Objective

To embed: convert a token (or sequence) to a vector that represents meaning.

Objective

To embed: convert a token (or sequence) to a vector that is useful to perform NLP tasks.

Objective

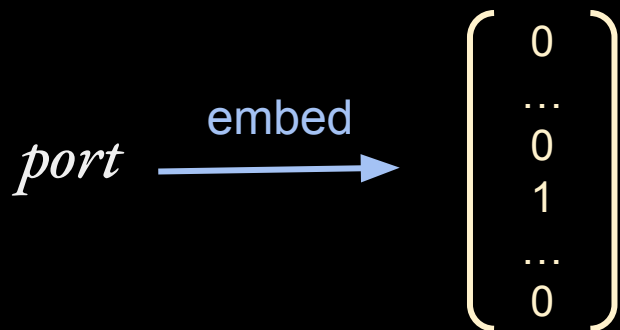


Objective

$$\textit{port} \xrightarrow{\text{embed}} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}$$

Objective

one-hot is sparse vector



Prefer dense vectors

- Less parameters (weights) for machine learning model.
- May generalize better implicitly.
- May capture synonyms

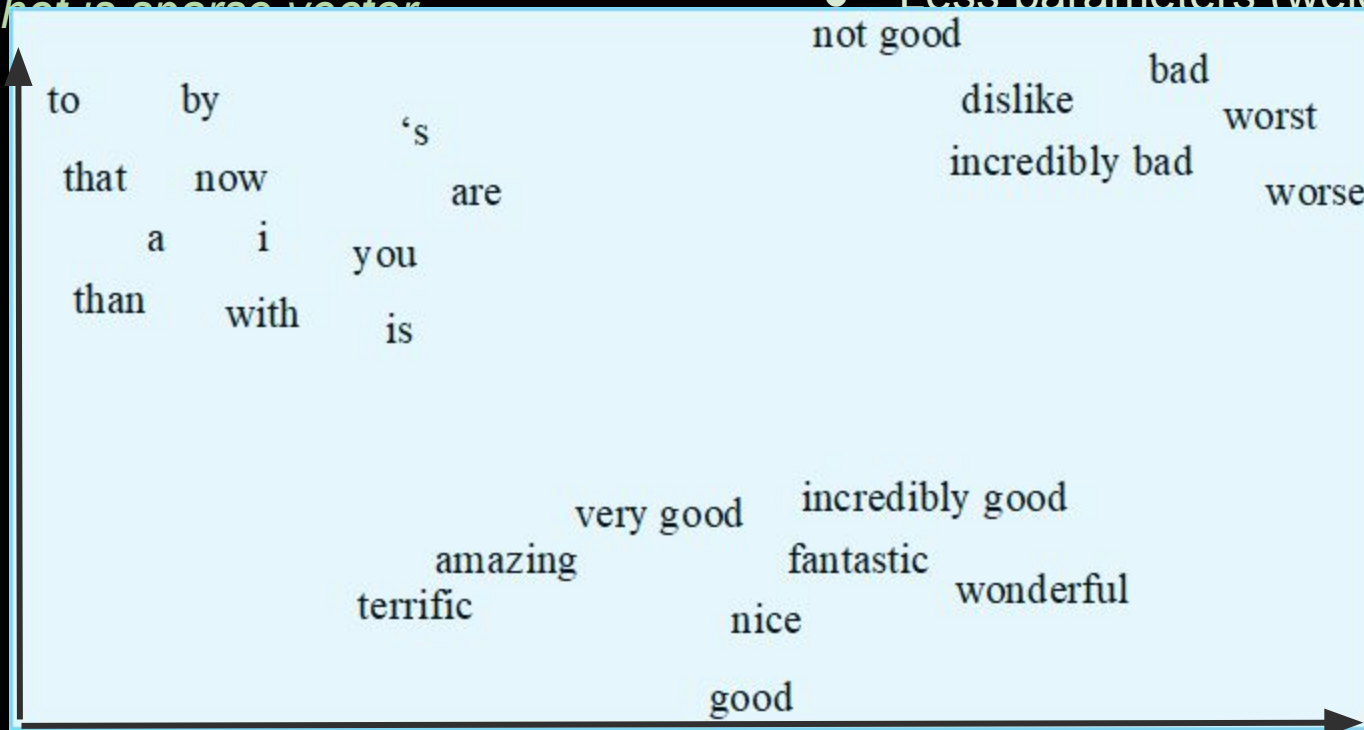
For deep learning, in practice, they work better. Why? Roughly, less parameters becomes increasingly important when you are learning multiple layers of weights rather than just a single layer.

Objective

Prefer dense vectors

- Less parameters (weights) for

one-hot is sparse vector

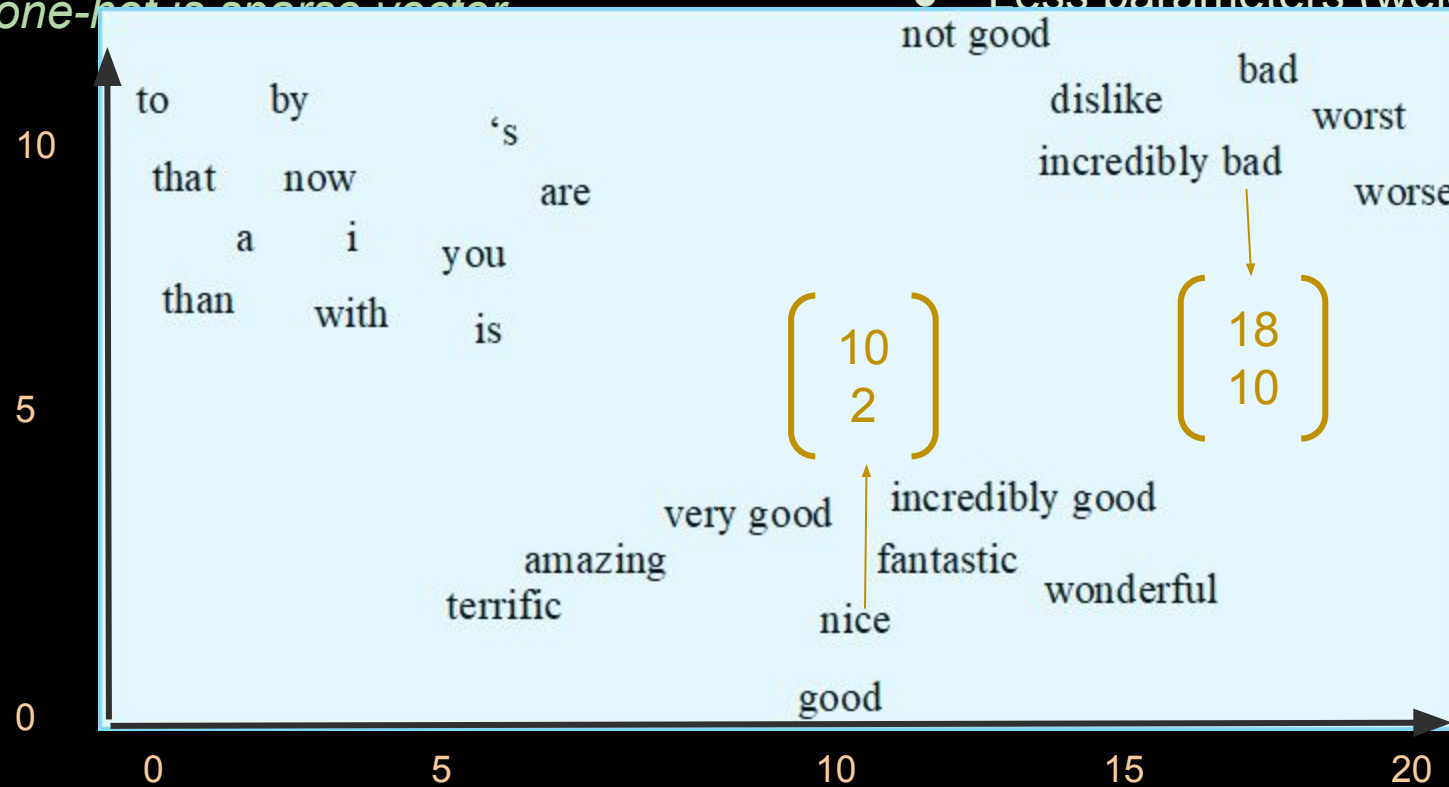


el.
implicitly.

ce, they work
rameters
nt when you are
ghts rather than

Objective

one-hot is sparse vector



Prefer dense vectors

- Less parameters (weights) for

el.
implicitly.

ce, they work
rameters
nt when you are
ghts rather than

(Jurafsky, 2012)

Objective

To embed: convert a token (or sequence) to a vector that represents **meaning**.

Objective

To embed: convert a token (or sequence) to a vector that represents **meaning**.

Wittgenstein, 1945: “*The meaning of a word is its use in the language*”

Distributional hypothesis -- A word's meaning is defined by all the different contexts it appears in (i.e. how it is “distributed” in natural language).

Firth, 1957: “*You shall know a word by the company it keeps*”

Objective

To embed: convert a token (or sequence) to a vector that represents **meaning**.

Wittgenstein, 1945: “*The meaning of a word is its use in the language*”

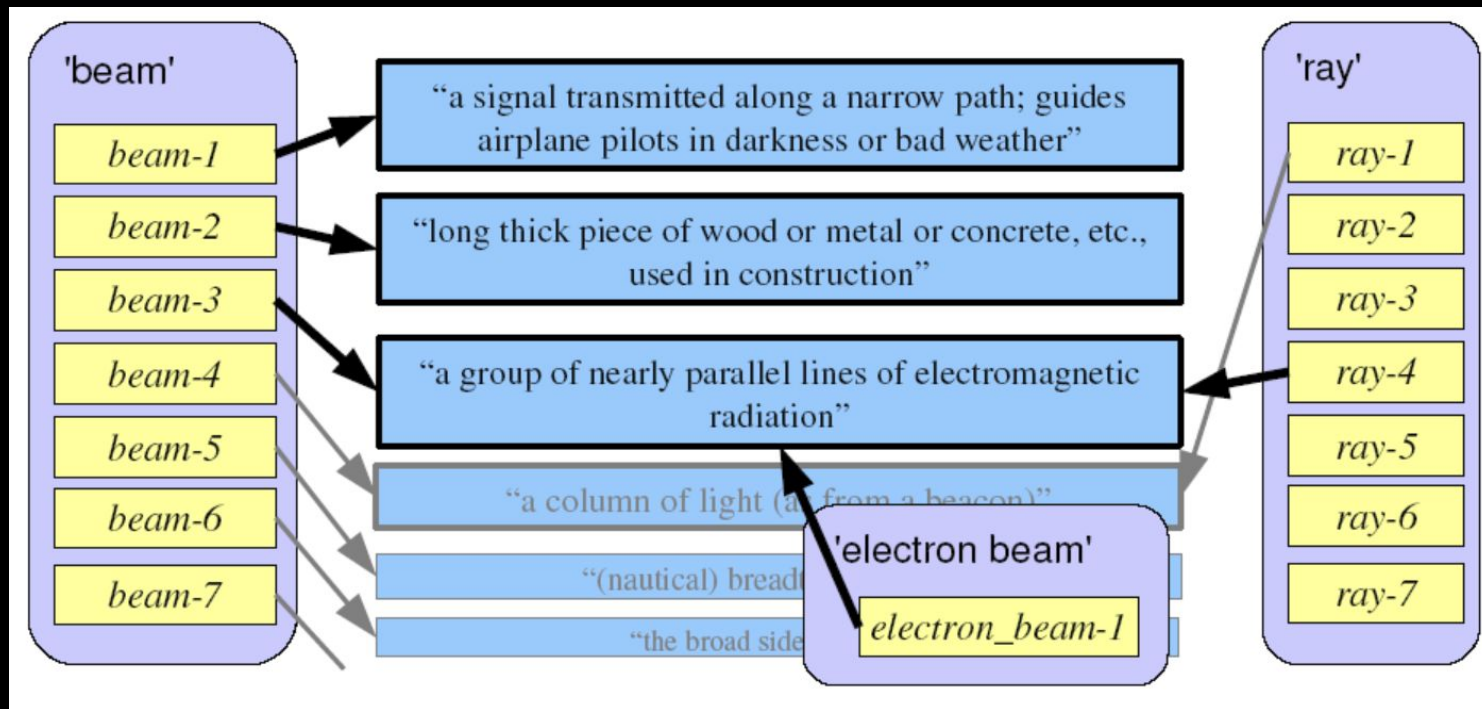
Distributional hypothesis -- A word's meaning is defined by all the different contexts it appears in (i.e. how it is “distributed” in natural language).

Firth, 1957: “*You shall know a word by the company it keeps*”

The nail hit the beam behind the wall.

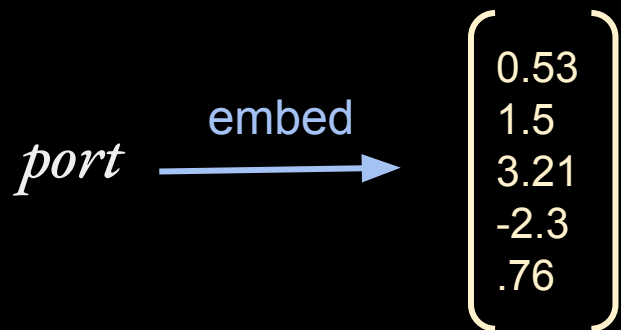


Distributional Hypothesis

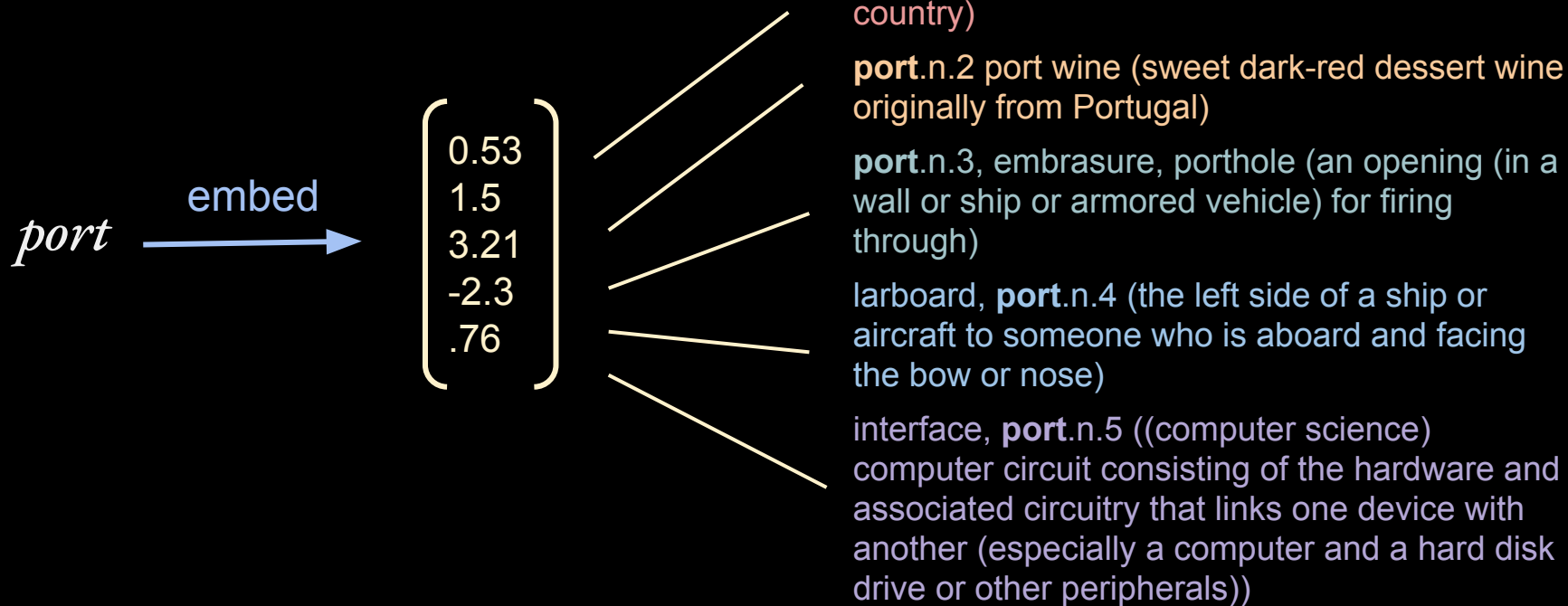


The nail hit the beam behind the wall.

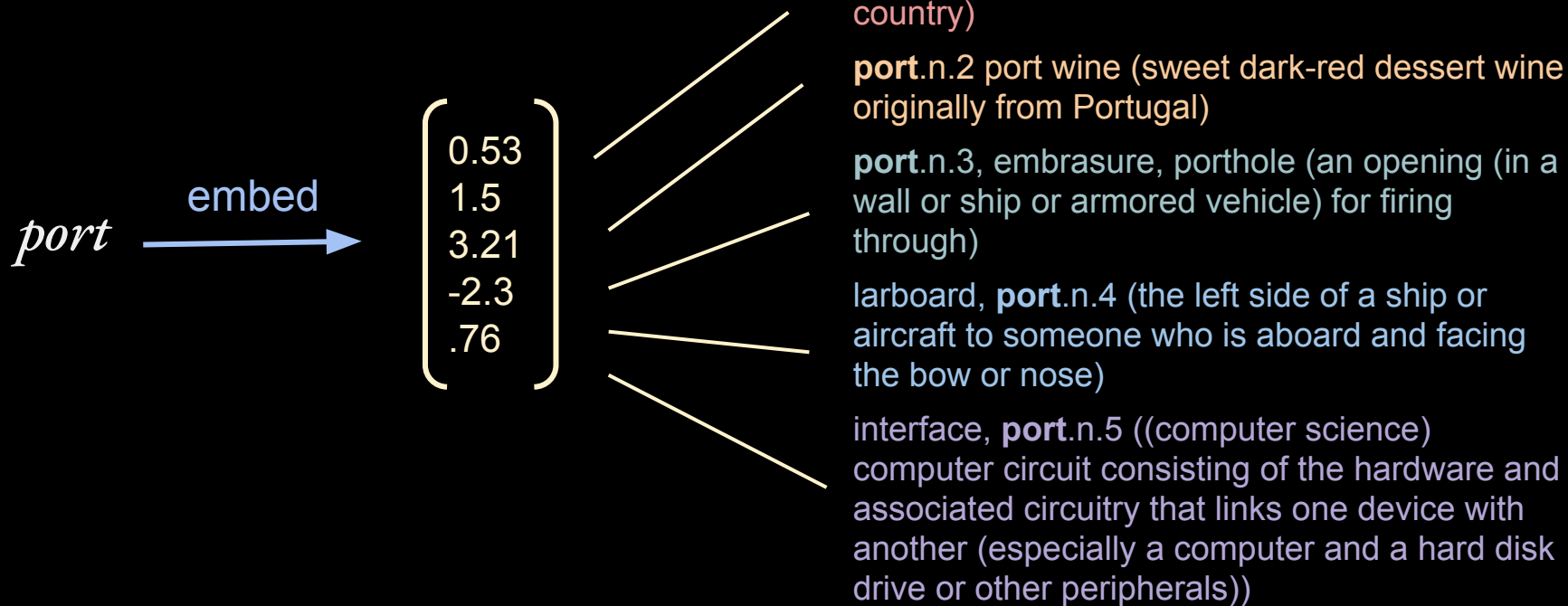
Objective



Objective



Objective



How?

1. One-hot representation
2. Selectors (represent context by “multi-hot” representation)
3. From PCA/Singular Value Decomposition
(Know as “Latent Semantic Analysis” in some circumstances)

Tf-IDF: Term Frequency, Inverse Document Frequency,

PMI: Point-wise mutual information, ...etc...

How?

1. One-hot representation
2. Selectors (represent context by “multi-hot” representation)
3. From PCA/Singular Value Decomposition
(Know as “Latent Semantic Analysis” in some circumstances)

“Neural Embeddings”:

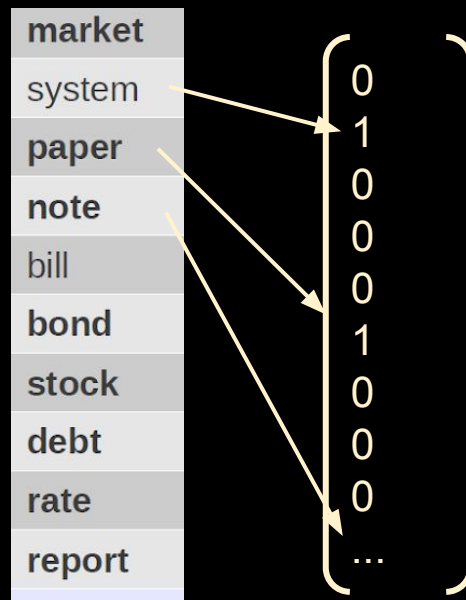
4. Word2vec
5. Fasttext
6. Glove
7. Bert

How?

1. One-hot represent
2. **Selectors (represent context by “multi-hot” representation)**
3. From PCA/Singular Value Decomposition
(Know as “Latent Semantic Analysis” in some circumstances)

..., word1, word2, **bill**, word3, word4, ...

4. Word2vec
5. Fasttext
6. Glove
7. Bert



How?

1. One-hot represent
2. Selectors (represent context by “multi-hot” representation)
3. From PCA/Singular Value Decomposition
(Know as “Latent Semantic Analysis” in some circumstances)

“Neural Embeddings”:

4. Word2vec
5. Fasttext
6. Glove
7. Bert

How?

1. One-hot represent
2. Selectors (represent context by “multi-hot” representation)
3. From PCA/Singular Value Decomposition
(Know as “Latent Semantic Analysis” in some circumstances)

“Neural Embeddings”:

4. Word2vec
5. Fasttext
6. Glove
7. Bert

How?

1. One-hot represent
2. Selectors (represent context by “multi-hot” representation)
3. From PCA/Singular Value Decomposition
(Know as “Latent Semantic Analysis” in some circumstances)

“Neural Embeddings”:

4. Word2vec
5. Fasttext
6. Glove
7. Bert

SVD-Based Embeddings

Singular Value Decomposition...

Concept, In Matrix Form:

f1, f2, f3, f4, ...

o1
o2
o3
...

oN

rows:
N observations

columns:
p features

f_p

SVD-Based Embeddings

context words are features

$f_1, f_2, f_3, f_4, \dots$

f_p

o_1
 o_2
 o_3
...

o_N

co-occurrence counts
are cells.

target words are
observations

SVD-Based Embeddings

Dimensionality reduction
-- try to represent with only p' dimensions

context words are features

$f_1, f_2, f_3, f_4, \dots$

o_1
 o_2
 o_3
...

o_N

co-occurrence counts
are cells.

f_p

$c_1, c_2, c_3, c_4, \dots$

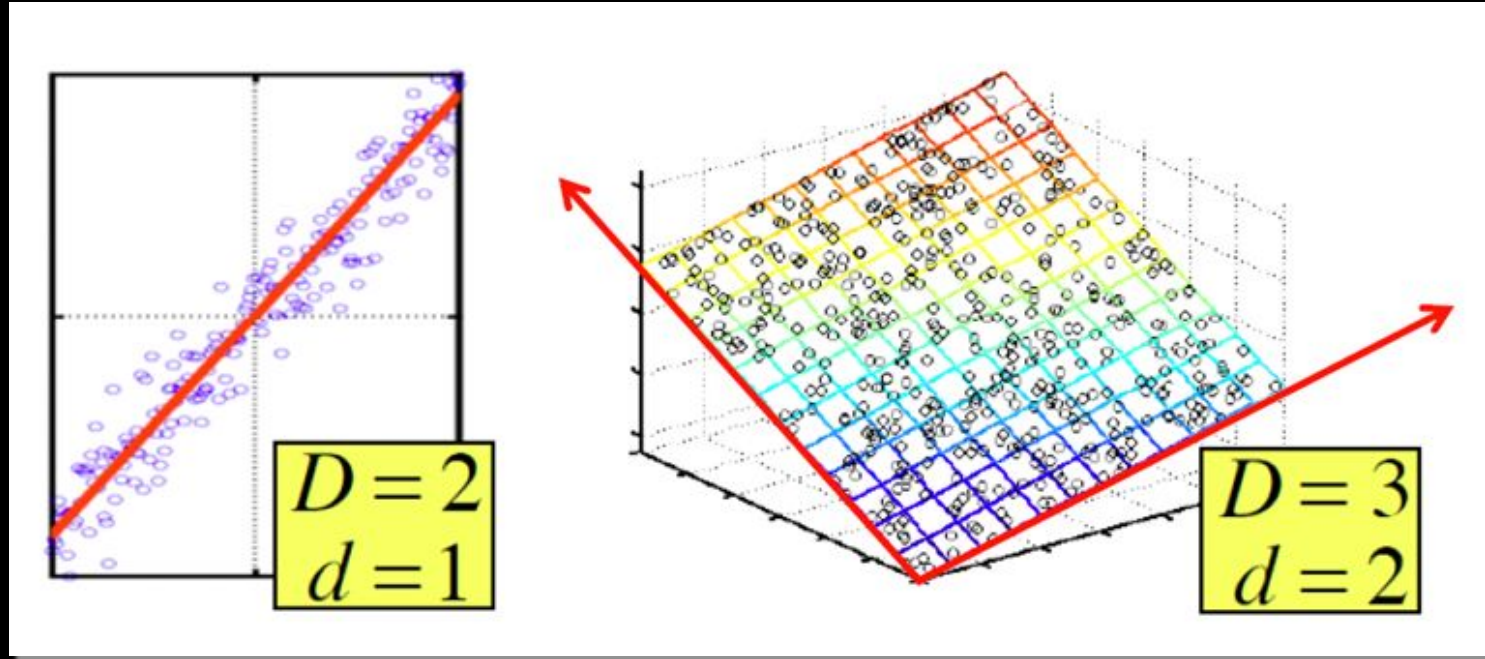
$c_{p'}$

o_1
 o_2
 o_3
...

o_N

target words are
observations

Concept: Dimensionality Reduction in 3-D, 2-D, and 1-D



Data (or, at least, what we want from the data) may be accurately represented with less dimensions.

Dimensionality Reduction

Rank: Number of linearly independent columns of A.

(i.e. columns that can't be derived from the other columns through addition).

Q: What is the rank of this matrix?

$$\begin{pmatrix} 1 & -2 & 3 \\ 2 & -3 & 5 \\ 1 & 1 & 0 \end{pmatrix}$$

Dimensionality Reduction

Rank: Number of linearly independent columns of A.
(i.e. columns that can't be derived from the other columns).

Q: What is the rank of this matrix?

$$\begin{pmatrix} 1 & -2 & 3 \\ 2 & -3 & 5 \\ 1 & 1 & 0 \end{pmatrix}$$

A: 2. The 1st is just the sum of the second two columns

... we can represent as linear combination of 2 vectors:

$$\begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix} \begin{pmatrix} -2 \\ -3 \\ 1 \end{pmatrix}$$

Dimensionality Reduction - PCA

Linear approximates of data in r dimensions.

Found via *Singular Value Decomposition*:

$$X_{[n \times p]} = U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

X: original matrix,

U: “left singular vectors”,

D: “singular values” (diagonal),

V: “right singular vectors”

Dimensionality Reduction - PCA

Linear approximates of data in r dimensions.

Found via *Singular Value Decomposition*:

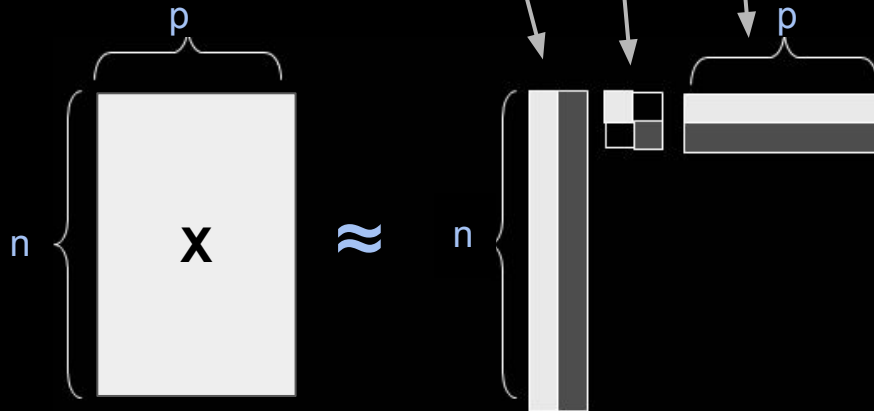
$$X_{[n \times p]} = U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

X : original matrix,

D : “singular values” (diagonal),

U : “left singular vectors”,

V : “right singular vectors”



Dimensionality Reduction - PCA - Example

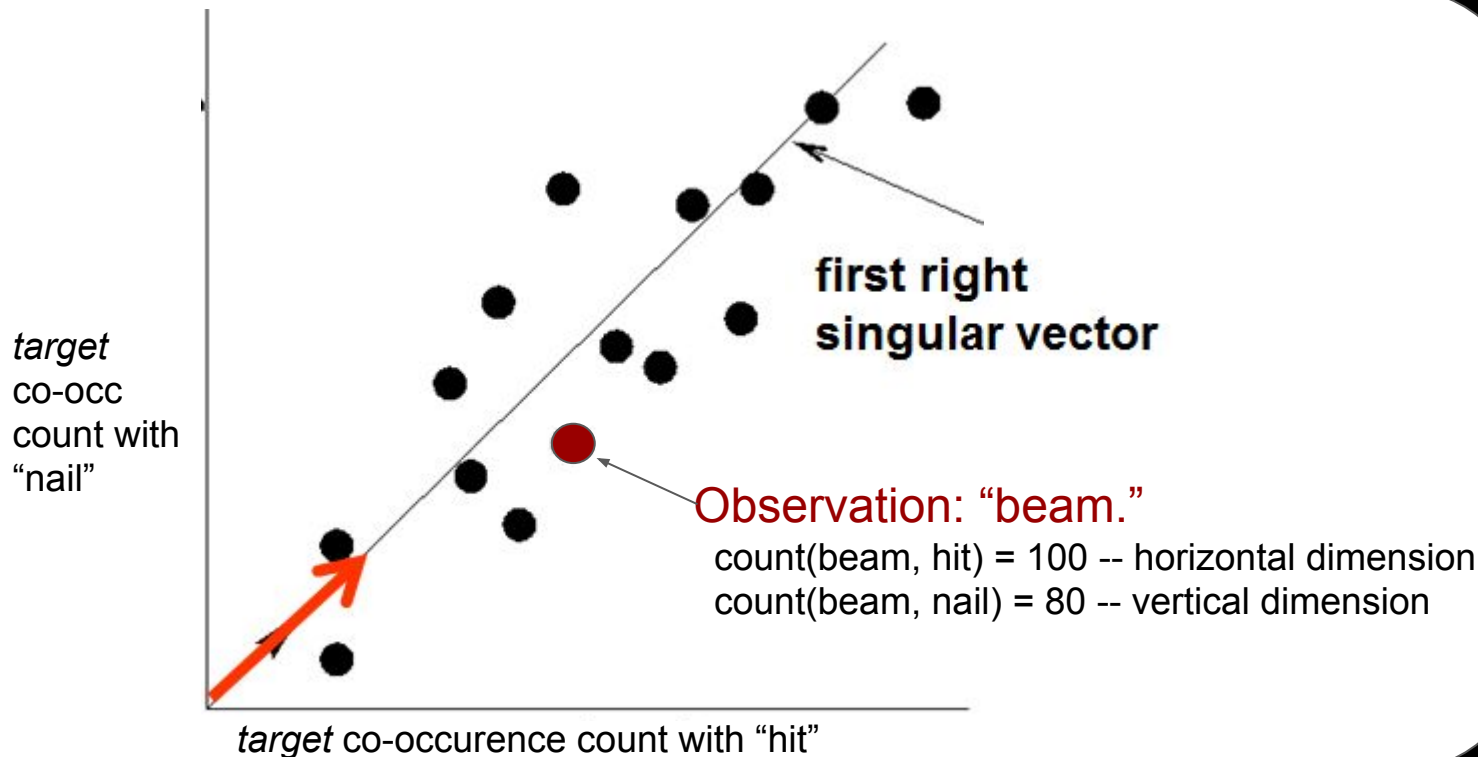
$$X_{[n \times p]} = U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

Word co-occurrence
counts:

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

Dimensionality Reduction - PCA - Example

$$X_{[n \times p]} \approx U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$



Dimensionality Reduction - PCA

Linear approximates of data in r dimensions.

Found via *Singular Value Decomposition*:

$$X_{[n \times p]} \cong U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

X: original matrix,

U: “left singular vectors”,

D: “singular values” (diagonal),

V: “right singular vectors”

Projection (dimensionality reduced space) in 3 dimensions:

$$(U_{[n \times 3]} D_{[3 \times 3]} V_{[p \times 3]}^T)$$

To reduce features in new dataset, A:

$$A_{[m \times p]} V D = A_{\text{small}[m \times 3]}$$

Dimensionality Reduction - PCA

Linear approximates of data in r dimensions.

Found via *Singular Value Decomposition*:

$$X_{[n \times p]} \cong U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

X: original matrix,

U: “left singular vectors”,

D: “singular values” (diagonal),

V: “right singular vectors”

To check how well the original matrix can be reproduced:

$$Z_{[n \times p]} = U D V^T, \text{ How does } Z \text{ compare to } X?$$

To reduce features in new dataset, A:

$$A_{[m \times p]} V D = A_{\text{small}[m \times 3]}$$

Dimensionality Reduction - PCA

Linear approximates of data in r dimensions.

Found via *Singular Value Decomposition*:

$$X_{[n \times p]} \cong U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

X: original matrix,

U: “left singular vectors”,

D: “singular values” (diagonal),

V: “right singular vectors”

To check how well the original matrix can be reproduced:

$$Z_{[n \times p]} = U D V^T, \text{ How does } Z \text{ compare to original } X?$$

To reduce features in new dataset:

$$A_{[m \times p]} V D = A_{\text{small}[m \times 3]}$$

Dimensionality Reduction - PCA

- **Goal: Minimize the sum of reconstruction errors:**

$$\sum_{i=1}^N \sum_{j=1}^D \|x_{ij} - z_{ij}\|^2$$

- where x_{ij} are the “old” and z_{ij} are the “new” coordinates

X: original matrix

D: “singular values”

“singular vectors”,

“singular vectors”

To check how well the original matrix can be reproduced:

$$Z_{[n \times p]} = U D V^T, \text{ How does } Z \text{ compare to original } X?$$

To reduce features in new dataset:

$$A_{[m \times p]} V D = A_{\text{small}[m \times 3]}$$

Dimensionality Reduction - PCA

This is the objective that SVD Solves

- Goal: Minimize the sum of reconstruction errors:

$$\sum_{i=1}^N \sum_{j=1}^D \|x_{ij} - z_{ij}\|^2$$

- where x_{ij} are the “old” and z_{ij} are the “new” coordinates

X: original matrix

D: “singular values”

“singular vectors”,
“singular vectors”

To check how well the original matrix can be reproduced:

$$Z_{[n \times p]} = U D V^T, \text{ How does } Z \text{ compare to original } X?$$

To reduce features in new dataset:

$$A_{[m \times p]} V D = A_{\text{small}[m \times 3]}$$

Dimensionality Reduction - PCA

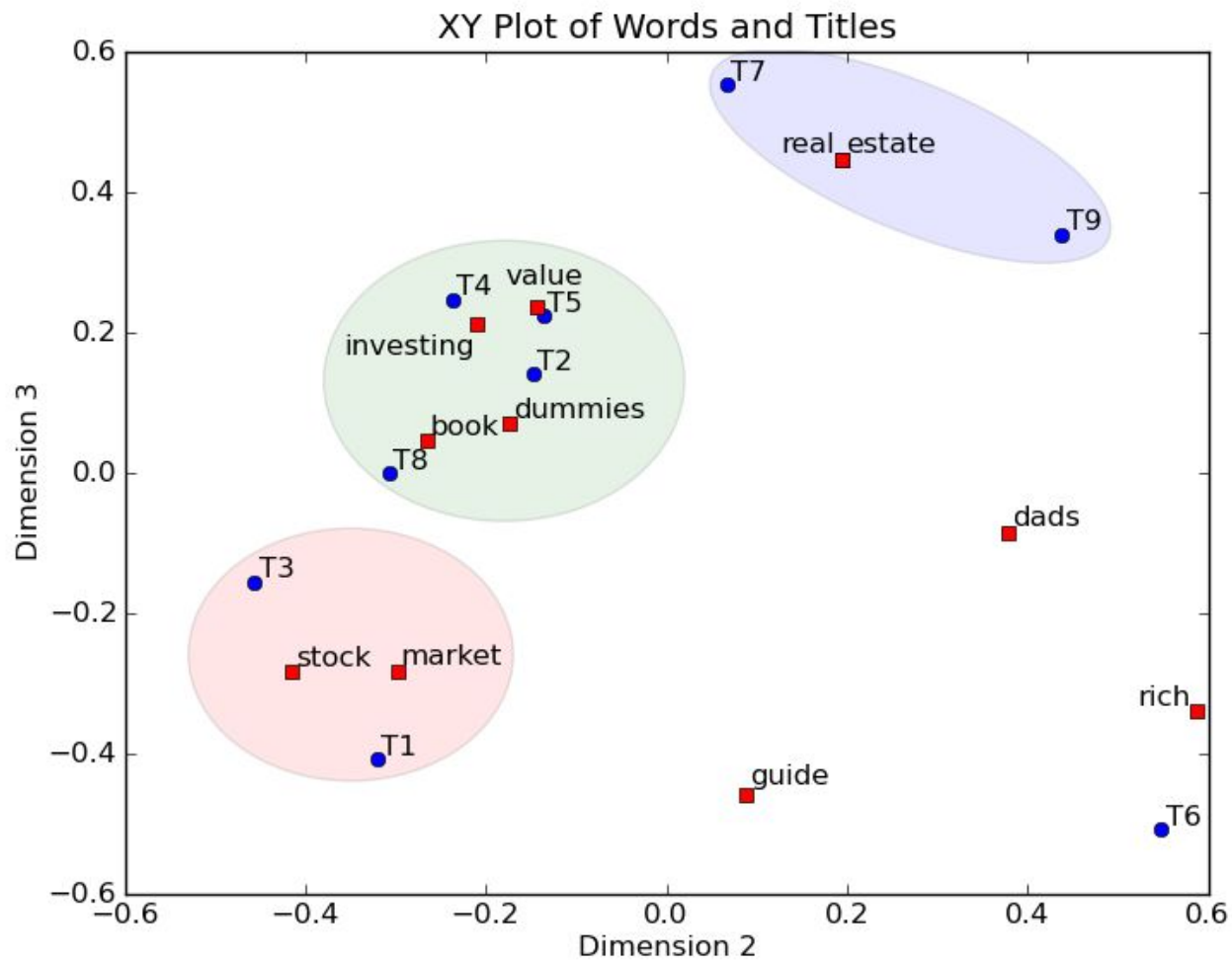
Linear approximates of data in r dimensions.

Found via *Singular Value Decomposition*:

$$X_{[n \times p]} \cong U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

U , D , and V are unique

D : always positive



How?

1. One-hot represent
2. Selectors (represent context by “multi-hot” representation)
3. From PCA/Singular Value Decomposition
(Know as “Latent Semantic Analysis” in some circumstances)

“Neural Embeddings”:

4. Word2vec
5. Fasttext
6. Glove
7. Bert

How?

1. One-hot represent
2. Selectors (represent context by “multi-hot” representation)
3. From PCA/Singular Value Decomposition
(Know as “Latent Semantic Analysis” in some circumstances)

“Neural Embeddings”:

4. **Word2vec**
5. Fasttext
6. Glove
7. Bert

Word2Vec

Principal: Predict missing word.

Similar to language modeling but predicting context, rather than next word.

$$p(\text{context} \mid \text{word})$$

Word2Vec

Principal: Predict missing word.

Similar to language modeling but predicting context, rather than next word.

To learn, maximize

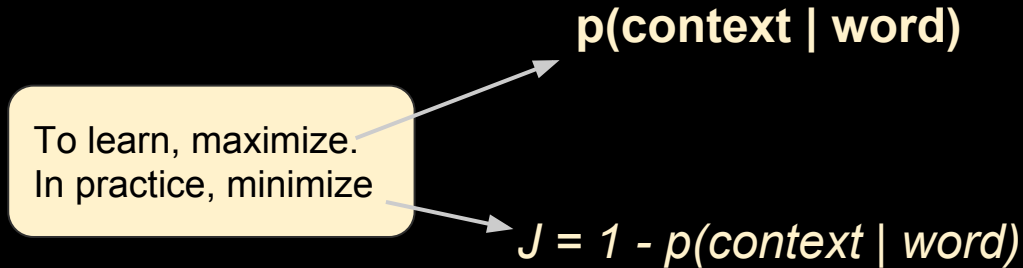
$p(\text{context} \mid \text{word})$

A yellow arrow points from the text 'To learn, maximize' to the mathematical expression $p(\text{context} \mid \text{word})$.

Word2Vec

Principal: Predict missing word.

Similar to language modeling but predicting context, rather than next word.



Word2Vec: Context

$$p(\text{context} \mid \text{word})$$

2 Versions of Context:

- Continuous bag of words (CBOW): Predict word from context
- Skip-Grams (SG): predict context words from target

Word2Vec: Context

$$p(\text{context} \mid \text{word})$$

2 Versions of Context:

- Continuous bag of words (CBOW): Predict word from context
- **Skip-Grams (SG): predict context words from target**

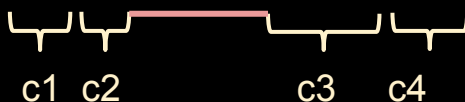
1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

Word2Vec: Context

$$p(\text{context} \mid \text{word})$$

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

The nail hit the beam behind the wall.



Word2Vec: Context

$p(\text{context} \mid \text{word})$

$x = (\text{hit}, \text{beam}), y = 1$

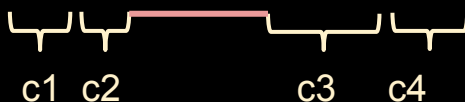
$x = (\text{the}, \text{beam}), y = 1$

$x = (\text{behind}, \text{beam}), y = 1$

...

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

The nail hit the beam behind the wall.



(Jurafsky, 2017)

Word2Vec: Context

$p(\text{context} \mid \text{word})$

$x = (\text{hit}, \text{beam}), y = 1$

$x = (\text{the}, \text{beam}), y = 1$

$x = (\text{behind}, \text{beam}), y = 1$

...

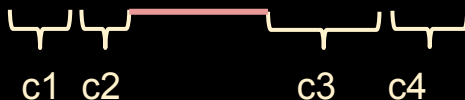
$x = (\text{happy}, \text{beam}), y = 0$

$x = (\text{think}, \text{beam}), y = 0$

...

1. Treat the target word and a neighboring context word as positive examples.
- 2. Randomly sample other words in the lexicon to get negative samples**
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

The nail hit the beam behind the wall.



(Jurafsky, 2017)

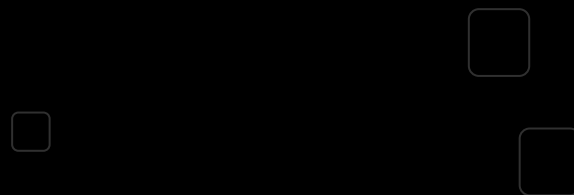
Word2Vec: Context

$p(\text{context} \mid \text{word})$

$x = (\text{hit}, \text{beam}), y = 1$
 $x = (\text{the}, \text{beam}), y = 1$
 $x = (\text{behind}, \text{beam}), y = 1$
...
 $x = (\text{happy}, \text{beam}), y = 0$
 $x = (\text{think}, \text{beam}), y = 0$
...

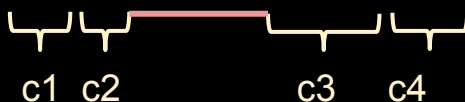
k negative example ($y=0$) for every positive.

How?



1. Treat the target word and a neighboring context word as positive examples.
- 2. Randomly sample other words in the lexicon to get negative samples**
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

The nail hit the beam behind the wall.



(Jurafsky, 2017)

Word2Vec: Context

$p(\text{context} \mid \text{word})$

$x = (\text{hit}, \text{beam}), y = 1$
 $x = (\text{the}, \text{beam}), y = 1$
 $x = (\text{behind}, \text{beam}), y = 1$
...
 $x = (\text{happy}, \text{beam}), y = 0$
 $x = (\text{think}, \text{beam}), y = 0$
...

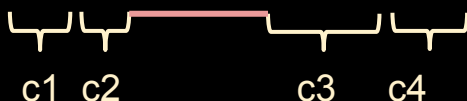
k negative example ($y=0$) for every positive.

How? Randomly draw from unigram distribution

$$P_{\square}(w) = \frac{\text{count}(w)^{\square}}{\sum_w \text{count}(w)^{\square}}$$

1. Treat the target word and a neighboring context word as positive examples.
- 2. Randomly sample other words in the lexicon to get negative samples**
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

The nail hit the beam behind the wall.


c1 c2 c3 c4

(Jurafsky, 2017)

Word2Vec: Context

$p(\text{context} \mid \text{word})$

$x = (\text{hit}, \text{beam}), y = 1$
 $x = (\text{the}, \text{beam}), y = 1$
 $x = (\text{behind}, \text{beam}), y = 1$
...
 $x = (\text{happy}, \text{beam}), y = 0$
 $x = (\text{think}, \text{beam}), y = 0$
...

k negative example ($y=0$) for every positive.

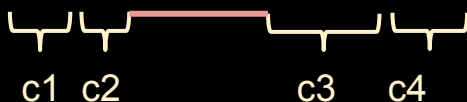
How? Randomly draw from unigram distribution adjusted:

$$P_{\alpha}(w) = \frac{\text{count}(w)^{\alpha}}{\sum_w \text{count}(w)^{\alpha}}$$

$\alpha = 0.75$

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

The nail hit the beam behind the wall.



(Jurafsky, 2017)

Word2Vec: Context

$p(\text{context} \mid \text{word})$

$x = (\text{hit}, \text{beam}), y = 1$
 $x = (\text{the}, \text{beam}), y = 1$
 $x = (\text{behind}, \text{beam}), y = 1$
...
 $x = (\text{happy}, \text{beam}), y = 0$
 $x = (\text{think}, \text{beam}), y = 0$
...

k negative example ($y=0$) for every positive.

How? Randomly draw from unigram distribution adjusted:

$$P_{\alpha}(w) = \frac{\text{count}(w)^{\alpha}}{\sum_w \text{count}(w)^{\alpha}}$$

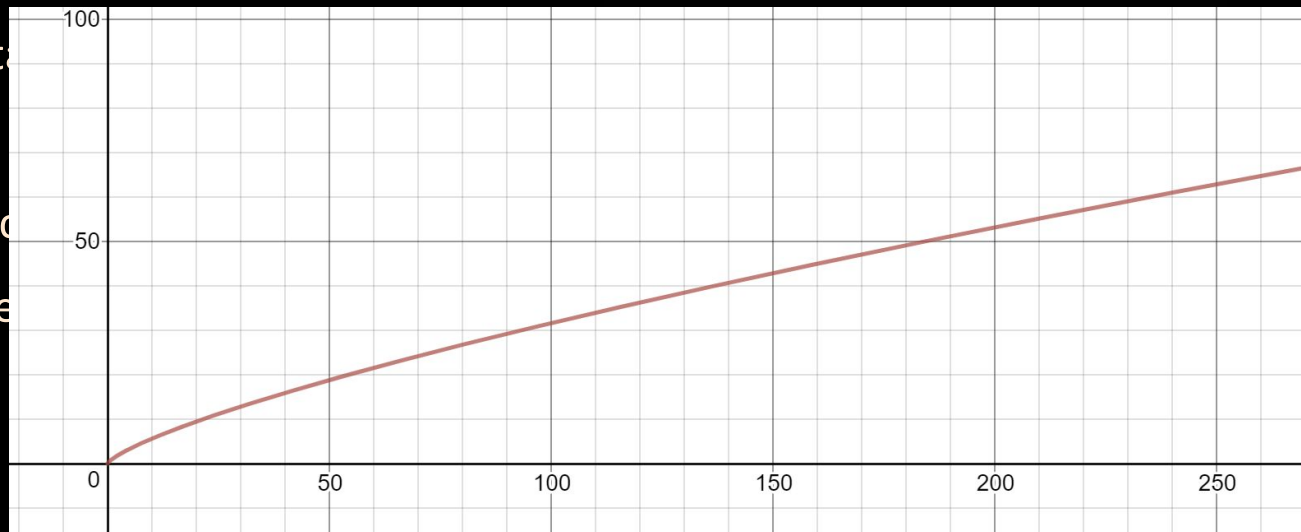
$\alpha = 0.75$

1. Treat the task as

2. Randomly

3. Use logistic

4. Use the weights



c1 c2

c3 c4

(Jurafsky, 2017)

Word2Vec: Context

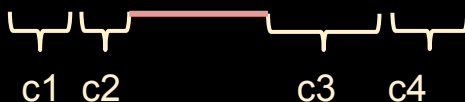
x = (hit, beam), y = 1
x = (the, beam), y = 1
x = (behind, beam), y = 1
...
x = (happy, beam), y = 0
x = (think, beam), y = 0
...

single context:

$$P(y=1 | c, t) = \frac{1}{1 + e^{-t \cdot c}}$$

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
- 3. Use logistic regression to train a classifier to distinguish those two cases**
4. Use the weights as the embeddings

The nail hit the beam behind the wall.



Word2Vec: Context

Logistic: $\sigma(z) = 1 / (1 + e^{-z})$

x = (hit, beam), y = 1
x = (the, beam), y = 1
x = (behind, beam), y = 1
...
x = (happy, beam), y = 0
x = (think, beam), y = 0
...

single context:

$$P(y=1 | c, t) = \frac{1}{1 + e^{-t \cdot c}}$$

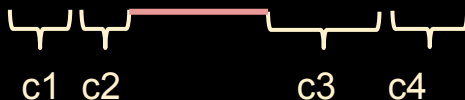
All Contexts

$$P(y=1 | c, t) = \prod_{i=1}^n \frac{1}{1 + e^{-t \cdot c_i}}$$



1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
- 3. Use logistic regression to train a classifier to distinguish those two cases**
4. Use the weights as the embeddings

The nail hit the beam behind the wall.



Word2Vec: Context

x = (hit, beam), y = 1
x = (the, beam), y = 1
x = (behind, beam), y = 1
...
x = (happy, beam), y = 0
x = (think, beam), y = 0
...

single context:

$$P(y=1 | c, t) = \frac{1}{1 + e^{-t \cdot c}}$$

Intuition: $t \cdot c$ is a measure of similarity:

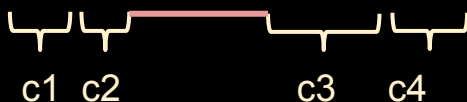
$$a \cdot b = \|a\| \|b\| \cos \theta$$

But, it is not a probability! To make it one, apply logistic activation:

$$\sigma(z) = 1 / (1 + e^{-z})$$

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
- 3. Use logistic regression to train a classifier to distinguish those two cases**
4. Use the weights as the embeddings

The nail hit the beam behind the wall.



Word2Vec: Context

x = (hit, beam), y = 1
x = (the, beam), y = 1
x = (behind, beam), y = 1
...
x = (happy, beam), y = 0
x = (think, beam), y = 0
...

single context:

$$P(y=1 | c, t) = \frac{1}{1 + e^{-t \cdot c}}$$

all contexts

$$P(y=1 | c, t) = \prod_{i=1}^n \frac{1}{1 + e^{-t \cdot c_i}}$$

Intuition: $t \cdot c$ is a measure of similarity:

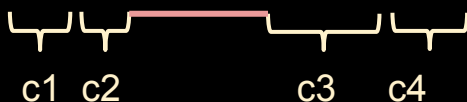
$$a \cdot b = \|a\| \|b\| \cos \theta$$

But, it is not a probability! To make it one, apply logistic activation:

$$\sigma(z) = 1 / (1 + e^{-z})$$

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

The nail hit the beam behind the wall.

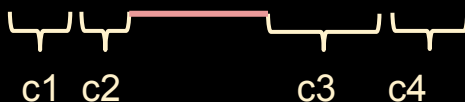


Word2Vec: How to Learn?

$$P(y=1 | c, t)$$

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
- 4. Use the weights as the embeddings**

The nail hit the beam behind the wall.



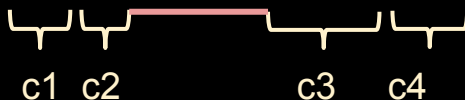
Word2Vec: How to Learn?

$$P(y=1 | c, t)$$

Assume $300 * |\text{vocab}|$ weights (parameters) for each of c and t

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
- 4. Use the weights as the embeddings**

The nail hit the beam behind the wall.



Word2Vec: How to Learn?

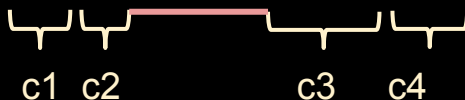
$$P(y=1 | c, t)$$



Assume $300 * |\text{vocab}|$ weights (parameters) for each of c and t
Start with random vectors (or all 0s)

1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
- 4. Use the weights as the embeddings**

The nail hit the beam behind the wall.



Word2Vec: How to Learn?

$$P(y=1 | c, t)$$



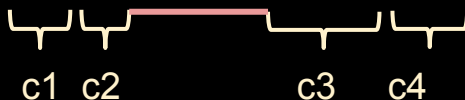
Assume $300 * |\text{vocab}|$ weights (parameters) for each of c and t
Start with random vectors (or all 0s)

Goal:


Maximize similarity of (c, t) in positive data ($y = 1$)

Minimize similarity of (c, t) in negative data ($y = 0$)

The nail hit the beam behind the wall.



Word2Vec: How to Learn?

$$P(y=1|c, t)$$


Assume $300 * |\text{vocab}|$ weights (parameters) for each of c and t
Start with random vectors (or all 0s)


Goal:

Maximize similarity of (c, t) in positive data ($y = 1$)

Minimize similarity of (c, t) in negative data ($y = 0$)

$$\sum_{(c,t)} (y) \log P(y = 1|c, t) + (y - 1) \log P(y = 0|c, t)$$

Word2Vec: How to Learn?

$$P(y=1|c, t)$$


Assume $300 * |\text{vocab}|$ weights (parameters) for each of c and t
Start with random vectors (or all 0s)

Goal:

Maximize similarity of (c, t) in positive data ($y = 1$)

Minimize similarity of (c, t) in negative data ($y = 0$)

$$\sum_{(c,t)} (y) \log P(y = 1|c, t) + (y - 1) \log P(y = 0|c, t)$$

$$1 - P(y = 1|c, t) = \frac{e^{-t \cdot c}}{1 + e^{-t \cdot c}}$$

Word2Vec: How to Learn?

$$P(y=1|c, t)$$

Assume
Start with

Optimized using gradient descent type methods.

for each of c and t

Goal:

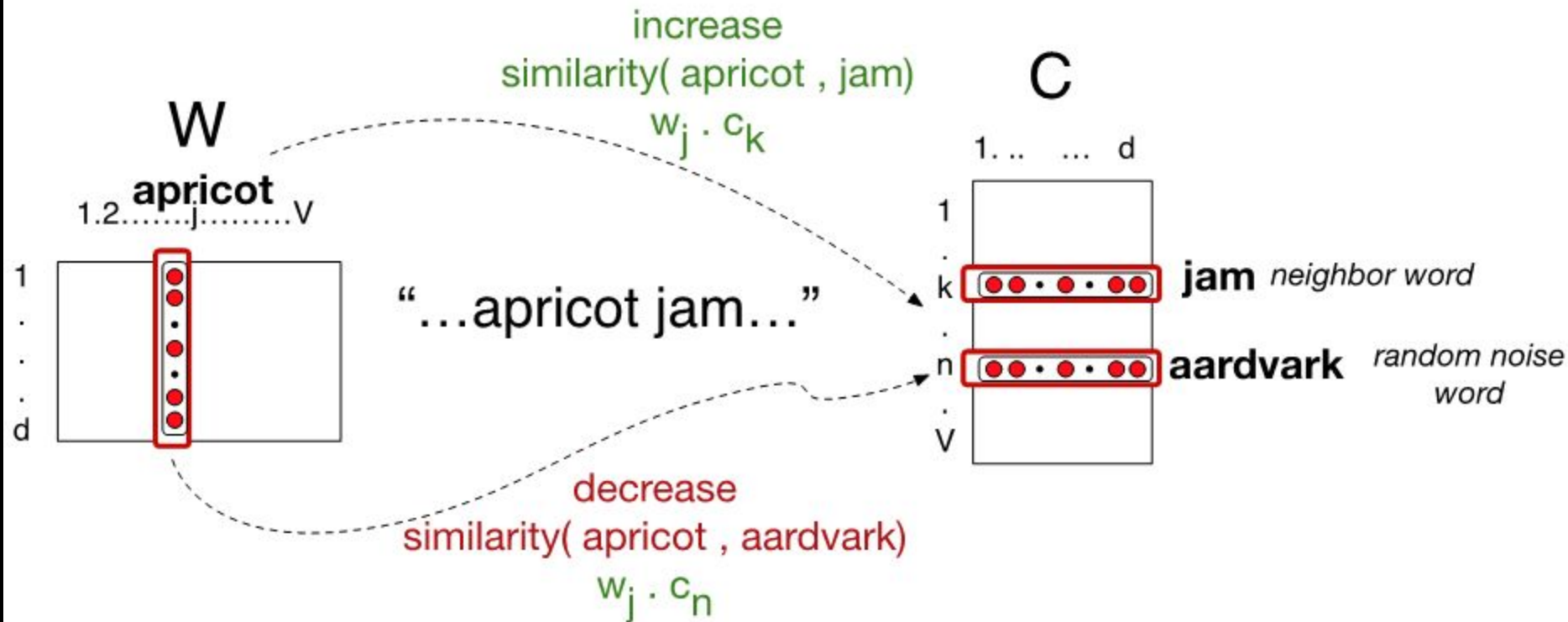
Maximize similarity of (c, t) in positive data ($y = 1$)

Minimize similarity of (c, t) in negative data ($y = 0$)

$$\sum_{(c,t)} (y) \log P(y = 1|c, t) + (y - 1) \log P(y = 0|c, t)$$

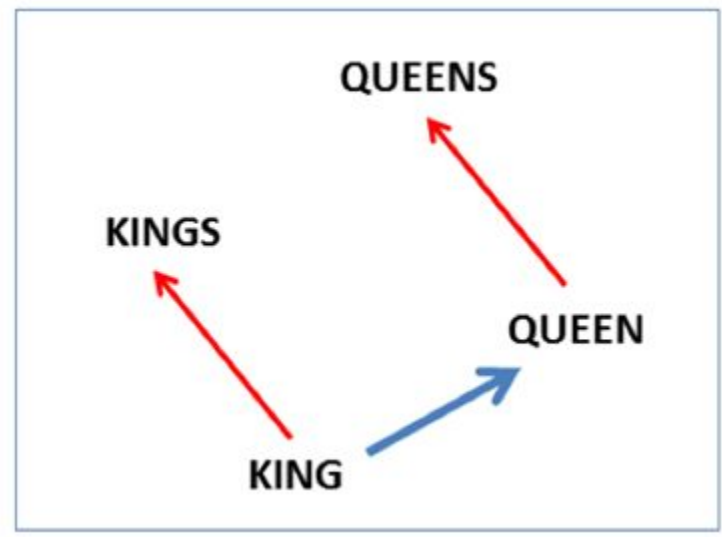
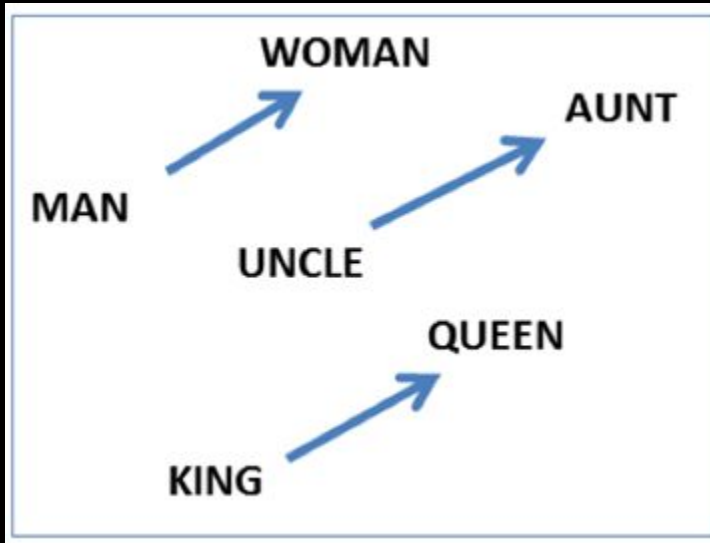
$$1 - P(y = 1|c, t) = \frac{e^{-t \cdot c}}{1 + e^{-t \cdot c}}$$

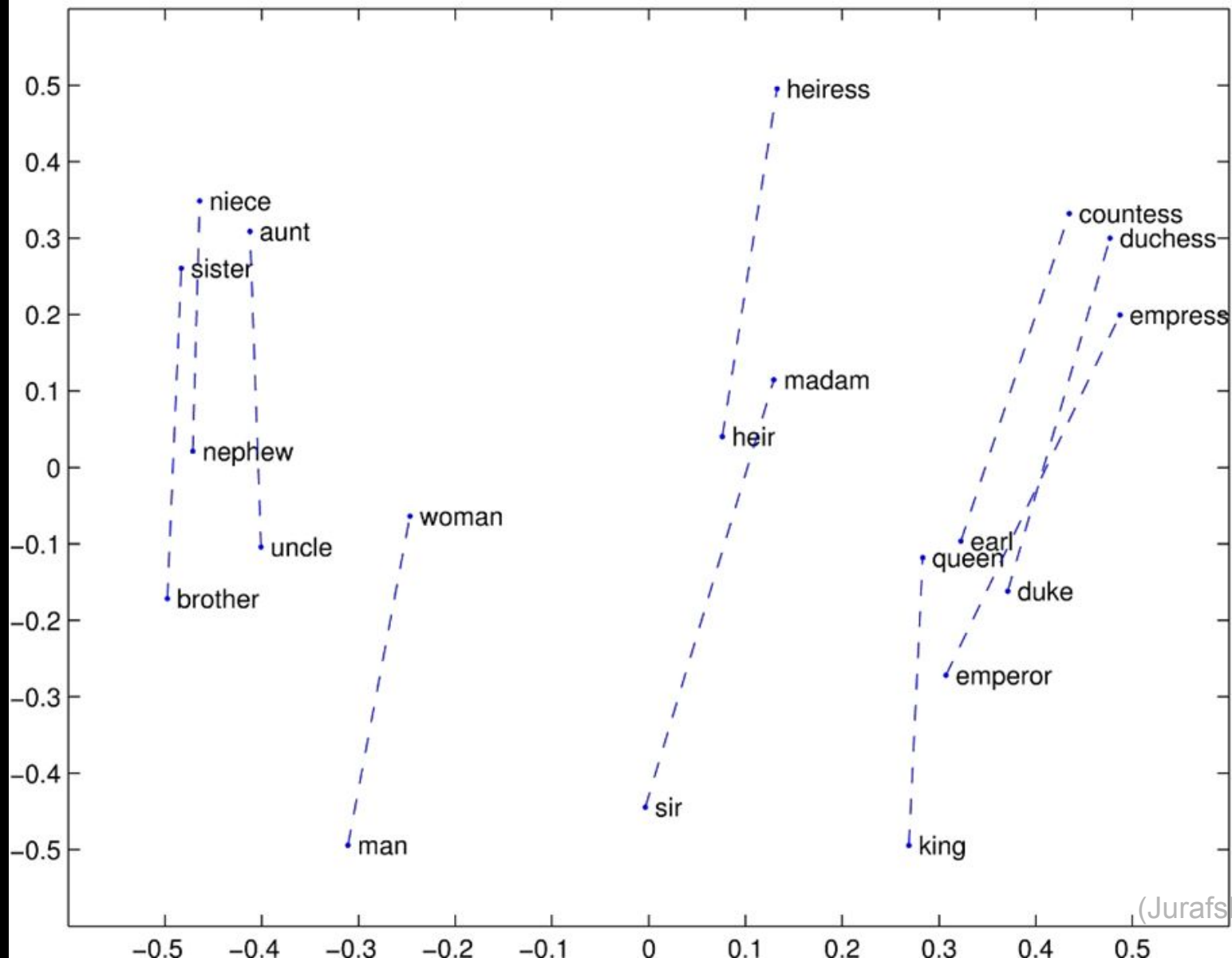
Word 2 Vec



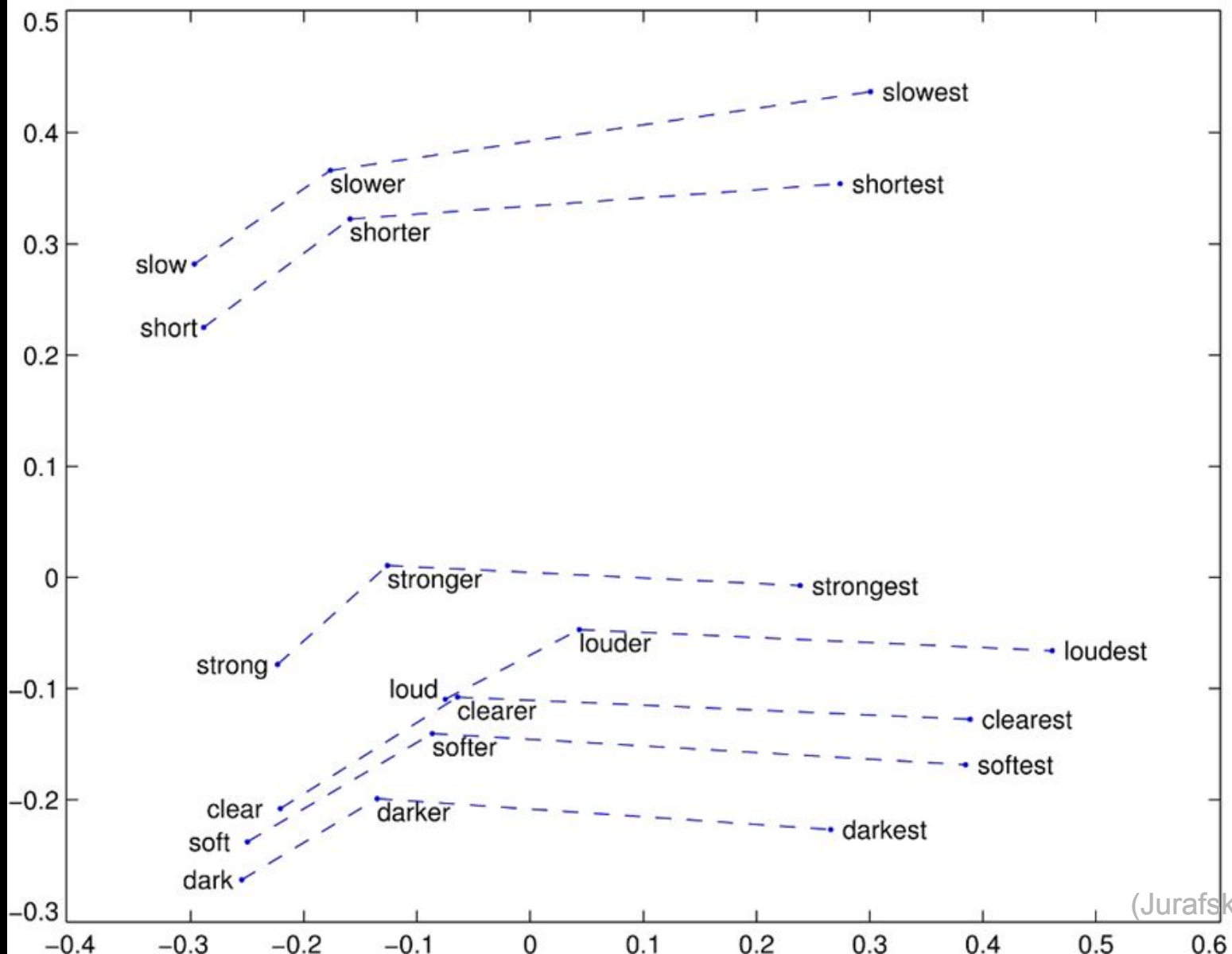
$$\sum_{(c,t)} (y) \log P(y = 1|c, t) + (y - 1) \log P(y = 0|c, t)$$

Word2Vec captures analogies (kind of)





(Jurafsky, 2017)



(Jurafsky, 2017)

Word2Vec: Quantitative Evaluations

Compare to manually annotated pairs of words: WordSim-353 (Finkelstein et al., 2002)

Compare to words in context (Huang et al., 2012)

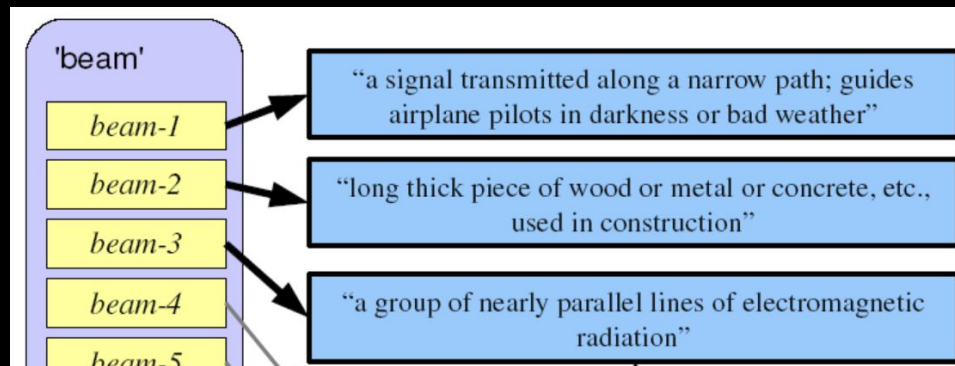
Answer [TOEFL synonym questions](#).

Current Trends in Embeddings

1. Contextual word embeddings (a different embedding depending on context):

The nail hit the beam behind the wall.

They reflected a beam off the moon.



Current Trends in Embeddings

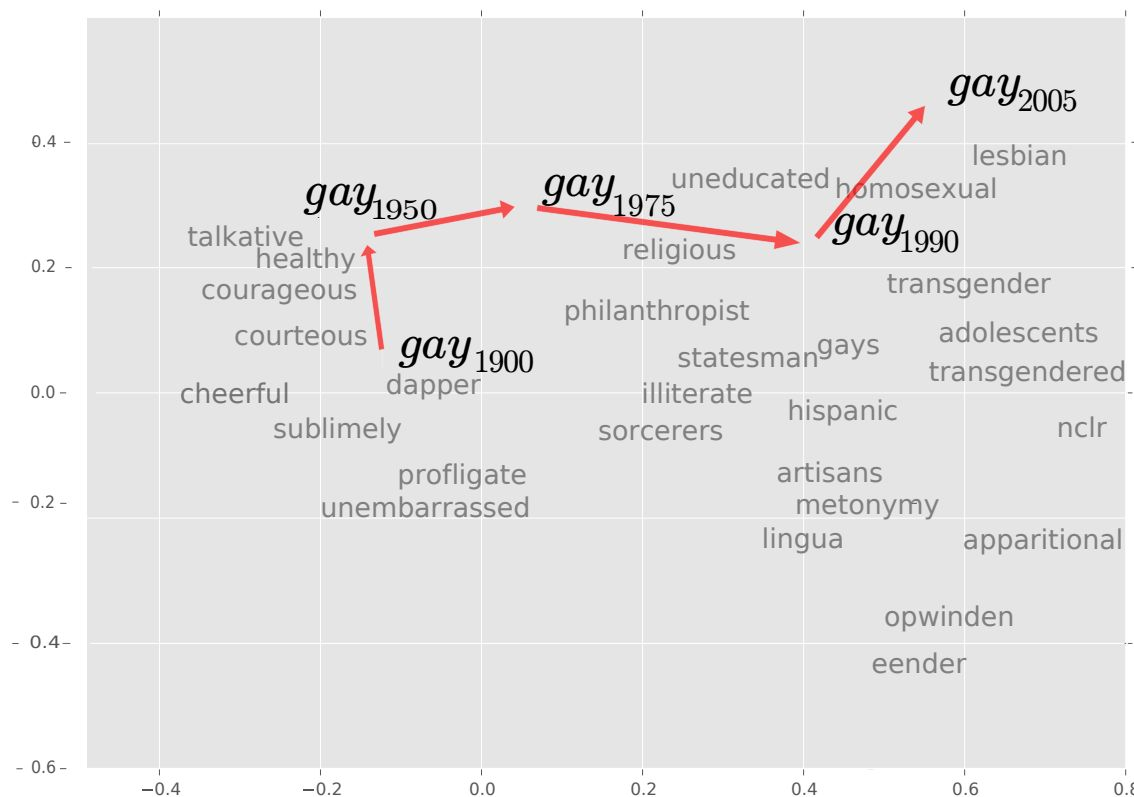
1. Contextual word embeddings (a different embedding depending on context):

The nail hit the beam behind the door.

They reflected a beam off the mirror.

2. Embeddings can capture changes in word meaning.

(Kulkarni et al., 2015)



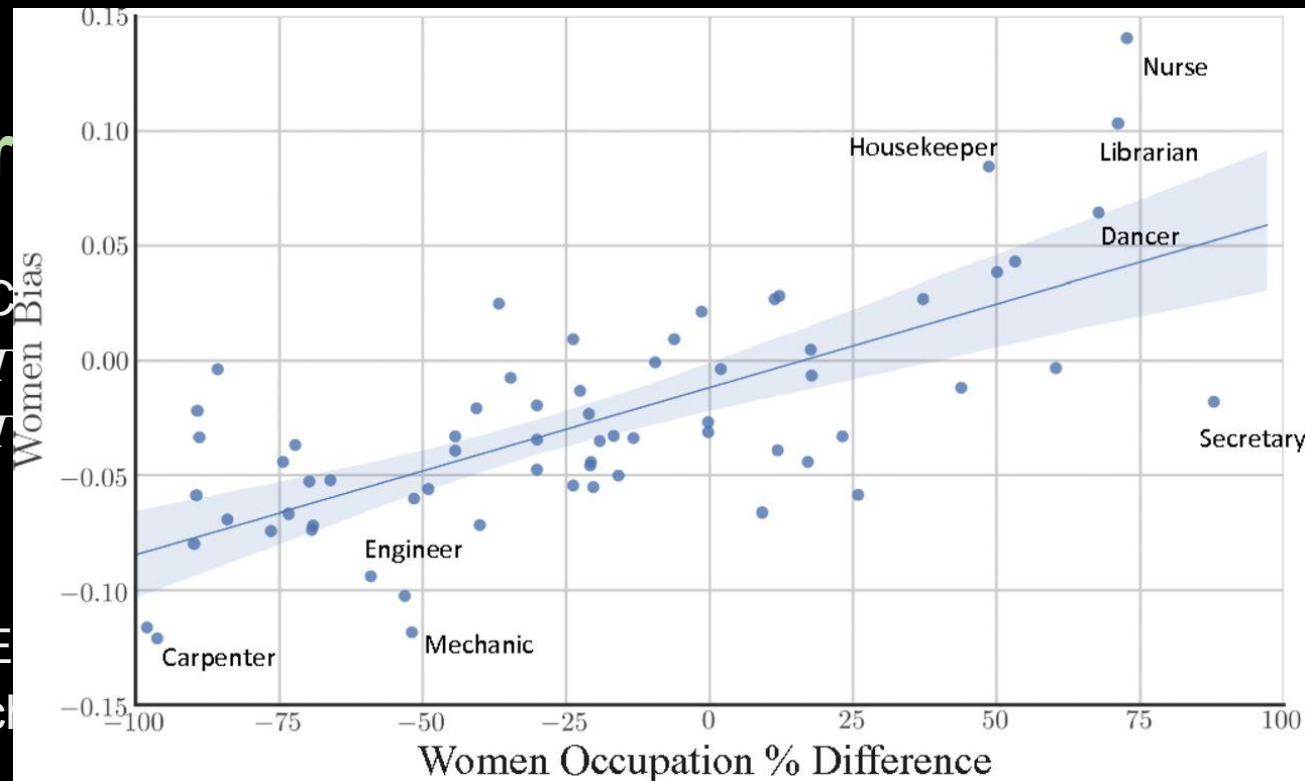
Current

1. Context

Women Bias

2. Embeddings

Context



g on context):

(Garg et al., 2018)

3. Embeddings capture demographic biases in data.

Current Trends

1. Contextual word embeddings

The nail hit the board

They reflected a board

2. Embeddings capture changes in word meaning

3. Embeddings capture demographic biases in data.

a. Efforts to debias

b. Useful for tracking bias over time.

(Garg et al., 2018)

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

Tolga Bolukbasi¹, Kai-Wei Chang², James Zou², Venkatesh Saligrama^{1,2}, Adam Kalai²
¹Boston University, 8 Saint Mary's Street, Boston, MA
²Microsoft Research New England, 1 Memorial Drive, Cambridge, MA
tolgab@bu.edu, kw@kwchang.net, jameszou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

Abstract

The blind application of machine learning runs the risk of amplifying biases present in data. Such a danger is facing us with *word embedding*, a popular framework to represent text data as vectors which has been used in many machine learning and natural language processing tasks. We show that even word embeddings trained on Google News articles exhibit female/male gender stereotypes to a disturbing extent. This raises concerns because their widespread use, as we describe, often tends to amplify these biases. Geometrically, gender bias is first shown to be captured by a direction in the word embedding. Second, gender neutral words are shown to be linearly separable from gender definition words in the word embedding. Using these properties, we provide a methodology for modifying an embedding to remove gender stereotypes, such as the association between the words *receptionist* and

Current Trends

1. Contextual word embeddings

The nail hit the bear

They reflected a bear

2. Embeddings can capture

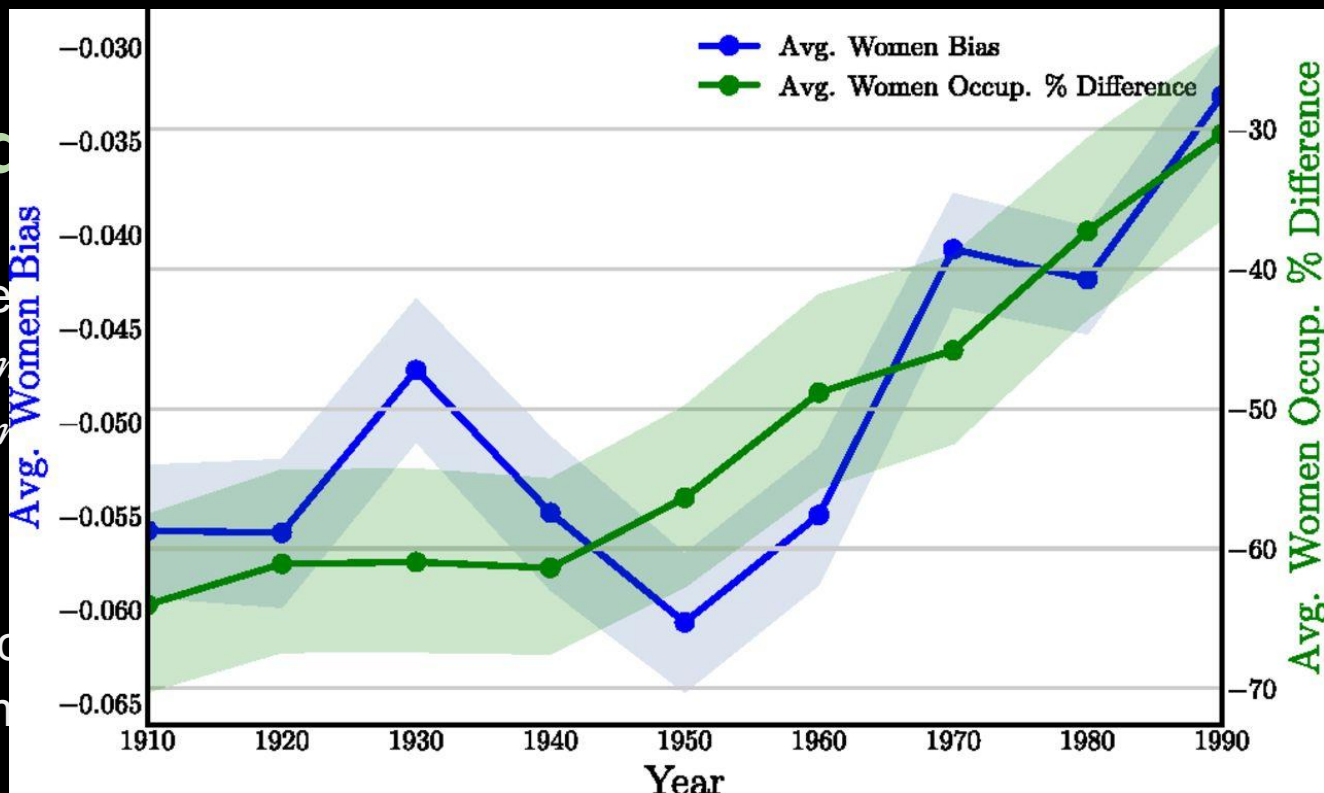
changes in word meaning

3. Embeddings capture demographic biases in data.

a. Efforts to debias

b. Useful for tracking bias over time.

(Garg et al., 2018)



Vector Semantics and Embeddings

Take-Aways

- Dense representation of meaning is desirable.
- Approach 1: Dimensionality reduction techniques
- Approach 2: Learning representations by trying to predict held-out words.
- Word2Vec skipgram model attempts to solve by predicting target word from context word:
maximize similarity between true pairs; minimize similarity between random pairs.
- Embeddings do in fact seem to capture meaning in applications
- Dimensionality reduction techniques just as good by some evaluations.
- Current Trends: Integrating context, Tracking changes in meaning.