

# Lexical and Vector Semantics

CSE538 - Spring 2024  
Natural Language Processing

# Topics

- Lexical Ambiguity (why word sense disambiguation)
- Word Vectors
- Topic Modeling

# Objectives

- Define common semantic tasks in NLP and learn some approaches to solve.
- Understand linguistic information necessary for semantic processing
- Motivate deep learning models necessary to capture language semantics.
- Learn word embeddings (the starting point for modern large language models)

# Terminology: lemma and wordform

- A **lemma** or **citation form**
  - Same stem, part of speech, rough semantics
- A **wordform**
  - The inflected word as it appears in text

Wordform	Lemma
banks	bank
sung	sing
duermes	dormir

(Jurafsky & Martin, SLP, 2019)

## Lemmas have senses

- One lemma “bank” can have many meanings:

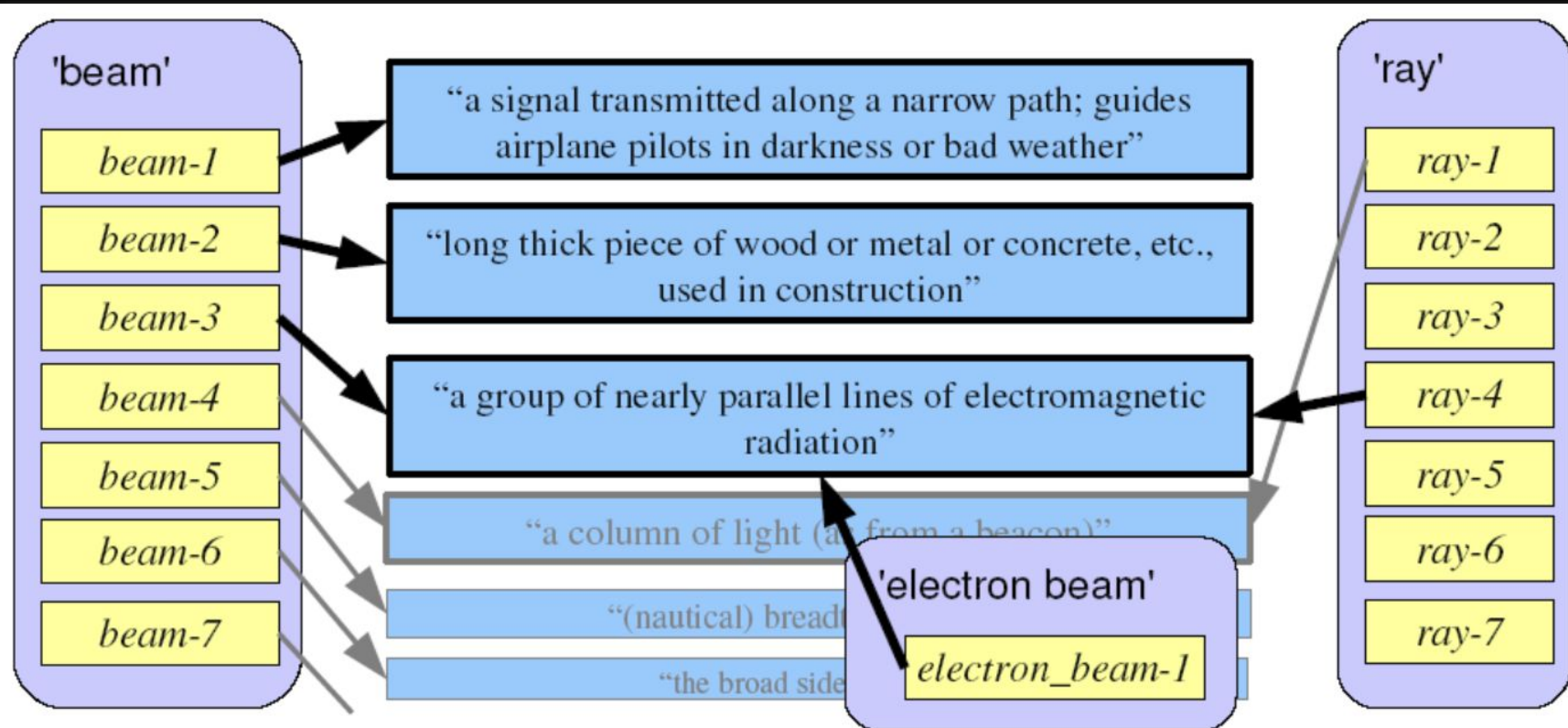
Sense 1: • ...a **bank**<sub>1</sub> can hold the investments in a custodial account...

Sense 2: • “...as agriculture burgeons on the east **bank**<sub>2</sub> the river will shrink even more”

- **Sense (or word sense)**
  - A discrete representation of an aspect of a word’s meaning.
- The lemma **bank** here has two senses

(Jurafsky & Martin, SLP, 2019)

# Lemmas have senses



(Schwartz, 2011)

# Homonymy

**Homonyms:** words that share a form but have unrelated, distinct meanings:

- bank<sub>1</sub>: financial institution, bank<sub>2</sub>: sloping land
- bat<sub>1</sub>: club for hitting a ball, bat<sub>2</sub>: nocturnal flying mammal

1. Homographs (bank/bank, bat/bat)

2. Homophones:

1. Write and right
2. Piece and peace

(Jurafsky & Martin, SLP, 2019)

# Homonymy causes problems for NLP applications

- Information retrieval
  - “bat care”
- Machine Translation
  - bat: murciélago (animal) or bate (for baseball)
- Text-to-Speech
  - bass (stringed instrument) vs. bass (fish)

(Jurafsky & Martin, SLP, 2019)

# Word Sense Disambiguation

He put the **port** on the ship.

He walked along the **port** of the steamer.

He walked along the **port** next to the steamer.

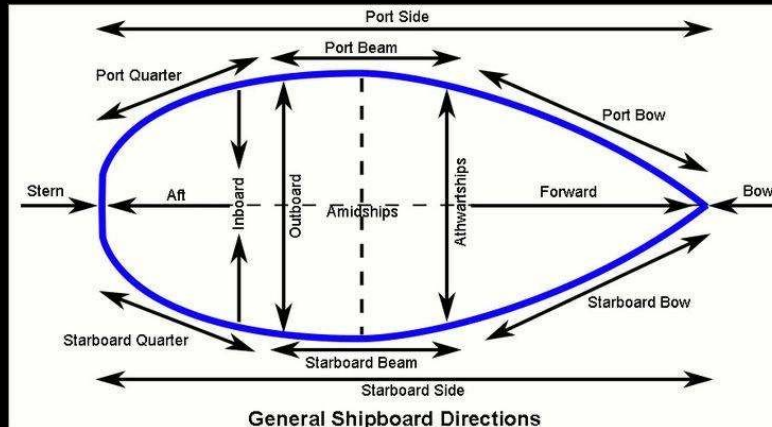


# Word Sense Disambiguation

He put the **port** on the ship.

He walked along the **port** of the steamer.

He walked along the **port** next to the steamer.

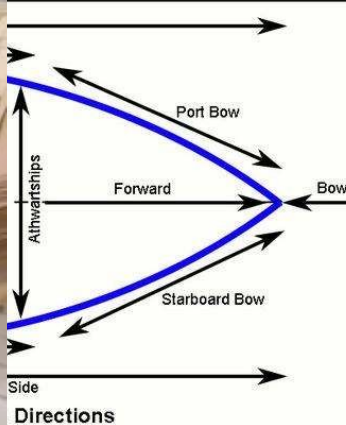


# Word Sense Disambiguation

He put the **port** on the ship.

He walked along the **port** of the steamer.

He walked along the **port** next to the steamer.



# Word Sense Disambiguation

He put the **port** on the ship.

He walked along the **port** of the steamer.

He walked along the **port** next to the steamer.

**port.n.1** (a place (seaport or airport) where people and merchandise can enter or leave a country)

**port.n.2** port wine (sweet dark-red dessert wine originally from Portugal)

# Word Sense Disambiguation

He put the **port** on the ship.

He walked along the **port** of the steamer.

He walked along the **port** next to the steamer.

**port.n.1** (a place (seaport or airport) where people and merchandise can enter or leave a country)

**port.n.2** port wine (sweet dark-red dessert wine originally from Portugal)

**port.n.3**, embrasure, porthole (an opening (in a wall or ship or armored vehicle) for firing through)

# Word Sense Disambiguation

He put the **port** on the ship.

He walked along the **port** of the steamer.

He walked along the **port** next to the steamer.

**port.n.1** (a place (seaport or airport) where people and merchandise can enter or leave a country)

**port.n.2** port wine (sweet dark-red dessert wine originally from Portugal)

**port.n.3**, embrasure, porthole (an opening (in a wall or ship or armored vehicle) for firing through)

larboard, **port.n.4** (the left side of a ship or aircraft to someone who is aboard and facing the bow or nose)

# Word Sense Disambiguation

He put the **port** on the ship.

He walked along the **port** of the steamer.

He walked along the **port** next to the steamer.

**port.n.1** (a place (seaport or airport) where people and merchandise can enter or leave a country)

**port.n.2** port wine (sweet dark-red dessert wine originally from Portugal)

**port.n.3**, embrasure, porthole (an opening (in a wall or ship or armored vehicle) for firing through)

larboard, **port.n.4** (the left side of a ship or aircraft to someone who is aboard and facing the bow or nose)

interface, **port.n.5** ((computer science) computer circuit consisting of the hardware and associated circuitry that links one device with another (especially a computer and a hard disk drive or other peripherals))

# Word Sense Disambiguation

He put the **port** on the ship.

He walked along the **port** of the steamer.

He walked along the **port** next to the steamer.

## As a verb...

1. **port** (put or turn on the left side, of a ship) "*port the helm*"
2. **port** (bring to port) "*the captain ported the ship at night*"
3. **port** (land at or reach a port) "*The ship finally ported*"
4. **port** (turn or go to the port or left side, of a ship) "*The big ship was slowly porting*"
5. **port** (carry, bear, convey, or bring) "*The small canoe could be ported easily*"
6. **port** (carry or hold with both hands diagonally across the body, especially of weapons) "*port a rifle*"
7. **port** (drink port) "*We were porting all in the club after dinner*"
8. **port** (modify (software) for use on a different machine or platform)

**port.n.1** (a place (seaport or airport) where people and merchandise can enter or leave a country)

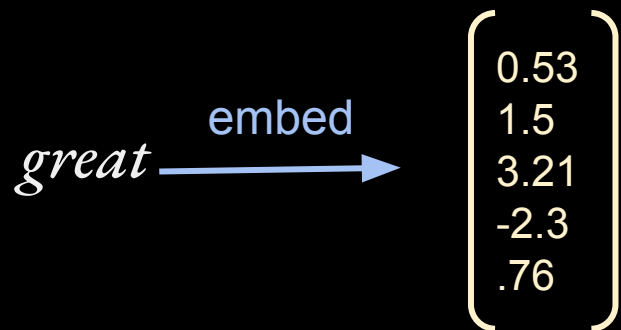
**port.n.2** port wine (sweet dark-red dessert wine originally from Portugal)

**port.n.3**, embrasure, porthole (an opening (in a wall or ship or armored vehicle) for firing through)

larboard, **port.n.4** (the left side of a ship or aircraft to someone who is aboard and facing the bow or nose)

interface, **port.n.5** ((computer science) computer circuit consisting of the hardware and associated circuitry that links one device with another (especially a computer and a hard disk drive or other peripherals))

# Objective





# Objective

*great* → embed

$$\begin{bmatrix} 0.53 \\ 1.5 \\ 3.21 \\ -2.3 \\ .76 \end{bmatrix}$$

?

**great.a.1** (relatively large in size or number or extent; larger than others of its kind)

**great.a.2**, outstanding (of major significance or importance)

**great.a.3** (remarkable or out of the ordinary in degree or magnitude or effect)

bang-up, bully, corking, cracking, dandy, **great.a.4**, groovy, keen, neat, nifty, not bad, peachy, slap-up, swell, smashing, old (very good)

capital, **great.a.5**, majuscule (uppercase)

big, enceinte, expectant, gravid, **great.a.6**, large, heavy, with child (in an advanced stage of pregnancy)

# Objective

*great* → embed

$$\begin{bmatrix} 0.53 \\ 1.5 \\ 3.21 \\ -2.3 \\ .76 \end{bmatrix}$$

?

**great.a.1** (relatively large in size or number or extent; larger than others of its kind)

**great.a.2**, outstanding (of major significance or importance)

**great.a.3** (remarkable or out of the ordinary in degree or magnitude or effect)

bang-up, bully, corking, cracking, dandy, **great.a.4**, groovy, keen, neat, nifty, not bad, peachy, slap-up, swell, smashing, old (very good)

capital, **great.a.5**, majuscule (uppercase)

big, enceinte, expectant, gravid, **great.a.6**, large, heavy, with child (in an advanced stage of pregnancy)

**great.n.1** (a person who has achieved distinction and honor in some field)

# Word Sense Disambiguation

A classification problem:

General Form:

$f(\text{sent\_tokens}, (\text{target\_index}, \text{lemma}, \text{POS})) \rightarrow \text{word\_sense}$

port.n.1  
port.n.2  
port.n.3,  
port.n.4  
port.n.5

He walked along the **port** next to the steamer.

# Word Sense Disambiguation

A classification problem:

General Form:

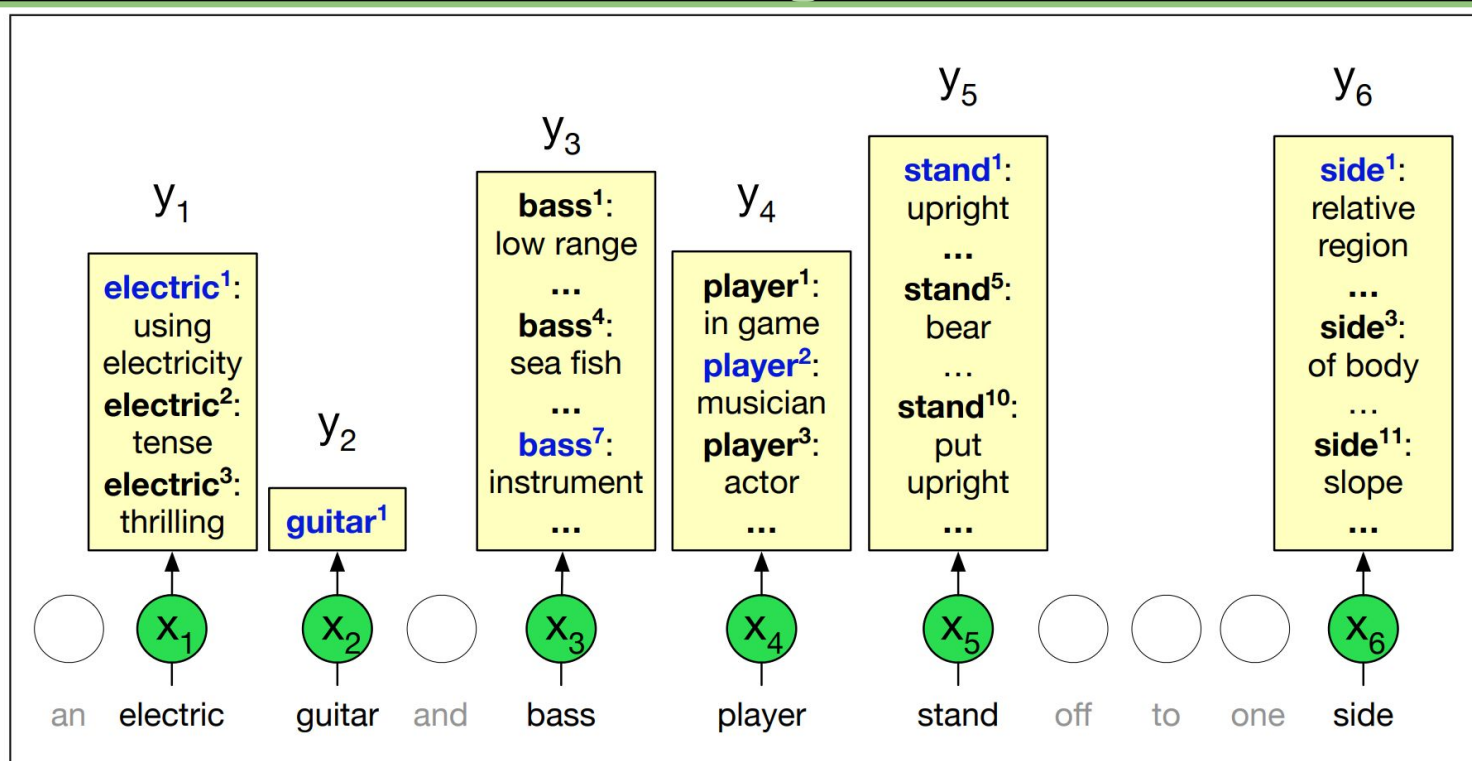
$$f(\text{sent\_tokens}, (\text{target\_index}, \text{lemma}, \text{POS})) \rightarrow \text{word\_sense}$$

Logistic Regression (or any discriminative classifier):

$$P_{\text{lemma}, \text{POS}}(\text{sense} = s \mid \text{features})$$

He walked along the **port** next to the steamer.

# Word Sense Disambiguation



**Figure 19.8** The all-words WSD task, mapping from input words ( $x$ ) to WordNet senses ( $y$ ). Only nouns, verbs, adjectives, and adverbs are mapped, and note that some words (like *guitar* in the example) only have one sense in WordNet. Figure inspired by [Chaplot and Salakhutdinov \(2018\)](#).

# Distributional Hypothesis:

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

# Distributional Hypothesis:

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

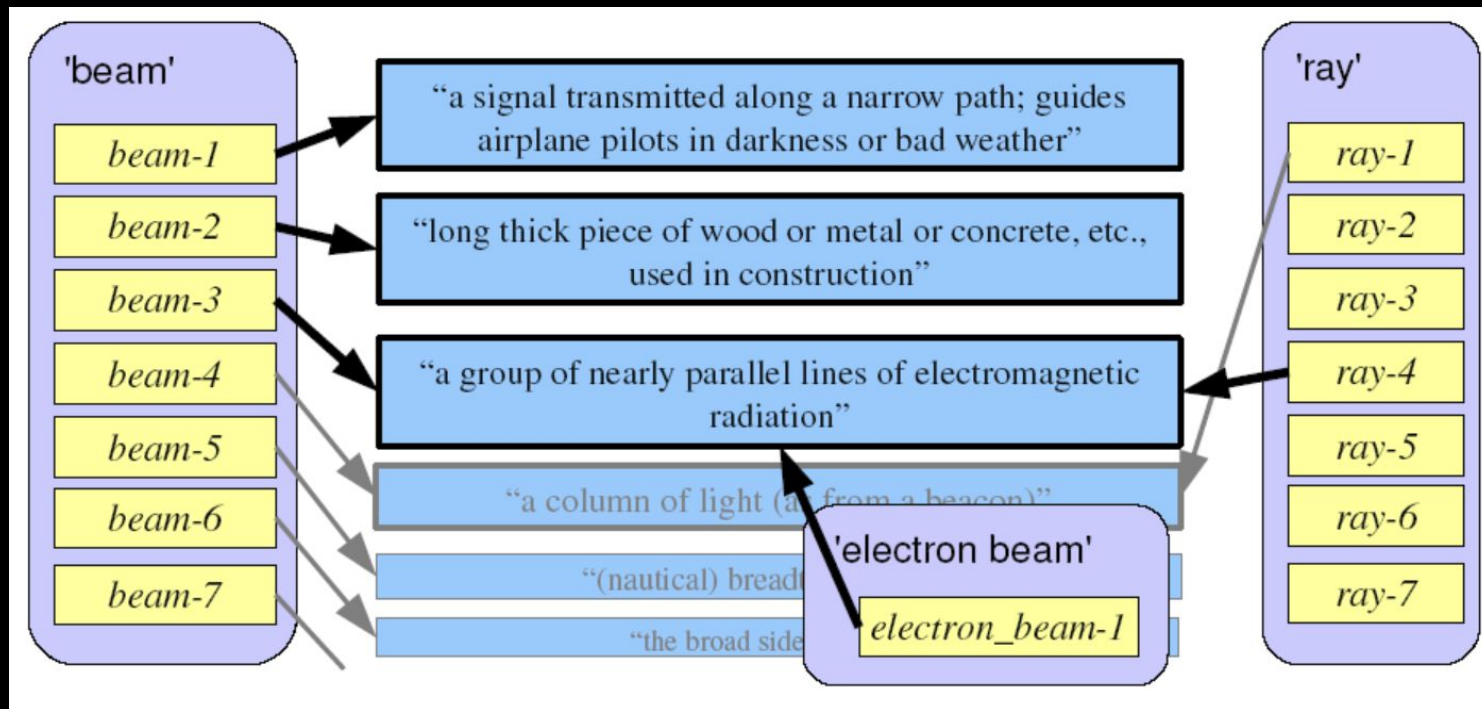
Distributional hypothesis -- A word's meaning is defined by all the different contexts it appears in (i.e. how it is “distributed” in natural language).

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.*



# Approaches to WSD

*I.e. how to operationalize the distributional hypothesis.*

1. Bag of words for context  
*E.g. multi-hot for any word in a defined “context”.*
2. Surrounding window with positions  
*E.g. one-hot per position relative to word).*
3. Lesk algorithm  
*E.g. compare context to sense definitions.*
4. Selectors -- other *target* words that appear with same context  
*E.g. counts for any selector.*
5. Contextual Embeddings  
*E.g. real valued vectors that “encode” the context (TBD).*

# Approaches to WSD

*I.e. how to operationalize the distributional hypothesis.*

1. Bag of words for context  
*E.g. multi-hot for any word in a defined “context”.*
2. Surrounding window with positions  
*E.g. one-hot per position relative to word).*
3. **Lesk algorithm**  
*E.g. compare context to sense definitions.*
4. Selectors -- other *target* words that appear with same context  
*E.g. counts for any selector.*
5. Contextual Embeddings  
*E.g. real valued vectors that “encode” the context (TBD).*

# Lesk Algorithm for WSD

**function** SIMPLIFIED LESK(*word*, *sentence*) **returns** best sense of *word*

*best-sense*  $\leftarrow$  most frequent sense for *word*

*max-overlap*  $\leftarrow 0$

*context*  $\leftarrow$  set of words in *sentence*

**for each** *sense* **in** senses of *word* **do**

*signature*  $\leftarrow$  set of words in the gloss and examples of *sense*

*overlap*  $\leftarrow$  COMPUTEOVERLAP(*signature*, *context*)

**if** *overlap*  $>$  *max-overlap* **then**

*max-overlap*  $\leftarrow$  *overlap*

*best-sense*  $\leftarrow$  *sense*

**end**

**return**(*best-sense*)

**Figure 19.10** The Simplified Lesk algorithm. The COMPUTEOVERLAP function returns the number of words in common between two sets, ignoring function words or other words on a stop list. The original Lesk algorithm defines the *context* in a more complex way.

# Lesk Algorithm for WSD

- bank.n.1 (sloping land (especially the slope beside a body of water)) "they pulled the canoe up on the bank"; "he sat on the bank of the river and watched the currents"
- bank.n.2 (a financial institution that accepts deposits and channels the money into lending activities) "he cashed a check at the bank"; "that bank holds the mortgage on my home"

```
overlap  $\leftarrow$  COMPUTEOVERLAP(signature, context)  
if overlap > max-overlap then  
    max-overlap  $\leftarrow$  overlap  
    best-sense  $\leftarrow$  sense  
end  
return(best-sense)
```

*The bank can guarantee deposits will cover future tuition costs, ...*

- bank.n.1 (sloping land (especially the slope beside a body of water)) "they pulled the canoe up on the bank"; "he sat on the bank of the river and watched the currents"
- bank.n.2 (a financial institution that accepts deposits and channels the money into lending activities) "he cashed a check at the bank"; "that bank holds the mortgage on my home"
- ...
- bank.n.4 (an arrangement of similar objects in a row or in tiers) "he operated a bank of switches"
- ...
- bank.n.8 (a building in which the business of banking transacted) "the bank is on the corner of Nassau and Witherspoon"
- bank.n.9 (a flight maneuver; aircraft tips laterally about its longitudinal axis (especially in turning)) "the plane went into a steep bank"

**end**

**return**(*best-sense*)

*The bank can guarantee deposits will cover future tuition costs, ...*

## Look-Algorithm for WSD

- **striker.n.1** (a forward on a soccer team)
- **striker.n.2** (someone receiving intensive training for a naval technical rating)
- **striker.n.3** (an employee on strike against an employer)
- **striker.n.4** (someone who hits) *"a hard hitter"; "a fine striker of the ball"; "blacksmiths are good hitters"*
- **striker.n.5** (the part of a mechanical device that strikes something)

*overlap*  $\leftarrow$  COMPUTEOVERLAP(*signature*, *context*)

**if** *overlap* > *max-overlap* **then**

*max-overlap*  $\leftarrow$  *overlap*

*best-sense*  $\leftarrow$  *sense*

**end**

**return**(*best-sense*)

*He addressed the strikers at the rally.*

# Approaches to WSD

*I.e. how to operationalize the distributional hypothesis.*

1. Bag of words for context  
*E.g. multi-hot for any word in a defined “context”.*
2. Surrounding window with positions  
*E.g. one-hot per position relative to word).*
3. **Lesk algorithm**  
*E.g. compare context to sense definitions.*
4. Selectors -- other *target* words that appear with same context  
*E.g. counts for any selector.*
5. Contextual Embeddings  
*E.g. real valued vectors that “encode” the context (TBD).*

# Approaches to WSD

*I.e. how to operationalize the distributional hypothesis.*

1. Bag of words for context  
*E.g. multi-hot for any word in a defined “context”.*
2. Surrounding window with positions  
*E.g. one-hot per position relative to word).*
3. Lesk algorithm  
*E.g. compare context to sense definitions.*
4. **Selectors -- other target words that appear with same context** *E.g. counts for any selector.*
5. Contextual Embeddings  
*E.g. real valued vectors that “encode” the context (TBD).*



# Selectors

... a word which can take the place of another given word within the same local context (Lin, 1997)

Original version: Local context defined by dependency parse

# Selectors

... a word which can take the place of another given word within the same local context (Lin, 1997)

Original version: Local context defined by dependency parse



*He addressed the strikers at the rally.*

# Selectors

... a word which can take the place of another given word within the same local context (Lin, 1997)

Original version: Local context defined by dependency parse (Lin, 1997)

Web version: Local context defined by lexical patterns matched on the Web (Schwartz, 2008).

*“He addressed the \* at the rally.”*

## Selectors

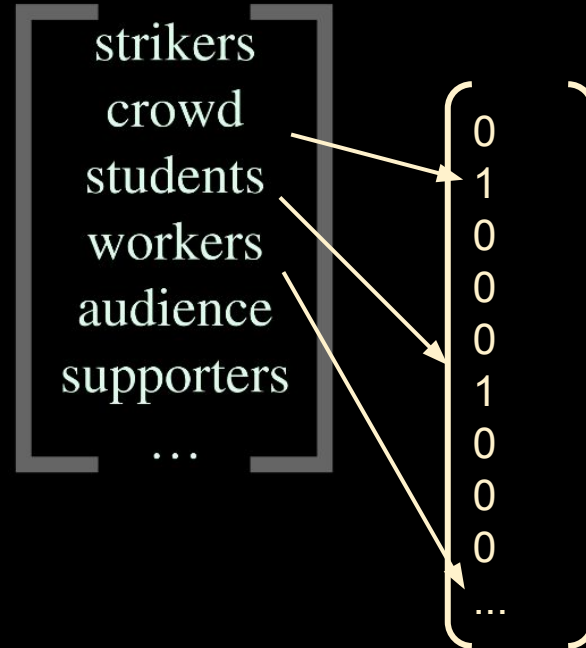
"He addressed the strikers at the rally."

strikers  
crowd  
students  
workers  
audience  
supporters

...

# Selectors

"He addressed the strikers at the rally."



# Selectors

"He addressed the strikers at the rally."

he  
man  
owners  
Mary  
...

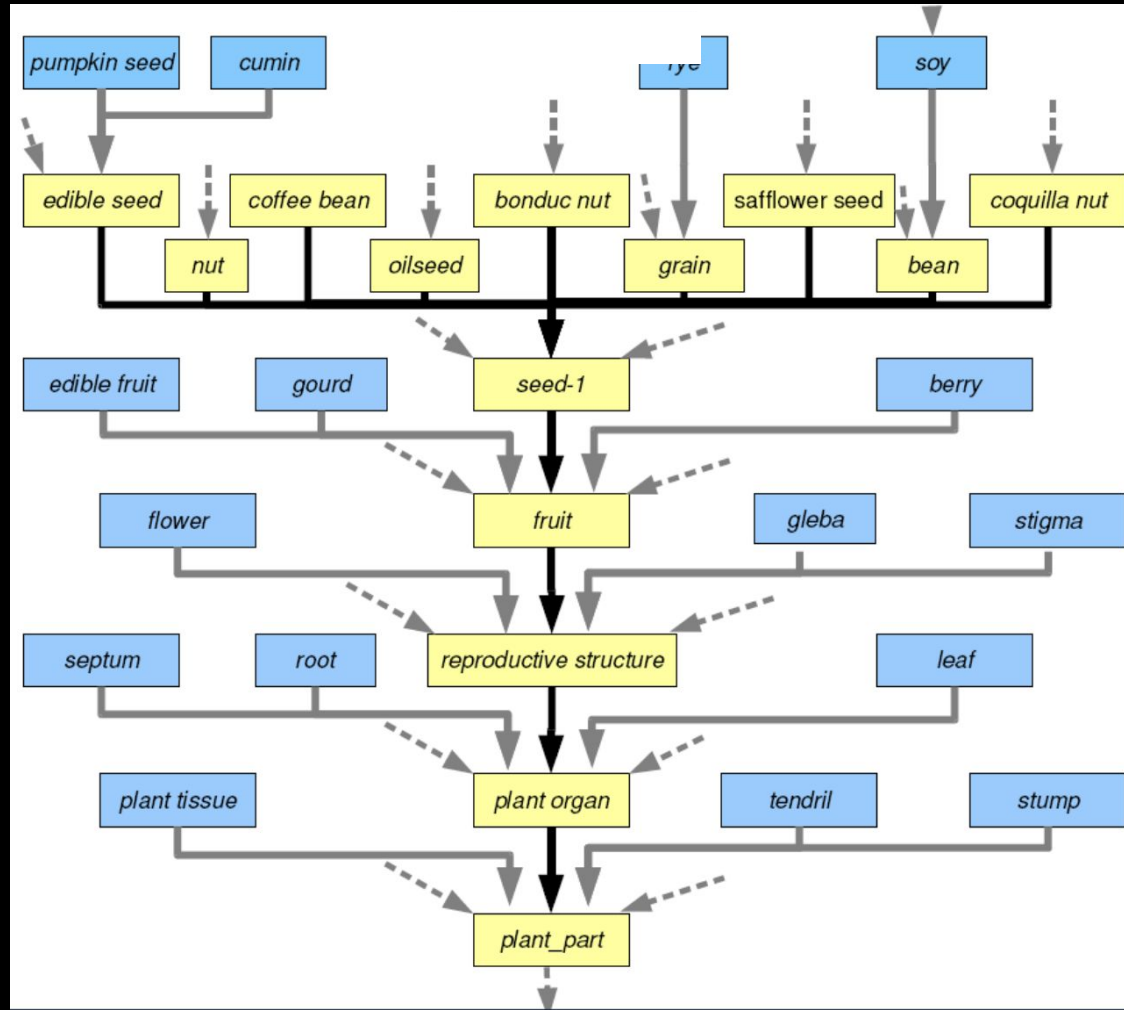
addressed  
scolded  
rallyied  
kept  
...

strikers  
crowd  
students  
workers  
audience  
supporters  
...

rally  
protest  
demonstration  
work  
stadium  
...

# Selectors

Leverages *hyponymy*:  
*concept1 <is-a> concept2*



# Why Are Selectors Effective?

Sets of selectors tend to vary extensively by word sense:

<i>bill-n.1</i>	<i>bill-n.2</i>	<i>bill-n.3</i>
bill	bill	<b>market</b>
it	<b>staff</b>	system
<b>legislation</b>	system	<b>paper</b>
system	<b>money</b>	<b>note</b>
<b>program</b>	<b>time</b>	bill
<b>law</b>	it	<b>bond</b>
<b>plan</b>	<b>tax</b>	<b>stock</b>
<b>you</b>	<b>work</b>	<b>debt</b>
<b>measure</b>	<b>rent</b>	<b>rate</b>
<b>project</b>	<b>tuition</b>	<b>report</b>

<i>occur-v.1</i>	<i>occur-v.2</i>	<i>occur-v.3</i>
be	go	go
<b>happen</b>	get	<b>look</b>
occur	Come	<b>break</b>
go	have	<b>remove</b>
<b>take</b>	<b>try</b>	<b>find</b>
work	<b>lead</b>	get
come	<b>listen</b>	<b>place</b>
<b>see</b>	work	<b>keep</b>
have	be	<b>stick</b>
<b>change</b>	<b>belong</b>	<b>stop</b>



- Polls show wide, generalized support for some vague concept of service, but the **bill** now under discussion lacks any passionate public backing.  
training set never contained: “but the \_ now under”
- ... in his lecture, refers to the “startling experience which almost every person confesses, that particular passages of conversation and action have **occurred** to him in the same order before, whether dreaming or waking ...  
small context is contradictory:  
“action have occurred” => occur-v.1 (“to happen or take place”)  
“occurred to him” => occur-v.2 (“to come to mind”)

<i>bill-n.1</i>	<i>bill-n.2</i>	<i>bill-n.3</i>
bill	bill	<b>market</b>
it	<b>staff</b>	system
<b>legislation</b>	system	<b>paper</b>
system	<b>money</b>	<b>note</b>
<b>program</b>	<b>time</b>	bill
law	it	<b>bond</b>
<b>plan</b>	<b>tax</b>	<b>stock</b>
you	work	<b>debt</b>
<b>measure</b>	<b>rent</b>	<b>rate</b>
<b>project</b>	<b>tuition</b>	<b>report</b>

<i>occur-v.1</i>	<i>occur-v.2</i>	<i>occur-v.3</i>
be	go	go
<b>happen</b>	get	<b>look</b>
occur	Come	<b>break</b>
go	have	<b>remove</b>
<b>take</b>	<b>try</b>	<b>find</b>
work	<b>lead</b>	get
come	<b>listen</b>	<b>place</b>
<b>see</b>	work	<b>keep</b>
have	be	<b>stick</b>
<b>change</b>	<b>belong</b>	<b>stop</b>

# Supervised Selectors

	<b>base</b>	<b>w/ sels</b>	<i>mfs</i>	<i>tests</i>
noun	87.9	<b>91.7</b>	80.9	2559
verb	83.3	<b>83.7</b>	76.5	2292
both	85.7	<b>87.9</b>	78.8	4851

**Accuracy over SemEval-2007: Task 17.**

# Vector Semantics

1. Word2vec
2. Topic Modeling - Latent Dirichlet Allocation (LDA)

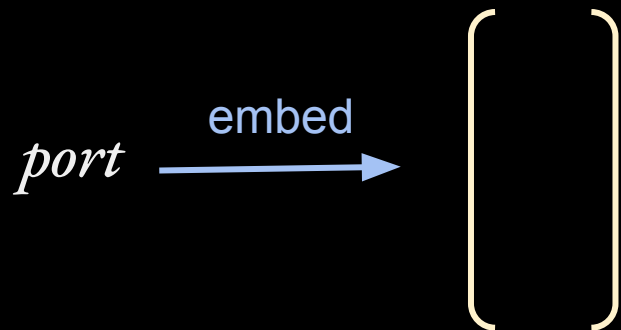
# Objective

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

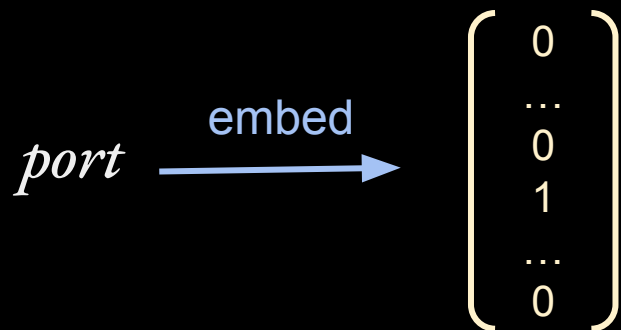
# Objective

To embed: convert a token (or sequence) to a vector that represents meaning, or is useful to perform downstream NLP application.

# Objective

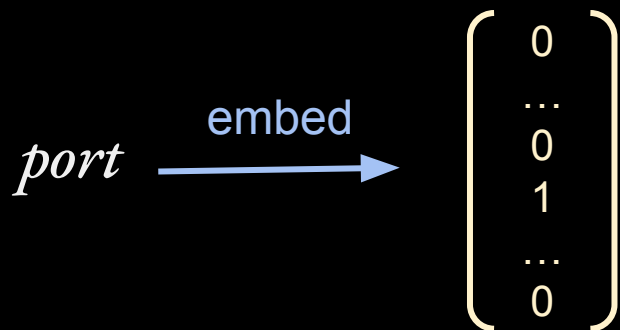


# Objective



# 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.

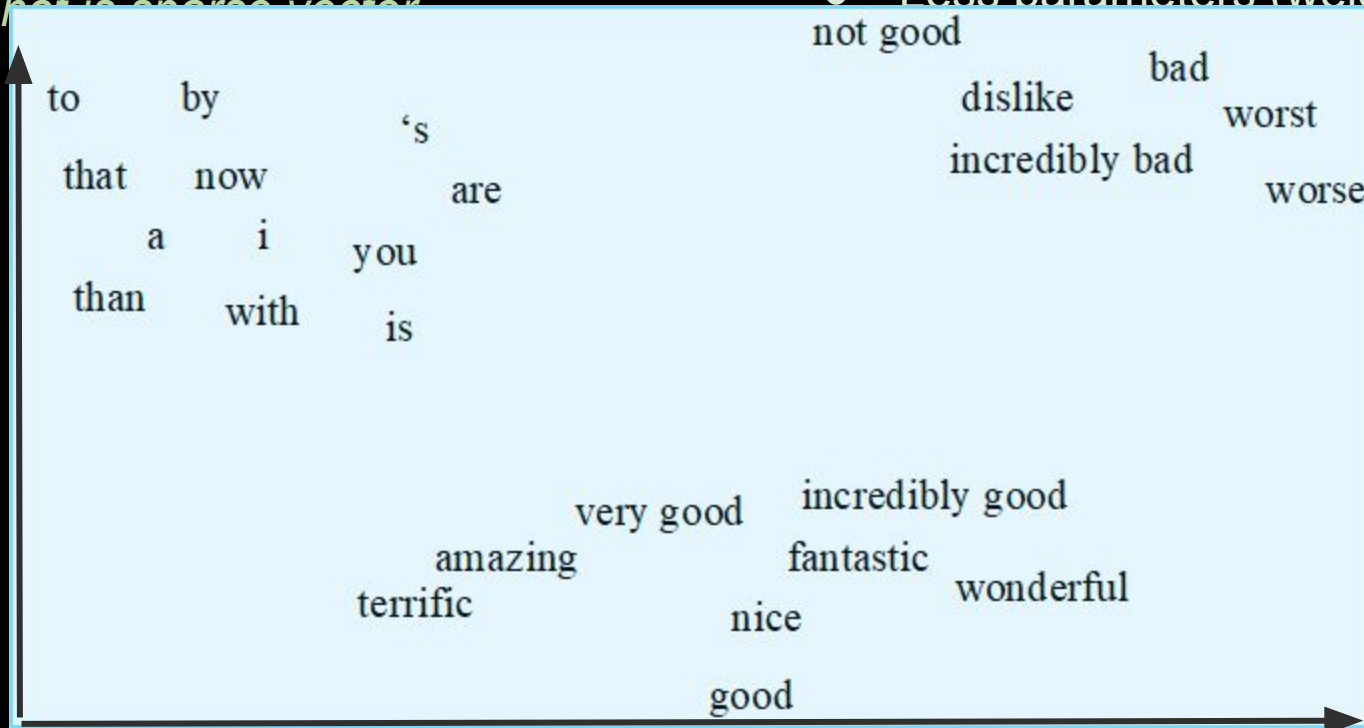


# Objective

## Prefer dense vectors

- Less parameters (weights) for

one-hot is sparse vector

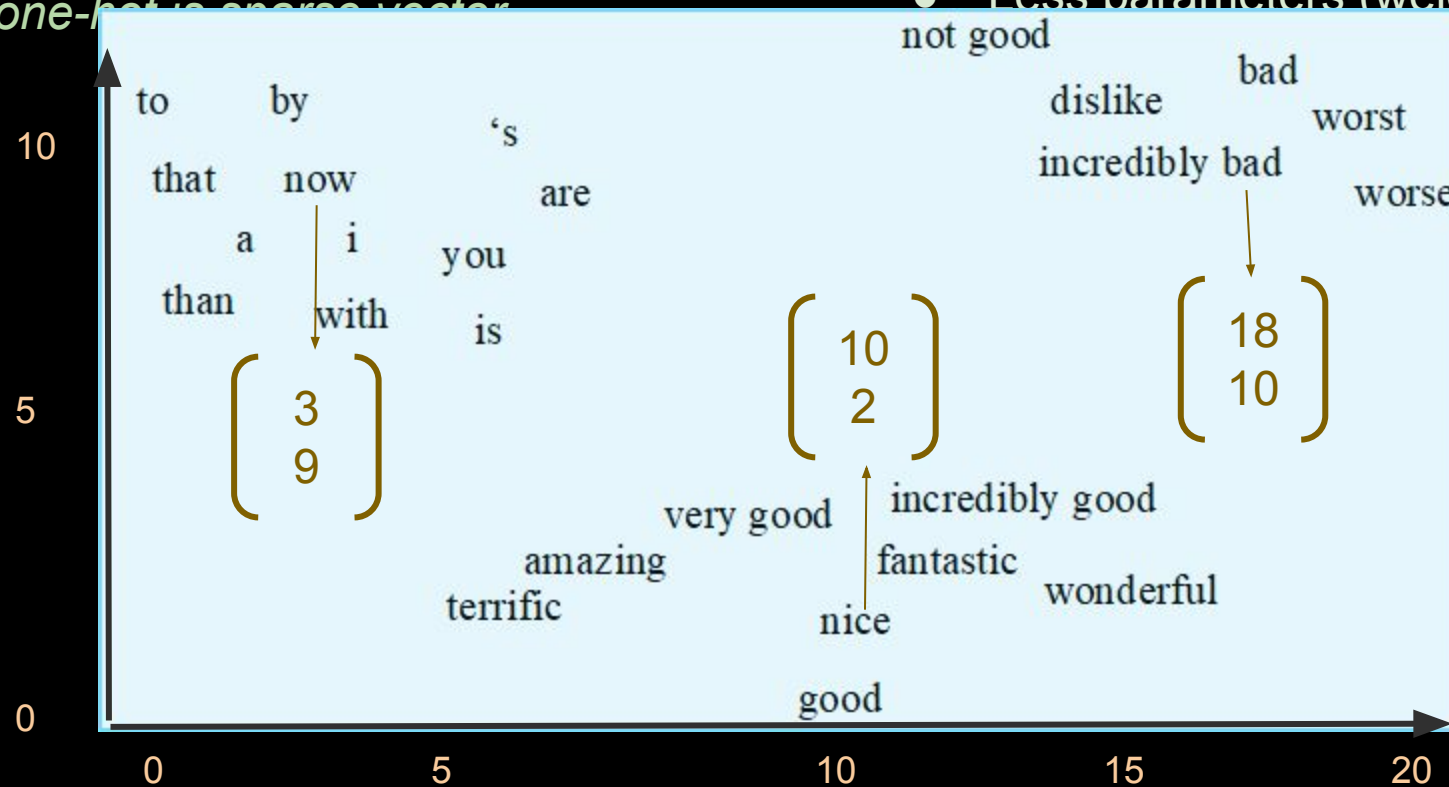


el.  
implicitly.

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).

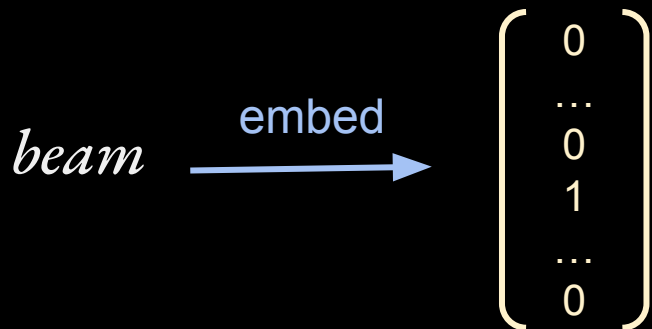
Firth, 1957: “*You shall know a word by the company it keeps*”

*The nail hit the beam behind the wall.*



# Word Vectors

*"one-hot encoding"*

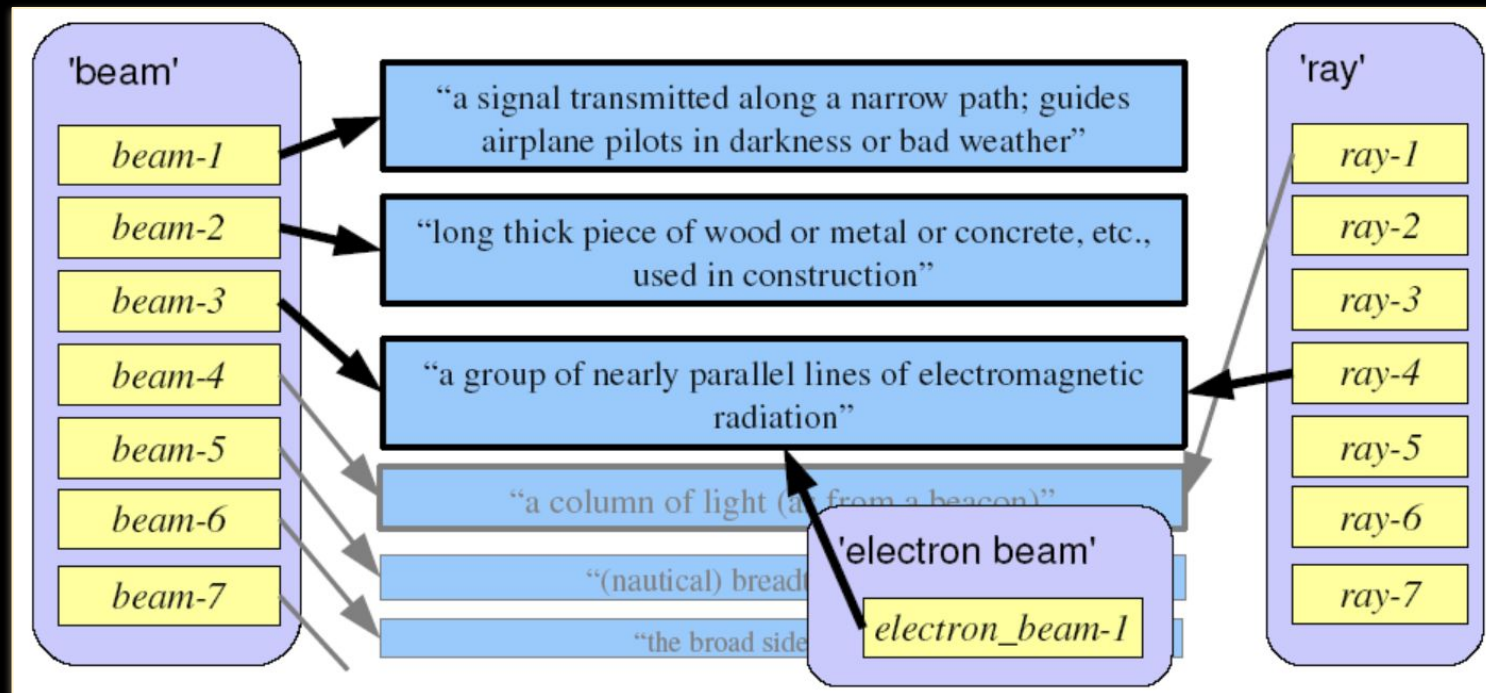


## Prefer dense vectors

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

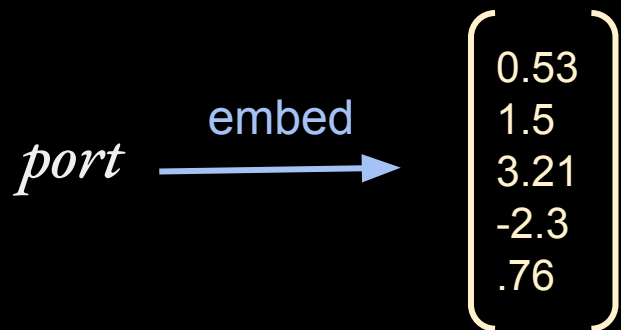
For deep learning, in practice, they work better. Why? Roughly, less parameters becomes increasingly important when you are learning multiple layers of weights rather than just a single layer.

# Distributional Hypothesis



*The nail hit the beam behind the wall.*

# Objective





# Objective

*port*

embed


$$\begin{bmatrix} 0.53 \\ 1.5 \\ 3.21 \\ -2.3 \\ .76 \end{bmatrix}$$

**port.n.1** (a place (seaport or airport) where people and merchandise can enter or leave a country)

**port.n.2** port wine (sweet dark-red dessert wine originally from Portugal)

**port.n.3**, embrasure, porthole (an opening (in a wall or ship or armored vehicle) for firing through)

larboard, **port.n.4** (the left side of a ship or aircraft to someone who is aboard and facing the bow or nose)

interface, **port.n.5** ((computer science) computer circuit consisting of the hardware and associated circuitry that links one device with another (especially a computer and a hard disk drive or other peripherals))

# Objective

*great* → embed

$$\begin{bmatrix} -0.2 \\ 0.3 \\ -1.1 \\ -2.1 \\ .26 \end{bmatrix}$$

?

**great.a.1** (relatively large in size or number or extent; larger than others of its kind)

**great.a.2**, outstanding (of major significance or importance)

**great.a.3** (remarkable or out of the ordinary in degree or magnitude or effect)

bang-up, bully, corking, cracking, dandy, **great.a.4**, groovy, keen, neat, nifty, not bad, peachy, slap-up, swell, smashing, old (very good)

capital, **great.a.5**, majuscule (uppercase)

big, enceinte, expectant, gravid, **great.a.6**, large, heavy, with child (in an advanced stage of pregnancy)

**great.n.1** (a person who has achieved distinction and honor in some field)

# Word2Vec

Principal: Predict missing word.

Similar to classification where  $y = \text{context}$  and  $x = \text{word}$ .

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

# Word2Vec

Principal: Predict missing word.

Similar to classification where  $y = \text{context}$  and  $x = \text{word}$ .

To learn, maximize

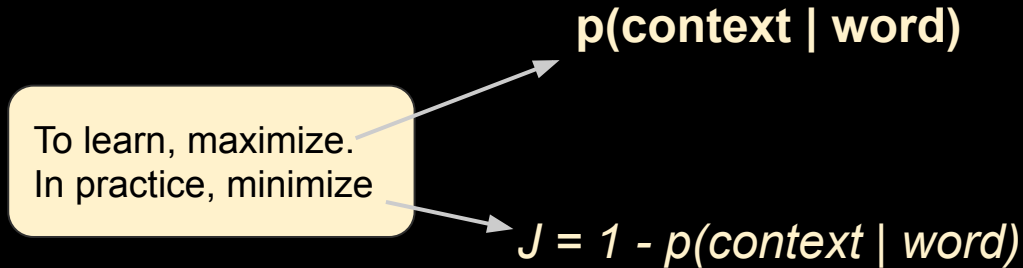
$p(\text{context} \mid \text{word})$

A yellow rectangular box containing the text "To learn, maximize" has a yellow arrow pointing from its right side to the expression  $p(\text{context} \mid \text{word})$ .

# Word2Vec

Principal: Predict missing word.

Similar to classification where  $y = \text{context}$  and  $x = \text{word}$ .



# Word2Vec: Context

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

2 Versions of Context:

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

# Word2Vec: Context

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

2 Versions of Context:

1. Continuous bag of words (CBOW): Predict word from context
2. **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

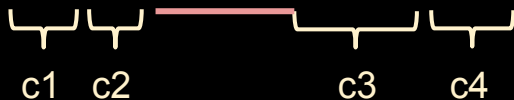
# Skip-Grams (SG): predict context words from target

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

## Steps:

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.*





# Skip-Grams (SG): predict context words from target

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

## Steps:

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

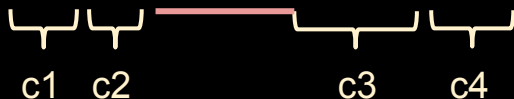
$x = (\text{hit}, \text{beam}), y = 1$

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

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

...

*The nail hit the beam behind the wall.*



# Skip-Grams (SG): predict context words from target

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

## Steps:

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

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

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

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

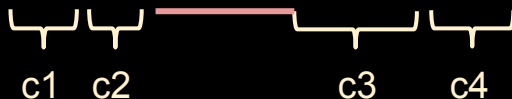
...

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

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

...

*The nail hit the beam behind the wall.*



# Skip-Grams (SG): predict context words from target

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

## Steps:

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

$x = (\text{hit}, \text{beam}), y = 1$        $k$  negative samples ( $y=0$ ) for every positive.

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

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

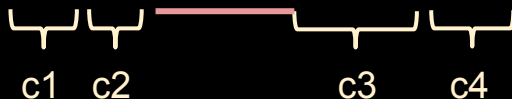
...

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

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

...

*The nail hit the beam behind the wall.*



# Skip-Grams (SG): predict context words from target

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

## Steps:

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

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

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

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

...

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

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

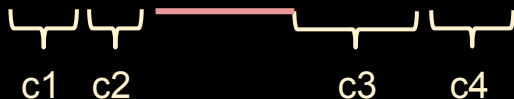
...

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

**How?** Randomly draw from unigram distribution

$$P(w) = \frac{\text{count}(w)}{\sum_w \text{count}(w)}$$

*The nail hit the beam behind the wall.*



# Skip-Grams

## Steps:

1. Treat the target word as a point in a vector space.
2. Randomly sample words from the vocabulary.
3. Use logistic regression to predict the target word from the sampled words.
4. Use the weights to update the word embeddings.



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

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

**How?** Randomly draw from unigram distribution,  $\alpha_{\text{adjusted}}$

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

*The nail hit the beam behind the wall.*

$\underbrace{\quad\quad}_{c1} \underbrace{\quad\quad}_{c2} \underbrace{\quad\quad}_{c3} \underbrace{\quad\quad}_{c4}$

where  
 $\alpha = 0.75$

# Skip-Grams (SG): predict context words from target

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

## Steps:

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

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

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

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

...

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

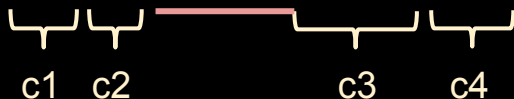
$x = (\text{think}, \text{beam}), y = 0$

...

single context:

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

*The nail hit the beam behind the wall.*



# Skip-Grams (SG): predict context words from target

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

## Steps:

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

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

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

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

...

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

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

...

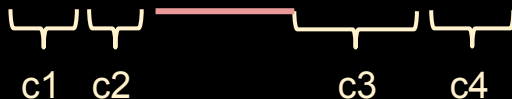
single context:

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

all contexts

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

*The nail hit the beam behind the wall.*



# Skip-Grams (SG): predict context words from target

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

## Steps:

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

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

single context:

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

all contexts

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

*nail hit the beam behind the wall.*

$\underbrace{\quad} \underbrace{\quad} \underbrace{\quad} \underbrace{\quad}$   
c1 c2 c3 c4



# Skip-Grams (SG): predict context words from target

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

## Steps:

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

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

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

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

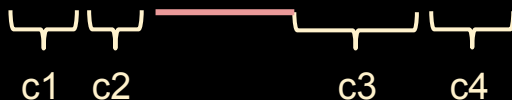
single context:

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

all contexts

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

*nail hit the beam behind the wall.*



# Skip-Grams (SG): predict context words from target

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

## Steps:

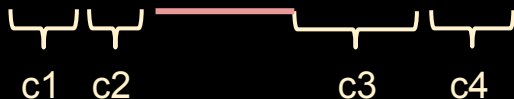
1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

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

all contexts

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

*The nail hit the beam behind the wall.*



# Skip-Grams (SG): predict context words from target

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

## Steps:

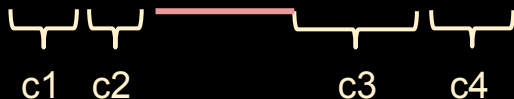
1. Treat the target word and a neighboring context word as positive examples.
2. Randomly sample other words in the lexicon to get negative samples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the **weights** as the embeddings

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

all contexts

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

*The nail hit the beam behind the wall.*



# Skip-Grams (SG): predict context words from target

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

## Steps:

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

all contexts

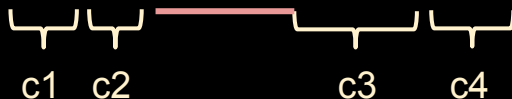
$$P(y=1 \mid \mathbf{c}, \mathbf{t}) = \prod_{i=1}^{\kappa} \frac{1}{1 + e^{-\mathbf{t} \cdot \mathbf{c}_i}}$$

x = (hit, beam), y = 1  
x = (the, beam), y = 1  
x = (behind, beam), y = 1

3a. assume  $dim * |vocab|$  weights for each of c and t,

...  
x = (happy, beam), y = 0  
x = (think, beam), y = 0

*The nail hit the beam behind the wall.*



# Skip-Grams (SG): predict context words from target

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

## Steps:

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

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  
...

all contexts

$$P(y=1 \mid \mathbf{c}, \mathbf{t}) = \prod_{i=1}^{\kappa} \frac{1}{1 + e^{-\mathbf{t} \cdot \mathbf{c}_i}}$$

- 3a. assume  $\text{dim} * |\text{vocab}|$  weights for each of c and t,  
initialized to random values (e.g.  $\text{dim} = 50$  or  $\text{dim} = 300$ )
- 3b.

# Skip-Grams (SG): predict context words from target

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

## Steps:

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

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  
...

all contexts

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

3a. assume  $dim * |\text{vocab}|$  weights for each of c and t,  
initialized to random values (e.g.  $dim = 50$  or  $dim = 300$ )

3b. optimize loss:

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

# Skip-Grams (SG): predict context words from target

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

## Steps:

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

all contexts

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

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  
...

3a. assume  $dim * |\text{vocab}|$  weights for each of c and t,  
initialized to random values (e.g.  $dim = 50$  or  $dim = 300$ )

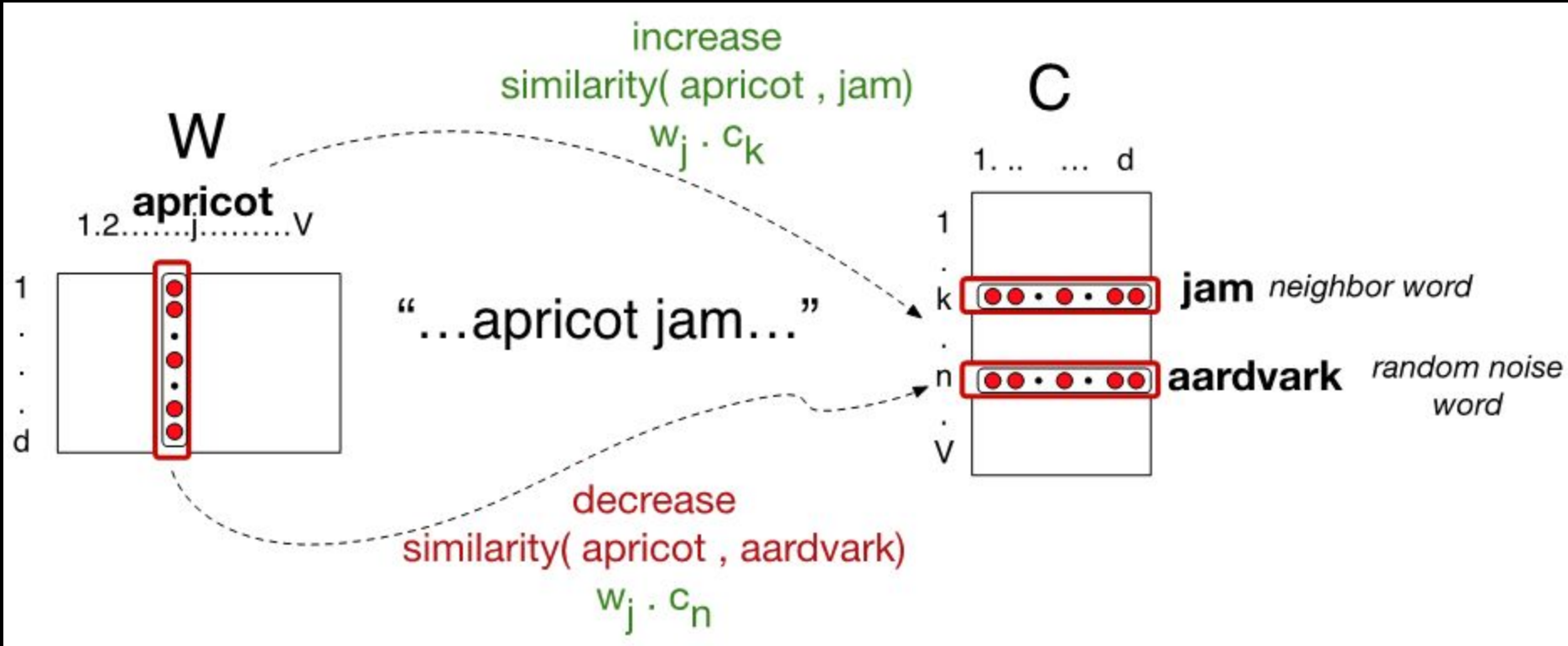
3b. optimize loss:

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

*Maximizes similarity of (c, t) in positive data (y = 1)*

*Minimizes similarity of (c, t) in negative data (y = 0)*

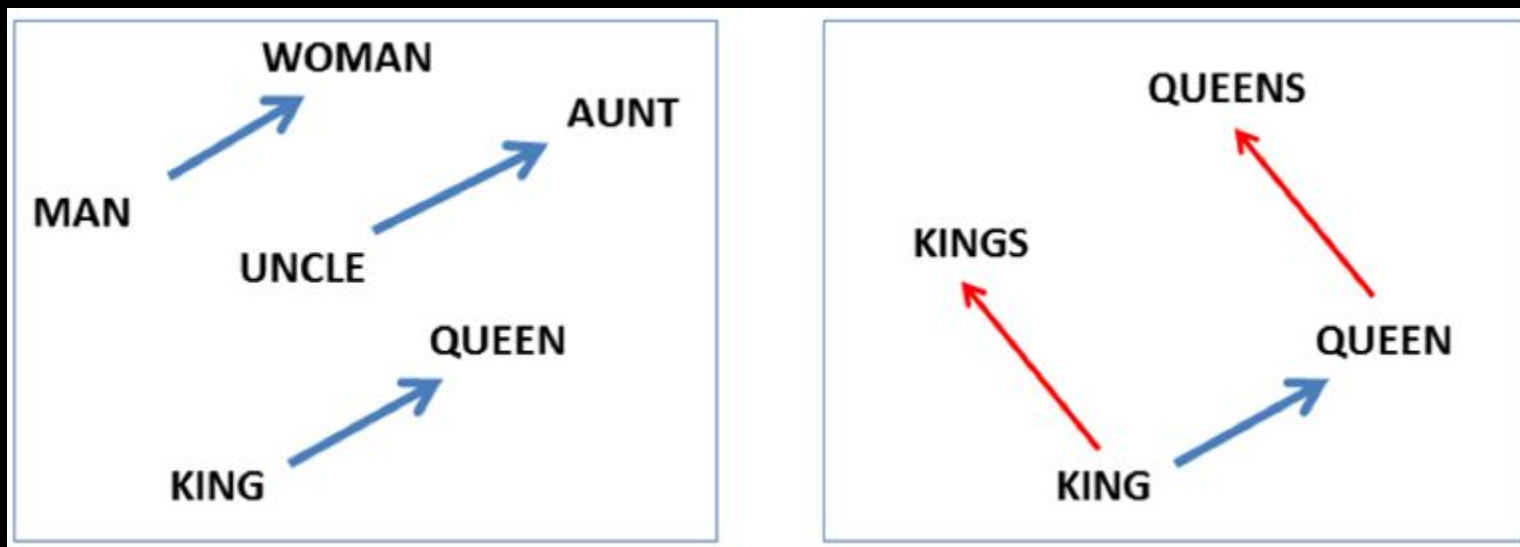
# Word 2 Vec

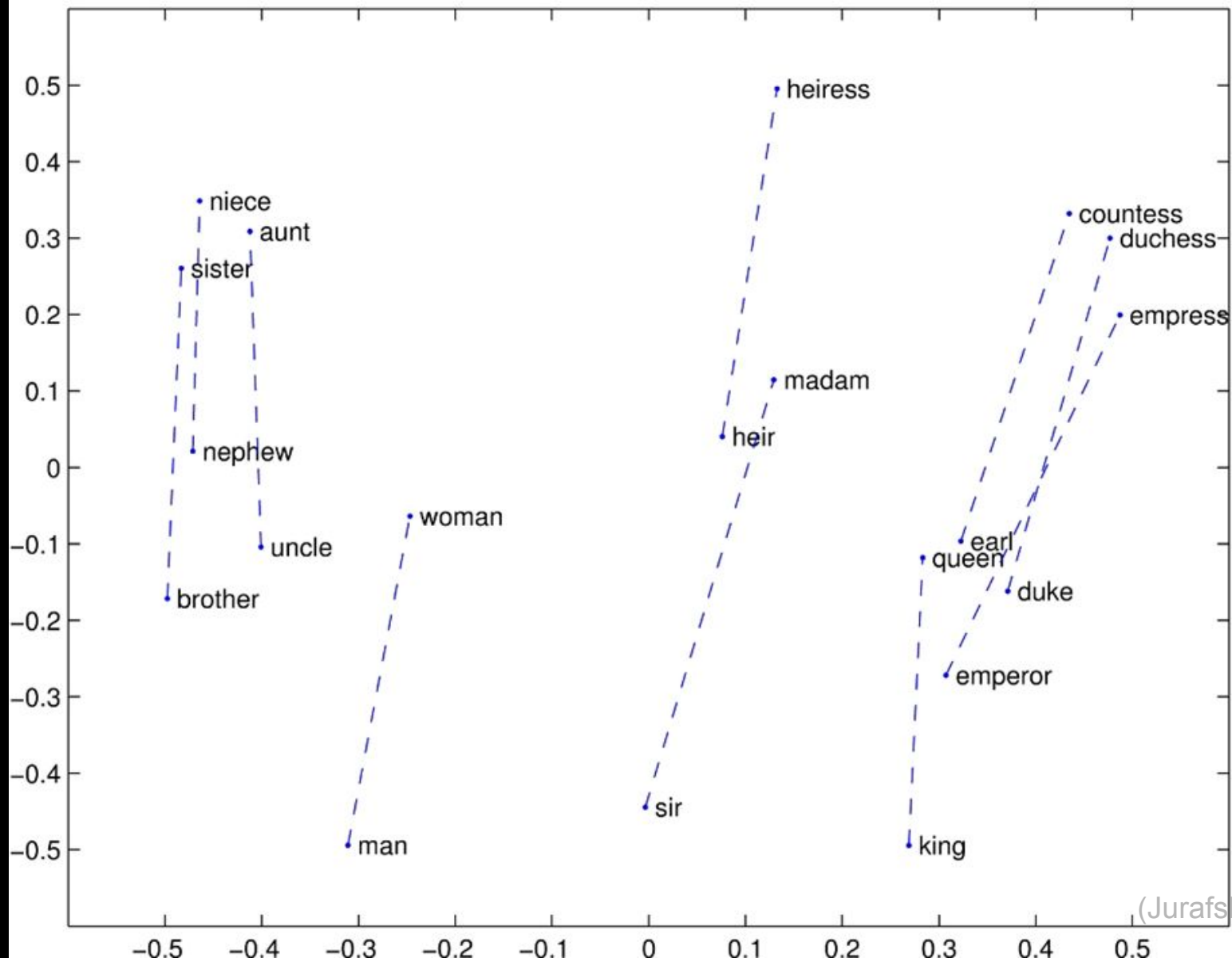


$$\sum_{(c,t)} (y) \log P(y = 1|c, t) + (y - 1) \log P(y = 0|c, t)$$

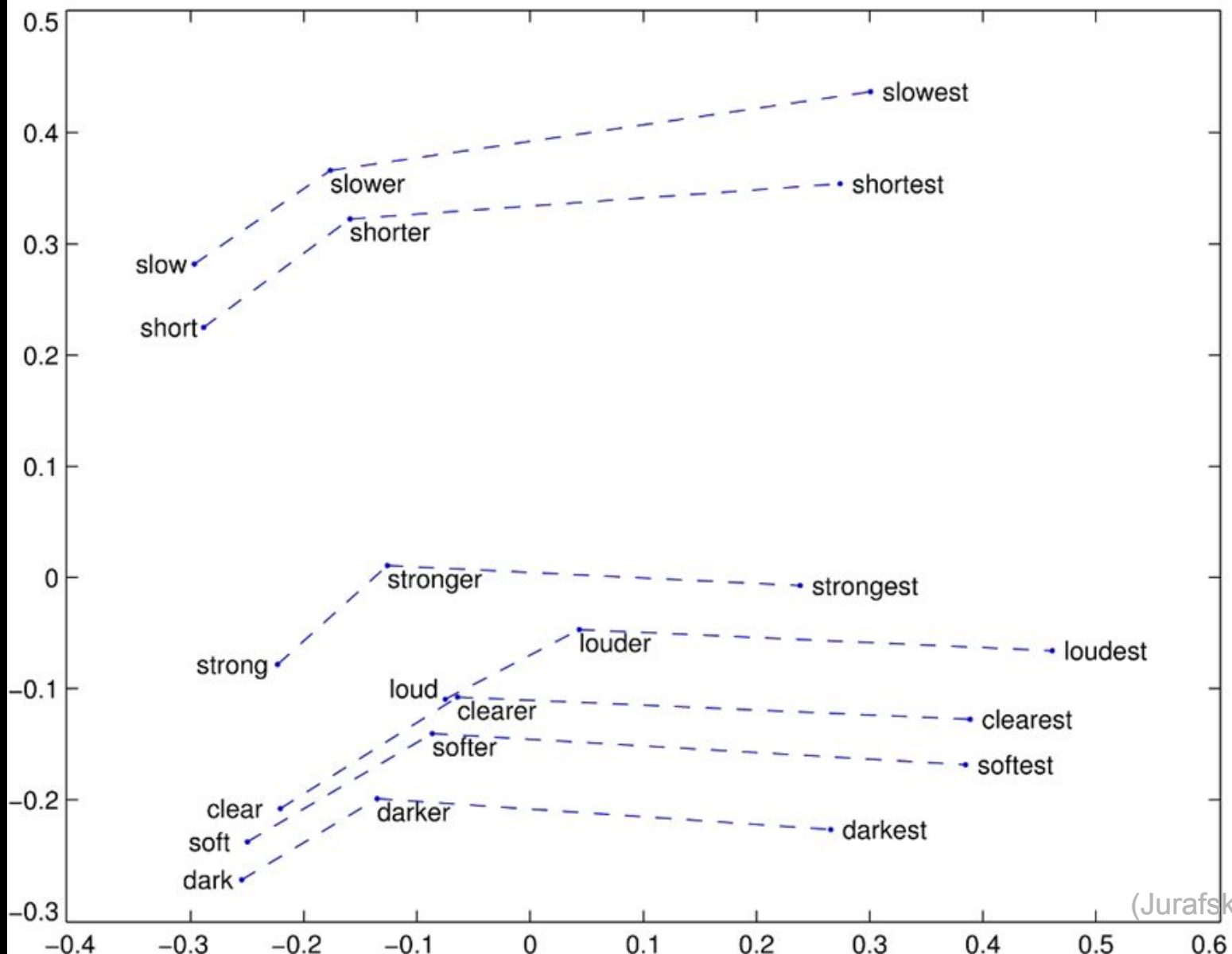


# Word2Vec captures analogies (kind of)





(Jurafsky, 2017)



(Jurafsky, 2017)

# Word2Vec: Quantitative Evaluations

1. Compare to manually annotated pairs of words: WordSim-353 (Finkelstein et al., 2002)
2. Compare to words in context (Huang et al., 2012)
3. Answer [TOEFL synonym questions](#).

# Word2Vec: Quantitative Evaluations

1. Compare to manually annotated pairs of words: WordSim-353 (Finkelstein et al., 2002)
2. Compare to words in context (Huang et al., 2012)
3. Answer [TOEFL synonym questions](#).

**What have we learned since Word2vec?** (a lot, but here are 2 important points)

1. Improved loss function: GLOVE embeddings (Pennington et al., 2014)
2. Word2Vec itself performs very similarly to PCA on a co-occurrence matrix ("LSA" Deerwater et al., 1988 – a much much older techniques!).

# Same multi-class loss function as LogReg!

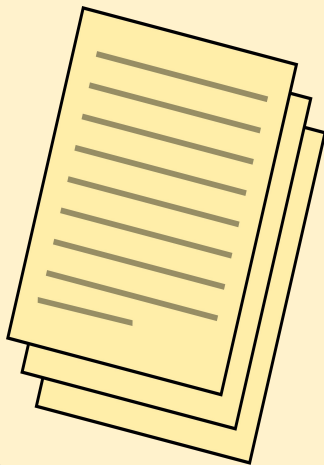
Logistic Regression Likelihood:  $L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^N p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$

Log Likelihood:  $\ell(\beta) = \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p(x_i))$

Log Loss:  $J(\beta) = -\frac{1}{N} \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p(x_i))$

Cross-Entropy Cost:  $J = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j} \log p(x_{i,j})$  (a “multiclass” log loss)

In vector algebra form:  $-\text{mean}(\text{sum}(y * \log(y\_pred)))$



2. Imagine that you are now 15 years old. Write about the life you are leading: your interests, your home life and your work in the age of 15. (You have 10 minutes to do this).

I am 15 years of age I work in [redacted]  
[redacted] that is a garage in [redacted]  
roads I live in [redacted] road number  
90, we do not have a lot of  
work getting there. I go by my mind  
copper S. with all my home made  
breads and lighted notes. when I  
get to work it is about 8 o'clock, I  
get out of bed about 5 to 8 clock, get  
dressed and have a glass of milk and  
go to work. after work I go  
to the car, and have a salad and  
a tomato. then I go down the



# Topic Modeling

**Topic:** A group of highly related words and phrases. (aka "semantic field")

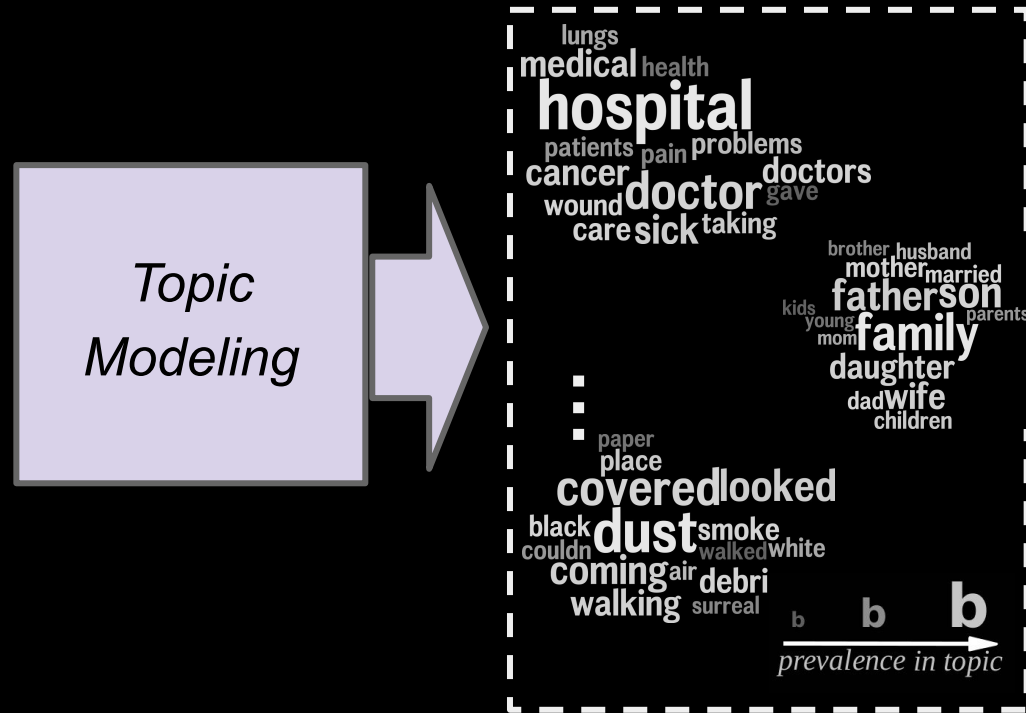
example: from WTC responder interviews  
(Son et al., 2021)





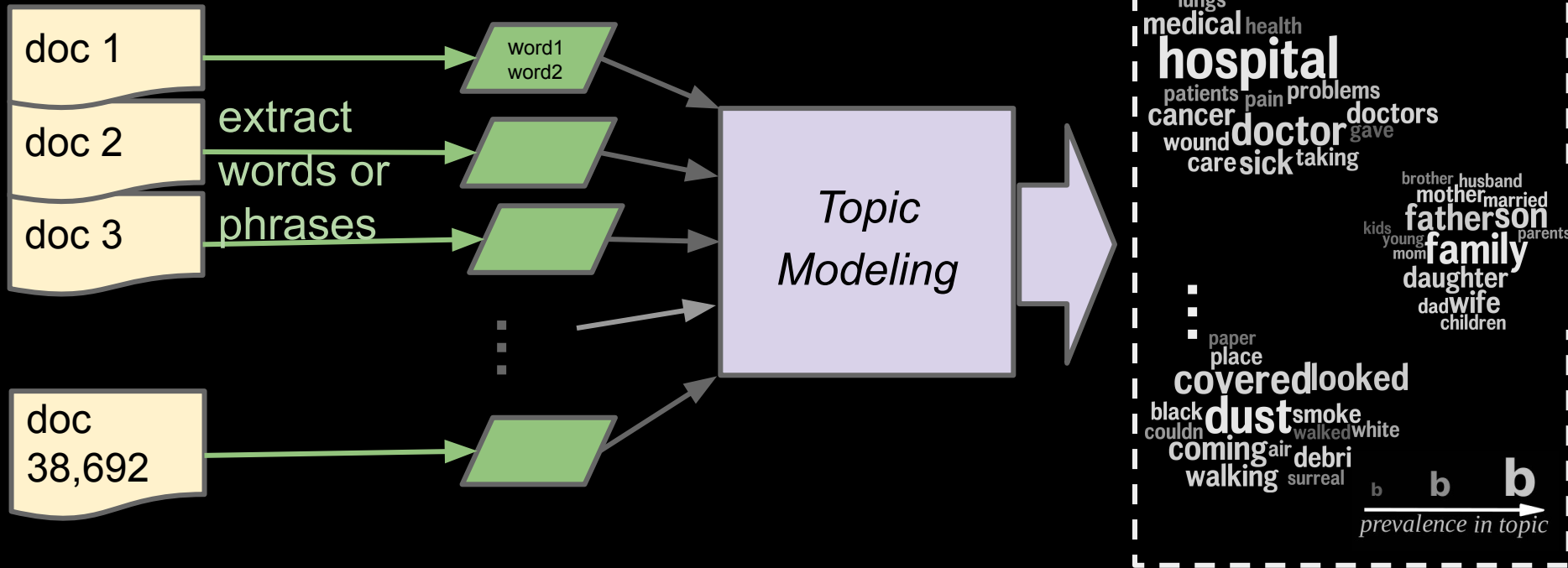
# Topic Modeling

*Topic*: A group of highly related words and phrases. (aka "semantic field")



# Topic Modeling

*Topic*: A group of highly related words and phrases. (aka "semantic field")



# Select Example Topics

**b**      **b**      **b**  
→  
*prevalence in topic*

fear  
sense  
thought  
feeling  
felt  
difficult  
overwhelming  
moment  
little\_bit  
happening  
kind  
thinking  
hard  
feel  
mind

big  
fires  
firemen  
fire  
guys  
burned  
site  
truck  
truck  
equipment  
firehouse  
vehicles  
burning  
smoke  
smell

lived  
moved  
queens  
place  
new\_york  
manhattan  
living  
live  
grew  
long\_island  
bronx  
city  
house  
brooklyn  
born

lungs  
health  
medical  
hospital  
patients  
pain  
problems  
cancer  
doctors  
wound  
doctor  
care  
sick  
taking

brother  
husband  
mother  
married  
father  
son  
kids  
young  
mom  
family  
daughter  
dad  
wife  
children

house  
call  
wife  
calls  
calling  
phone  
working  
called  
cell\_phone  
number  
home  
office  
contact  
told  
touch

paper  
place  
covered  
looked  
black  
dust  
smoke  
couldn  
coming  
air  
debris  
walking  
surreal

years  
months  
money  
retired  
half  
end  
year  
ten  
ve  
9/11

lights  
bridge  
driving  
traffic  
stopped  
cars  
manhattan  
precinct  
road  
city  
drive  
drove  
police  
car  
bus

running  
collapsed  
falling  
fall  
floor  
lobby  
buildings  
tower  
ran  
fell  
standing  
coming  
collapse  
towers  
building

# Generating Topics from Documents

- *Latent Dirichlet Allocation* -- a Bayesian probabilistic model where by words which appear in similar *contexts* (i.e. in essays that have similar sets of words) will be clustered into a prespecified number of topics.
- Rule of thumb:  $|\text{topics}| = \frac{|\text{observations}|}{100}$
- Each document receives a score per topic -- a probability:  $p(\text{topic}|\text{doc})$ .

## Doc 1

topic 1: .05

topic 2: .02

topic 3: .01

...

topic 100: .07

## Doc 2

topic 1: .03

topic 2: .01

topic 3: .03

...

topic 100: .05

## Doc 3

topic 1: .04

topic 2: .03

topic 3: .03

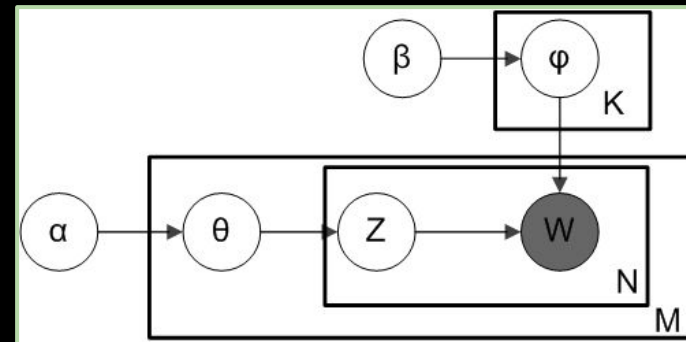
...

topic 100: .06

# Latent Dirichlet Allocation

(Blei et al., 2003)

- LDA specifies a Bayesian probabilistic model where by documents are viewed as a distribution of topics, and topics are a distribution of words.



## Observed:

$W$  -- observed word in document  $m$

## Inferred:

$\theta$  -- topic distribution for document  $m$ ,

$Z$  -- topic for word  $n$  in document  $m$

$\phi$  -- word distribution for topic  $k$

## Priors

$\alpha$  -- parameter for Dirichlet prior on the topics per document.

$\beta$  -- parameter for Dirichlet prior on the words per topic.

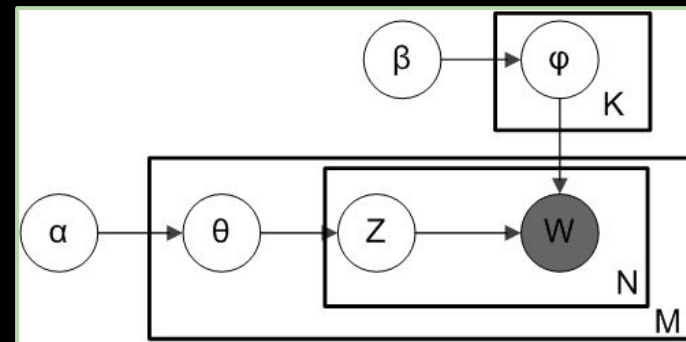
# Latent Dirichlet Allocation

(Blei et al., 2003)

- LDA specifies a Bayesian probabilistic model where by documents are viewed as a distribution of topics, and topics are a distribution of words.
- How to estimate (i.e. fit) the model parameters given data and priors? Common choices:
  - Gibb's Sampling (best)
  - variational Bayesian Inference (fastest).
- Key Output: the "posterior"  $\varphi = p(\text{word} \mid \text{topic})$ , the probability of a word given a topic.

From this and  $p(\text{topic})$ , we can get:  $p(\text{topic} \mid \text{word})$

$$p(\text{topic} \mid \text{doc}) = \sum_{\text{word} \in \text{topic}} p(\text{topic} \mid \text{word}) p(\text{word} \mid \text{doc})$$



## Observed:

$\mathbf{W}$  -- observed word in document  $m$

## Inferred:

$\theta$  -- topic distribution for document  $m$ ,

$\mathbf{Z}$  -- topic for word  $n$  in document  $m$

$\varphi$  -- word distribution for topic  $k$

## Priors

$\alpha$  -- parameter for Dirichlet prior on the topics per document.

$\beta$  -- parameter for Dirichlet prior on the words per topic.

# Example

Most prevalent words for 4 topics are listed at the top and words associated with them from a Yelp review are colored accordingly below.

Ranard, B.L., Werner, R.M., Antanavicius, T., Schwartz, H.A., Smith, R.J., Meisel, Z.F., Asch, D.A., Ungar, L.H. & Merchant, R.M. (2016). Yelp Reviews Of Hospital Care Can Supplement And Inform Traditional Surveys Of The Patient Experience Of Care. Health Affairs, 35(4), 697-705.

Labor and Delivery	Patient treatment	Surgery/ procedure and peri-op	Insurance and Billing
Baby Birth	Care Staff	Surgery Procedure	Insurance Billing
Nurses	Nurses	Surgeon	Bill
Labor	Hospital	Recovery	Hospital
Delivery	Doctors	Day	Department
Experience	Great	Staff	Company
Nurse	Caring	Experience	Paid

It depends what you look for in a **hospital**. Remember that this is a teaching **hospital** so you must adjust your **expectations** accordingly. This means many students who, **bleed** their **hearts**, may ask you the same **questions** again and again. I waited for hours on standby to **deliver** my **baby** by emergency **c-section**. The kind **nurses** who **served** me during **recovery** and the **anesthesiologist** on **duty** during my **surgery** **deserve** praise. My **OB** was very competent, but I wish he were willing to do an extraversion or at least given me an **epidural**. Im grateful they ultimately did what was best for my kid. However, I think **things** could have happened a lot more **smoothly** with better **pain** control. The only other thing to watch out for is your bills. This is the only institution I have been to that bills me **prior** to **billing insurance**. I **fought** two years to **claim** a **credit** through a **database** system **change**. The cafeteria gets flack for being all vegetarian but you just have to know what to order. **Stay** there for 1-2 **weeks** and you get the **hang** of whats good and whats not.

# Topic Modeling Packages

Most Reliable: [Mallet](#) (Java; uses Gibb's Sampling),  
pymallet (slower than Mallet but high quality results)

Ease of use: [Gensim](#) (python; uses variational inference;  
implements word2vec as well)



# Topic Modeling

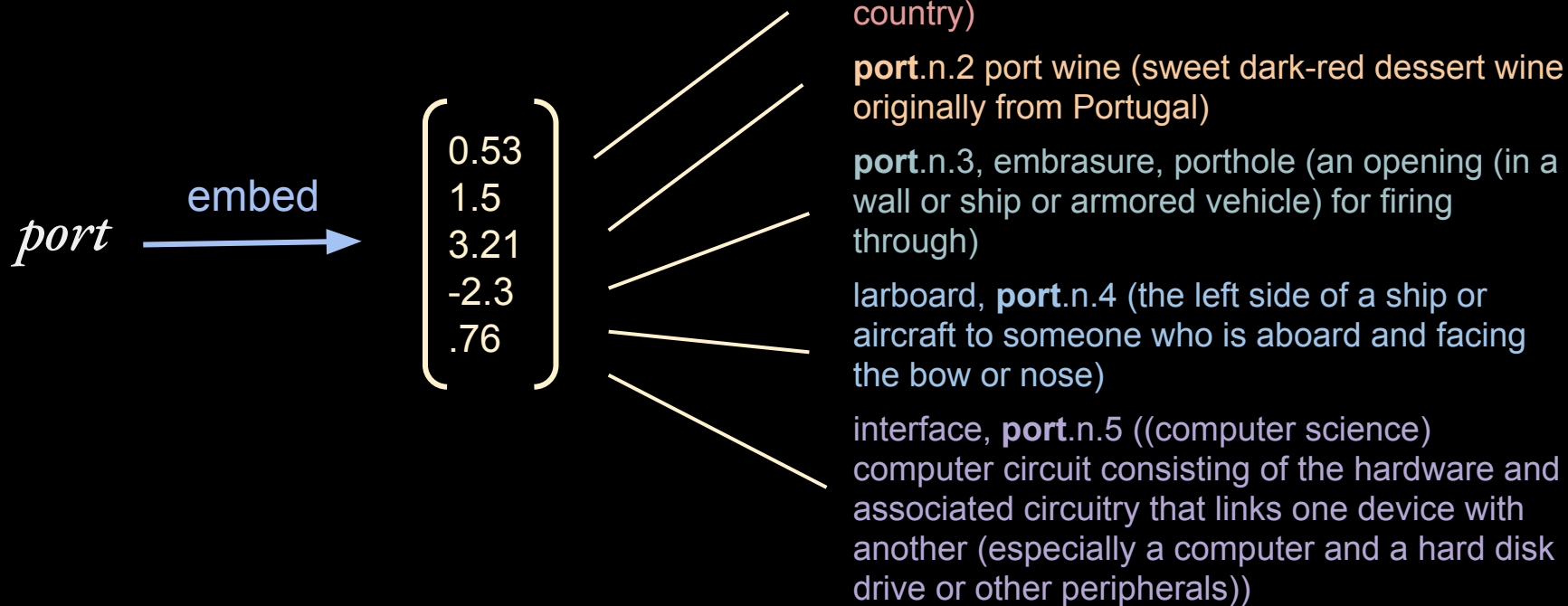
Common applications:

- **Open vocabulary content analysis:** Describing the latent semantic categories of words or phrases present across a set of documents
- **Embeddings for predictive task:** for all topics, use  $p(\text{topic}|\text{document})$  as score. Feed to predictive model (e.g. classifier).





# Objective





# PCA-Based Embeddings

*also known as "Latent Semantic Analysis"*

Dimensionality reduction

-- try to represent with only  $p'$  dimensions

Supplement (details not on test)



# PCA-Based Embeddings

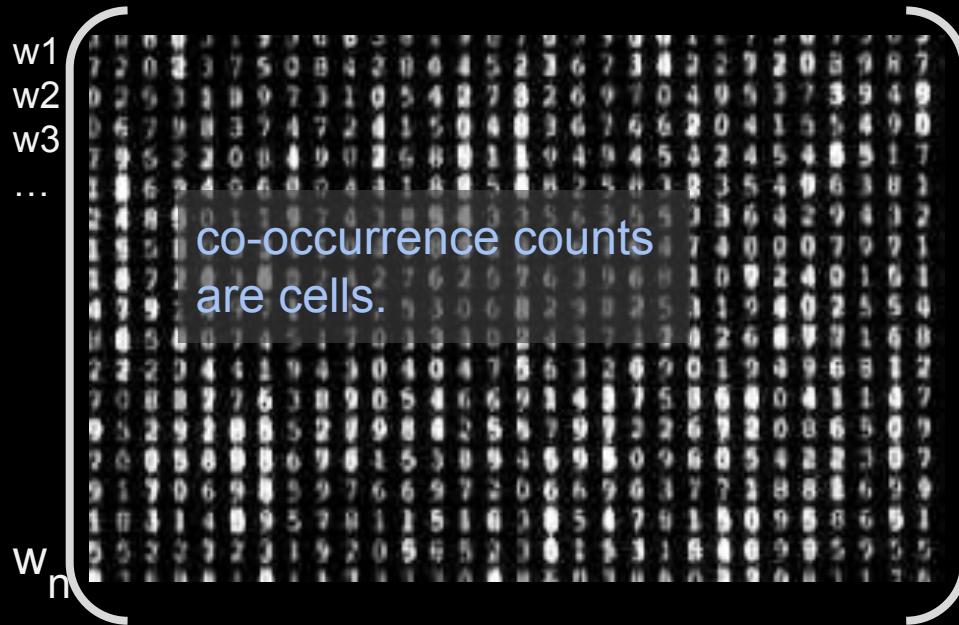
Dimensionality reduction

-- try to represent with only  $p'$  dimensions

$$p' < p$$

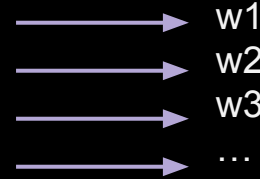
context words are features

$w_1, w_2, w_3, w_4, \dots$



co-occurrence counts  
are cells.

$w_p$



$w_1$   
 $w_2$   
 $w_3$   
 $\dots$



$w_n$

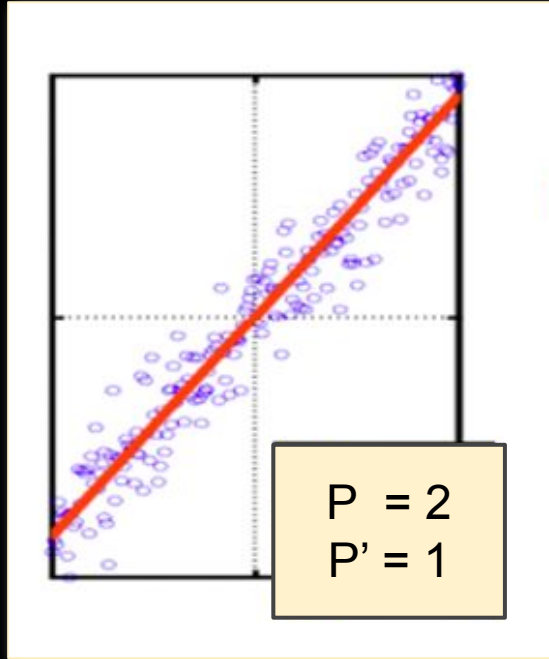
$c_1, c_2, c_3, c_4, \dots$

$cp'$

target words are  
observations

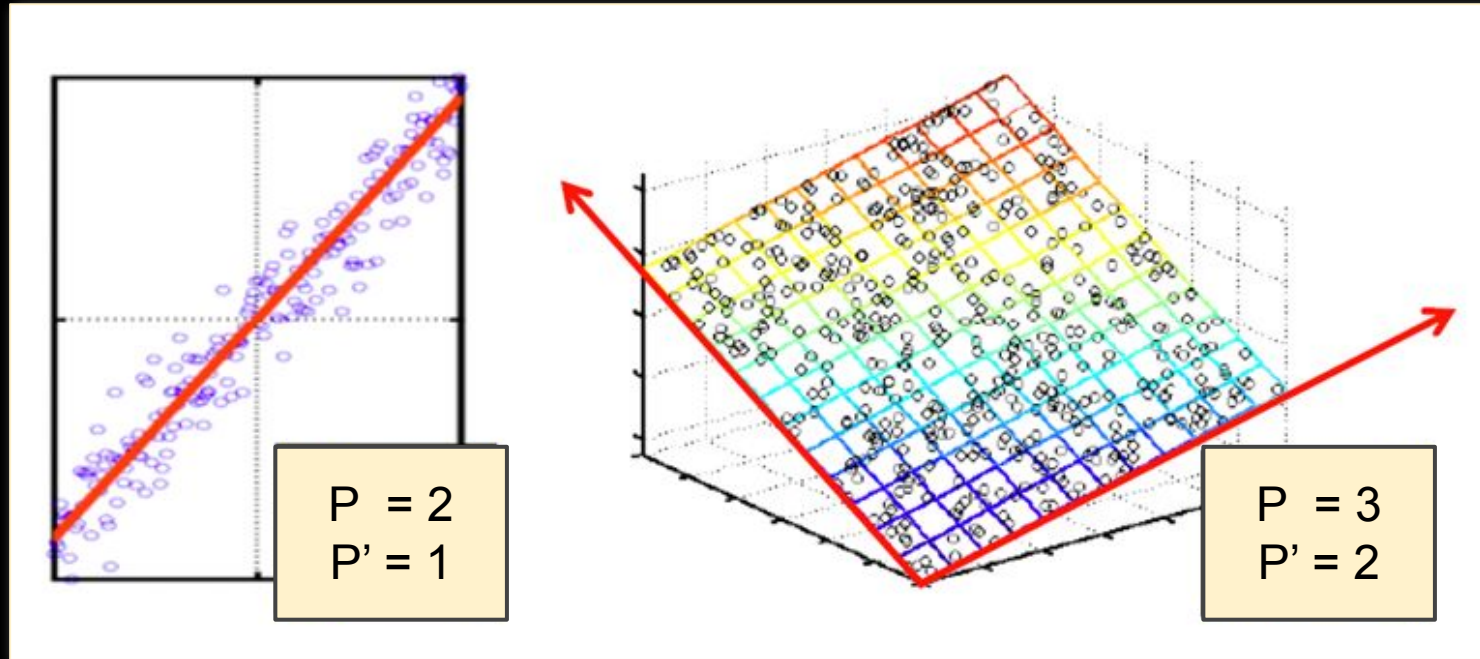


# Concept: Dimensionality Reduction in 3-D, 2-D, and 1-D



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

# 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.

# Concept: 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: How many columns do we really need?

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

# Concept: 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: How many columns do we really need?

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

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

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

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

# SVD-Based Embeddings

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

context words are features

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

$f_p$

$o_1$   
 $o_2$   
 $o_3$   
...

$o_n$

co-occurrence counts  
are cells.

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

$cp'$

$o_1$   
 $o_2$   
 $o_3$   
...

$o_n$

target words are  
observations

# 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*:

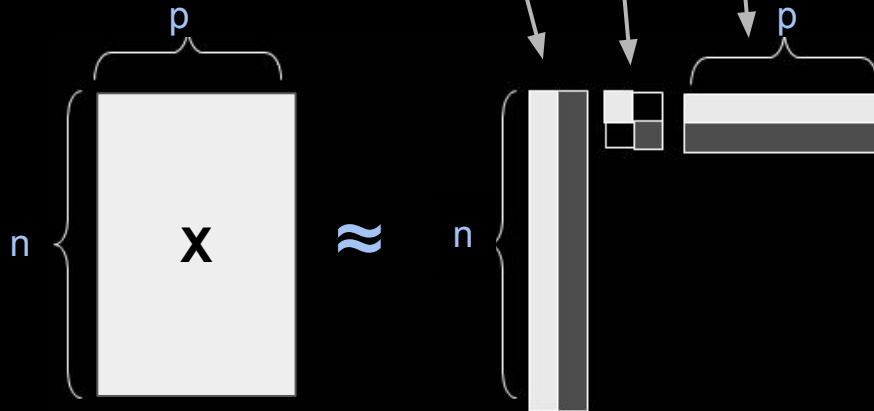
$$X_{[n \times p]} = U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

$X$ : original matrix,

$D$ : “singular values” (diagonal),

$U$ : “left singular vectors”,

$V$ : “right singular vectors”



# Dimensionality Reduction - PCA - Example

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

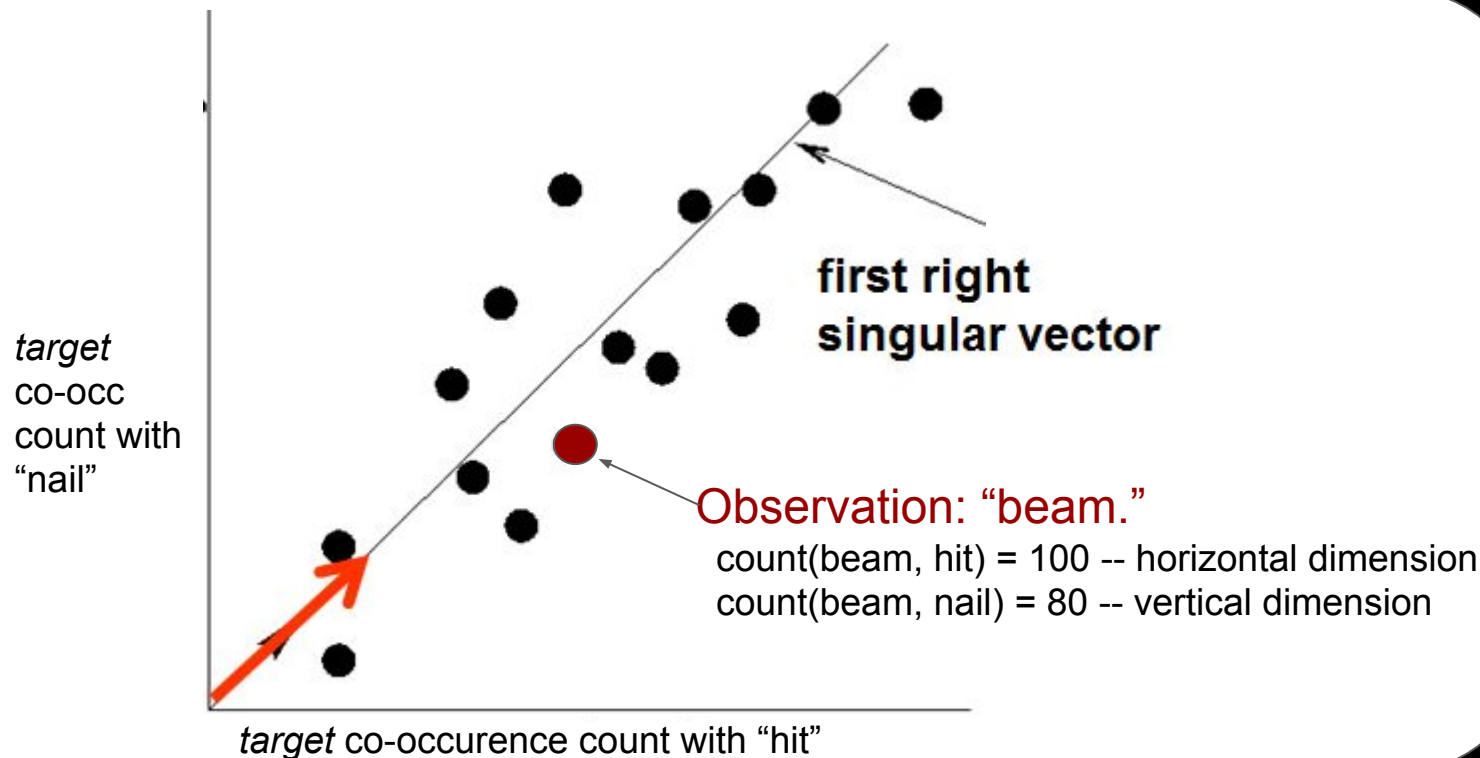
Word co-occurrence  
counts:

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



# 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:

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

# 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:

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

# Dimensionality Reduction - PCA

Li 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:

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

# Dimensionality Reduction - PCA

The loss function 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:

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

# 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