# TTIC 31190: Natural Language Processing

## Lecture 3: Word Representations

Fall 2023

# Announcements

- TA (Jiamin Yang) Tutorial Sessions & Office Hours
  - Fridays 3 pm – 4 pm; TTIC Room 530
  - This week and next: tutorials on Python programming (numpy, PyTorch, etc.)
  - Office hour 4 pm – 5 pm

- Assignment 1 to be released today; due in two weeks

# Recap

- Linguistic Morphology

- Lexical Semantics

- Word Tokenization

# Linguistic Morphology

- **morphology**: study of how words are built from morphemes

- **morphemes**: meaning-bearing units in a language, often classified into **stems** and **affixes**

- type/token ratio correlated with morphological richness of a language

- types of word formation: **inflection**, **derivation**, **compounding**

- morphological decomposition is sometimes hierarchical (`unlockable`)

# Linguistic Morphology

- **lemmatization**: convert wordform to lemma (may depend on context)

- **stemming:** removing affixes from words to get stems (simple, rule-based)

# Lexical Semantics

- **word sense**: discrete representation of an aspect of a word's meaning

- most common words have multiple senses
  - though some sense distinctions are subtle

- semantic relationships among senses:
  - **synonymy**: senses have same meanings, can be used interchangeably
  - **antonymy**: senses are opposites in some dimension of meaning, otherwise are similar
  - **hyponymy** (and **hypernymy**): subclass (or superclass) relationship

# Lexical Semantics

- **word sense disambiguation (WSD)**: NLP task of determining intended sense of a word based on its context
  - methods use words from context of the ambiguous word
  - unclear if useful for downstream tasks
  - today often done implicitly as part of another task

# Word Tokenization

- to do NLP on some text, we need to preprocess it:
  - tokenize documents into sentences
  - tokenize sentences into tokens


- rule-based tokenizers exist for many languages


- for writing systems without whitespace, tokenization becomes complex (often treated as an NLP problem)

# Word Tokenization

- useful terms: **type**, **token**, **type/token ratio**

- when adding data, number of types keeps increasing

- most types are extremely rare (**Zipf's law**)

- Data-driven tokenizers: **Byte Pair Encoding (BPE)**
  - splits words based on data, very common in deep learning

# Question

How does ChatGPT (GPT-2 etc.) tokenize texts from different languages, with a unified tokenizer and fixed vocabulary size?

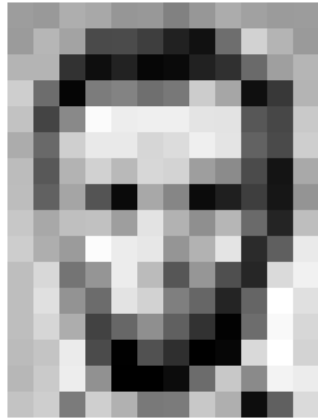Byte-level BPE (BBPE)

GPT-2 vocabulary size: 50257

"That's great 👍"

| | |
|---|---|
| T | 54 |
| h | 68 |
| a | 61 |
| t | 74 |
| ' | 2019 |
| s | 73 |
| | 20 |
| g | 67 |
| r | 72 |
| e | 65 |
| a | 61 |
| t | 74 |
| | 20 |
| 👍 | 1F44D |

| | |
|---|---|
| T | 54 |
| h | 68 |
| a | 61 |
| t | 74 |
| ◆ | e2 |
| ◆ | 80 |
| ◆ | 99 |
| s | 73 |
| | 20 |
| g | 67 |
| r | 72 |
| e | 65 |
| a | 61 |
| t | 74 |
| | 20 |
| ◆ | f0 |
| ◆ | 9f |
| ◆ | 91 |
| ◆ | 8d |

# Digital Representations

How does a computer see?

# Digital Representations

How does a computer see?

# Digital Representations

How does a computer see?

# Digital Representations

How does a computer read?

Birds  are  n't  real  .

# Digital Representations

How does a computer read?

Birds    are   n't    real    .

515      834   45    3435   9

# Digital Representations

How does a computer read?

Birds   are   n't   real   .

515     834   45    3435   9

"Raw" input is often uninteresting/unwieldy to work with.

# Word Representations

Representing words in **vector** space that captures meaningful structure

# How to represent a word

- Stems and affixes
- Dictionary definition
- Lemma and wordforms
- Senses
- Relationships between words and senses

# Annotated Database for Lexical Semantics

- WordNet (Fellbaum, 1998): https://wordnet.princeton.edu/

# All-Words WSD



**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).

# WordNet

- hierarchically organized lexical database
- fine-grained sense inventories, relationships among senses
- originally developed for English; other languages now available
- English WordNet version 3.0 contains:

| Category | Unique Strings |
|---|---|
| Noun | 117,798 |
| Verb | 11,529 |
| Adjective | 22,479 |
| Adverb | 4,481 |

# How is "sense" defined in WordNet?

- **synset** (**synonym set**):
  - set of near-synonyms, instantiates a sense or concept
  - has a **gloss** (roughly, a definition)

- example: $\text{chump}_1$ has gloss "*a person who is gullible and easy to take advantage of*"

- $\text{chump}_1$ belongs to a synset with 8 other senses:

  $\text{fool}_2$, $\text{gull}_1$, $\text{mark}_9$, $\text{patsy}_1$, $\text{fall guy}_1$, $\text{sucker}_1$, $\text{soft touch}_1$, $\text{mug}_2$

- each of **these** senses has this same gloss
  - not **every** sense of these words; $\text{gull}_2$ is the aquatic bird

WordNet has three synsets for the noun `fool`:

- S: (n) **fool**, sap, saphead, muggins, tomfool (a person who lacks good judgment)
- S: (n) chump, **fool**, gull, mark, patsy, fall guy, sucker, soft touch, mug (a person who is gullible and easy to take advantage of)
- S: (n) jester, **fool**, motley fool (a professional clown employed to entertain a king or nobleman in the Middle Ages)

**Ambiguity**

one form, multiple meanings → split form
- the three senses of `fool` belong to different synsets

**Variability**

multiple forms, one meaning → merge forms
- each synset contains senses of several different words

# Hypernyms in WordNet

1. (n) {chump, fool, gull, mark, patsy, fall guy, sucker, soft touch, mug} (a person who is gullible and easy to take advantage of)

2. (n) {victim, dupe} (a person who is tricked or swindled)

3. (n) {person, individual, someone, somebody, mortal, soul} (a human being)

4. (n) {organism, being} (a living thing that has (or can develop) the ability to act or function independently)

5. (n) {living thing, animate thing} (a living (or once living) entity)

6. (n) {whole, unit} (an assemblage of parts that is regarded as a single entity)

7. (n) {object, physical object} (a tangible and visible entity; an entity that can cast a shadow)

8. (n) {physical entity} (an entity that has physical existence)

9. (n) {entity} (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

# How to represent a word

- Until the ~2010s, in NLP, words == atomic symbols
- Nowadays, **vector representations**, **word == vectors**



man = [-0.3, 0.85, 0.45]
king = [0.2, 0.72, 0.35]
woman = [0.3, 0.48, 0.29]
queen = [0.89, 0.41, 0.2]

Similar words are "**nearby in the vector space**"

# How to represent a word

man = [-0.3, 0.85, 0.45]

king = [0.2, 0.72, 0.35]

woman = [0.3, 0.48, 0.29]

queen = [0.89, 0.41, 0.2]

**CHAPTER**

**6**

# Vector Semantics and Embeddings

荃者所以在鱼，得鱼而忘荃　Nets are for fish;
　　　　　　　　　　　　　Once you get the fish, you can forget the net.
言者所以在意，得意而忘言　Words are for meaning;
　　　　　　　　　　　　　Once you get the meaning, you can forget the words
　　　　　　　　　　　　　　　　　　庄子(Zhuangzi), Chapter 26

[SLP3, Chapter 6]

cat    chef    chicken    civic    cooked    council ...

17    91    253    104    5    6001    ...

cat | chef | chicken | civic | cooked | council ...

| cat | chef | chicken | civic | cooked | council |
|-----|------|---------|-------|--------|---------|
| 0.1 | -0.1 | -0.4 | 0.1 | -0.5 | 0.6 |
| 7.9 | 2.1 | 2.4 | 0 | -1.1 | -1.3 |
| 2.4 | 3.8 | 9.7 | -1.5 | 7.6 | 0 |
| -1.3 | -0.1 | -1.0 | 2.4 | -3.1 | 3.4 |
| 0.5 | 5.3 | 3.2 | 0.2 | 4.2 | -0.6 |

**"embeddings"**

# Motivations

**Variability**

multiple forms,
similar meaning

really    reallly    realllly

$$\begin{bmatrix} 2.1 \\ -7.9 \\ 8.4 \\ -1.3 \end{bmatrix} \quad \begin{bmatrix} 2.3 \\ -6.1 \\ 7.8 \\ -0.8 \end{bmatrix} \quad \begin{bmatrix} 1.9 \\ -6.8 \\ 7.7 \\ -1.0 \end{bmatrix}$$

# Representation Learning for Engineering

- Engineering: these representations are often useful for downstream tasks!

- Transfer learning:

Images ⇒      ⇒ Object Classification

Text   ⇒      ⇒ Context Prediction

# Representation Learning for Engineering

- Engineering: these representations are often useful for downstream tasks!

- Transfer learning:

# How to represent a word

$|V|$

| | | | |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 |

$w =$
$\in V$

*better*    *winner* ... *cat*    *champion*

$|V|$

# How to represent a word

| 1 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 |

$|V|$

$w =$
$\in V$

*better*　*winner* … *cat*　*champion*

$|V|$

- "One-hot" representation of words

$$\mathbf{Rep}(w) \in \{0,1\}^{|V|}$$

# How to represent a word

|V| $\left[\begin{array}{cccc}\end{array}\right.$

| | | | |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 |

$w$ = *better*   *winner*  ...  *cat*   *champion*
$\in V$

|V|

- "One-hot" representation of words

$$\mathrm{Rep}(w) \in \{0,1\}^{|V|}$$

|V| could be very large! (e.g. 50K)

# How to represent a word

$|V|$

| better | winner | ... cat | champion |
|:---:|:---:|:---:|:---:|
| 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 |

$w = \in V$

$|V|$



winner of world cup 2022

Images   News   Videos   Team   Groups   Qatar   Woman   Golden Boot   Fifa 22

About 3,820,000,000 results (0.47 seconds)

See results about — World Cup — Soccer competition

2022 World Cup / Champion

## Argentina national football team

Recent News. 2022 FIFA World Cup, international football (soccer) tournament that took place in Qatar from November 20 to December 18, 2022, and was contested by the men's national teams of 32 countries. Argentina won its third World Cup victory in the tournament after defeating France in the final match.

# How to represent a word



$|V|$

$$w = \in V$$

better | winner ... cat | champion

$|V|$

winner of world cup 2022

Images   News   Videos   Team   Groups   Qatar   Woman   Golden Boot   Fifa 22

About 3,820,000,000 results (0.47 seconds)

See results about   World Cup   Soccer competition

2022 World Cup / Champion

## Argentina national football team

Recent News. 2022 FIFA World Cup, international football (soccer) tournament that took place in Qatar from November 20 to December 18, 2022, and was contested by the men's national teams of 32 countries. Argentina won its third World Cup victory in the tournament after defeating France in the final match.

Vectors are orthogonal!

# Word Representation

- What is an ideal word representation?

- It should probably capture information about usage and meaning:
  - Part of speech tags (noun, verb, adj., adv., etc.)
  - The intended sense
  - Semantic similarities (winner **vs.** champion)
  - Semantic relationships (antonyms, hypernyms, etc.)

# Features?

| | better | winner | ... | cat | champion |
|---|---|---|---|---|---|
| Is noun? | 0 | 1 | | 1 | 1 |
| Is verb? | 0 | 0 | | 0 | 1 |
| Is adj.? | 1 | 0 | | 0 | 0 |
| Is animal? | 0 | 0 | | 1 | 0 |
| ... | 0 | 0 | | 0 | 0 |
| | 0 | 0 | | 0 | 0 |
| | 0 | 0 | | 1 | 0 |
| | 0 | 0 | | 0 | 1 |

?

$|V|$

# Features?

| | 0 | 1 | 1 | 1 | |
|---|---|---|---|---|---|
| Is noun? | 0 | 1 | 1 | 1 | |
| Is verb? | 0 | 0 | 0 | 1 | |
| Is adj.? | 1 | 0 | 0 | 0 | |
| Is animal? | 0 | 0 | 1 | 0 | |
| ... | 0 | 0 | 0 | 0 | |
| | 0 | 0 | 0 | 0 | |
| | 0 | 0 | 1 | 0 | ? |
| | 0 | 0 | 0 | 1 | |
| | 1 | 1 | 1 | 1 | |
| | 0 | 0 | 0 | 0 | |
| | 0 | 0 | 0 | 0 | |

# Features?

## WordNet

# Word Representation

- What is an ideal word representation?

- It should probably capture information about usage and meaning:
  - Part of speech tags (noun, verb, adj., adv., etc.)
  - The intended sense
  - Semantic similarities (winner **vs.** champion)
  - Semantic relationships (antonyms, hypernyms, etc.)

**Distributional Semantics:**
How much of this can we capture from context/data alone?

# Main Idea

"The meaning of a word is its use in the language."

[Ludwig Wittgenstein 1943]

"Usage":

Words are defined by their environments

(the words around them)

# Main Idea

Consider encountering a new word: tezgüino.

1. A bottle of tezgüino is on the table.

2. Everybody likes tezgüino.

3. Don't have tezgüino before you drive.

4. We make tezgüino out of corn.

What do you think the tezgüino is?

loud
motor oil
tortillas
choices
wine

# Main Idea

Consider encountering a new word: tezgüino.

1. A bottle of tezgüino **is on the table.**
2. **Everybody likes** tezgüino.
3. **Don't have** tezgüino **before you drive.**
4. **We make** tezgüino **out of corn.**

**What do you think the** tezgüino **is?**

|  | context | | | |
|---|---|---|---|---|
| | 1 | 2 | 3 | 4 |
| tezgüino | 1 | 1 | 1 | 1 |
| loud | 0 | 0 | 0 | 0 |
| motor oil | 1 | 0 | 0 | 1 |
| tortillas | 0 | 1 | 0 | 1 |
| choices | 0 | 1 | 0 | 0 |
| wine | 1 | 1 | 1 | 0 |

term

# Main Idea

Consider encountering a new word: tezgüino.

1. A bottle of tezgüino **is on the table.**

2. **Everybody likes** tezgüino.

3. **Don't have** tezgüino **before you drive.**

4. **We make** tezgüino **out of corn.**

What do you think the tezgüino is?

| | **context** | | | | similarity? |
|---|---|---|---|---|---|
| **term** | **1** | **2** | **3** | **4** | |
| tezgüino | 1 | 1 | 1 | 1 | |
| loud | 0 | 0 | 0 | 0 | 0 |
| motor oil | 1 | 0 | 0 | 1 | 2 |
| tortillas | 0 | 1 | 0 | 1 | 2 |
| choices | 0 | 1 | 0 | 0 | 1 |
| wine | 1 | 1 | 1 | 0 | 3 |

# Distributional Hypothesis

- These representations encode distributional properties of each word
- The distributional hypothesis: words with similar meaning are used in similar contexts.

*"You shall know a word by the company it keeps."*

J.R. Firth, *A Synopsis of Linguistic Theory*, 1957

*"The meaning of a word is its use in the language."*

[Ludwig Wittgenstein 1943]

*"If A and B have almost identical environments we say that they are synonyms."*

[Harris 1954]

# Distributional Hypothesis

- How can we automate the process of constructing representations of word meaning from its "company"?

- First solution: word-word co-occurrence counts

words we are computing vectors for:

cat  chef  chicken  civic  cooked  council

context
words:

the
cat
chicken
city
cook

… , the club may also employ a **chef** to prepare and cook food items .

… is up to remy , linguini , and the **chef** colette to cook for many people …

… cooking program the cook and the **chef** with simon bryant , who is …

chef

|        | chef |
|--------|------|
| the    | 0    |
| cat    | 0    |
| chicken| 0    |
| city   | 0    |
| cook   | 0    |

… , the club may also employ a **chef** to prepare and cook food items .

… is up to remy , linguini , and the **chef** colette to cook for many people …

… cooking program the cook and the **chef** with simon bryant , who is …

window size: $w = 1$

$$\begin{array}{c} \text{chef} \\ \begin{array}{r} \text{the} \\ \text{cat} \\ \text{chicken} \\ \text{city} \\ \text{cook} \end{array} \begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \end{array}$$

… , the club may also employ a **chef** to prepare and cook food items .

… is up to remy , linguini , and the **chef** colette to cook for many people …

… cooking program the cook and the **chef** with simon bryant , who is …

window size: $w = 4$

$$
\begin{array}{c}
\text{chef} \\
\begin{array}{r}
\text{the} \\
\text{cat} \\
\text{chicken} \\
\text{city} \\
\text{cook}
\end{array}
\begin{bmatrix}
3 \\
0 \\
0 \\
0 \\
3
\end{bmatrix}
\end{array}
$$

words we are computing vectors for:

|  | cat | chef | chicken | civic | cooked | council |
|---|---|---|---|---|---|---|
| the | 24708 | 7410 | 7853 | 16486 | 3463 | 316380 |
| cat | 2336 | 14 | 23 | 0 | 1 | 36 |
| chicken | 23 | 21 | 1640 | 1 | 181 | 7 |
| city | 116 | 89 | 62 | 943 | 7 | 27033 |
| cook | 12 | 113 | 34 | 6 | 34 | 51 |

context words:

- once we have word vectors, we can compute similarities!
- many ways to define similarity of two vectors
- a simple way: **dot product** (also called inner product):

$$\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^\top \mathbf{v} = \sum_i u_i v_i$$

$$\mathbf{u} = \text{a vector}$$

$$u_i = \text{entry i in the vector}$$

- dot product is large when the vectors have very large (or very negative) values in the same dimensions

with dot product as similarity function, let's find the most similar words ("nearest neighbors") to each word:

| | cat | chef | chicken | civic | cooked | council |
|---|---|---|---|---|---|---|
| nearest neighbors | council | council | council | council | council | council |
| | cat | cat | cat | cat | cat | cat |
| | civic | civic | civic | civic | civic | civic |
| | chicken | chicken | chicken | chicken | chicken | chicken |
| | chef | chef | chef | chef | chef | chef |
| | cooked | cooked | cooked | cooked | cooked | cooked |

with dot product as similarity function, let's find the most similar words ("nearest neighbors") to each word:

| | cat | chef | chicken | civic | cooked | council |
|---|---|---|---|---|---|---|
| | council | council | council | council | council | council |
| nearest neighbors | cat | cat | cat | cat | cat | cat |
| | civic | civic | civic | civic | civic | civic |
| | chicken | chicken | chicken | chicken | chicken | chicken |
| | chef | chef | chef | chef | chef | chef |
| | cooked | cooked | cooked | cooked | cooked | cooked |

with dot product as similarity function, let's find the rds ("nearest neighbors") to each word:

**council**

| | council |
|---|---|
| the | 316380 |
| cat | 36 |
| chicken | 7 |
| city | 27033 |
| cook | 51 |

| | chef | chicken | civic | cooked | council |
|---|---|---|---|---|---|
| council | council | council | council | council | council |
| cat | cat | cat | cat | cat | cat |
| civic | civic | civic | civic | civic | civic |
| chicken | chicken | chicken | chicken | chicken | chicken |
| chef | chef | chef | chef | chef | chef |
| cooked | cooked | cooked | cooked | cooked | cooked |

- dot product is large when vectors have large values in same dimensions, doesn't control for vector length

- vector length: $$\|\mathbf{u}\| = \sqrt{\sum_i u_i^2}$$

$\mathbf{u} =$ a vector

- **cosine similarity**: $$\frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \, \|\mathbf{v}\|}$$

$u_i =$ entry i in the vector

this is the cosine of the angle between the two vectors!

now using cosine similarity:

| | cat | chef | chicken | civic | cooked | council |
|---|---|---|---|---|---|---|
| | cat | chef | chicken | civic | cooked | council |
| nearest neighbors | chef | civic | cooked | council | chef | civic |
| | cooked | cooked | chef | chef | civic | chef |
| | civic | council | civic | cooked | council | cooked |
| | council | cat | council | cat | cat | cat |
| | chicken | chicken | cat | chicken | chicken | chicken |

# Any issue?

Raw frequency count is probably a bad representation!

Counts of common words are very large, but not very useful
- "the", "it", "they"
- Not very informative

Many ways proposed for improving raw counts

# Any issue?

Raw frequency count is probably a bad representation!

Counts of common words are very large, but not very useful

- "the", "it", "they"

- Not very informative

Many ways proposed for improving raw counts

- **TF-IDF**
- **PMI**
- **word2vec**

# TF-IDF

TF (Term Frequency) - IDF (Inverse Document Frequency)

# TF-IDF

TF (Term Frequency) - IDF (Inverse Document Frequency)

- Information Retrieval (IR) workhorse!
- A common baseline model
- Sparse vectors
- Words are represented by (a simple function of) the counts of nearby words

# TF-IDF Term-Document Matrix

- Consider a matrix of word counts across documents: term-document matrix

document

$$\mathtt{tf}(w, d) = \# \text{ of times word } w \text{ appears in document } d$$

$$\mathrm{tf}_{t,d} = \mathrm{count}(t, d)$$

word

$$\mathtt{idf}(w) = \log\left(\frac{\# \text{ of documents}}{\# \text{ of documents in which word } w \text{ occurs}}\right)$$

$$\mathrm{idf}_t = \log_{10}\left(\frac{N}{\mathrm{df}_t}\right)$$

$$\mathtt{tf-idf}(w, d) = \mathtt{tf}(w, d) \cdot \mathtt{idf}(w)$$

$$w_{t,d} = \mathrm{tf}_{t,d} \times \mathrm{idf}_t$$

# TF-IDF Term-Document Matrix

- Consider a matrix of word counts across documents: term-document matrix

$$\text{tf}_{t,d} = \text{count}(t,d)$$

|  | As You Like It | Twelfth Night | Julius Caesar | Henry V |
|---|---|---|---|---|
| battle | 1 | 0 | 7 | 13 |
| good | 114 | 80 | 62 | 89 |
| fool | 36 | 58 | 1 | 4 |
| wit | 20 | 15 | 2 | 3 |

# TF-IDF Term-Document Matrix

- Consider a matrix of word counts across documents: term-document matrix

$$\text{tf}_{t,d} = \text{count}(t,d)$$

| | As You Like It | Twelfth Night | Julius Caesar | Henry V | |
|---|---|---|---|---|---|
| battle | 1 | 0 | 7 | 13 | word vector |
| good | 114 | 80 | 62 | 89 | |
| fool | 36 | 58 | 1 | 4 | |
| wit | 20 | 15 | 2 | 3 | |

bag-of-words
(document representation)

# TF-IDF Term-Document Matrix

- Consider a matrix of word counts across documents: term-document matrix

document

word

$$\mathtt{tf}(w, d) = \# \text{ of times word } w \text{ appears in document } d$$

$$\mathtt{idf}(w) = \log\left(\frac{\# \text{ of documents}}{\# \text{ of documents in which word } w \text{ occurs}}\right)$$

$$\mathtt{tf\text{-}idf}(w, d) = \mathtt{tf}(w, d) \cdot \mathtt{idf}(w)$$

$$\mathrm{tf}_{t,d} = \mathrm{count}(t, d)$$

$$\mathrm{idf}_t = \log_{10}\left(\frac{N}{\mathrm{df}_t}\right)$$

$$w_{t,d} = \mathrm{tf}_{t,d} \times \mathrm{idf}_t$$

≈ 0 for words like "the"

# TF-IDF Term-Document Matrix

- IDF from 37 Shakespeare plays

$$idf_t = \log_{10}\left(\frac{N}{df_t}\right)$$

$$tf_{t,d} = count(t,d)$$

| word | df | idf |
|------|-----|-------|
| Romeo | 1 | 1.57 |
| salad | 2 | 1.27 |
| Falstaff | 4 | 0.967 |
| forest | 12 | 0.489 |
| battle | 21 | 0.246 |
| wit | 34 | 0.037 |
| fool | 36 | 0.012 |
| good | 37 | 0 |
| sweet | 37 | 0 |

| | As You Like It | Twelfth Night | Julius Caesar | Henry V |
|------|------|------|------|------|
| battle | 1 | 0 | 7 | 13 |
| good | 114 | 80 | 62 | 89 |
| fool | 36 | 58 | 1 | 4 |
| wit | 20 | 15 | 2 | 3 |

# TF-IDF Term-Document Matrix

- IDF from 37 Shakespeare plays

$$\mathrm{idf}_t = \log_{10}\left(\frac{N}{\mathrm{df}_t}\right)$$

$$w_{t,d} = \mathrm{tf}_{t,d} \times \mathrm{idf}_t$$

| word | df | idf |
|------|-----|-------|
| Romeo | 1 | 1.57 |
| salad | 2 | 1.27 |
| Falstaff | 4 | 0.967 |
| forest | 12 | 0.489 |
| battle | 21 | 0.246 |
| wit | 34 | 0.037 |
| fool | 36 | 0.012 |
| good | 37 | 0 |
| sweet | 37 | 0 |

| | As You Like It | Twelfth Night | Julius Caesar | Henry V |
|-------|-------|-------|-------|-------|
| battle | 0.074 | 0 | 0.22 | 0.28 |
| good | 0 | 0 | 0 | 0 |
| fool | 0.019 | 0.021 | 0.0036 | 0.0083 |
| wit | 0.049 | 0.044 | 0.018 | 0.022 |

# TF-IDF Variations

**Variants of term frequency (tf) weight**

| weighting scheme | tf weight |
|---|---|
| binary | $0, 1$ |
| raw count | $f_{t,d}$ |
| term frequency | $f_{t,d} \bigg/ \sum_{t' \in d} f_{t',d}$ |
| log normalization | $\log(1 + f_{t,d})$ |
| double normalization 0.5 | $0.5 + 0.5 \cdot \dfrac{f_{t,d}}{\max_{\{t' \in d\}} f_{t',d}}$ |
| double normalization K | $K + (1 - K) \dfrac{f_{t,d}}{\max_{\{t' \in d\}} f_{t',d}}$ |

**Variants of inverse document frequency (idf) weight**

| weighting scheme | idf weight $(n_t = \lvert \{d \in D : t \in d\} \rvert)$ |
|---|---|
| unary | $1$ |
| inverse document frequency | $\log \dfrac{N}{n_t} = -\log \dfrac{n_t}{N}$ |
| inverse document frequency smooth | $\log \left( \dfrac{N}{1 + n_t} \right) + 1$ |
| inverse document frequency max | $\log \left( \dfrac{\max_{\{t' \in d\}} n_{t'}}{1 + n_t} \right)$ |
| probabilistic inverse document frequency | $\log \dfrac{N - n_t}{n_t}$ |

[Image: Wikipedia]

# TF-IDF Usage

- TF-IDF was designed for and still excels at document retrieval
- The BM25 model (very similar to TF-IDF) is still a strong document retrieval baseline!

$$BM25(d, \mathbf{q}) = \sum_{q_i \in \mathbf{q}} \text{idf}_{q_i} \cdot \frac{tf_{q_i,d} \cdot (k_1 + 1)}{tf_{q_i,d} + k_1(1 - b + b\frac{|D|}{avg|D|})}$$

Recent history in NLP might suggest that learned dense representations should always outperform sparse features, but this is not necessarily true: as shown in Figure 1, the BM25 model (Robertson et al., 2009) can outperform a dual encoder based on BERT, particularly on longer documents (See § 7).

**Containing Passage Retrieval**

recall@1

- BM25-uni
- BM25-bi
- DE-BERT-768
- BERT-init

passage length: 50, 100, 200, 400

Fancy, computationally expensive neural network models!

Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. 2021. Sparse, Dense, and Attentional Representations for Text Retrieval. *Transactions of the Association for Computational Linguistics*, 9:329–345.

# Pointwise Mutual Information (PMI)

# Pointwise Mutual Information (PMI)

- consider two random variables, $X$ and $Y$

- do two events $X = x$ and $Y = y$ occur together more often than if they were independent?

$$\text{pmi}(x, y) = \log_2 \frac{p_{X,Y}(x, y)}{p_X(x)\, p_Y(y)}$$

- if they are independent, PMI = 0

# PMI for Word Vectors

- for word vectors,
  $X$ is the **center word**
  $Y$ is the **context word**

- each probability can be estimated using counts we already computed!

$\#(x, y)$ = co-occurrence count of x and y
$N$ = total count

$$\mathrm{pmi}(x, y) = \log_2 \frac{p_{X,Y}(x, y)}{p_X(x)\, p_Y(y)}$$

$$p_{X,Y}(x, y) = \frac{\#(x, y)}{N}$$

$$p_X(x) = \frac{\sum_y \#(x, y)}{N}$$

$$p_Y(y) = \frac{\sum_x \#(x, y)}{N}$$

# Top co-occurrence counts with "chicken"

| | | | | | |
|---|---|---|---|---|---|
| 14464 | , | 1525 | or | 508 | pork |
| 7853 | the | 1225 | for | 500 | meat |
| 6276 | and | 1061 | 's | 481 | be |
| 5931 | . | 940 | fried | 479 | he |
| 5213 | a | 906 | on | 452 | such |
| 3963 | of | 889 | was | 445 | his |
| 3282 | in | 869 | that | 417 | at |
| 2520 | to | 828 | are | 405 | soup |
| 2438 | " | 777 | by | 389 | made |
| 2339 | is | 746 | from | 384 | rice |
| 2127 | with | 710 | it | 375 | but |
| 1818 | ( | 600 | beef | 350 | has |
| 1745 | ) | 590 | which | 330 | fish |
| 1640 | chicken | 557 | also | 325 | other |
| 1594 | as | 531 | an | 318 | this |

# Words with largest PMI with "chicken"

| | | | | | |
|---|---|---|---|---|---|
| 10.2 | fried | 7.0 | robot | 6.1 | pig |
| 9.7 | chicken | 6.9 | burger | 6.0 | breeds |
| 9.3 | pork | 6.8 | recipe | 6.0 | vegetable |
| 9.0 | beef | 6.6 | vegetables | 6.0 | potato |
| 8.7 | soup | 6.6 | potatoes | 5.9 | goose |
| 7.8 | sauce | 6.6 | goat | 5.9 | dixie |
| 7.7 | curry | 6.5 | eggs | 5.9 | kung |
| 7.6 | cooked | 6.4 | cow | 5.9 | pie |
| 7.5 | lamb | 6.4 | pizza | 5.8 | menu |
| 7.4 | dish | 6.4 | rice | 5.8 | steamed |
| 7.3 | shrimp | 6.3 | ribs | 5.8 | tastes |
| 7.3 | egg | 6.3 | tomatoes | 5.7 | beans |
| 7.2 | sandwich | 6.2 | cheese | 5.7 | butter |
| 7.2 | dishes | 6.2 | duck | 5.7 | barn |
| 7.2 | meat | 6.1 | chili | 5.7 | breed |

words we are computing vectors for:

|  | cat | chef | chicken | civic | cooked | council |
|---|---|---|---|---|---|---|
| the | 24708 | 7410 | 7853 | 16486 | 3463 | 316380 |
| cat | 2336 | 14 | 23 | 0 | 1 | 36 |
| chicken | 23 | 21 | 1640 | 1 | 181 | 7 |
| city | 116 | 89 | 62 | 943 | 7 | 27033 |
| cook | 12 | 113 | 34 | 6 | 34 | 51 |

context words:

words we are computing vectors for:

|  | cat | chef | chicken | civic | cooked | council |
|---|---|---|---|---|---|---|
| the | 0.1 | -0.1 | -0.4 | 0.1 | -0.5 | 0.6 |
| cat | 7.9 | 2.1 | 2.4 | 0 | -1.1 | -1.3 |
| chicken | 2.4 | 3.8 | 9.7 | -1.5 | 7.6 | 0 |
| city | -1.3 | -0.1 | -1.0 | 2.4 | -3.1 | 3.4 |
| cook | 0.5 | 5.3 | 3.2 | 0.2 | 4.2 | -0.6 |

context words:

$\mathrm{pmi(cat, the)}$ $\mathrm{pmi(chef, the)}$ ...

using counts:

| cat | chef | chicken | civic | cooked | council |
|-----|------|---------|-------|--------|---------|
| cat | chef | chicken | civic | cooked | council |
| chef | civic | cooked | council | chef | civic |
| cooked | cooked | chef | chef | civic | chef |
| civic | council | civic | cooked | council | cooked |
| council | cat | council | cat | cat | cat |
| chicken | chicken | cat | chicken | chicken | chicken |

nearest neighbors

using PMIs:

| cat | chef | chicken | civic | cooked | council |
|------|------|---------|-------|--------|---------|
| cat | chef | chicken | civic | cooked | council |
| chicken | chicken | cooked | council | chicken | civic |
| chef | cooked | chef | chef | chef | chicken |
| cooked | cat | cat | cat | cat | chef |
| civic | council | council | chicken | council | cooked |
| council | civic | civic | cooked | civic | cat |

nearest neighbors

# Positive PMI (PPMI)

- some have found benefit by truncating PMI at 0 ("positive PMI")

$$\text{PPMI}(u, v) = \max\{0, \text{PMI}(u, v)\}$$

- negative PMI: words occur together less than we would expect, i.e., they are anticorrelated

- these anticorrelations may need more data to reliably estimate

- however, negative PMIs do seem reasonable!

| Largest PMIs: | PMIs close to zero: | Smallest PMIs: |
|---|---|---|
| 10.2 fried | 0.003 climbed | -4.6 users |
| 9.7 chicken | 0.003 detailing | -4.6 data |
| 9.3 pork | 0.002 turkish | -4.7 discussion |
| 9.0 beef | 0.002 oaks | -4.7 museum |
| 8.7 soup | 0.001 productivity | -4.7 below |
| 7.8 sauce | 0.000 swing | -4.8 editors |
| 7.7 curry | -0.001 structures | -4.8 railway |
| 7.6 cooked | -0.001 thirteenth | -4.8 committee |
| 7.5 lamb | -0.001 commentators | -4.8 elected |
| 7.4 dish | -0.001 palmer | -4.9 championship |
| 7.3 shrimp | -0.002 obstacles | -5.0 archive |
| 7.3 egg | -0.003 horns | -5.3 edits |
| 7.2 sandwich | -0.003 burning | -6.1 deletion |

# words we are computing vectors for:

- downside: large context word vocabulary needed for good vectors (1,000 to 10,000)

- hard to work with high-dimensional vectors

- we can reduce dimensionality (SVD, etc.), but this is difficult to scale to large vocabularies

cooked  council

$$\begin{bmatrix} -0.5 \\ -1.1 \\ 7.6 \\ -3.1 \\ 4.2 \end{bmatrix} \quad \begin{bmatrix} 0.6 \\ -1.3 \\ 0 \\ 3.4 \\ -0.6 \end{bmatrix}$$

# word2vec

# word2vec

- Learning representations with neural networks



$$\overrightarrow{king} - \overrightarrow{man} + \overrightarrow{woman} = \overrightarrow{queen}$$

# word2vec

- Learning representations with neural networks

## Efficient Estimation of Word Representations in Vector Space

**Tomas Mikolov**
Google Inc., Mountain View, CA
tmikolov@google.com

**Kai Chen**
Google Inc., Mountain View, CA
kaichen@google.com

**Greg Corrado**
Google Inc., Mountain View, CA
gcorrado@google.com

**Jeffrey Dean**
Google Inc., Mountain View, CA
jeff@google.com

## Distributed Representations of Words and Phrases and their Compositionality

**Tomas Mikolov**
Google Inc.
Mountain View
mikolov@google.com

**Ilya Sutskever**
Google Inc.
Mountain View
ilyasu@google.com

**Kai Chen**
Google Inc.
Mountain View
kai@google.com

**Greg Corrado**
Google Inc.
Mountain View
gcorrado@google.com

**Jeffrey Dean**
Google Inc.
Mountain View
jeff@google.com

[Mikolov et al., 2013]

# word2vec

- Learning representations with neural networks

- Instead of counting, train a classifier (neural network) to **predict** context (e.g. neighboring words)

*Count-based Distributional Semantics* ➡ *Neural Distributional Semantics*

- Training is **self-supervised**: no annotated data required, just raw text

- Word embeddings learned via backpropagation

# Neural Word Embeddings

cat
dog
paw
great
good
printer
zoom
stonks
red
bandaid
cash
jumped
scintillating

The excited dog jumped over the annoyed cat

The fluffy Samoyed jumped over the backyard fence

The quick brown fox jumped over the lazy dog



Intuition: word embedding for "jumped" should be learned (from random initialization) such that it can well-predict surrounding context.

# word2vec

- CBOW (Continuous Bag-of-Words): learn representations that predict a word given context

$$P(w_t \mid w_{t+1}, \ldots, w_{t+k}, w_{t-1}, \ldots, w_{t-k})$$

- Skipgram: learn representations that predict the context given a word
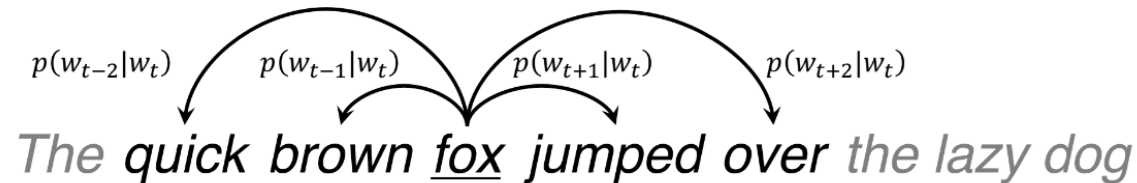
$$P(w_{t+1}, \ldots, w_{t+k}, w_{t-1}, \ldots, w_{t-k} \mid w_t)$$

# word2vec

- CBOW (Continuous Bag-of-Words): learn representations that predict a word given context



- Skipgram: learn representations that predict the context given a word

# skipgram

Randomly initialized.
(To be learned via backprop)

$$
W = \begin{array}{c} a \\ aardvark \\ able \\ are \\ \vdots \\ zyzzyva \end{array}
\begin{bmatrix}
1.2 & -0.1 & 0.3 & \ldots & 0.1 \\
0.2 & 0.7 & -0.4 & \ldots & 1.1 \\
-0.7 & 0.5 & 0.6 & \ldots & -0.8 \\
0.1 & 0.9 & 0.8 & \ldots & 0.7 \\
& & \vdots & & \\
0.3 & -0.2 & 0.7 & \ldots & 0.4
\end{bmatrix}
\qquad
U = \begin{array}{c} a \\ aardvark \\ able \\ are \\ \vdots \\ zyzzyva \end{array}
\begin{bmatrix}
2.1 & -0.5 & 1.3 & \ldots & 1.4 \\
-0.4 & -0.7 & 0.5 & \ldots & 0.1 \\
0.2 & 0.1 & 0.4 & \ldots & -0.7 \\
0.5 & 0.8 & 0.1 & \ldots & 0.4 \\
& & \vdots & & \\
-0.3 & 0.3 & 0.2 & \ldots & 0.6
\end{bmatrix}
$$

$$\theta = \{W, U\}$$

$W : V \times d$ input embedding matrix

$U : V \times d$ output embedding matrix

# skipgram

Randomly initialized.
(To be learned via backprop)

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_{\text{out}} \cdot w_{\text{input}})}{\sum_{v \in V} \exp(u_v \cdot w_{\text{input}})}$$

Just a (log) linear model!

**softmax**

|  | | | | | |
|---|---|---|---|---|---|
| a | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| aardvark | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| able | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| are | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| ⋮ | | | ⋮ | | |
| zyzzyva | 0.3 | -0.2 | 0.7 | ... | 0.4 |

|  | | | | | |
|---|---|---|---|---|---|
| a | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| aardvark | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| able | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| are | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| ⋮ | | | ⋮ | | |
| zyzzyva | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$$W \qquad\qquad U$$

$$\theta = \{W, U\}$$

$W : V \times d$ input embedding matrix

$U : V \times d$ output embedding matrix

# skipgram

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_{\text{out}} \cdot w_{\text{input}})}{\sum_{v \in V} \exp(u_v \cdot w_{\text{input}})}$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| *aardvark* | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | ... | 0.4 |

$$W$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| *aardvark* | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$$U$$

it is a far , far better rest that I go to , than I have ever known

# skipgram

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_\text{out} \cdot w_\text{input})}{\sum_{v \in V} \exp(u_v \cdot w_\text{input})}$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| *aardvark* | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | ... | 0.4 |

$W$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| *aardvark* | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$U$

it is a far , far better **rest** that I go to , than I have ever known

Pick a window centered at a word and predict the context (window size is a hyperparameter)

# skipgram

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_{\text{out}} \cdot w_{\text{input}})}{\sum_{v \in V} \exp(u_v \cdot w_{\text{input}})}$$

| | | | | | |
|---|---|---|---|---|---|
| $a$ | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| $aardvark$ | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| $able$ | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| $are$ | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| $\vdots$ | | | $\vdots$ | | |
| $zyzzyva$ | 0.3 | -0.2 | 0.7 | ... | 0.4 |

$$W$$

| | | | | | |
|---|---|---|---|---|---|
| $a$ | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| $aardvark$ | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| $able$ | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| $are$ | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| $\vdots$ | | | $\vdots$ | | |
| $zyzzyva$ | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$$U$$

it is a far , far better rest that I go to , than I have ever known

$$L_t = -\log p_\theta(x_{t-2} \mid x_t) - \log p_\theta(x_{t-1} \mid x_t)$$
$$- \log p_\theta(x_{t+1} \mid x_t) - \log p_\theta(x_{t+2} \mid x_t)$$

# skipgram

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_{\text{out}} \cdot w_{\text{input}})}{\sum_{v \in V} \exp(u_v \cdot w_{\text{input}})}$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| *aardvark* | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | ... | 0.4 |

$$W$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| *aardvark* | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$$U$$

it is a far , far better rest that I go to , than I have ever known

$$L_t = -\log p_\theta(x_{t-2} \mid x_t) - \log p_\theta(x_{t-1} \mid x_t)$$
$$- \log p_\theta(x_{t+1} \mid x_t) - \log p_\theta(x_{t+2} \mid x_t)$$

# skipgram

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_\text{out} \cdot w_\text{input})}{\sum_{v \in V} \exp(u_v \cdot w_\text{input})}$$

|  | | | | |
|---|---|---|---|---|
| *a* | 1.2 | -0.1 | 0.3 | ... 0.1 |
| *aardvark* | 0.2 | 0.7 | -0.4 | ... 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | ... -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | ... 0.7 |
| ⋮ | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | ... 0.4 |

$W$

|  | | | | |
|---|---|---|---|---|
| *a* | 2.1 | -0.5 | 1.3 | ... 1.4 |
| *aardvark* | -0.4 | -0.7 | 0.5 | ... 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | ... -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | ... 0.4 |
| ⋮ | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | ... 0.6 |

$U$

it is a far , far better rest that I go to , than I have ever known

$$L_t = -\log p_\theta(x_{t-2} \mid x_t) - \log p_\theta(x_{t-1} \mid x_t)$$
$$- \log p_\theta(x_{t+1} \mid x_t) - \log p_\theta(x_{t+2} \mid x_t)$$

# skipgram

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_\text{out} \cdot w_\text{input})}{\sum_{v \in V} \exp(u_v \cdot w_\text{input})}$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| *aardvark* | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | ... | 0.4 |

$W$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| *aardvark* | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$U$

it is a far , far better rest that I go to , than I have ever known

$$L_t = -\log p_\theta(x_{t-2} \mid x_t) - \log p_\theta(x_{t-1} \mid x_t)$$
$$- \log p_\theta(x_{t+1} \mid x_t) - \log p_\theta(x_{t+2} \mid x_t)$$

# skipgram

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_{\text{out}} \cdot w_{\text{input}})}{\sum_{v \in V} \exp(u_v \cdot w_{\text{input}})}$$

|  | | | | | |
| --- | --- | --- | --- | --- | --- |
| *a* | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| *aardvark* | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | ... | 0.4 |

$$W$$

|  | | | | | |
| --- | --- | --- | --- | --- | --- |
| *a* | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| *aardvark* | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$$U$$

it is a far , far better rest that I go to , than I have ever known

$$L_t = -\log p_\theta(x_{t-2} \mid x_t) - \log p_\theta(x_{t-1} \mid x_t)$$
$$- \log p_\theta(x_{t+1} \mid x_t) - \log p_\theta(x_{t+2} \mid x_t)$$

# skipgram

$$p_\theta(\text{out} \,|\, \text{input}) = \frac{\exp(u_{\text{out}} \cdot w_{\text{input}})}{\sum_{v \in V} \exp(u_v \cdot w_{\text{input}})}$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| *aardvark* | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | ... | 0.4 |

$$W$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| *aardvark* | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$$U$$

it is a far , far better rest that I go to , than I have ever known

# skipgram

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_{\text{out}} \cdot w_{\text{input}})}{\sum_{v \in V} \exp(u_v \cdot w_{\text{input}})}$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| *aardvark* | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | ... | 0.4 |

$$W$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| *aardvark* | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$$U$$

it is a far , far better rest that I go to , than I have ever known

# CBOW (Continuous Bag-of-Words)

Use the context to predict the center word

| | a | 1.2 | -0.1 | 0.3 | ... | 0.1 |
| | aardvark | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| | able | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| | are | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| | ⋮ | | | ⋮ | | |
| | zyzzyva | 0.3 | -0.2 | 0.7 | ... | 0.4 |

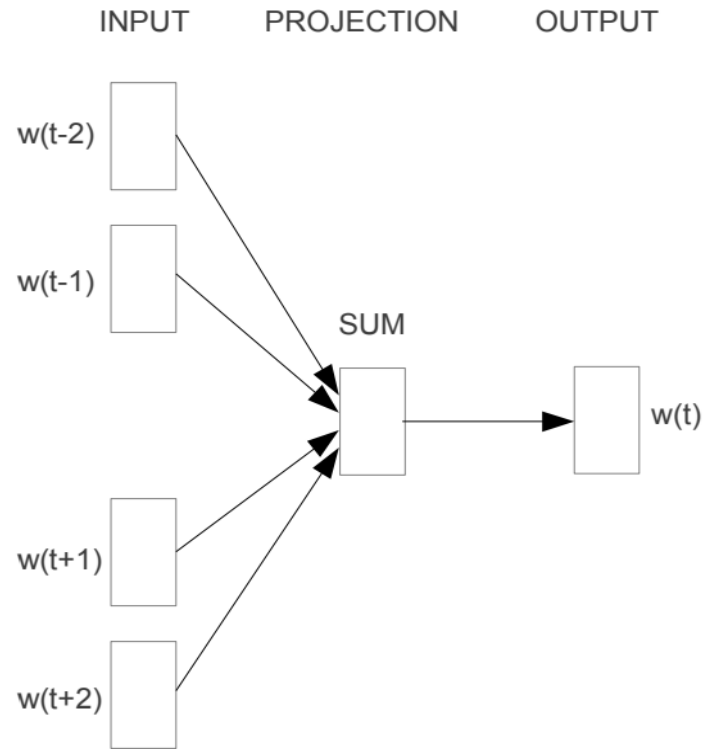| | a | 2.1 | -0.5 | 1.3 | ... | 1.4 |
| | aardvark | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| | able | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| | are | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| | ⋮ | | | ⋮ | | |
| | zyzzyva | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$$p_\theta(x_t \mid x_{t-w}, \ldots, x_{t+w}) \propto \exp\left(u_{x_t} \cdot \frac{1}{2w} \sum_{k \in \{-w, \ldots, -1, 1, w\}} w_{x_{t+k}}\right)$$
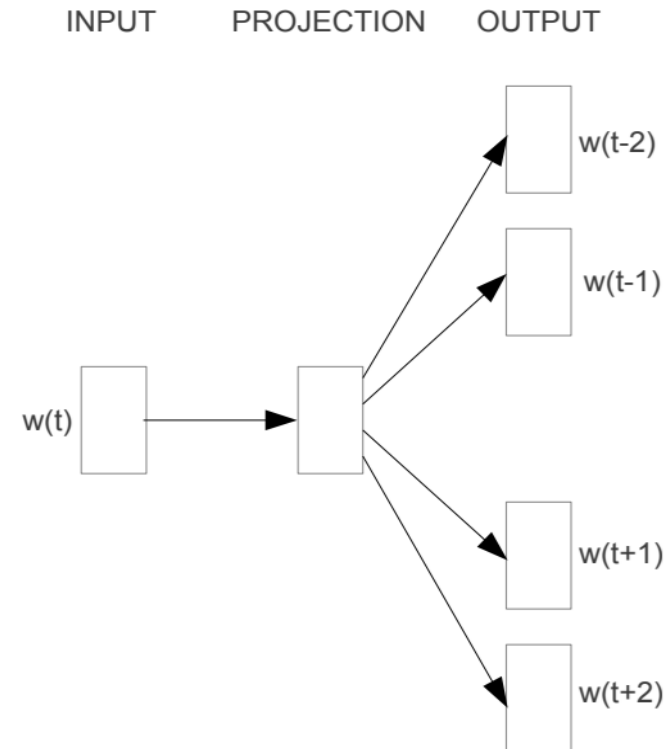
$$W$$

$$U$$

it is a far , far better rest that I go to , than I have ever known

$$L_t = -\log p_\theta(x_t \mid x_{t-2}, x_{t-1}, x_{t+1}, x_{t+2})$$

# CBOW (Continuous Bag-of-Words)

Use the context to predict the center word

| | 1.2 | -0.1 | 0.3 | ... | 0.1 |
|---|---|---|---|---|---|
| *a* | | | | | |
| *aardvark* | 0.2 | 0.7 | -0.4 | ... | 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | ... | -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | ... | 0.7 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | ... | 0.4 |

| | 2.1 | -0.5 | 1.3 | ... | 1.4 |
|---|---|---|---|---|---|
| *a* | | | | | |
| *aardvark* | -0.4 | -0.7 | 0.5 | ... | 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | ... | -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | ... | 0.4 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | ... | 0.6 |

$$p_\theta(x_t \mid x_{t-w}, \ldots, x_{t+w}) \propto \exp\left(u_{x_t} \cdot \frac{1}{2w} \sum_{k\in\{-w,\ldots,-1,1,w\}} w_{x_{t+k}}\right) \qquad W \qquad\qquad U$$

it is a far , far better rest that I go to , than I have ever known

$$L_t = -\log p_\theta(x_t \mid x_{t-2}, x_{t-1}, x_{t+1}, x_{t+2})$$

# word2vec



[Mikolov et al. 2013]

# skipgram w/ Negative Sampling

$$p_\theta(\text{out} \mid \text{input}) = \frac{\exp(u_{\text{out}} \cdot w_{\text{input}})}{\sum_{v \in V} \exp(u_v \cdot w_{\text{input}})}$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 1.2 | -0.1 | 0.3 | . . . | 0.1 |
| *aardvark* | 0.2 | 0.7 | -0.4 | . . . | 1.1 |
| *able* | -0.7 | 0.5 | 0.6 | . . . | -0.8 |
| *are* | 0.1 | 0.9 | 0.8 | . . . | 0.7 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | 0.3 | -0.2 | 0.7 | . . . | 0.4 |

$$W$$

|  | | | | | |
|---|---|---|---|---|---|
| *a* | 2.1 | -0.5 | 1.3 | . . . | 1.4 |
| *aardvark* | -0.4 | -0.7 | 0.5 | . . . | 0.1 |
| *able* | 0.2 | 0.1 | 0.4 | . . . | -0.7 |
| *are* | 0.5 | 0.8 | 0.1 | . . . | 0.4 |
| ⋮ | | | ⋮ | | |
| *zyzzyva* | -0.3 | 0.3 | 0.2 | . . . | 0.6 |

$$U$$

- Vocabulary size V: 50K – 30M
- Very expensive O(|V|)

# skipgram w/ Negative Sampling

- Treat the target word and a neighboring context word as positive examples $(x, y)$

- Randomly sample other words outside of context to get negative samples $(x, v)$

- learn to distinguish between true pair $(x, y)$ and negative samples $(x, v)$ with a binary classifier

- New objective

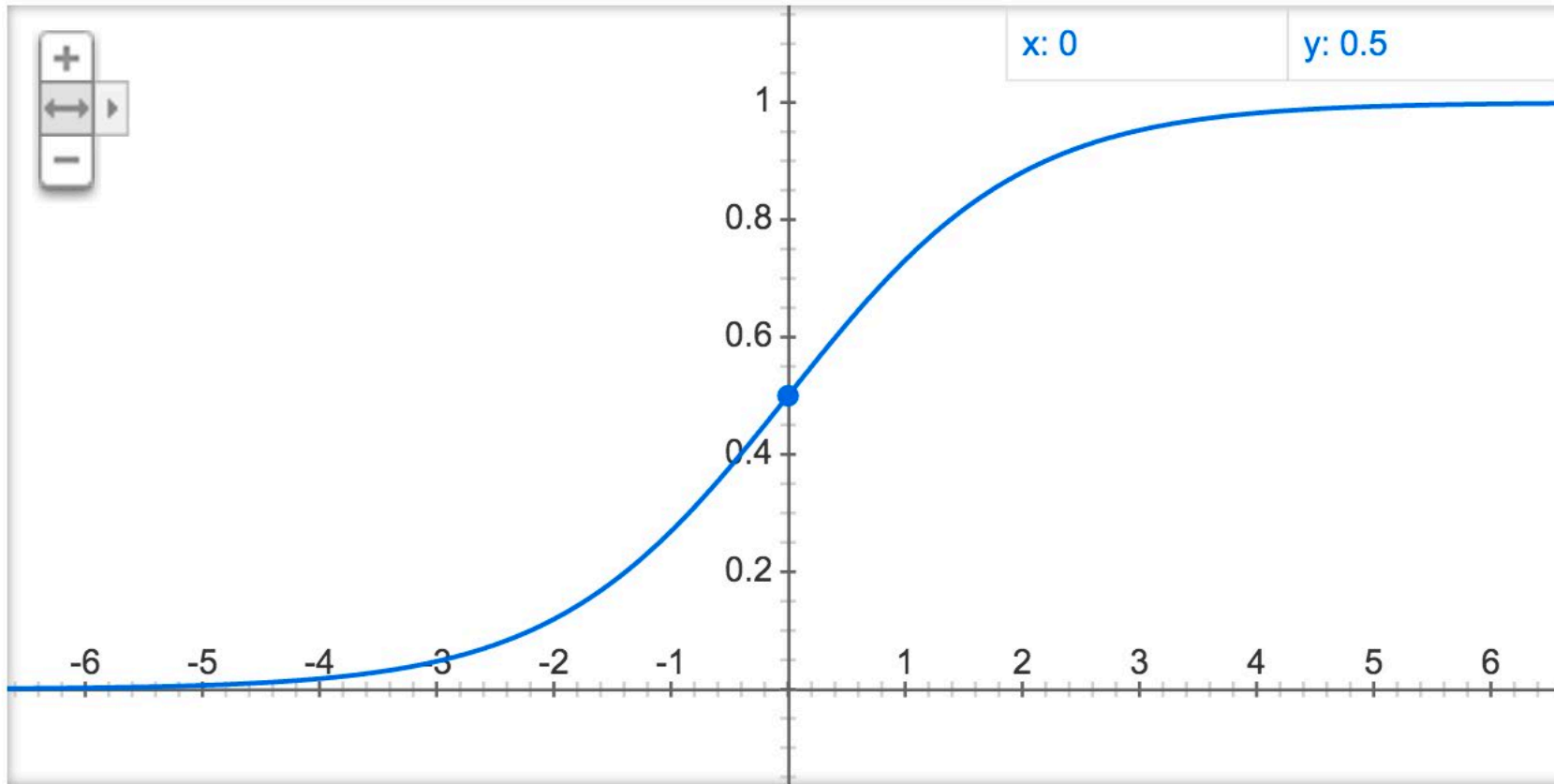$$\log p((x, y) \text{ is a true pair}) + \sum_{k \in C} \log p((x, k) \text{ is a negative pair})$$

C = Negative Samples

# skipgram w/ Negative Sampling

- Treat the target word and a neighboring context word as positive examples $(x, y)$

- Randomly sample other words outside of context to get negative samples $(x, v)$

- learn to distinguish between true pair $(x, y)$ and negative samples $(x, v)$ with a binary classifier

- New objective

$$\log p((x, y) \text{ is a true pair}) + \sum_{k \in C} \log p((x, k) \text{ is a negative pair})$$

$$p((x, c) \text{ is a true pair}) = \sigma(u_c \cdot w_x) = \frac{1}{1 + \exp(-u_c \cdot w_x)}$$

# skipgram w/ Negative Sampling

- Treat the target word and a neighboring context word as positive examples $(x, y)$

- Randomly sample other words outside of context to get negative samples $(x, v)$

- learn to distinguish between true pair $(x, y)$ and negative samples $(x, v)$ with a binary classifier

- New objective

$$\log p((x, y) \text{ is a true pair}) + \sum_{k \in C} \log p((x, k) \text{ is a negative pair})$$

$$= \log \sigma(u_y \cdot w_x) + \sum_{k \in C} \log(\sigma(-u_k \cdot w_x))$$

(logistic) sigmoid: $\sigma(x) = \dfrac{1}{1 + \exp\{-x\}}$

# skipgram w/ Negative Sampling

- Treat the target word and a neighboring context word as positive examples $(x, y)$

- Randomly sample other words outside of context to get negative samples $(x, v)$

- learn to distinguish between true pair $(x, y)$ and negative samples $(x, v)$ with a binary classifier

- New objective

$$\log p((x, y) \text{ is a true pair}) + \sum_{k \in C} \log p((x, k) \text{ is a negative pair})$$

**Much cheaper to compute: O(|C|)**

# Choosing negative samples

- According to unigram probabilities  *P(w)*
- More common to choose from a flattened version

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

- From Mikolov et al. (2013):
  - α=0.75 works well empirically (why?)
  - Usually 2-20 sampled negative words

# Contrastive Learning

- Learning to contrast positive vs. negative samples is a very powerful idea!
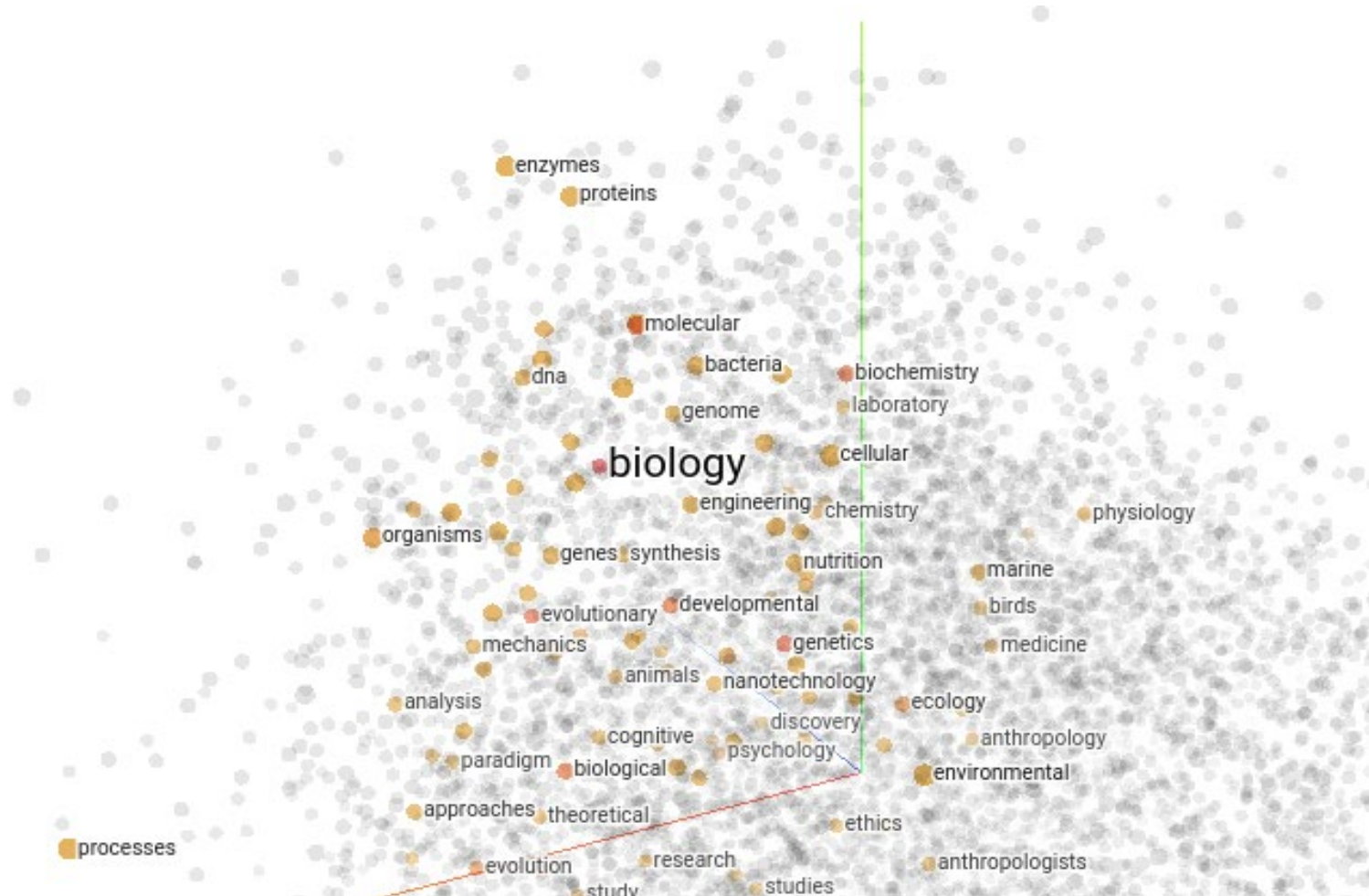
o Representation learning



(a) Original  (b) Crop and resize  (c) Crop, resize (and flip)  (d) Color distort. (drop)  (e) Color distort. (jitter)

(f) Rotate $\{90°, 180°, 270°\}$  (g) Cutout  (h) Gaussian noise  (i) Gaussian blur  (j) Sobel filtering

[SimCLR; Chen et al. 2020]

# Contrastive Learning

- Learning to contrast positive vs. negative samples is a very powerful idea!

o Density estimation: Generative Adversarial Networks (GANs)



[Goodfellow et al. 2014]

# Word2vec Embeddings



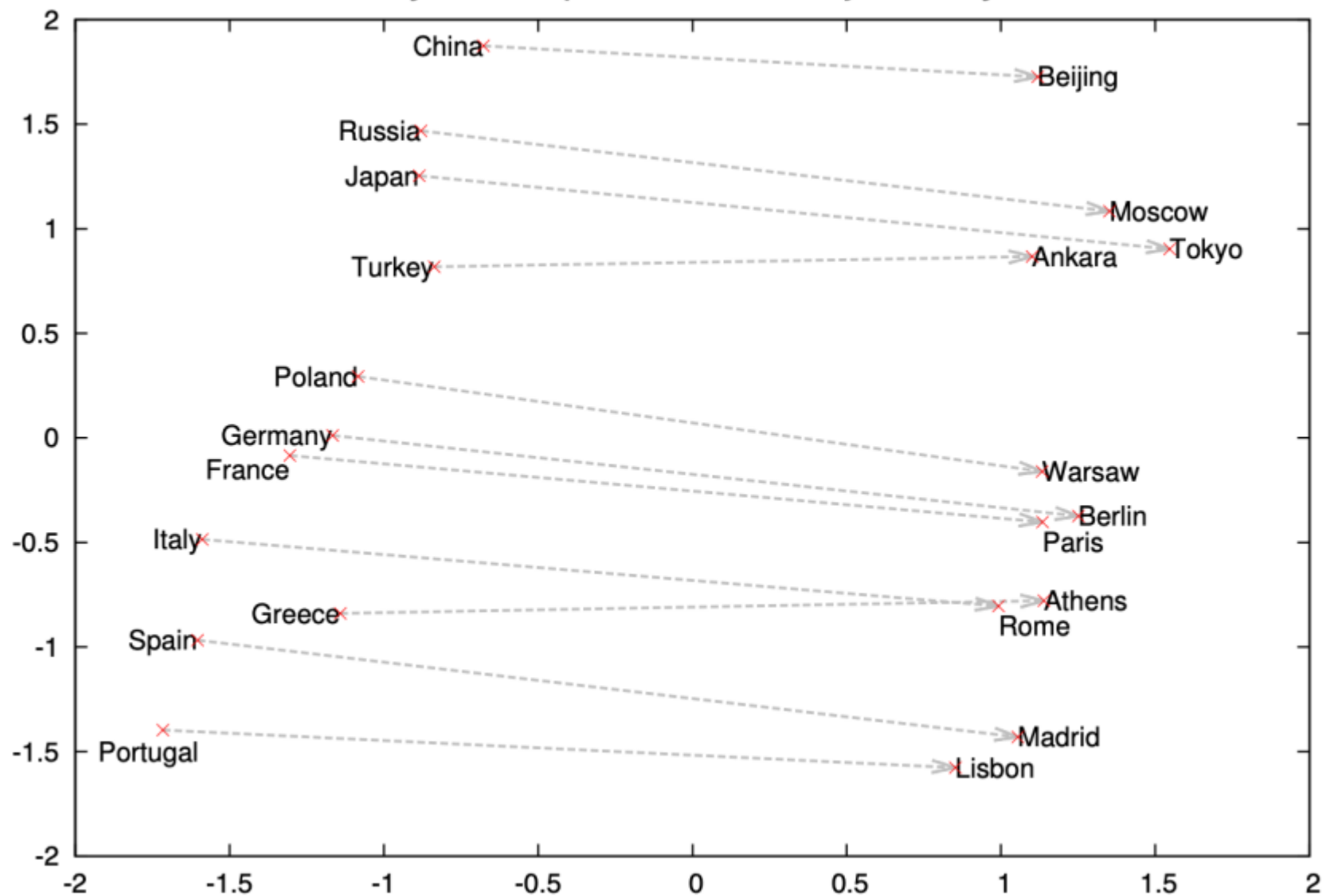[https://projector.tensorflow.org/]

# Word2vec Embeddings

- Regularities in the vector space correspond to regularities in language space!



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

$$w_{\text{apple}} - w_{\text{apples}} \approx w_{\text{car}} - w_{\text{cars}}$$

Country and Capital Vectors Projected by PCA

[Mikolov et al. 2013]

# Word Embeddings You can Download

- word2vec [Mikolov et al. 2013]:

  https://code.google.com/archive/p/word2vec/


- GloVe [Pennington et al. 2014]:

  https://nlp.stanford.edu/projects/glove/
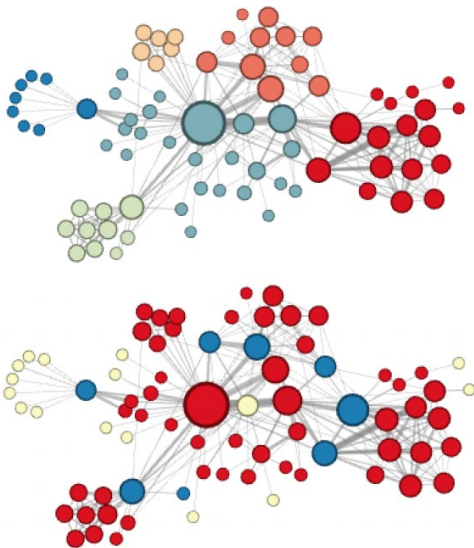

- fasttext [Bojanowsi et al. 2017]:
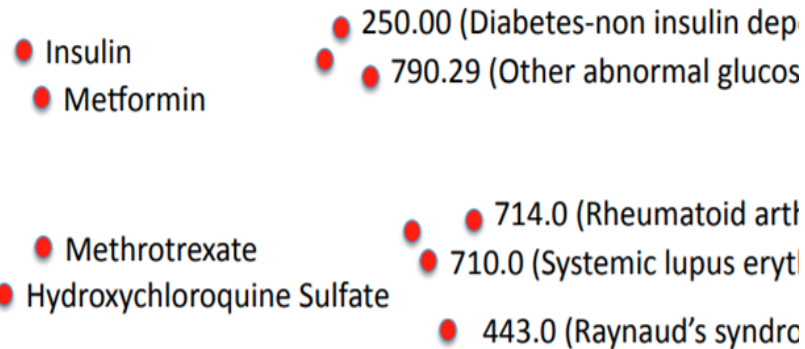
  https://fasttext.cc/

# Extensions

- Neural word embeddings

  - Multilingual?

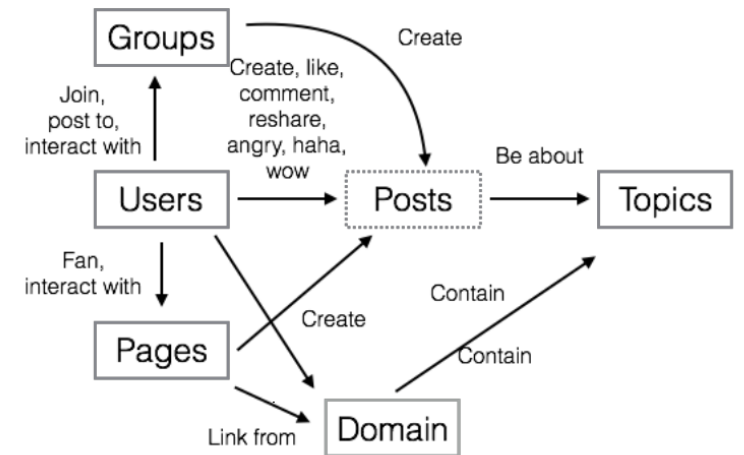  - Social biases?

  - Contextualized?

# Extensions

"You shall know ~~a word~~ by the company it keeps" anything?



Node2Vec
[Grover and Leskovec 2016]

Concept2Vec
[Choi et al. 2016]

World2Vec
[Facebook AI Research]

# Summary (1/2)

- Annotated Database for Lexical Semantics -- **WordNet**

- **Word vectors**: the use of a vector of numbers to represent a word

- **Distributional word vectors**: the use of distributional statistics (e.g., word co-occurrence counts) in defining word vectors

- Computing similarity of two vectors:
  - **dot product** is a simple starting point
  - **cosine similarity** accounts for vector length and works better for word vectors

# Summary (2/2)

- Simple distributional word vectors: word-word co-occurrence counts

- Improving by reducing the influence of common context words
- **TF-IDF (Term Frequency – Inverse Document Frequency)**
- **PMI (Pointwise Mutual Information)**

- Learning representations from data: **word2vec**
  - **skipgram** (w/ negative sampling)
  - **CBOW (Continuous Bag-of-Words)**