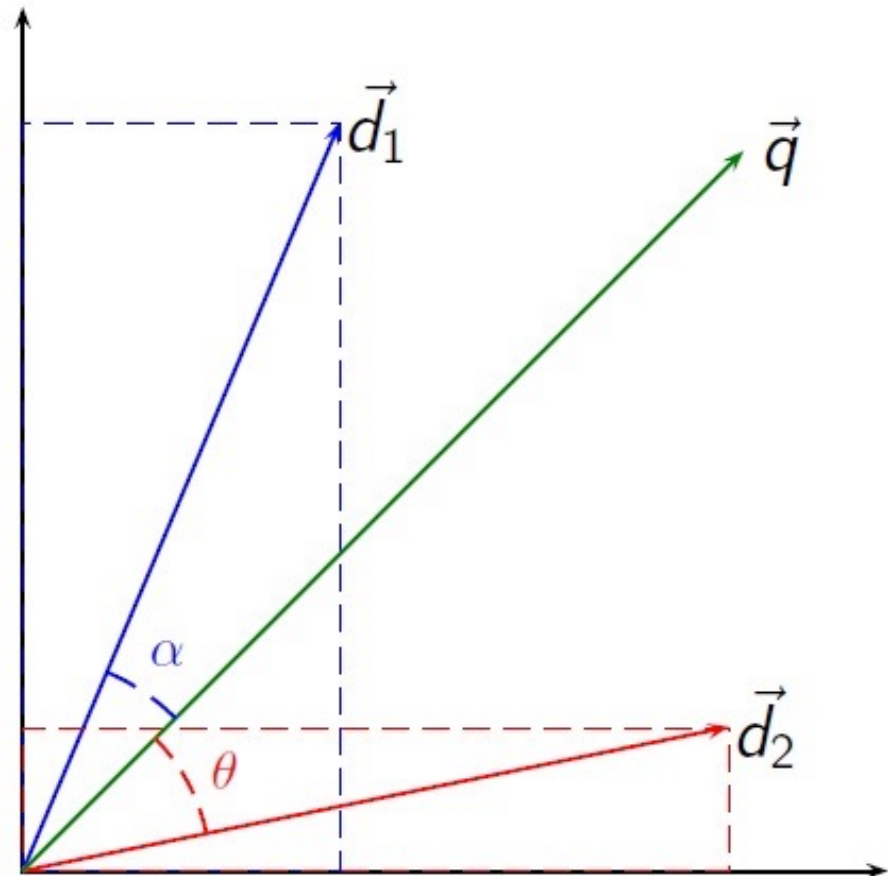


Lecture 28: Vector Semantics

Mark Hasegawa-Johnson
4/2024

CC0: Public domain. Re-
use, Remix, Redistribute at
will.



By Riclas - Own work, CC BY 3.0,
<https://commons.wikimedia.org/w/index.php?curid=9076846>

Outline

- What is a word? wordforms vs. lemmas vs. word senses
- What is meaning? synonymy, similarity, and relatedness
- Vector semantics: CBOW and skip-gram
- Generative training vs. Contrastive training
- Visualizations
- Bias

What is a word?

[w] word - Wiktionary

https://en.wiktionary.org/wiki/word

visibility

- Show translations
- Show declension
- Show quotations
- Show derived terms

In other languages

- Deutsch
- Español
- Français
- 한국어
- Italiano
- Русский
- GWY
- Tiếng Việt
- 中文

🗨️ 78 more

Print/export

- Create a book
- Download as PDF
- Printable version

If you have time, leave us a note.

Noun [edit]

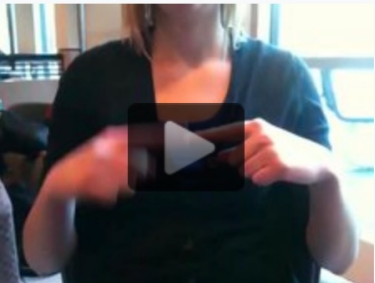
word (*countable and uncountable, plural words*)

- The smallest unit of language that has a particular meaning and can be expressed by itself; the smallest **discrete, meaningful** unit of language. (*contrast morpheme.*) [quotations ▼]
 - The smallest discrete unit of spoken language with a particular meaning, composed of one or more **phonemes** and one or more **morphemes** [quotations ▼]
 - The smallest discrete unit of written language with a particular meaning, composed of one or more **letters** or **symbols** and one or more **morphemes** [quotations ▼]
 - A discrete, meaningful unit of language approved by an **authority** or **native speaker** (*compare non-word*). [quotations ▼]
- Something like such a unit of language:
 - A **sequence** of **letters**, characters, or **sounds**, considered as a **discrete entity**, though it does not necessarily belong to a language or have a meaning [quotations ▼]

Examples

The word *inventory* may be pronounced with four syllables (/ɪn.vən.tɔ.ɪ/) or only three (/ɪn'ven.tɪ/).

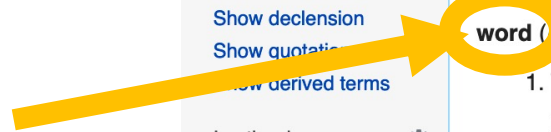
The word *island* is six letters long; the *s* has never been pronounced but was added under the influence of *isle*.



The word *about* signed in American Sign Language.

What is a word?

Is this a word?



The screenshot shows the Wiktionary page for the word "word". The browser address bar shows "https://en.wiktionary.org/wiki/word". The page title is "[w] word - Wiktionary". The main content area is titled "Noun [edit]" and lists the word "word" as "countable and uncountable, plural words". It provides three numbered definitions:

1. The smallest unit of language that has a particular meaning and can be expressed by itself; the smallest **discrete, meaningful** unit of language. (*contrast morpheme*.) [quotations ▼]
 1. The smallest discrete unit of spoken language with a particular meaning, composed of one or more **phonemes** and one or more **morphemes** [quotations ▼]
 2. The smallest discrete unit of written language with a particular meaning, composed of one or more **letters** or **symbols** and one or more **morphemes** [quotations ▼]
 3. A discrete, meaningful unit of language approved by an **authority** or **native speaker** (*compare non-word*). [quotations ▼]
2. Something like such a unit of language:
 1. A **sequence** of **letters**, characters, or **sounds**, considered as a **discrete entity**, though it does not necessarily belong to a language or have a meaning [quotations ▼]

The left sidebar contains navigation options: "visibility" (Show translations, Show declension, Show quotations, Show derived terms), "In other languages" (Deutsch, Español, Français, 한국어, Italiano, Русский, GŬŬ, Tiếng Việt, 中文, 78 more), and "Print/export" (Create a book, Download as PDF, Printable version, If you have time, leave us a note).

The right sidebar is titled "Examples" and contains two paragraphs and a video:

The word *inventory* may be pronounced with four syllables (/ɪn.vən.tɔ.ɪ/) or only three (/ɪn'ven.tɪ/).

The word *island* is six letters long; the *s* has never been pronounced but was added under the influence of *isle*.

The word *about* signed in American Sign Language.

What is a word?

Is this a word?

The image shows a browser window displaying the Wiktionary page for the word "word". The page title is "[w] word - Wiktionary" and the URL is "https://en.wiktionary.org/wiki/word". The main content is under the "Noun" section, with the word "word" circled in yellow. A yellow arrow points from the text "Is this a word?" to the circled word. Another yellow circle highlights the phrase "plural words" in the definition, with a yellow arrow pointing from the text "Is this a different word, or the same word?". The definition lists three points: 1. The smallest unit of language that has a particular meaning and can be expressed by itself; the smallest discrete, meaningful unit of language. (contrast morpheme.) [quotations ▼] 2. Something like such a unit of language: 1. A sequence of letters, characters, or sounds, considered as a discrete entity, though it does not necessarily belong to a language or have a meaning [quotations ▼] 2. The smallest discrete unit of spoken language with a particular meaning, composed of one or more phonemes and one or more morphemes [quotations ▼] 3. The smallest discrete unit of written language with a particular meaning, composed of one or more letters or symbols and one or more morphemes [quotations ▼] 3. A discrete, meaningful unit of language approved by an authority or native speaker (compare non-word). [quotations ▼]. To the right of the main content, there are two text boxes. The top one discusses the pronunciation of "inventory" with four or three syllables. The bottom one discusses the word "island" and its pronunciation, mentioning that the 's' has never been pronounced but was added under the influence of "isle". Below this text is a video player showing a person signing, with a caption: "The word about signed in American Sign Language." The left sidebar of the Wiktionary page includes options for translations, declension, quotations, and derived terms, as well as a list of other languages (Deutsch, Español, Français, 한국어, Italiano, Русский, GŬY, Tiếng Việt, 中文) and a button for "78 more".

Wordform

A wordform is a unique sequence of characters.

- Wordforms are much easier for computers to find than lemmas, therefore most automatic processing deals with wordforms.
- ...however, we lose something. “dog” and “dogs” become completely unrelated – as unrelated as “dog” and “exaggerate.”

word (countable and uncountable, plural **words**)

1. The smallest unit of language that has a particular meaning and can be expressed by itself; the smallest **discrete, meaningful** unit of language. (*contrast morpheme*.) [quotations ▼]
 1. The smallest discrete unit of spoken language with a particular meaning, composed of one or more **phonemes** and one or more **morphemes** [quotations ▼]
 2. The smallest discrete unit of written language with a particular meaning, composed of one or more **letters** or **symbols** and one or more **morphemes** [quotations ▼]
 3. A discrete, meaningful unit of language approved by an **authority** or **native speaker** (*compare non-word*). [quotations ▼]
2. Something like such a unit of language:
 1. A **sequence** of **letters**, characters, or **sounds**, considered as a **discrete entity**, though it does not necessarily belong to a language or have a meaning [quotations ▼]

Lemma

A lemma is what humans usually think of as a “word.” It is defined to be the form of the word that appears in a dictionary.

- In dictionaries designed for human beings,...
- other wordforms that can be easily predicted from the lemma are not listed.

word (*countable and uncountable, plural words*)

1. The smallest unit of language that has a particular meaning and can be expressed by itself; the smallest **discrete, meaningful** unit of language. (*contrast morpheme.*) [quotations ▼]
 1. The smallest discrete unit of spoken language with a particular meaning, composed of one or more **phonemes** and one or more **morphemes** [quotations ▼]
 2. The smallest discrete unit of written language with a particular meaning, composed of one or more **letters** or **symbols** and one or more **morphemes** [quotations ▼]
 3. A discrete, meaningful unit of language approved by an **authority** or **native speaker** (*compare non-word*). [quotations ▼]
2. Something like such a unit of language:
 1. A **sequence** of **letters**, characters, or **sounds**, considered as a **discrete entity**, though it does not necessarily belong to a language or have a meaning [quotations ▼]

What is a word?

Is this a word?

Are these the same word, or different words?

The image shows a browser window displaying the Wiktionary page for the word "word". The page title is "[w] word - Wiktionary" and the URL is "https://en.wiktionary.org/wiki/word". The page content includes a sidebar on the left with options like "Show translations", "Show declension", and "Show quotations". The main content area is titled "Noun [edit]" and lists three definitions of "word".

Annotations in yellow highlight the following elements:

- The word "word" in the title "Noun [edit] word".
- The phrase "countable and uncountable, plural words" in the first definition.
- The first definition: "1. The smallest unit of language that has a particular meaning and can be expressed by itself; the smallest discrete, meaningful unit of language. (contrast morpheme.) [quotations ▼]"
- The second definition: "2. The smallest discrete unit of spoken language with a particular meaning, composed of one or more phonemes and one or more morphemes. [quotations ▼]"
- The third definition: "3. A sequence of letters, characters, or sounds, considered as a discrete entity, though it does not necessarily belong to a language or have a meaning [quotations ▼]"

Text on the right side of the image asks: "Is this a different word, or the same word?". Below this text is a video player showing a person signing, with a caption: "The word *about* signed in American Sign Language."

Word sense

Often, a word has different meanings that are completely unrelated. We think of them as different words, that just happen to be spelled and pronounced the same way.

We say that these are different “senses” of the same word.



The Bank of England. By Diliff - Own work, CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=40912212>



The Bank of the Thames. By Diliff - Own work, CC BY 3.0, <https://commons.wikimedia.org/w/index.php?curid=3639626>

Wordform, lemma, and word sense

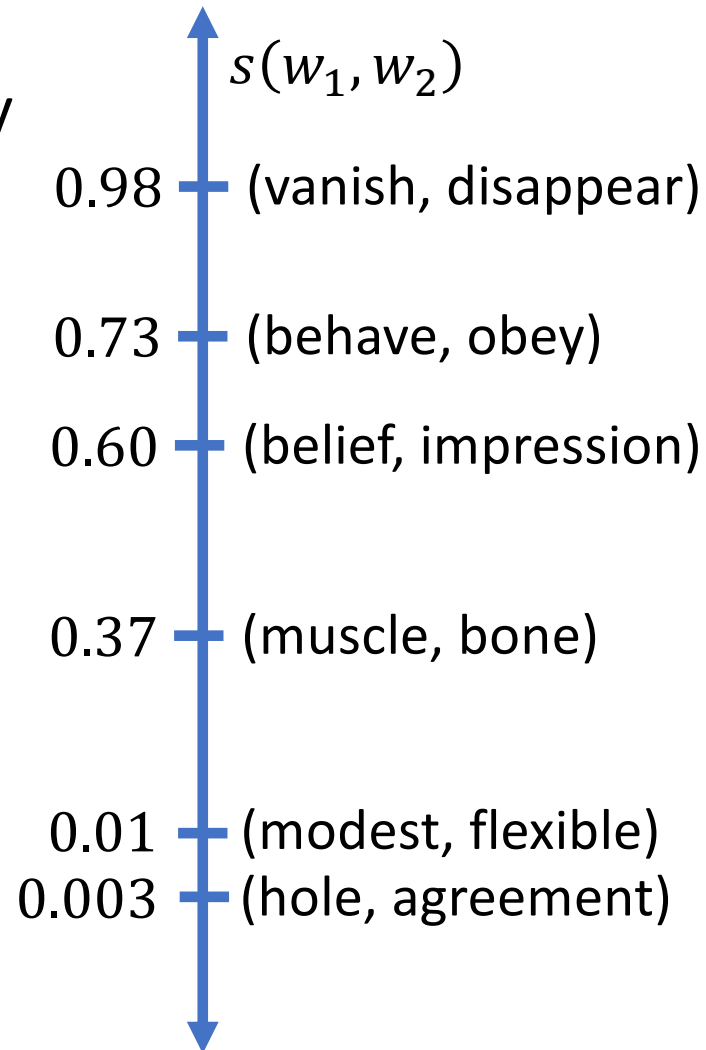
- wordform
 - easy for a computer to work with: just look for space-bounded sequences of characters
- lemma
 - This is what humans think of as a word. A set of wordforms whose spellings, pronunciations, and meanings can all be derived from one another by applying simple rules.
- word sense
 - A meaning so distinct from the other meanings of the word that it's hard to consider them the same word.

Outline

- What is a word? wordforms vs. lemmas vs. word senses
- What is meaning? synonymy, similarity, and relatedness
- Vector semantics: CBOW and skip-gram
- Generative training vs. Contrastive training
- Visualizations
- Bias

Synonymy and similarity

- Words are “synonyms” if they have exactly the same meaning.
- No words ever have ***exactly*** the same meaning, so no two words are ever exactly synonyms.
- We prefer to talk about word similarity, $0 \leq s(w_1, w_2) \leq 1$
 - $s(w_1, w_2) = 1$: w_1 and w_2 are perfect synonyms. Never happens in practice, but sometimes close.
 - $s(w_1, w_2) = 0$: w_1 and w_2 are completely different.



SimLex-999

SimLex-999 is a gold standard resource for the evaluation of models that learn the meaning of words and concepts.

SimLex-999 provides a way of measuring how well models capture *similarity*, rather than *relatedness* or *association*. The scores in SimLex-999 therefore differ from other well-known evaluation datasets such as *WordSim-353* (Finkelstein et al. 2002). The following two example pairs illustrate the difference - note that *clothes* are not similar to *closets* (different materials, function etc.), even though they are very much related:

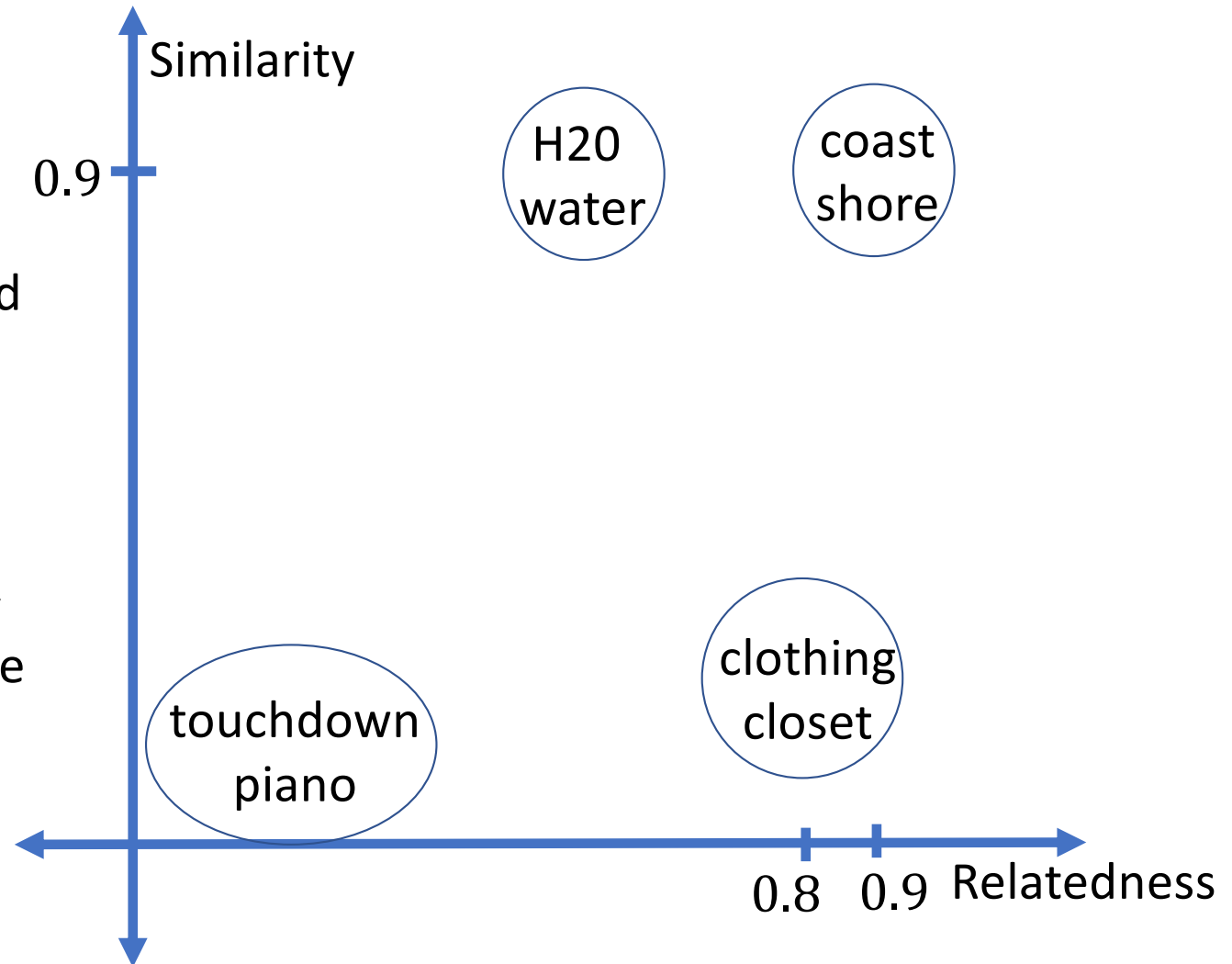
Pair	Simlex-999 rating	WordSim-353 rating
<i>coast - shore</i>	9.00	9.10
<i>clothes - closet</i>	1.96	8.00

- Algorithms that try to estimate the similarity of two wordforms can be tested on databases such as SimLex-999.
- Humans rated the similarity of each word pair on a 10-point scale.

Similarity vs. Relatedness

Similar: words can be used interchangeably in most contexts

Related: there is some connection between the two words, such that they tend to appear in the same documents.



Outline

- What is a word? wordforms vs. lemmas vs. word senses
- What is meaning? synonymy, similarity, and relatedness
- Vector semantics: CBOW and skip-gram
- Generative training vs. Contrastive training
- Visualizations
- Bias

Review: Naïve Bayes: the “Bag-of-words” model

We can estimate the likelihood of an e-mail by pretending that the e-mail is just a bag of words (order doesn't matter).

With only a few thousand spam e-mails, we can get a pretty good estimate of these things:

- $P(W = \text{“hi”}|Y = \text{spam}), P(W = \text{“hi”}|Y = \text{ham})$
- $P(W = \text{“vitality”}|Y = \text{spam}), P(W = \text{“vitality”}|Y = \text{ham})$
- $P(W = \text{“production”}|Y = \text{spam}), P(W = \text{“production”}|Y = \text{ham})$

Then we can approximate $P(X|Y)$ by assuming that the words, W , are **conditionally independent of one another given the category label**:

$$P(X = x|Y = y) \approx \prod_{i=1}^n P(W = w_i|Y = y)$$



Similarity: The Internet is the database

Similarity = words can be used interchangeably in most contexts

How do we measure that in practice?

Answer: extract examples of word w_1 , +/- C words (C=2 or 3):

...hot, although iced coffee is a popular...

...indicate that moderate coffee consumption is benign...

...and of w_2 :

...consumed as iced tea. Sweet tea is...

...national average of tea consumption in Ireland...

The words “iced” and “consumption” appear in both contexts, so we can conclude that $s(\text{coffee}, \text{tea}) > 0$. No other words are shared, so we can conclude $s(\text{coffee}, \text{tea}) < 1$.

skip-gram context probability

Consider the “...hot although iced coffee is a popular...”.

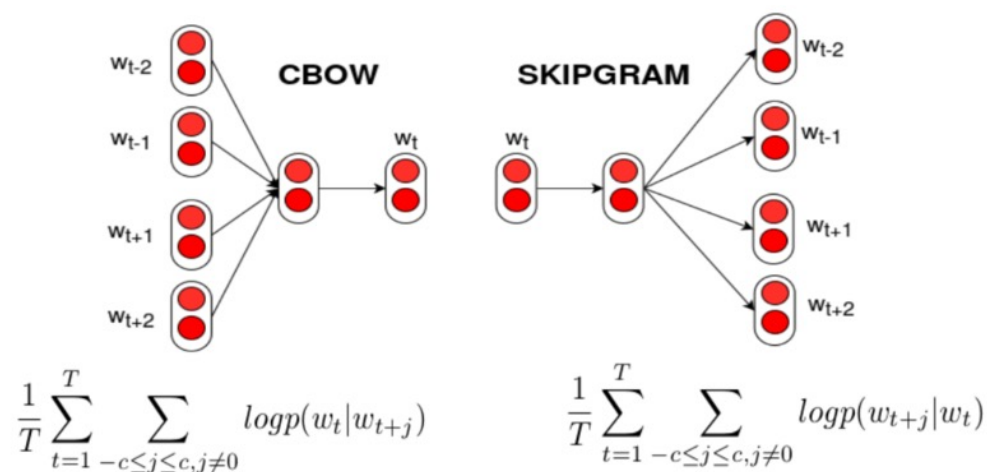
Define the target word to be $w_t = \text{coffee}$.

Define the context words $w_{t-3} = \text{hot}$, $w_{t-2} = \text{although}$, ..., $w_{t+3} = \text{popular}$.

The skip-gram probability is a naïve Bayes model of the context:

$$p(w_{t-3}, \dots, w_{t+3} | w_t) = \prod_{\substack{i \neq 0 \\ i=-3 \\ i=3}}^3 p(w_{t+i} | w_t)$$

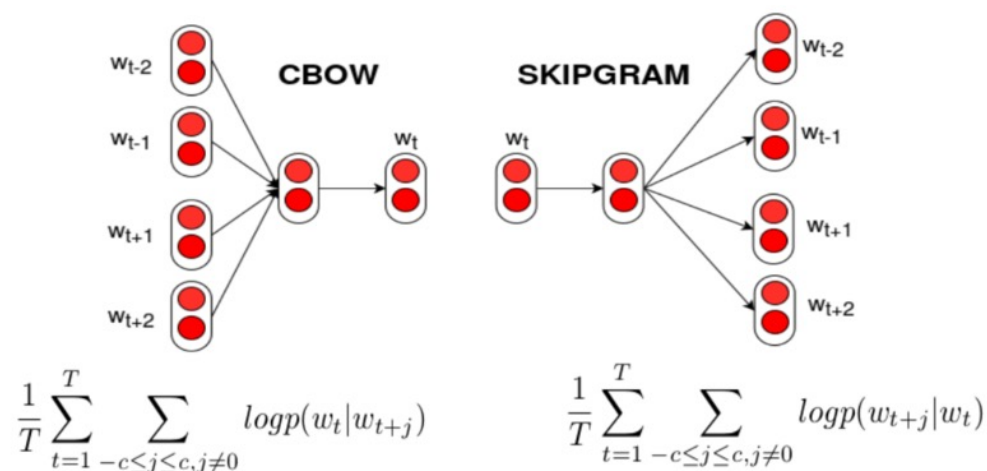
The skip-gram model



- Skip-gram is a model of word meaning:
- The meaning of a word is defined to be the distribution of context words that it can predict.
- We find out which words w_t can predict by learning neural nets that predict its context words w_{t+j} :

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln P(w_{t+j} | w_t)$$

The “continuous bag of words” model (CBOW)



- CBOW is a similar model of word meaning:
- The meaning of a word is defined to be the distribution of context words that predict it the best.
- We find out which words predict w_t by learning neural nets that predict w_t given its context words, w_{t+j} , for $-c \leq j \leq c$:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln P(w_t | w_{t+j})$$

“Probability,” for a NN, means softmax

- What does it mean that we train a neural net to compute $P(w_t|w_{t+j})$?
- It's a probability, so it must mean a softmax:

$$P(W_t = w_t | W_{t+j} = w_{t+j}) = \frac{\exp(e_{t,t+j})}{\sum_{t'} \exp(e_{t',t+j})}$$

- But what are the inputs to the neural net? What is $e_{t,t+j}$?

Vector Semantics

- The simplest useful assumption is this: a word is a vector.

$$P(W_t = w_t | W_{t+j} = w_{t+j}) = \frac{\exp(\mathbf{v}_t^T \mathbf{v}_{t+j})}{\sum_{v \in \mathcal{V}} \exp(\mathbf{v}^T \mathbf{v}_{t+j})}$$

- ...where \mathbf{v}_t is a d-dimensional vector, $\mathbf{v}_t = [v_{t,1}, \dots, v_{t,d}]^T$, and \mathcal{V} is the set of all such vectors (the “vocabulary”)
- The only trainable parameters in this model are the word vectors!
- The dictionary, \mathcal{V} , is a matrix, with as many columns as there are words in the vocabulary:

$$\mathcal{V} = [\text{vec}("a"), \dots, \text{vec}("zzz")] = \begin{bmatrix} \text{vec}("a")_1 & \cdots & \text{vec}("zzz")_1 \\ \vdots & \ddots & \vdots \\ \text{vec}("a")_d & \cdots & \text{vec}("zzz")_d \end{bmatrix}$$

cosine similarity

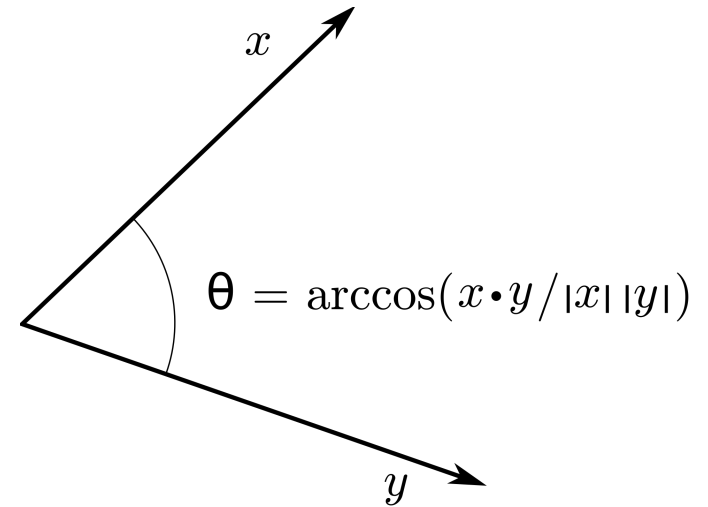
If words w_1 and w_2 are similar, w_1 is represented by vector \mathbf{v}_1 , and w_2 by vector \mathbf{v}_2 , then the angle between the two vectors should be small.

Angle between two vectors can be measured by their dot product:

$$\cos \theta = \frac{\mathbf{v}_1^T \mathbf{v}_2}{|\mathbf{v}_1| |\mathbf{v}_2|}$$

where

$$\mathbf{v}_1^T \mathbf{v}_2 = \sum_{i=1}^d v_{1,i} v_{2,i}, \quad |\mathbf{v}_1| = \sqrt{\sum_{i=1}^d v_{1,i}^2}$$



By BenFrantzDale at the English Wikipedia, CC BY-SA 3.0,
<https://commons.wikimedia.org/w/index.php?curid=49972362>

Vector Semantics: Variations

There are many ways to make this model more flexible. For example:

- Every word could have two different vectors: one (\mathbf{v}) for when it's being predicted, one (\mathbf{v}') for when it is predicting, thus $e_{t,t+j} = \mathbf{v}_t^T \mathbf{v}'_{t+j}$.
- We could weight the dot-product depending on the delay between the word vectors, thus $e_{t,j,t+j} = \mathbf{v}_t^T \mathbf{W}(j) \mathbf{v}'_{t+j}$.
- We could use a two-layer network to calculate the similarity, for example, $e_{t,t+j} = \mathbf{w}_2^T \max\left(\mathbf{0}, \mathbf{W}_1(j) \begin{bmatrix} \mathbf{v}_t \\ \mathbf{v}'_{t+j} \end{bmatrix}\right)$.

Vector Semantics

The basic CBOW probability is:

$$P(W_t = w_t | W_{t+j} = w_{t+j}) = \frac{\exp(\mathbf{v}_t^T \mathbf{v}_{t+j})}{\sum_{v \in \mathcal{V}} \exp(\mathbf{v}^T \mathbf{v}_{t+j})}$$

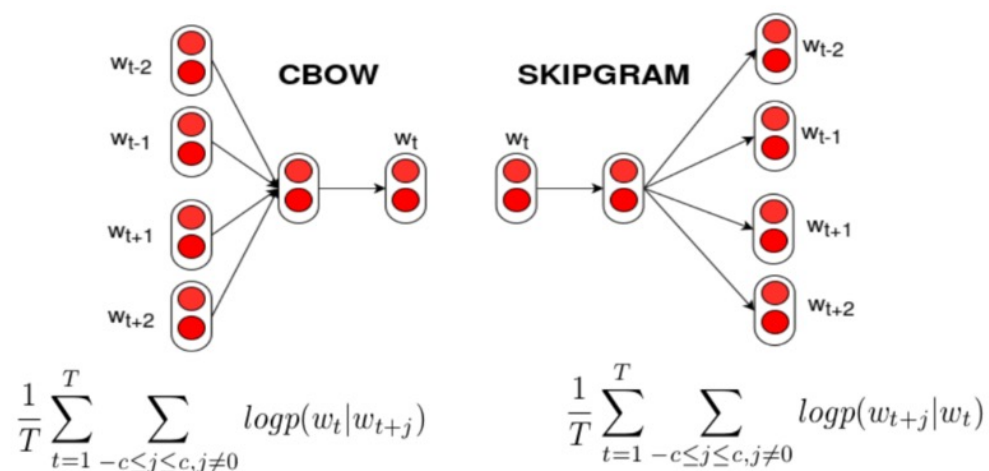
Its logarithm is:

$$\ln P(W_t = w_t | W_{t+j} = w_{t+j}) = \mathbf{v}_t^T \mathbf{v}_{t+j} - \ln \sum_{v \in \mathcal{V}} \exp(\mathbf{v}^T \mathbf{v}_{t+j})$$

The gradient of the log softmax is:

$$\nabla_{\mathbf{v}_t} \ln P(W_t = w_t | W_{t+j} = w_{t+j}) = \mathbf{v}_{t+j} \left(1 - \frac{\exp(\mathbf{v}_t^T \mathbf{v}_{t+j})}{\sum_{v \in \mathcal{V}} \exp(\mathbf{v}^T \mathbf{v}_{t+j})} \right)$$

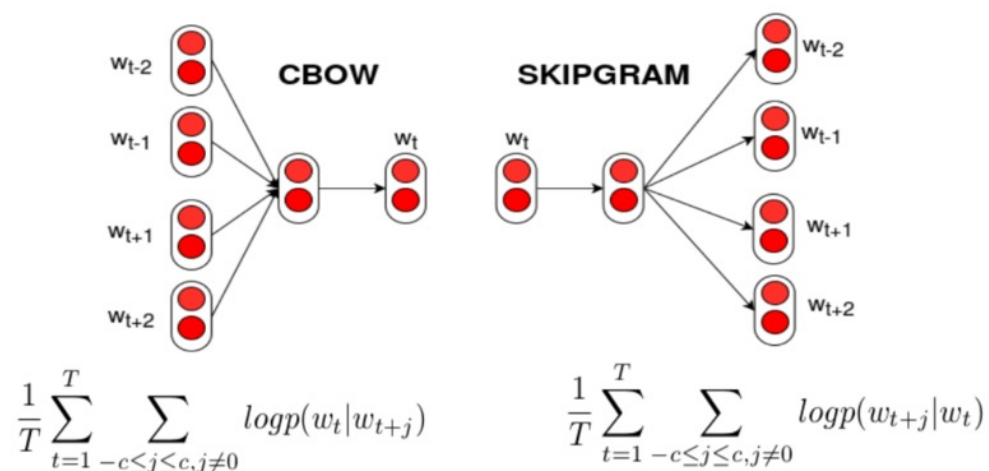
Training a CBOW model



To find the parameters, we use gradient descent:

$$\begin{aligned}
 \nabla_{\mathbf{v}_t} \mathcal{L} &= -\frac{1}{T} \sum_{t: W_t = w_t} \sum_{j=-c, j \neq 0}^c \nabla_{\mathbf{v}_t} \ln P(W_t = w_t | W_{t+j} = w_{t+j}) \\
 &= -\frac{1}{T} \sum_{t: W_t = w_t} \sum_{j=-c, j \neq 0}^c \mathbf{v}_{t+j} \left(1 - \frac{\exp(\mathbf{v}_t^T \mathbf{v}_{t+j})}{\sum_{\mathbf{v} \in \mathcal{V}} \exp(\mathbf{v}^T \mathbf{v}_{t+j})} \right)
 \end{aligned}$$

Training a CBOW model



The CBOW model is trained by setting every vector equal to a weighted average of the words that occurred near it!

$$\mathbf{v}_t \leftarrow \mathbf{v}_t - \eta \nabla_{\mathbf{v}_t} \mathcal{L} = \mathbf{v}_t + \frac{\eta}{T} \sum_{t: W_t = w_t} \sum_{j=-c, j \neq 0}^c \mathbf{v}_{t+j} \left(1 - P(W_t = w_t | W_{t+j} = w_{t+j}) \right)$$

- There is more weight on words that don't yet predict well ($P(w_t | w_{t+j})$ small).
- The weight is accumulated over the corpus, so if a word occurs near w_t often, then it gets more total weight.

Try the quiz!

Try the quiz:

https://us.prairielearn.com/pl/course_instance/147925/assessment/2411691

Outline

- What is a word? wordforms vs. lemmas vs. word senses
- Synonymy, similarity, and relatedness
- Vector semantics: CBOW and skip-gram
- **Generative training vs. Contrastive training**
- **Visualizations**
- **Bias**

Contrastive loss vs. Generative loss

- A generative loss is one like this:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=1}^T \sum_{j=-c, j \neq 0}^c \ln \frac{\exp(\mathbf{v}_t^T \mathbf{v}_{t+j})}{\sum_{\mathbf{v} \in \mathcal{V}} \exp(\mathbf{v}^T \mathbf{v}_{t+j})}$$

- Notice that this loss term compares each word, w_t , to every other word in the dictionary.
- Sometimes, generative training can take a very long time to converge.
- Sometimes, we get faster training using contrastive loss.

Contrastive loss

We train the neural network by listing, as positive examples, the words that occur in the context of “ $w_t = \text{coffee}$,” e.g.,

$$\mathcal{D}_+(w_t) = \{w_{t-3}, w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2}, w_{t+3}\}$$

Create a contrastive database by choosing the same number of words, at random, from outside that context window

$$\mathcal{D}_-(w_t) = \{\text{aardvark, dog, gazebo, actor, precipitates, iceberg}\}$$

Contrastive loss

The generative loss is based on the probability of generating w_t , which, for CBOW, is:

$$P(W_t = w_t | W_{t+j} = w_{t+j}) = \frac{\exp(\mathbf{v}_t^T \mathbf{v}_{t+j})}{\sum_{\mathbf{v} \in \mathcal{V}} \exp(\mathbf{v}^T \mathbf{v}_{t+j})}$$

The contrastive loss is based on the probability of the sets $\mathcal{D}_+(w_t)$ and $\mathcal{D}_-(w_t)$. For skip-gram, we could write:

$$\begin{aligned} \Pr(\mathcal{D}_+(w_t), \mathcal{D}_-(w_t) | w_t) &= \prod_{w' \in \mathcal{D}_+(w_t)} \Pr(w' \in \mathcal{D}_+(w_t) | w_t) \prod_{w' \in \mathcal{D}_-(w_t)} \Pr(w' \in \mathcal{D}_-(w_t) | w_t) \\ &= \prod_{\mathbf{v}' \in \mathcal{D}_+(w_t)} \frac{1}{1 + e^{-\mathbf{v}'^T \mathbf{v}_t}} \prod_{\mathbf{v}' \in \mathcal{D}_-(w_t)} \left(1 - \frac{1}{1 + e^{-\mathbf{v}'^T \mathbf{v}_t}}\right) = \prod_{\mathbf{v}' \in \mathcal{D}_+(w_t)} \frac{1}{1 + e^{-\mathbf{v}'^T \mathbf{v}_t}} \prod_{\mathbf{v}' \in \mathcal{D}_-(w_t)} \frac{1}{1 + e^{\mathbf{v}'^T \mathbf{v}_t}} \end{aligned}$$

Training with contrastive loss

The coefficients $\mathbf{v}_t = [v_{t,1}, \dots, v_{t,d}]^T$ for each vector are chosen to maximize the log probability of the dataset:

$$\begin{aligned}\mathcal{L} &= -\ln p(\text{Data}) = -\frac{1}{T} \sum_{t=1}^T (\ln p(\mathcal{D}_+(w_t)|w_t) + \ln p(\mathcal{D}_-(w_t)|w_t)) \\ &= -\frac{1}{T} \sum_{t=1}^T \left(\sum_{\mathbf{v}' \in \mathcal{D}_+(w_t)} \ln \frac{1}{1 + e^{-\mathbf{v}'^T \mathbf{v}_t}} + \sum_{\mathbf{v}' \in \mathcal{D}_-(w_t)} \ln \frac{1}{1 + e^{\mathbf{v}'^T \mathbf{v}_t}} \right)\end{aligned}$$

Outline

- What is a word? wordforms vs. lemmas vs. word senses
- Synonymy, similarity, and relatedness
- Vector semantics: CBOW and skip-gram
- Generative training vs. Contrastive training
- **Visualizations**
- **Bias**

Visualizations: Similarity

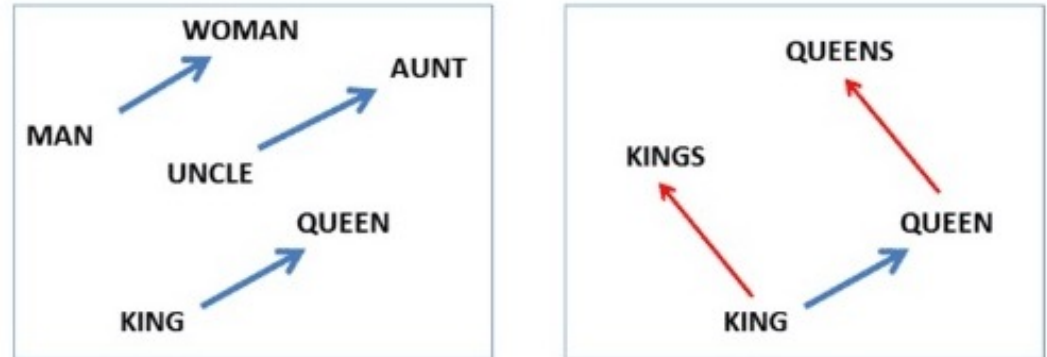
Mikolov et al. (2013) tested word2vec on SimLex-999, and had better results than previously published baselines. Here are some examples from their paper.

Model (training time)	Redmond	Havel	ninjutsu	graffiti	capitulate
Collobert (50d) (2 months)	conyers lubbock keene	plauen dzerzhinsky osterreich	reiki kohona karate	cheesecake gossip dioramas	abdicate accede rearm
Turian (200d) (few weeks)	McCarthy Alston Cousins	Jewell Arzu Ovitz	- - -	gunfire emotion impunity	- - -
Mnih (100d) (7 days)	Podhurst Harlang Agarwal	Pontiff Pinochet Rodionov	- - -	anaesthetics monkeys Jews	Mavericks planning hesitated
Skip-Phrase (1000d, 1 day)	Redmond Wash. Redmond Washington Microsoft	Vaclav Havel president Vaclav Havel Velvet Revolution	ninja martial arts swordsmanship	spray paint grafitti taggers	capitulation capitulated capitulating

Table 6: Examples of the closest tokens given various well known models and the Skip-gram model trained on phrases using over 30 billion training words. An empty cell means that the word was not in the vocabulary.

Visualizations: Relatedness

$$\text{vec}(\text{"woman"}) - \text{vec}(\text{"man"}) + \text{vec}(\text{"king"}) = \text{vec}(\text{"queen"})$$



Christian S. Perone, "Voynich Manuscript: word vectors and t-SNE visualization of some patterns," in *Terra Incognita*, 16/01/2016, <http://blog.christianperone.com/2016/01/voynich-manuscript-word-vectors-and-t-sne-visualization-of-some-patterns/>.

Mikolov (2013) showed that word2vec captures similarity relationships among words. For example, the difference between the vectors for "woman" and "man" is roughly the same as the difference between the vectors for "queen" and "king." Perone (2016) showed that this effect works differently depending on the training corpus: in his blog post, he looks at word relatedness in the 15th century Voynich manuscript.

Outline

- What is a word? wordforms vs. lemmas vs. word senses
- Synonymy, similarity, and relatedness
- Vector semantics: CBOW and skip-gram
- Generative training vs. Contrastive training
- Visualizations
- **Bias**

Learning biased analogies from data

- It's useful that algorithms like word2vec learn appropriate analogies, like "Paris → France as Tokyo → Japan" and "kings → king as queens → queen."
- Unfortunately, it also learns other analogies that were implied in the training corpus, but that are invalid analogies.
- The paper that first demonstrated that problem was named after one of the worst such discovered analogies:

"Man is to Computer Programmer as Woman is to Homemaker?
Debiasing Word Embeddings," Bolukbasi et al., 2016

Biased analogies

- Bolukbasi et al. defined a “male-female” continuum by subtracting $\text{vec}(\text{“female”}) - \text{vec}(\text{“male”})$, $\text{vec}(\text{“woman”}) - \text{vec}(\text{“man”})$, and so on, then averaging these difference vectors.
- They then created a “neutral-specific” continuum by averaging gender-specific words, averaging gender-neutral words, and subtracting.
 - Gender-specific: dictionary definition includes gender-specific language
 - Gender-neutral: all other words

The Male-Female vs. Neutral-Specific Space

Here's the resulting 2D space, from Bolukbasi et al., 2016:



Summary

- What is a word? Lemmas, wordforms, and word sense
- Synonymy, similarity, and relatedness
- Context bag-of-words (CBOW), generative loss:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=1}^T \sum_{j=-c, j \neq 0}^c \ln \frac{\exp(\mathbf{v}_t^T \mathbf{v}_{t+j})}{\sum_{\mathbf{v} \in \mathcal{V}} \exp(\mathbf{v}^T \mathbf{v}_{t+j})}$$

- Skip-gram, contrastive loss:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=1}^T \left(\sum_{\mathbf{v}' \in \mathcal{D}_+(w_t)} \ln \frac{1}{1 + e^{-\mathbf{v}'^T \mathbf{v}_t}} + \sum_{\mathbf{v}' \in \mathcal{D}_-(w_t)} \ln \frac{1}{1 + e^{\mathbf{v}'^T \mathbf{v}_t}} \right)$$

- Visualizations: $\text{vec}(\text{"aunt"}) - \text{vec}(\text{"uncle"}) + \text{vec}(\text{"king"}) = \text{vec}(\text{"queen"})$
- $\text{vec}(\text{"male"}) - \text{vec}(\text{"female"})$: differences OK for words whose dictionary definition includes gender, but not for other words