

Vector Semantics and Embeddings, Part 2

Kenneth Lai

Brandeis University

November 14, 2022

Announcements

- ▶ For Wednesday
 - ▶ HW4 due
- ▶ Course registration for spring 2023 begins tomorrow!

Today's Plan

- ▶ Vector Semantics and Embeddings
 - ▶ TF-IDF
 - ▶ word2vec
 - ▶ Properties of Embeddings

Vector Semantics & Embeddings

TF-IDF

But raw frequency is a bad representation

- The co-occurrence matrices we have seen represent each cell by word frequencies.
- Frequency is clearly useful; if *sugar* appears a lot near *apricot*, that's useful information.
- But overly frequent words like *the*, *it*, or *they* are not very informative about the context
- It's a paradox! How can we balance these two conflicting constraints?

Two common solutions for word weighting

tf-idf: tf-idf value for word t in document d :

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

Words like "the" or "it" have very low idf

PMI: (Pointwise mutual information)

- $\text{PMI}(w_1, w_2) = \log \frac{p(w_1, w_2)}{p(w_1)p(w_2)}$

See if words like "good" appear more often with "great" than we would expect by chance

Term frequency (tf)

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

Instead of using raw count, we squash a bit:

$$tf_{t,d} = \log_{10}(\text{count}(t,d)+1)$$

Document frequency (df)

df_t is the number of documents t occurs in.

(note this is not collection frequency: total count across all documents)

"*Romeo*" is very distinctive for one Shakespeare play:

	Collection Frequency	Document Frequency
Romeo	113	1
action	113	31

Inverse document frequency (idf)

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

N is the total number of documents
in the collection

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

What is a document?

Could be a play or a Wikipedia article

But for the purposes of tf-idf, documents can be **anything**; we often call each paragraph a document!

Final tf-idf weighted value for a word

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

Raw counts:

	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:

	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

Vector Semantics & Embeddings

TF-IDF

Vector Semantics & Embeddings

PPMI

Pointwise Mutual Information

Pointwise mutual information:

Do events x and y co-occur more than if they were independent?

$$\text{PMI}(X, Y) = \log_2 \frac{P(x, y)}{P(x)P(y)}$$

PMI between two words: (Church & Hanks 1989)

Do words x and y co-occur more than if they were independent?

$$\text{PMI}(\text{word}_1, \text{word}_2) = \log_2 \frac{P(\text{word}_1, \text{word}_2)}{P(\text{word}_1)P(\text{word}_2)}$$

Positive Pointwise Mutual Information

- PMI ranges from $-\infty$ to $+\infty$
- But the negative values are problematic
 - Things are co-occurring **less than** we expect by chance
 - Unreliable without enormous corpora
 - Imagine w_1 and w_2 whose probability is each 10^{-6}
 - Hard to be sure $p(w_1, w_2)$ is significantly different than 10^{-12}
 - Plus it's not clear people are good at "unrelatedness"
- So we just replace negative PMI values by 0
- Positive PMI (**PPMI**) between word1 and word2:

$$\text{PPMI}(word_1, word_2) = \max\left(\log_2 \frac{P(word_1, word_2)}{P(word_1)P(word_2)}, 0\right)$$

Computing PPMI on a term-context matrix

Matrix F with W rows (words) and C columns (contexts)

f_{ij} is # of times w_i occurs in context c_j

$$P_{ij} = \frac{f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}} \quad P_{i*} = \frac{\sum_{j=1}^C f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}} \quad P_{*j} = \frac{\sum_{i=1}^W f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

$$pmi_{ij} = \log_2 \frac{P_{ij}}{P_{i*} P_{*j}} \quad ppmi_{ij} = \begin{cases} pmi_{ij} & \text{if } pmi_{ij} > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$p_{ij} = \frac{f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

$$p(w=\text{information}, c=\text{data}) = 3982/11716 = .3399$$

$$p(w=\text{information}) = 7703/11716 = .6575$$

$$p(c=\text{data}) = 5673/11716 = .4842$$

$$p(w_i) = \frac{\sum_{j=1}^C f_{ij}}{N} \quad p(c_j) = \frac{\sum_{i=1}^W f_{ij}}{N}$$

	p(w,context)					p(w)
	computer	data	result	pie	sugar	p(w)
cherry	0.0002	0.0007	0.0008	0.0377	0.0021	0.0415
strawberry	0.0000	0.0000	0.0001	0.0051	0.0016	0.0068
digital	0.1425	0.1436	0.0073	0.0004	0.0003	0.2942
information	0.2838	0.3399	0.0323	0.0004	0.0011	0.6575
p(context)	0.4265	0.4842	0.0404	0.0437	0.0052	

$$pmi_{ij} = \log_2 \frac{P_{ij}}{P_i * P_j}$$

	p(w,context)					p(w)
	computer	data	result	pie	sugar	p(w)
cherry	0.0002	0.0007	0.0008	0.0377	0.0021	0.0415
strawberry	0.0000	0.0000	0.0001	0.0051	0.0016	0.0068
digital	0.1425	0.1436	0.0073	0.0004	0.0003	0.2942
information	0.2838	0.3399	0.0323	0.0004	0.0011	0.6575
p(context)	0.4265	0.4842	0.0404	0.0437	0.0052	

$pmi(\text{information}, \text{data}) = \log_2 (.3399 / (.6575 * .4842)) = .0944$

Resulting PPMI matrix (negatives replaced by 0)

	computer	data	result	pie	sugar
cherry	0	0	0	4.38	3.30
strawberry	0	0	0	4.10	5.51
digital	0.18	0.01	0	0	0
information	0.02	0.09	0.28	0	0

Weighting PMI

PMI is biased toward infrequent events

- Very rare words have very high PMI values

Two solutions:

- Give rare words slightly higher probabilities
- Use add-one smoothing (which has a similar effect)

Weighting PMI: Giving rare context words slightly higher probability

Raise the context probabilities to $\alpha = 0.75$:

$$\text{PPMI}_\alpha(w, c) = \max(\log_2 \frac{P(w, c)}{P(w)P_\alpha(c)}, 0)$$

$$P_\alpha(c) = \frac{\text{count}(c)^\alpha}{\sum_c \text{count}(c)^\alpha}$$

This helps because $P_\alpha(c) > P(c)$ for rare c

Consider two events, $P(a) = .99$ and $P(b) = .01$

$$P_\alpha(a) = \frac{.99^{.75}}{.99^{.75} + .01^{.75}} = .97 \quad P_\alpha(b) = \frac{.01^{.75}}{.01^{.75} + .99^{.75}} = .03$$

Distributed Representations of Words

- ▶ More generally, two approaches to distributed, distributional representations (Baroni et al. 2014):
 - ▶ Count-based
 - ▶ Count occurrences of words in contexts, optionally followed by some mathematical transformation (e.g., tf-idf, PPMI, SVD)
 - ▶ Prediction-based
 - ▶ Given some context vector(s) \mathbf{c} , predict some word \mathbf{x} (or vice versa)
 - ▶ a.k.a. **language modeling**-based

(e.g., word2vec,  , )

Language Models

- ▶ Given some context vector(s) \mathbf{c} , predict some word \mathbf{x} (or vice versa)
- ▶ Two approaches to language models:
 - ▶ **Generative** models
 - ▶ Model the joint probability distribution $P(\mathbf{x}, \mathbf{c})$
 - ▶ Examples: n-gram language models
 - ▶ Unigram: predict $P(\mathbf{x}_i)$
 - ▶ Bigram: predict $P(\mathbf{x}_i | \mathbf{x}_{i-1})$
 - ▶ Trigram: predict $P(\mathbf{x}_i | \mathbf{x}_{i-2}, \mathbf{x}_{i-1})$

Language Models

- ▶ Given some context vector(s) \mathbf{c} , predict some word \mathbf{x} (or vice versa)
- ▶ Two approaches to language models:
 - ▶ **Discriminative** models
 - ▶ Predict the conditional probability $P(\mathbf{x}|\mathbf{c})$ (or $P(\mathbf{c}|\mathbf{x})$) directly
 - ▶ Examples: neural network language models
 - ▶ Feedforward: word2vec (Mikolov et al., 2013a, 2013b)

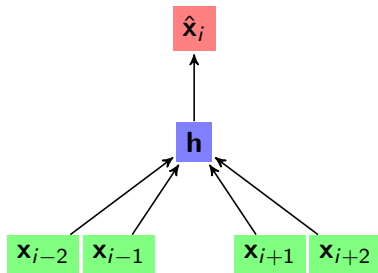


- ▶ Recurrent: (Peters et al., 2018)

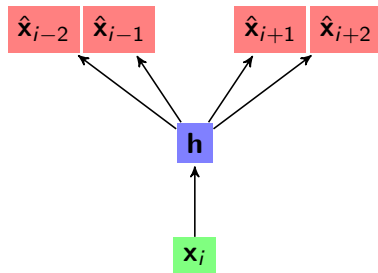


- ▶ Transformer: (Devlin et al., 2019)

- ▶ Based on a feedforward neural network language model

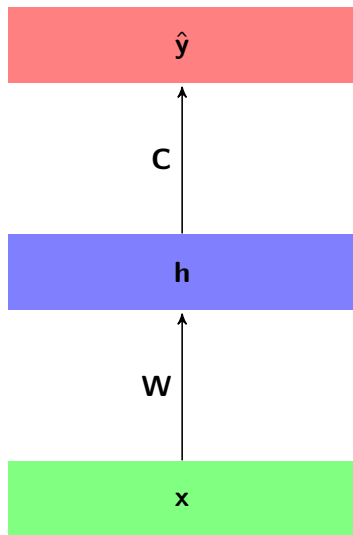


CBOW



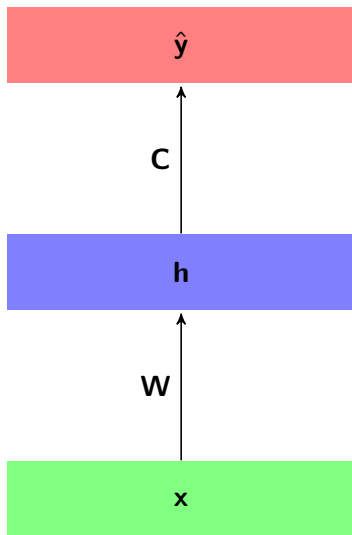
Skip-gram

Neural Networks



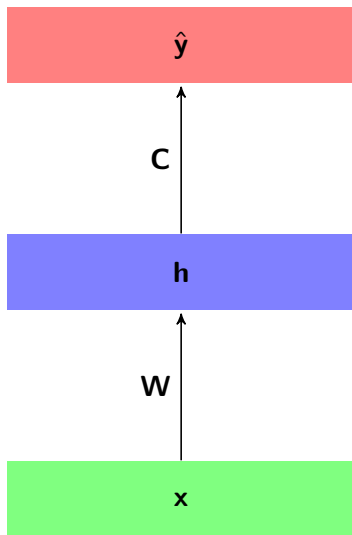
- ▶ Output layer
- ▶ Hidden layer(s)
- ▶ Input layer

Neural Networks



- ▶ Output layer
- ▶ Hidden layer(s)
- ▶ Input layer
- ▶ x is the input
- ▶ h is the hidden layer output
 - ▶ Can be seen as intermediate representation of the input
- ▶ \hat{y} is the predicted output
 - ▶ $\hat{=}$ predicted

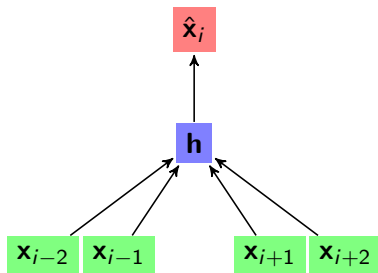
Neural Networks



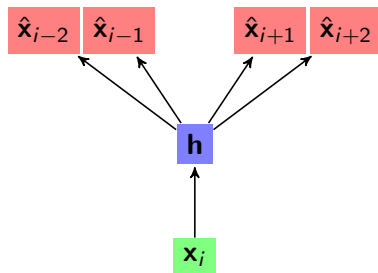
- ▶ **Output layer**
- ▶ **Hidden layer(s)**
- ▶ **Input layer**
- ▶ $\mathbf{h} = g(\mathbf{x} \cdot \mathbf{W})$
- ▶ $\hat{\mathbf{y}} = f(\mathbf{h} \cdot \mathbf{C})$
 - ▶ \mathbf{W} and \mathbf{C} are **weight** (or **parameter**) **matrices**
 - ▶ May or may not include a bias term
 - ▶ g and f are **activation functions**

word2vec

- ▶ Based on a feedforward neural network language model



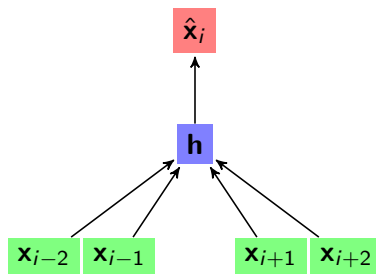
CBOW



Skip-gram

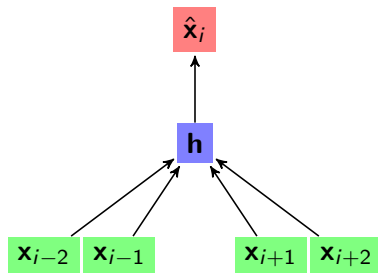
- ▶ Continuous bag of words (CBOW): use context to predict current word
- ▶ Skip-gram: use current word to predict context

CBOW



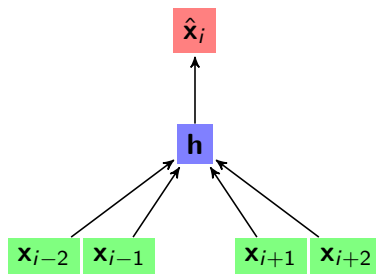
- ▶ Input layer: one-hot word vectors
 - ▶ $[0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]$
 - ▶ Context words within some window

CBOW



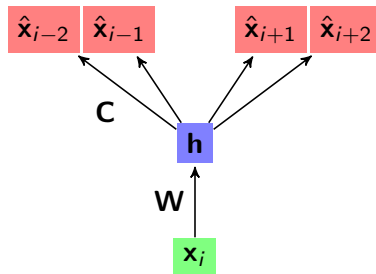
- ▶ Hidden (projection) layer: identity activation function, no bias
 - ▶ Weight matrix shared for all context words
 - ▶ Input \rightarrow hidden = table lookup (in weight matrix)
 - ▶ Context word vectors are averaged

CBOW



- ▶ Output layer: softmax activation function
 - ▶ Numbers \rightarrow probabilities

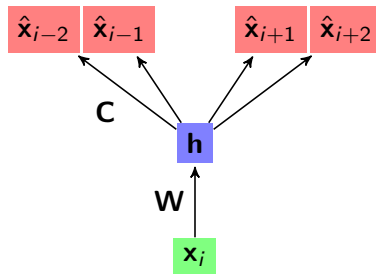
Skip-gram



► Input layer: one-hot word vectors

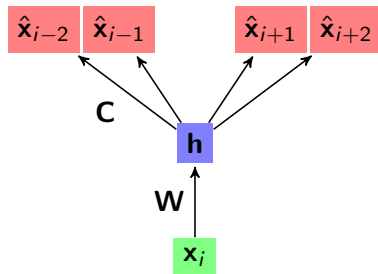
► $[0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]$

Skip-gram



- ▶ Hidden (projection) layer: identity activation function, no bias
 - ▶ Input \rightarrow hidden = table lookup (in weight matrix)

Skip-gram



- ▶ Output layer: softmax activation function
 - ▶ Predict context words within some window
 - ▶ Separate classification for each context word
 - ▶ Closer context words sampled more than distant context words

word2vec

- ▶ Skip-gram model: for each word, word2vec learns two word embeddings
 - ▶ Target word vector \mathbf{w} (row of \mathbf{W} , = output of hidden layer)
 - ▶ Context word vector \mathbf{c} (column of \mathbf{C})
- ▶ Common final word embeddings
 - ▶ Add $\mathbf{w} + \mathbf{c}$
 - ▶ Just \mathbf{w} (throw away \mathbf{c})

Vector
Semantics &
Embeddings

Properties of Embeddings

The kinds of neighbors depend on window size

Small windows ($C = +/- 2$) : nearest words are syntactically similar words in same taxonomy

- *Hogwarts* nearest neighbors are other fictional schools
- *Sunnydale, Evernight, Blandings*

Large windows ($C = +/- 5$) : nearest words are related words in same semantic field

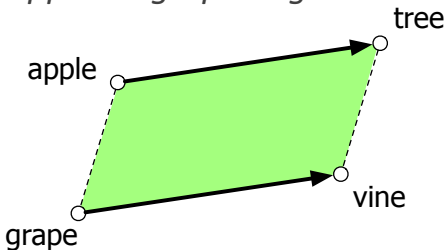
- *Hogwarts* nearest neighbors are Harry Potter world:
- *Dumbledore, half-blood, Malfoy*

Analogical relations

The classic parallelogram model of analogical reasoning
(Rumelhart and Abrahamson 1973)

To solve: "apple is to tree as grape is to _____"

Add $\overrightarrow{\text{tree}} - \overrightarrow{\text{apple}}$ to $\overrightarrow{\text{grape}}$ to get **vine**



Analogical relations via parallelogram

The parallelogram method can solve analogies with both sparse and dense embeddings (Turney and Littman 2005, Mikolov et al. 2013b)

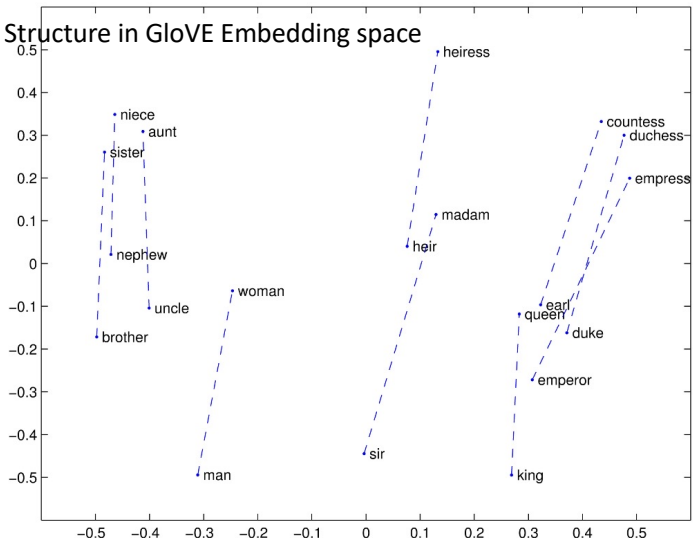
$\vec{\text{king}} - \vec{\text{man}} + \vec{\text{woman}}$ is close to $\vec{\text{queen}}$

$\vec{\text{Paris}} - \vec{\text{France}} + \vec{\text{Italy}}$ is close to $\vec{\text{Rome}}$

For a problem $a : a^* :: b : b^*$, the parallelogram method is:

$$\hat{b}^* = \underset{x}{\operatorname{argmin}} \operatorname{distance}(x, a^* - a + b)$$

Structure in GloVe Embedding space



Caveats with the parallelogram method

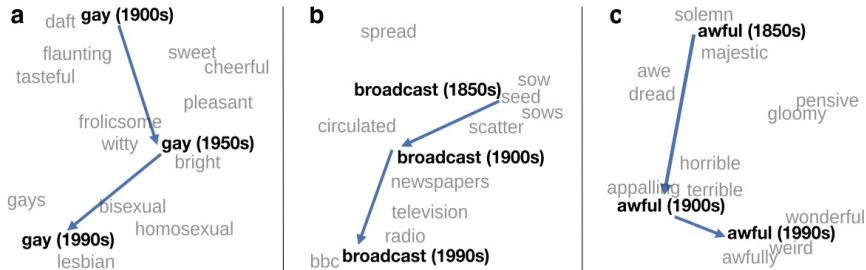
It only seems to work for frequent words, small distances and certain relations (relating countries to capitals, or parts of speech), but not others. (Linzen 2016, Gladkova et al. 2016, Ethayarajh et al. 2019a)

Understanding analogy is an open area of research (Peterson et al. 2020)

Embeddings as a window onto historical semantics

Train embeddings on different decades of historical text to see meanings shift

~30 million books, 1850-1990, Google Books data



William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. Proceedings of ACL.

Embeddings reflect cultural bias!

Bolukbasi, Tolga, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. "Man is to computer programmer as woman is to homemaker? debiasing word embeddings." In *NeurIPS*, pp. 4349-4357. 2016.

Ask "Paris : France :: Tokyo : x"

- x = Japan

Ask "father : doctor :: mother : x"

- x = nurse

Ask "man : computer programmer :: woman : x"

- x = homemaker

Algorithms that use embeddings as part of e.g., hiring searches for programmers, might lead to bias in hiring

Historical embedding as a tool to study cultural biases

Garg, N., Schiebinger, L., Jurafsky, D., and Zou, J. (2018). Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences* 115(16), E3635–E3644.

- Compute a **gender or ethnic bias** for each adjective: e.g., how much closer the adjective is to "woman" synonyms than "man" synonyms, or names of particular ethnicities
 - Embeddings for **competence** adjective (*smart, wise, brilliant, resourceful, thoughtful, logical*) are biased toward men, a bias slowly decreasing 1960-1990
 - Embeddings for **dehumanizing** adjectives (barbaric, monstrous, bizarre) were biased toward Asians in the 1930s, bias decreasing over the 20th century.
- These match the results of old surveys done in the 1930s

Vector
Semantics &
Embeddings

Properties of Embeddings