

AOS2– Deep learning

Lecture 03: Word embeddings

Sylvain Rousseau

Natural language processing

- What is natural language processing (NLP)?
 - Field of artificial intelligence
 - Read and understand human languages
- NLP is hard
 - Homographs like the word “set”
 - “I made her duck” (at least 5 different meanings)
 - No strict rules: “How R U?”

NLP tasks

Some NLP subtasks:

- Part-of-speech (POS) tagging
label tokens in a text based on their role: nouns, verbs, adjectives, adverbs, punctuation...
- Named Entity Recognition (NER)
identifying specific words or phrases (“entities”) and categorizing them as persons, locations...
- Topic modelling
identify topics in a text
- Sentiment analysis
interpret and classify emotions in subjective data.
- ...

Most of them rely on good representation of words !

Word vector representation

Notations

- A **word** as a symbol is denoted w or w_t if it is part of a sequence
- A **word vector** is in bold, $\boldsymbol{w} \in \mathbb{R}^N$

$$\boldsymbol{w} = (-0.3320, 0.3401, \dots, -1.6941)$$

- **Vocabulary**: list of known words of length V

$$\mathcal{V} = \{w_1, \dots, w_V\}$$

- **Embedding**: a vector for each word in the vocabulary

One-hot representation

- To each word w_k of the vocabulary \mathcal{V} , we associate the vector $\mathbf{w}_k \in \mathbb{R}^{\mathcal{V}}$ such that

$$\mathbf{w}_k = \overset{1}{(0, 0, \dots, 0, \overset{i(k)}{1}, 0, \dots, 0)}^{\mathcal{V}} \quad \mathbf{w}_{kj} = \begin{cases} 0 & \text{if } j \neq i(k) \\ 1 & \text{if } j = i(k) \end{cases}$$

- Properties
 - A text can be represented of a matrix of size $\mathcal{V} \times L$
 - Extremely sparse representation
- Drawbacks:
 - Very inefficient representation as \mathcal{V} is typically large $\mathcal{V} \sim 10^{5-7}$
 - All words are equidistant ($\|\mathbf{w}_i - \mathbf{w}_j\| = \delta_{ij}\sqrt{2}$)
 - Impossible to add a new word without changing the dimension of the representation

Word embedding

- Starts with the **distributional hypothesis** stating that words in similar contexts have similar meanings. Captured in a famous quote by Firth:
 “You shall know a word by the company it keeps”, Firth 1957
- Count-based approaches (counting co-occurrences, frequencies of words)
 - Co-occurrence + SVD
 - GloVe Pennington, Socher, and Manning 2014
- Prediction-based approaches (optimizing a loss function)
 - Word2vec Mikolov et al. 2013
 - FastText Bojanowski et al. 2017
 - ELMO Peters et al. 2018
 - Bert Devlin et al. 2019
 - ...

Window-based context

- The context of a word is defined as a window of specific size around that word
- Window of size 5 center at “jumps”

w_{t-2} w_{t-1} w_t w_{t+1} w_{t+2}

The brown fox jumps over the lazy dog

- Notations (window of size $2m + 1$)
 - **Center word:** w_t
 - **Context:** $w_{t-m}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+m}$

Co-occurrence matrix

- Co-occurrence of center word and words in its context
- Example corpus composed of 3 sentences:

- I like deep learning .

- I like NLP .

- I enjoy flying .

- Vocabulary (8 tokens)

$$\mathcal{V} = \left\{ \text{I}, \text{like}, \text{enjoy}, \text{deep}, \text{learning}, \text{NLP}, \text{flying}, \text{.} \right\}$$

Example

Window of size 3: $m = 1$

	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

- Properties
 - Similar words have similar representation
 - Co-occurrence matrix is symmetric
 - Co-occurrence matrix is sparse
- Drawbacks:
 - Static representation
 - Co-occurrence matrix is huge: size of C is $\mathcal{V} \times \mathcal{V}$

Dimension reduction

We want to reduce the number of columns of C and retain most of the information: PCA!

- To perform a PCA we do a SVD decomposition of C

$$C = USU^T$$

where $U^T U = I$, ($U = V$ because C is symmetric)

- Instead of C we choose N first columns of US

$$\tilde{C} = (US)_{[1:N]}$$

where N is typically 25–1000.

- \tilde{C} is of size $V \times N$ instead of $V \times V$
- \tilde{C} is the best rank- N matrix preserving correlations (CC^T)

word2vec

Prediction-based approach

Principle: learn to model words through a prediction task

- **Word2vec** from Mikolov et al. 2013
 - Simplest prediction-based approach
 - Two different architectures
 - Continuous Bag Of Words (CBOW)
 - Skip-gram
 - Variants:
 - Negative sampling
 - Hierarchical softmax

Continuous Bag of Words (CBOW)

- Principle: model the distribution of **center word conditional to its context**

$$p(w \mid w_{t-m}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+m})$$

Or $p(w \mid \mathcal{C})$ with $\mathcal{C} = \{w_{t-m}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+m}\}$

- Idea: use a neural network to learn that conditional distribution given given one-hot encoded context words

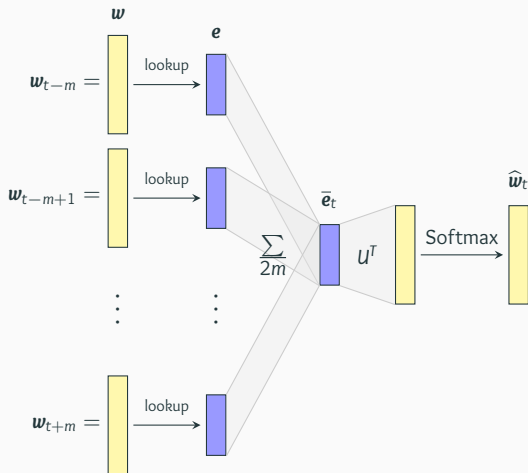
CBOW model

- Use a neural network to model $p(w \mid C)$:

- Inputs are one-hot encoded context words: w_{t+k} , $k \in \{-m, \dots, m\}$, $k \neq 0$
- First (shallow) layer: simple lookup
- First real layer: averaging
- Second layer: linear transform
- Output layer: softmax to get a distribution on the vocabulary: \hat{w}_t

- Parameters

- Embeddings: e_{t+k} for $k \in \{-m, \dots, m\}$, $k \neq 0$
- Matrix: U



Layers

- First (shallow) layer: lookup table $E = [\mathbf{e}_1, \dots, \mathbf{e}_V]$, (size of E is $N \times V$)

$$\mathbf{e}_{t-k} = E\mathbf{w}_{t-k}$$

- First real layer: averaging

$$\bar{\mathbf{e}}_t = \frac{1}{2m} \sum_{\substack{k=-m \\ k \neq 0}}^m \mathbf{e}_{t+k}$$

- Second layer: linear transform

$$U^T \bar{\mathbf{e}}_t$$

- Third layer: softmax

$$\hat{\mathbf{w}}_t = \text{softmax}(U^T \bar{\mathbf{e}}_t)$$

- Relationship between one-hot encoded context words and “one-hot” prediction of expected word

$$\hat{\mathbf{w}}_t = \text{softmax} \left(\mathbf{u}^T \cdot \left(\frac{1}{2m} \sum_{k=-m}^m E \mathbf{w}_{t+k} \right) \right)$$

- Matrix E is the embedding matrix: $\mathbf{e}_i = E \mathbf{w}_i$

$$\hat{\mathbf{w}}_t = \text{softmax} \left(\mathbf{u}^T \cdot \left(\frac{1}{2m} \sum_{k=-m}^m \mathbf{e}_{t+k} \right) \right), \quad \hat{\mathbf{w}}_t = \text{softmax} (\mathbf{u}^T \bar{\mathbf{e}}_t)$$

Expanded form

- Softmax is done on values:

$$\langle \mathbf{u}_i, \bar{\mathbf{e}}_t \rangle \quad \text{for } i = 1, \dots, \mathcal{V}$$

where $U = [\mathbf{u}_1, \dots, \mathbf{u}_n]$

- Per-word expanded version (probability that expected word w_t is w_i)

$$\hat{w}_{ti} = \frac{\exp(\langle \mathbf{u}_i, \bar{\mathbf{e}}_t \rangle)}{\sum_{k=1}^{\mathcal{V}} \exp(\langle \mathbf{u}_k, \bar{\mathbf{e}}_t \rangle)}$$

- Compare predicted distribution $\hat{\mathbf{w}}_t$ over vocabulary with one-hot encoded expected center word \mathbf{w}_t with cross-entropy:

$$H(\mathbf{w}_t, \hat{\mathbf{w}}_t) = - \sum_{j=1}^V \mathbf{w}_{tj} \log \hat{\mathbf{w}}_{tj}$$

- Since \mathbf{w}_t is a one-hot vector, it reduces to

$$H(\mathbf{w}_t, \hat{\mathbf{w}}_t) = - \log \hat{\mathbf{w}}_{t,i(t)} \quad (i(t) \text{ is the one-hot index of token } w_t)$$

- Final loss is

$$H(\mathbf{w}_t, \hat{\mathbf{w}}_t) = - \langle \mathbf{u}_{i(t)}, \bar{\mathbf{e}}_t \rangle + \log \left(\sum_{k=1}^V \exp(\langle \mathbf{u}_k, \bar{\mathbf{e}}_t \rangle) \right)$$

- Parameters to optimize are: \mathbf{e}_{t+k} for $k \in \{-m, \dots, m\}, k \neq 0$, \mathbf{u}_k for $k \in V$

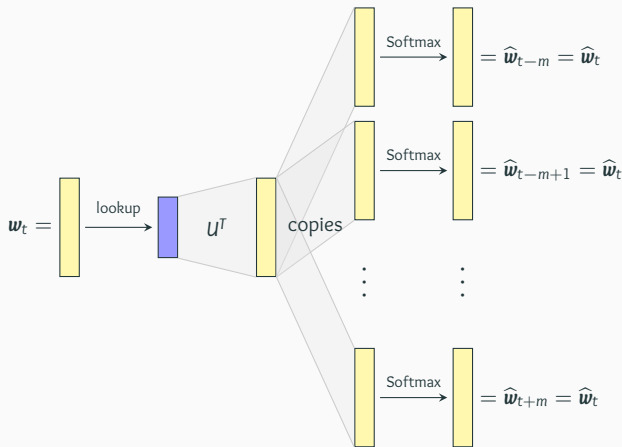
- Principle: model the distribution of **context words conditional to its center**

$$p(w_{t-m}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+m} \mid w_t)$$

- Idea: use a neural network to learn that conditional distribution given given one-hot encoded context words

Skip-gram model

- Reverse of CBOW model
 - Reverse arrows
 - Replace averaging by copies
 - Swap softmax/lookup
- First layer: center word lookup
- Second layer: linear transform
- Third (shallow layer): copies
- Third real layer: softmax to get distributions over \mathcal{V}
- Parameters
 - Embeddings: \mathbf{e}_{t+k} for $k \in \{-m, \dots, m\}, k \neq 0$
 - Matrix: U



- First layer: center word lookup

$$\mathbf{e}_t = E\mathbf{w}_t$$

- Second layer: linear transform

$$U^T \mathbf{e}_t$$

- Third (shallow layer): copies
- Third real layer: softmax to get distributions over \mathcal{V} (which are all the same)

$$\hat{\mathbf{w}}_t = \text{softmax}(U^T \mathbf{e}_t)$$

- Probability of word of index i in \mathcal{V} to appear in the context of w_t

$$\hat{w}_{ti} = \frac{\exp(\langle \mathbf{u}_i, \mathbf{e}_t \rangle)}{\sum_{j=1}^{\mathcal{V}} \exp(\langle \mathbf{u}_j, \mathbf{e}_t \rangle)}$$

- For context word w_k , corresponding probability is $\hat{w}_{t,i(k)}$

$$\hat{w}_{t,i(k)} = \frac{\exp(\langle \mathbf{u}_{i(k)}, \mathbf{e}_t \rangle)}{\sum_{j=1}^{\mathcal{V}} \exp(\langle \mathbf{u}_j, \mathbf{e}_t \rangle)}$$

where $i(k)$ is the index of token w_k in vocabulary \mathcal{V}

- Sum of cross-entropies on each output

$$\mathcal{L} = \sum_{\substack{k=t-m \\ k \neq t}}^{t+m} H(\mathbf{w}_{i(k)}, \hat{\mathbf{w}}_t)$$

- For word w_k , we have

$$\begin{aligned} H(\mathbf{w}_{i(k)}, \hat{\mathbf{w}}_t) &= \sum_{j=1}^V \mathbf{w}_{i(k)j} \log \hat{\mathbf{w}}_{tj} \\ &= \log \hat{\mathbf{w}}_{ti(k)} \end{aligned} \quad (\mathbf{w}_{i(k)} \text{ is the one-hot vector})$$

- Final loss is

$$\mathcal{L} = \sum_{\substack{k=t-m \\ k \neq t}}^{t+m} H(\mathbf{w}_k, \hat{\mathbf{w}}_k) = - \sum_{\substack{k=t-m \\ k \neq t}}^{t+m} \log \hat{\mathbf{w}}_{ti(k)}$$

- Expanded loss is

$$\mathcal{L} = - \sum_{\substack{k=t-m \\ k \neq t}}^{t+m} \langle \mathbf{u}_{i(k)}, \mathbf{e}_t \rangle + 2m \log \sum_{j=1}^{\mathcal{V}} \exp (\langle \mathbf{u}_j, \mathbf{e}_t \rangle)$$

- Parameters to optimize are: \mathbf{e}_t and \mathbf{u}_j for $j = 1, \dots, \mathcal{V}$

Cosine similarity

- Cosine similarity

$$S = \frac{\langle \mathbf{w}_i, \mathbf{w}_j \rangle}{\|\mathbf{w}_i\| \|\mathbf{w}_j\|} = \cos \theta$$

- $S \in (-1, 1)$
- Orthogonal if $S = 0$
- Same orientation if $S = 1$ ($\mathbf{w}_i = \lambda \mathbf{w}_j$, $\lambda \geq 0$)
- Two random tokens are \approx orthogonal

$$\mathbb{E}(S) = 0$$

- Non-zero is significant

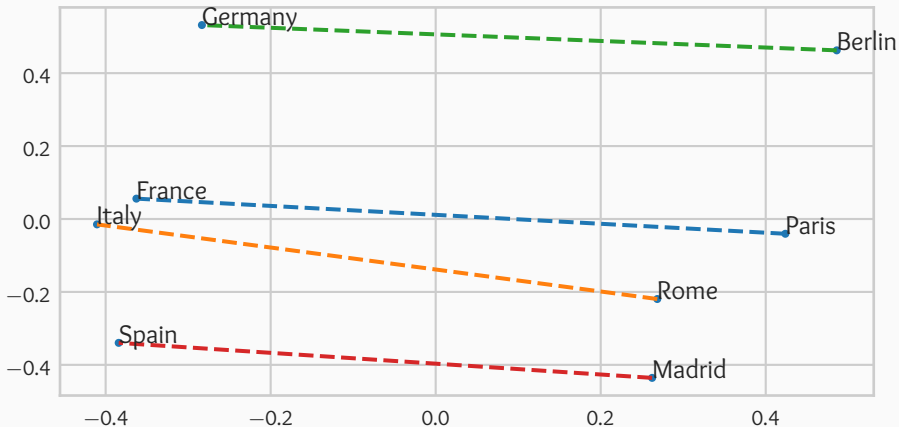
$$\text{Var}(S) = \mathcal{O}(1/N)$$

- Maximum cosine similarity

car	score	cat	score
cars	0.8139	cats	0.8064
automobile	0.7040	dog	0.7927
truck	0.6975	kitten	0.7527
vehicle	0.6854	rabbit	0.7517
roadster	0.6779	pet	0.7104
SUV	0.6734	fox	0.7061
motorbike	0.6674	feline	0.6898
motorcycle	0.6509	kittens	0.6837
motorcar	0.6486	puppy	0.6721
driver	0.6452	tomcat	0.6718

Capitals and countries

PCA in 2D



Concept of “capital of” as “Madrid” - “Spain”

- Most similar tokens to:

 (“Madrid” - “Spain”) + “France”

which translates to

 <capital of> “France”

- Same thing for

 (“Madrid” - “Spain”) + “Germany”

(“Madrid” - “Spain”) + “France”	score
Paris	0.7620
Lille	0.7085
Marseille	0.6987
Strasbourg	0.6825
France	0.6809

(“Madrid” - “Spain”) + “Germany”	score
Munich	0.7471
Stuttgart	0.7351
Berlin	0.7149
Dresden	0.7014
Germany	0.6990

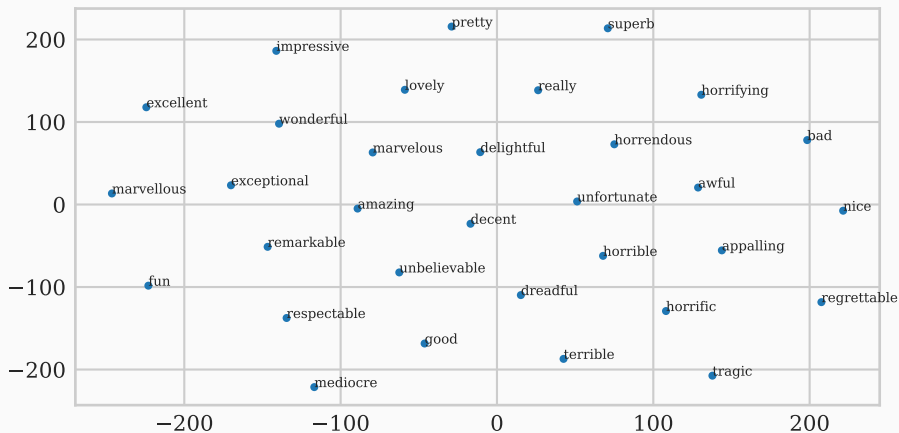
Additive compositionality

- 10 closest token of a sum of embeddings

Czech + currency		German + airlines	
	score		score
Czech	0.7478	airlines	0.7295
currency	0.7478	German	0.7295
koruna	0.6906	Interflug	0.7184
currencies	0.6406	airline	0.6487
Slovak	0.6321	TAROM	0.6281
Slovakia	0.6178	TUIfly	0.6279
forint	0.6016	Eurowings	0.6246
kroon	0.5836	Crossair	0.6223
króna	0.5812	Cargolux	0.6222
markka	0.5798	Germanwings	0.6167

Word2vec: illustrations

Depth 2 and 5—neighborhood of “good” and “bad” cosine similarity, t-SNE visualization

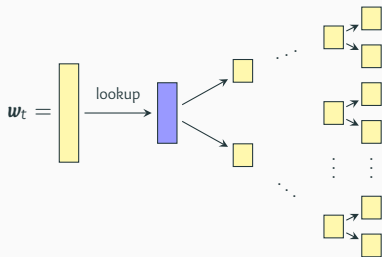


Properties

- Semantic information is encoded in a linear way (synonyms, analogy, additive compositionality)
 - Synonyms are close
 - Additive compositionality to simple question answering
 - Concepts as vector for analogy
- Shortcomings
 - Softmax operation is costly (size $V \sim 10^{5-7}$)
 - Don't use token frequencies: very frequent words carry less information than rare ones
 - Same embedding given to opposite words
- Improvements
 - Negative sampling
 - Hierarchical softmax

Hierarchical softmax

- Skip-gram model: $\hat{\mathbf{w}}_k = \text{softmax}(\mathbf{U}^T \mathbf{e}_t)$
- Replace linear transform + softmax layers by a **binary tree structure**
- If \mathcal{E} is the set of pairs of nodes from root to index i
$$\hat{\mathbf{w}}_{ki} = \prod_{(n,m) \in \mathcal{E}} \sigma(\varepsilon(m) \langle \mathbf{v}_n, \mathbf{e}_t \rangle)$$
- $\varepsilon(m)$ encodes “is node left/right”
$$\varepsilon(m) = \begin{cases} -1 & \text{if node } m \text{ is a left child} \\ 1 & \text{otherwise} \end{cases}$$
- \mathbf{v}_n attached to each (non-leaf) node including root



Negative sampling

- Based on the Skip-gram model but different loss function
- Skip-gram: model \mathcal{C} conditional to w_t
- Negative sampling: model the probability that for a token w , the token c is in the context ($D = 1$)

$$\Pr((w, c) \text{ in corpus}) = \Pr(D = 1 \mid w, c)$$

- Skip-gram model

$$\Pr(D = 1 \mid w, c) = \sigma(\langle \mathbf{e}_w, \mathbf{u}_c \rangle), \quad \Pr(D = 0 \mid w, c) = \sigma(-\langle \mathbf{e}_w, \mathbf{u}_c \rangle)$$

- Loss

$$\mathcal{L} = -\log \sigma(\langle \mathbf{e}_w, \mathbf{u}_c \rangle) - \sum_{r \in \mathcal{N}} \log \sigma(-\langle \mathbf{e}_w, \tilde{\mathbf{u}}_r \rangle)$$

where \mathcal{N} is a random sample of tokens most likely not part of context

GloVe

- Principle: Model co-occurrence matrix from words representation

$$\log C(w_i, w_j) = \langle \mathbf{w}_i, \tilde{\mathbf{w}}_j \rangle + b_i + \tilde{b}_j$$

where \mathbf{w}_i and $\tilde{\mathbf{w}}_j$ are two distinct word representations, b_i and \tilde{b}_i are biases

- Squared loss that weights all co-occurrences equally

$$J = \sum_{i,j=1}^V \left(\langle \mathbf{w}_i, \tilde{\mathbf{w}}_j \rangle + b_i + \tilde{b}_j - \log C(w_i, w_j) \right)^2$$

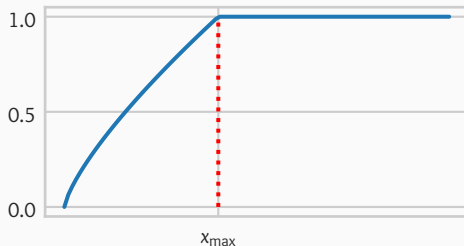
- Rare co-occurrence are given same importance as more frequent one

- Weighted squared loss

$$J = \sum_{i,j=1}^V f(C(w_i, w_j)) \left(\langle \mathbf{w}_i, \tilde{\mathbf{w}}_j \rangle + b_i + \tilde{b}_j - \log C(w_i, w_j) \right)^2$$

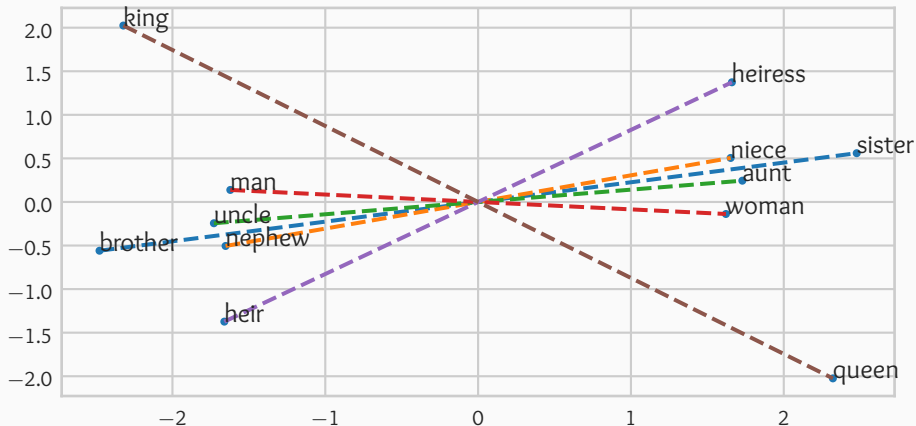
- Weighting function f with $\alpha = 3/4$
- No weight on non-existent words
- Weight rapidly increasing
- Don't put too much weight on frequent words

$$f(x) = \begin{cases} \left(\frac{x}{x_{\max}} \right)^\alpha & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases}$$

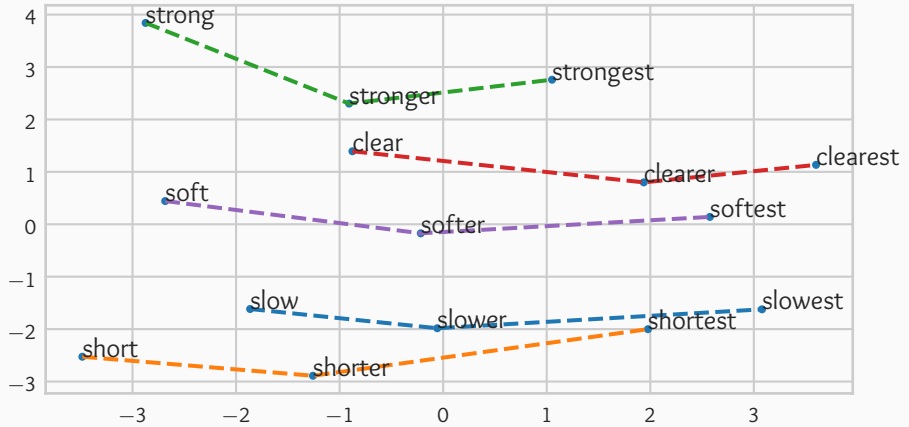


Illustrations

GloVe (recentered)

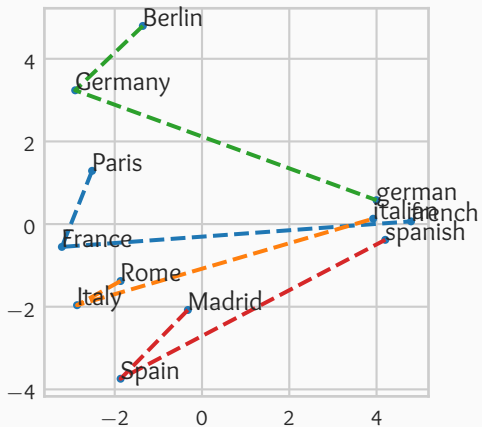


GloVe



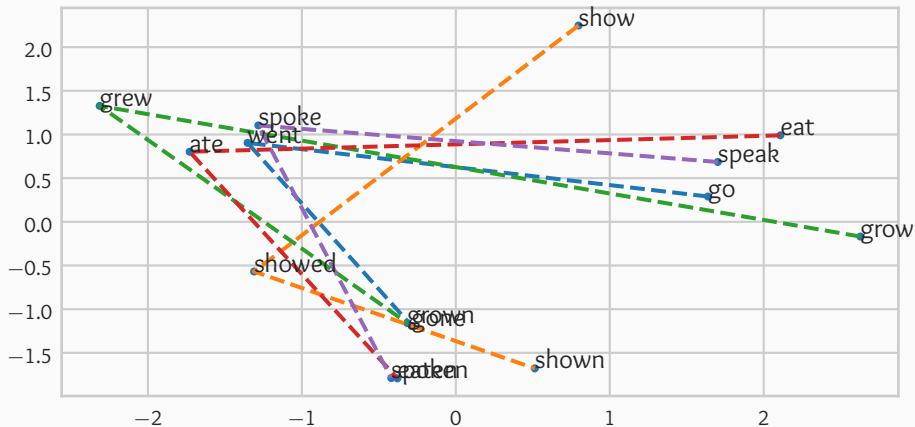
Illustrations

GloVe



Illustrations

GloVe (recentered)



- [1] John R. Firth. “**A Synopsis of Linguistic Theory, 1930-1955**”. In: *Studies in linguistic analysis* (1957).
- [2] Tomas Mikolov et al. “**Distributed Representations of Words and Phrases and Their Compositionality**”. In: *Advances in Neural Information Processing Systems*. 2013, pp. 3111–3119.
- [3] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. “**Glove: Global Vectors for Word Representation**”. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2014, pp. 1532–1543.
- [4] Piotr Bojanowski et al. “**Enriching Word Vectors with Subword Information**”. In: *Transactions of the Association for Computational Linguistics* 5 (2017), pp. 135–146.

- [5] Matthew E. Peters et al. **“Deep Contextualized Word Representations”**. 2018. arXiv: 1802.05365.
- [6] Jacob Devlin et al. **“BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”**. May 24, 2019. arXiv: 1810.04805 [cs]. URL: <http://arxiv.org/abs/1810.04805> (visited on 11/30/2021).