

# Statistical Natural Language Processing

## Distributed representations

Çağrı Çöltekin

University of Tübingen  
Seminar für Sprachwissenschaft

Summer Semester 2017

# Representations of linguistic units

- Most ML methods we use depend on how we represent the objects of interest, such as
  - words, morphemes
  - sentences, phrases
  - letters, phonemes
  - documents
  - speakers, authors
  - ...
- The way we represent these objects interacts with the ML methods
- We will mostly talk about word representations
  - They are also applicable any of the above

# Symbolic (one-hot) representations

A common way to represent words is one-hot vectors

$$\text{cat} = (0, \dots, 1, 0, 0, \dots, 0)$$

$$\text{dog} = (0, \dots, 0, 1, 0, \dots, 0)$$

$$\text{book} = (0, \dots, 0, 0, 1, \dots, 0)$$

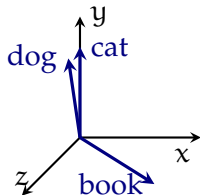
...

- No notion of similarity
- Large and sparse vectors

## More useful vector representations

- The idea is to represent similar words with similar vectors

$\text{cat} = (0, 3, 1, \dots, 4)$   
 $\text{dog} = (0, 3, 0, \dots, 3)$   
 $\text{book} = (4, 1, 4, \dots, 5)$   
 ...



- The similarity between the vectors may represent similarities based on
  - syntactic
  - semantic
  - topical
  - form
  - ... features useful in a particular task

# Where do the vector representations come from?

- The vectors are (almost certainly) learned from the data
- Typically using an unsupervised (or self-supervised) method
- The idea goes back to,

*You shall know a word by the company it keeps.*  
—Firth (1957)

- In practice, we make use of the contexts (company) of the words to determine their representations
- The words that appear in similar contexts are mapped to similar representations

# How to calculate word vectors?

count word in context

	$c_1$	$c_2$	$c_3$	$\dots$	$c_m$
cat	0	3	1	$\dots$	4
dog	0	3	0	$\dots$	3
book	4	1	4	$\dots$	5

- + Now words that appear in the same contexts will have similar vectors
  - The frequencies are often normalized (PMI, TF-IDF)
  - The data is highly correlated: lots of redundant information
  - Still large and sparse

# How to calculate word vectors?

count, factorize, truncate

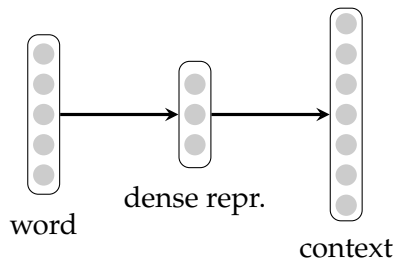
$$\begin{array}{c}
 \\
 w_1 \\
 w_2 \\
 w_3 \\
 \dots
 \end{array}
 \begin{array}{c}
 c_1 \quad c_2 \quad c_3 \quad \dots \quad c_m \\
 \left[ \begin{array}{ccccc}
 0 & 3 & 1 & \dots & 4 \\
 0 & 3 & 0 & \dots & 3 \\
 4 & 1 & 4 & \dots & 5 \\
 \dots & & & & 
 \end{array} \right] =
 \end{array}$$

$$\begin{array}{c}
 \\
 w_1 \\
 w_2 \\
 w_3 \\
 \dots
 \end{array}
 \begin{array}{c}
 z_1 \quad z_2 \quad z_3 \quad \dots \quad z_m \\
 \left[ \begin{array}{cccc}
 1 & 5 & 9 & \dots & 4 \\
 1 & 4 & 1 & \dots & 3 \\
 9 & 1 & 1 & \dots & 5 \\
 \dots & & & & 
 \end{array} \right]
 \end{array}
 \begin{array}{c}
 \left[ \begin{array}{ccc}
 \sigma_1 & \dots & 0 \\
 \vdots & \ddots & \vdots \\
 0 & \dots & \sigma_m
 \end{array} \right]
 \end{array}
 \begin{array}{c}
 c_1 \quad c_2 \quad c_3 \quad \dots \quad c_m \\
 \left[ \begin{array}{ccccc}
 0 & 3 & 1 & \dots & 4 \\
 0 & 3 & 0 & \dots & 3 \\
 9 & 1 & 8 & \dots & 0 \\
 \dots & & & & 
 \end{array} \right]
 \end{array}
 \begin{array}{c}
 u_1 \\
 u_2 \\
 u_3 \\
 \dots
 \end{array}$$

# How to calculate word vectors?

predict the context from the word, or word from the context

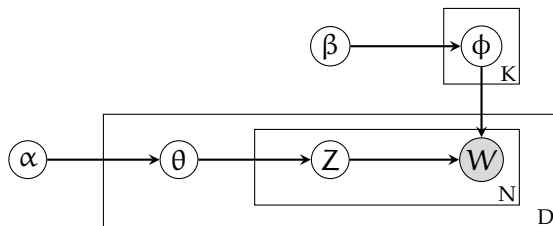
- The task is predicting
  - the context of the word from the word itself
  - or the word from its context
- Task itself is not interesting
- We are interested in the hidden layer representations learned





# How to calculate word vectors?

latent variable models (e.g., LDA)



- Assume that the each 'document' is generated based on a mixture of latent variables
- Learn the probability distributions
- Typically used for *topic modeling*
- Can model words too (as a mixture of latent variables)

# A toy example

A four-sentence corpus with *bag of words* (BOW) model.

The corpus:

S1: She likes cats and  
dogs

S2: He likes dogs and  
cats

S3: She likes books

S4: He reads books

Term-document (sentence) matrix

	S1	S2	S3	S4
she	1	0	1	0
he	0	1	0	1
likes	1	1	1	0
reads	0	0	0	1
cats	1	1	0	0
dogs	1	1	0	0
books	0	0	1	1
and	1	1	0	0

# A toy example

A four-sentence corpus with *bag of words* (BOW) model.

The corpus:

S1: She likes cats and  
dogs

S2: He likes dogs and  
cats

S3: She likes books

S4: He reads books

Term-term (left-context) matrix

	#	she	he	likes	reads	cats	dogs	books	and
she	2	0	0	0	0	0	0	0	0
he	2	0	0	0	0	0	0	0	0
likes	0	2	1	0	0	0	0	0	0
reads	0	0	1	0	0	0	0	0	0
cats	0	0	0	1	0	0	0	0	1
dogs	0	0	0	1	0	0	0	0	1
books	0	0	0	1	1	0	0	0	0
and	0	0	0	0	0	1	1	0	0

## Term-document matrices

- The rows are about the terms: similar terms appear in similar contexts
- The columns are about the context: similar contexts contain similar words
- The term-context matrices are typically sparse and large

Term-document (sentence) matrix

	S1	S2	S3	S4
she	1	0	1	0
he	0	1	0	1
likes	1	1	1	0
reads	0	0	0	1
cats	1	1	0	0
dogs	1	1	0	0
books	0	0	1	1
and	1	1	0	0

## SVD (again)

- Singular value decomposition is a well-known method in linear algebra
- An  $n \times m$  ( $n$  terms  $m$  documents) term-document matrix  $\mathbf{X}$  can be decomposed as

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

- $\mathbf{U}$  is a  $n \times r$  unitary matrix, where  $r$  is the rank of  $\mathbf{X}$  ( $r \leq \min(n, m)$ ). Columns of  $\mathbf{U}$  are the eigenvectors of  $\mathbf{X}\mathbf{X}^T$
  - $\mathbf{\Sigma}$  is a  $r \times r$  diagonal matrix of singular values (square root of eigenvalues of  $\mathbf{X}\mathbf{X}^T$  and  $\mathbf{X}^T\mathbf{X}$ )
  - $\mathbf{V}^T$  is a  $r \times m$  unitary matrix. Columns of  $\mathbf{V}$  are the eigenvectors of  $\mathbf{X}^T\mathbf{X}$
- One can consider  $\mathbf{U}$  and  $\mathbf{V}$  as PCA performed for reducing dimensionality of rows (terms) and columns (documents)

# Truncated SVD

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

- Using eigenvectors (from  $\mathbf{U}$  and  $\mathbf{V}$ ) that correspond to  $k$  largest singular values ( $k < r$ ), allows reducing dimensionality of the data with minimum loss
- The approximation,

$$\hat{\mathbf{X}} = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k$$

results in the best approximation of  $\mathbf{X}$ , such that  $\|\hat{\mathbf{X}} - \mathbf{X}\|_F$  is minimum

# Truncated SVD

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

- Using eigenvectors (from  $\mathbf{U}$  and  $\mathbf{V}$ ) that correspond to  $k$  largest singular values ( $k < r$ ), allows reducing dimensionality of the data with minimum loss
- The approximation,

$$\hat{\mathbf{X}} = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k$$

results in the best approximation of  $\mathbf{X}$ , such that  $\|\hat{\mathbf{X}} - \mathbf{X}\|_F$  is minimum

- Note that  $r$  may easily be millions (of words or contexts), while we choose  $k$  much smaller (a few hundreds)

## Truncated SVD (2)

$$\begin{bmatrix}
 x_{1,1} & x_{1,2} & x_{1,3} & \dots & x_{1,m} \\
 x_{1,1} & x_{1,2} & x_{1,3} & \dots & x_{1,m} \\
 x_{2,1} & x_{2,2} & x_{2,3} & \dots & x_{2,m} \\
 x_{3,1} & x_{3,2} & x_{3,3} & \dots & x_{3,m} \\
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 x_{n,1} & x_{n,2} & x_{n,3} & \dots & x_{n,m}
 \end{bmatrix} =$$

$$\begin{bmatrix}
 u_{1,1} & \dots & u_{1,k} \\
 u_{2,1} & \dots & u_{2,k} \\
 u_{3,1} & \dots & u_{3,k} \\
 \vdots & \ddots & \vdots \\
 u_{n,1} & \dots & u_{n,k}
 \end{bmatrix} \times
 \begin{bmatrix}
 \sigma_1 & \dots & 0 \\
 \vdots & \ddots & \vdots \\
 0 & \dots & \sigma_k
 \end{bmatrix} \times
 \begin{bmatrix}
 u_{1,1} & u_{1,2} & \dots & u_{1,m} \\
 \vdots & \vdots & \ddots & \vdots \\
 u_{k,1} & u_{k,2} & \dots & u_{n,m}
 \end{bmatrix}$$



## Truncated SVD (2)

$$\begin{bmatrix} x_{1,1} & x_{1,2} & x_{1,3} & \dots & x_{1,m} \\ x_{2,1} & x_{2,2} & x_{2,3} & \dots & x_{2,m} \\ x_{3,1} & x_{3,2} & x_{3,3} & \dots & x_{3,m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & x_{n,3} & \dots & x_{n,m} \end{bmatrix} =$$

$$\begin{bmatrix} u_{1,1} & \dots & u_{1,k} \\ u_{2,1} & \dots & u_{2,k} \\ u_{3,1} & \dots & u_{3,k} \\ \vdots & \ddots & \vdots \\ u_{n,1} & \dots & u_{n,k} \end{bmatrix} \times \begin{bmatrix} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_k \end{bmatrix} \times \begin{bmatrix} u_{1,1} & u_{1,2} & \dots & u_{1,m} \\ \vdots & \vdots & \ddots & \vdots \\ u_{k,1} & u_{k,2} & \dots & u_{n,m} \end{bmatrix}$$

The term<sub>1</sub> can be represented using the first row of  $\mathbf{U}_k$

## Truncated SVD (2)

$$\begin{bmatrix} x_{1,1} & x_{1,2} & x_{1,3} & \dots & x_{1,m} \\ x_{1,1} & x_{1,2} & x_{1,3} & \dots & x_{1,m} \\ x_{2,1} & x_{2,2} & x_{2,3} & \dots & x_{2,m} \\ x_{3,1} & x_{3,2} & x_{3,3} & \dots & x_{3,m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & x_{n,3} & \dots & x_{n,m} \end{bmatrix} =$$

$$\begin{bmatrix} u_{1,1} & \dots & u_{1,k} \\ u_{2,1} & \dots & u_{2,k} \\ u_{3,1} & \dots & u_{3,k} \\ \vdots & \ddots & \vdots \\ u_{n,1} & \dots & u_{n,k} \end{bmatrix} \times \begin{bmatrix} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_k \end{bmatrix} \times \begin{bmatrix} u_{1,1} & u_{1,2} & \dots & u_{1,m} \\ \vdots & \vdots & \ddots & \vdots \\ u_{k,1} & u_{k,2} & \dots & u_{n,m} \end{bmatrix}$$

The document<sub>1</sub> can be represented using the first column of  $\mathbf{V}_k^T$

# Truncated SVD example

The corpus:

(S1) She likes cats and dogs

(S2) He likes dogs and cats

(S3) She likes books

(S4) He reads books

	S1	S2	S3	S4
she	1	0	1	0
he	0	1	0	1
likes	1	1	1	0
reads	0	0	0	1
cats	1	1	0	0
dogs	1	1	0	0
books	0	0	1	1
and	1	1	0	0

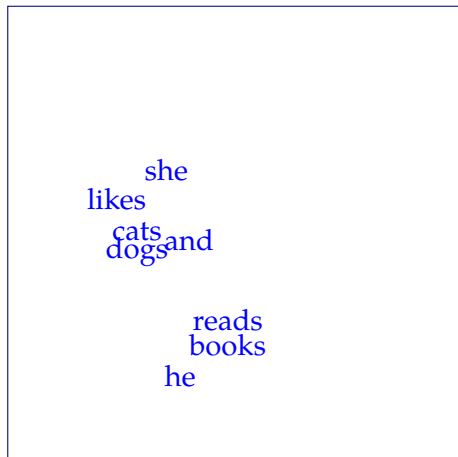
Truncated SVD ( $k = 2$ )

$$\mathbf{U} = \begin{bmatrix} -0.30 & 0.28 \\ -0.24 & -0.63 \\ -0.52 & 0.15 \\ -0.03 & -0.49 \\ -0.43 & 0.01 \\ -0.43 & 0.01 \\ -0.03 & -0.49 \\ -0.43 & 0.01 \end{bmatrix} \begin{array}{l} \text{she} \\ \text{he} \\ \text{likes} \\ \text{reads} \\ \text{cats} \\ \text{dogs} \\ \text{books} \\ \text{and} \end{array}$$

$$\mathbf{\Sigma} = \begin{bmatrix} 3.11 & 0 \\ 0 & 1.81 \end{bmatrix}$$

$$\mathbf{V}^T = \begin{array}{c} \text{S1} \quad \text{S2} \quad \text{S3} \quad \text{S4} \\ \begin{bmatrix} -0.68 & 0.26 & -0.11 & -0.66 \\ -0.66 & -0.23 & 0.48 & 0.50 \end{bmatrix} \end{array}$$

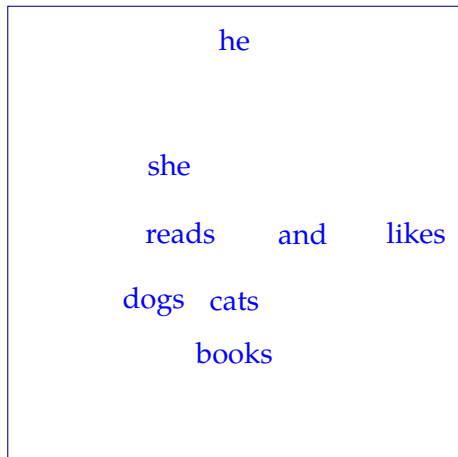
# Truncated SVD (with BOW sentence context)



The corpus:

- (S1) She likes cats and dogs
- (S2) He likes dogs and cats
- (S3) She likes books
- (S4) He reads books

# Truncated SVD (with single word context)



The corpus:

- (S1) She likes cats and dogs
- (S2) He likes dogs and cats
- (S3) She likes books
- (S4) He reads books

# SVD: LSI/LSA

- SVD applied to term-document matrices are called
  - *Latent semantic analysis* (LSA) if the aim is constructing *term* vectors
  - *Latent semantic indexing* (LSI) if the aim is constructing *document* vectors
- The well known Google *PageRank* algorithm is a variation of the SVD

## SVD based vectors: practical concerns

- In practice, instead of raw counts of terms within contexts, the term-document matrices typically contain
  - pointwise mutual information
  - tf-idfvalues.
- If the aim is finding latent (semantic) topics, frequent/syntactic words (*stopwords*) are often removed
- Depending on the measure used, it may also be important to normalize for the document length

## SVD-based vectors: applications

- The SVD-based methods is commonly used in information retrieval
  - The system builds document vectors using SVD
  - The search terms are also considered as a 'document'
  - System retrieves the documents whose vectors are similar to the search term
- The SVD-based methods for semantic similarity is also common
- It was shown that the vector space models outperform humans in TOEFL synonym questions and SAT analogy questions

the song



## Predictive models

- Instead of dimensionality reduction through SVD, we try to predict
  - either the target word from the context
  - or the context given the target word
- We assign each word to a fixed-size random vector
- We use a standard ML model and try to reduce the prediction error with a method like gradient descent
- During learning, the algorithm optimizes the vectors as well as the model parameters
- In this context, the word-vectors are called **embeddings**
- This types of models has been very popular during last few years

# Predictive models

- The idea is the ‘locally’ predict the context a particular word occurs
- Both the context and the words are represented as low dimensional dense vectors
- Typically, neural networks are used for the prediction
- The hidden layer representations are the vectors we are interested

# word2vec

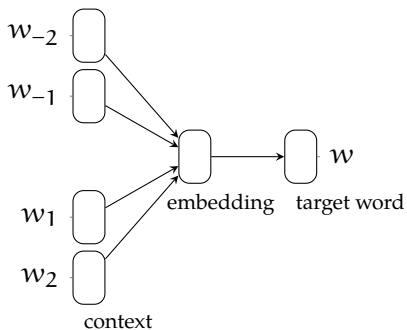
- **word2vec** is a popular algorithm and open source application for training word vectors (Mikolov et al. 2013)
- It has two modes of operation

**CBOW** or continuous bag of words predict the word using a window around the word

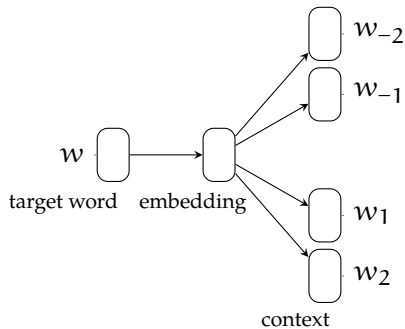
**Skip-gram** does the reverse, it predicts the words in the context of the target word using the target word as the predictor

## word2vec

## CBOW and skip-gram modes



CBOW



Skip-gram

# Learning in word2vec

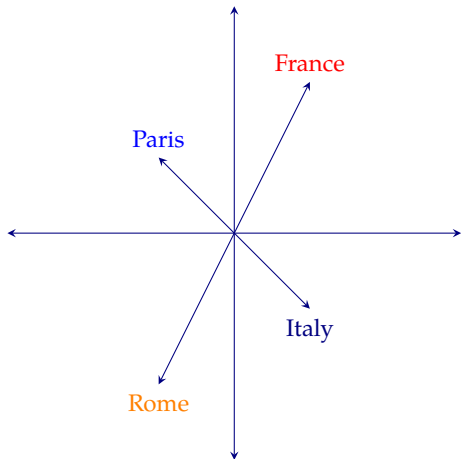
- The algorithm learns two sets of embeddings (one for context, one for target)
- The learning method is simply logistic regression, where word vectors are also updated (besides model parameters)
- A particular problem with predicting one-hot vectors at the output layer is computation of softmax for large vocabulary
- Negative examples are sampled from the larger corpus
- It performs well, and it is much faster than earlier (more complex) ANN architectures developed for this task

# GloVe

- GloVe is another popular method for obtaining word vectors (Pennington, Socher, and Manning 2014)
- It tries to combine intuitions from both SVD-like ‘counting’ methods, and prediction-based methods
- It typically performs better on smaller data sets

# Word vectors and syntactic/semantic relations

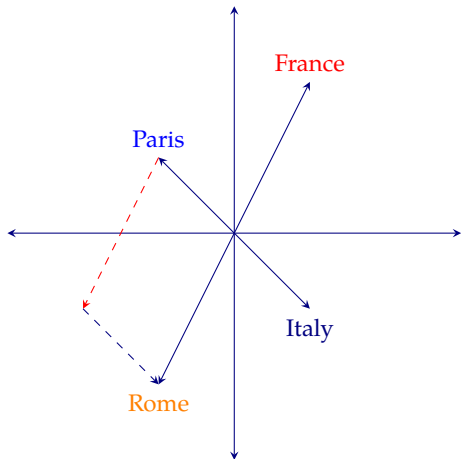
Word vectors map some syntactic/semantic relations to vector operations



# Word vectors and syntactic/semantic relations

Word vectors map some syntactic/semantic relations to vector operations

- Paris - France + Italy = Rome

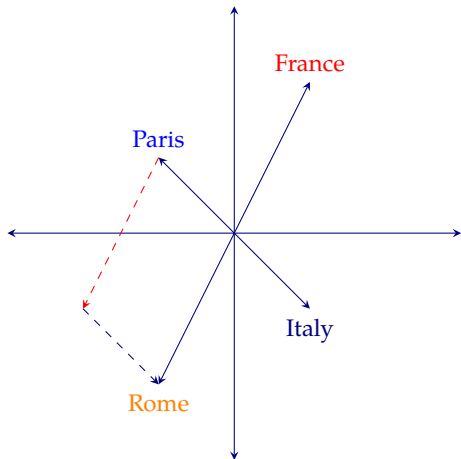




# Word vectors and syntactic/semantic relations

Word vectors map some syntactic/semantic relations to vector operations

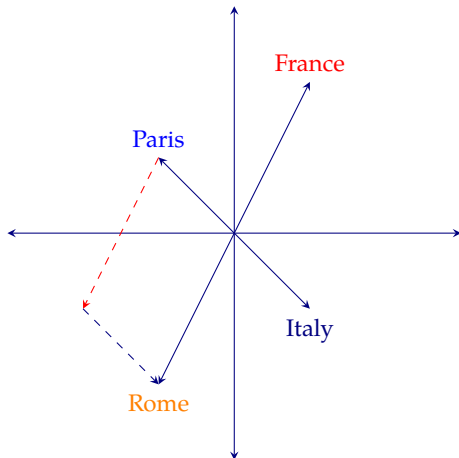
- Paris - France + Italy = Rome
- king - man + woman = queen



# Word vectors and syntactic/semantic relations

Word vectors map some syntactic/semantic relations to vector operations

- Paris - France + Italy = Rome
- king - man + woman = queen
- duck - ducks + mouse = mice



# Using vector representations

- Dense vector representations are useful for many ML methods
- They are particularly suitable for neural network models
- ‘General purpose’ vectors can be trained on unlabeled data
- They can also be trained for a particular purpose, resulting in ‘task specific’ vectors
- Dense vector representations are not specific to words, they can be obtained and used for any (linguistic) object of interest

# Evaluating vector representations

- Like other unsupervised methods, there are no ‘correct’ labels
- Evaluation can be based on
  - Intrinsic evaluation based on success on finding analogy/synonymy
  - Extrinsic evaluation, based on whether they improve a particular task (e.g., parsing, sentiment analysis) or not
  - Correlation with human judgments

# Summary

- Dense vector representations of linguistic units (as opposed to symbolic representations) allow calculating similarity / difference between the units
- They can be either based on counting (SVD), or predicting (word2vec, GloVe)
- They are particularly suitable for ANNs, deep learning architectures

# Summary

- Dense vector representations of linguistic units (as opposed to symbolic representations) allow calculating similarity / difference between the units
- They can be either based on counting (SVD), or predicting (word2vec, GloVe)
- They are particularly suitable for ANNs, deep learning architectures

Next:

Mon Text classification

Wed SMT (?)

# Assignment 1 scores

