Non-Standard Datenbanken und Data Mining

From Clustering to Embedding

Prof. Dr. Ralf Möller Universität zu Lübeck Institut für Informationssysteme



IM FOCUS DAS LEBEN

Acknowledgments

- Slides have been taken from
 - "Improving Distributional Similarity with Lessons Learned from Word Embeddings"
 - Omer Levy, Yoav Goldberg, Ido Dagan
 - Stanford CS224d: Deep Learning for NLP
 - Richard Socher
 - Rensselaer: Natural Language Processing
 - Heng Li



Recap bag-of-words approaches

• LSI: Documents as vectors, dimension reduction

Words are not independent of each other

- Word similarity measures
- Extend query with similar words automatically
- Extend query with most frequent followers/predecessors
- Insert words in anticipated gaps in a string query

Need to represent some aspects of word semantics



Beyond bags of words



Underlying Theory: The Distributional Hypothesis (Harris, '54; Firth, '57)

"Similar words occur in similar contexts"

https://www.tensorflow.org/tutorials/word2vec

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



References

• Harris 54

Harris, Zellig. Distributional structure. *Word* 10(23). 146–162. **1954**.

• Micholov et al. 13

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. In Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2 (NIPS'13). **2013**.

• Firth 57

Firth, John R. A synopsis of linguistic theory 1930–1955. In *Studies in linguistic analysis*, 1–32. Oxford: Blackwell. **1957**.

• Pennington et al. 14

Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global Vectors for Word Representation. **2014**.



Point(wise) Mutual Information: PMI

 Measure of association used in information theory and statistics

$$ext{pmi}(x;y) \equiv \log rac{p(x,y)}{p(x)p(y)} = \log rac{p(x|y)}{p(x)} = \log rac{p(y|x)}{p(y)}$$

- Positive PMI: PPMI(x, y) = max(pmi(x, y), 0)
- Quantifies the discrepancy between the probability of their coincidence given their joint distribution and their individual distributions, assuming independence
- Finding collocations and associations between words
- Countings of occurrences and co-occurrences of words in a text corpus can be used to approximate the probabilities p(x) or p(y) and p(x,y) respectively



PMI – Example

word 1	word 2	count word 1	count word 2	count of co-occurrences	PMI	
puerto	rico	1938	1311	1159	10.0349081703	
hong	kong	2438	2694	2205	9.72831972408	
los	angeles	3501	2808	2791	9.56067615065	
carbon	dioxide	4265	1353	1032	9.09852946116	
prize	laureate	5131	1676	1210	8.85870710982	
san	francisco	5237	2477	1779	8.83305176711	
nobel	prize	4098	5131	2498	8.68948811416	
ice	hockey	5607	3002	1933	8.6555759741	
star	trek	8264	1594	1489	8.63974676575	
car	driver	5578	2749	1384	8.41470768304	
it	the	283891	3293296	3347	-1.72037278119	
are	of	234458	1761436	1019	-2.09254205335	
this	the	199882	3293296	1211	-2.38612756961	
is	of	565679	1761436	1562	-2.54614706831	
and	of	1375396	1761436	2949	-2.79911817902	
а	and	984442	1375396	1457	-2.92239510038	
in	and	1187652	1375396	1537	-3.05660070757	
to	and	1025659	1375396	1286	-3.08825363041	
to	in	1025659	1187652	1066	-3.12911348956	
of	and	1761436	1375396	1190	-3.70663100173	

- Counts of pairs of words getting the most and the least PMI scores in the first 50 millions of words in Wikipedia (dump of October 2015)
- Filtering by 1,000 or more co-occurrences.
 - The frequency of each count can be obtained by dividing its value by 50,000,952. (Note: natural log is used to calculate the PMI values in this example, instead of log base 2)



Applications of PMI Data

- Extend query with most frequent followers/predecessors
- Insert words in anticipated gaps in a string query



PMI – Co-occurrence Matrix

	Add-2 Smoothed Count(w,context)							
	computer	data	pinch	result	sugar			
apricot	2	2	3	2	3			
pineapple	2	2	3	2	3			
digital	4	3	2	3	2			
information	3	8	2	6	2			
	F	PPMI(w,c	ontext)					
	computer	PPMI(w,c data	ontext) pinch	result	sugar			
apricot	computer -	PPMI(w,c data -	ontext) pinch 2.25	result -	sugar 2.25			
apricot pineapple	Computer - -	PPMI(w,c data - -	ontext) pinch 2.25 2.25	result -	sugar 2.25 2.25			
apricot pineapple digital	computer - - 1.66	PPMI(w,c data - - 0.00	ontext) pinch 2.25 2.25	result - - 0.00	sugar 2.25 2.25			



1. Clustering Approach to Word Semantics

Clustering vectors to visualize similarity in cooccurrence matrices (Rohde et al. 2005)

Use whatever clustering algorithm you prefer to determine "related" words

Application:

NIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

Extend query with related words automatically



WRIST ANKLE

SHOULDER ٨RM FG HAND

FOOT

HAWAII

COW

TURTLE OYSTER

MONTREAL NASHVILLE TOKYO

Apply SVD-based Dimension Reduction





2. Embedding Approaches to Word Semantics

- Represent each word with a low-dimensional vector
- Word similarity = vector similarity
- Key idea: Predict surrounding words of every word



Represent the meaning of **words** – word2vec

- 2 basic structural models:
 - Continuous Bag of Words (CBOW): use a window of words to predict the middle word
 - Skip-gram (SG): use a word to predict the surrounding ones in window.





Word2vec – Continuous Bag of Word

- E.g. "The cat <sat> on floor"
 - Window size = 2

























A **logistic function** or **logistic curve** is a common "S" shape (sigmoid curve), with equation:

$$f(x)=rac{L}{1+e^{-k(x-x_0)}}$$

where

- *e* = the natural logarithm base (also known as Euler's number),
- x₀ = the *x*-value of the sigmoid's midpoint,
- L = the curve's maximum value, and
- *k* = the steepness of the curve.^[1]





softmax(z)

The **softmax function**, or **normalized exponential function**, is a generalization of the logistic function that "squashes" a *K*-dimensional vector \mathbf{z} of arbitrary real values to a *K*-dimensional vector $\sigma(\mathbf{z})$ of real values in the range [0, 1] that add up to 1. The function is given by

$$egin{aligned} &\sigma: \mathbb{R}^K o [0,1]^K \ &\sigma(\mathbf{z})_j = rac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} & ext{for } j = 1, \, ..., \, K. \end{aligned}$$

In probability theory, the output of the softmax function can be used to represent a categorical distribution – that is, a probability distribution over K different possible outcomes.



[Wikipedia]







We can consider either W or W' as the word's representation.



Word Analogies







Word Analogies



25 IM FUCUS DAS LEBEN

What is word2vec?

- word2vec is not a single algorithm
- It is a **software package** for representing words as vectors, containing:
 - Two distinct models
 - CBoW
 - Skip-Gram (SG)
 - Various training methods
 - Negative Sampling (NS)
 - Hierarchical Softmax
 - A rich preprocessing pipeline
 - Dynamic Context Windows
 - Subsampling
 - Deleting Rare Words

UNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

Marco saw a furry little wampimuk hiding in the tree.



Marco saw a furry little wampimuk hiding in the tree.



Marco saw a furry little wampimuk hiding in the tree.

<u>words</u> wampimuk wampimuk wampimuk wampimuk





...

- SGNS finds a vector \vec{w} for each word w in our vocabulary V_W
- Each such vector has d latent dimensions (e.g. d = 100)
- Effectively, it learns a matrix W whose rows represent V_W
- **Key point:** it also derives a similar auxiliary matrix *C* of context vectors
- In fact, each word has two embeddings







• Maximize: $\sigma(\vec{w} \cdot \vec{c})$						
 c was observed with w 						
words contexts						
wampimuk	furry					
wampimuk	little					
wampimuk	hiding					
wampimuk in						



UNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

• Maximize: a	$\sigma(\vec{w}\cdot\vec{c})$	• Minimize: $\sigma(\vec{w} \cdot \vec{c}')$				
– <i>c</i> was obs	erved with w	 - c' was hallucinated with w 				
<u>words</u>	<u>contexts</u>					
wampimuk	furry	<u>words</u>	<u>contexts</u>			
wampimuk	little	wampimuk	Australia			
wampimuk	hiding	wampimuk	cyber			
wampimuk	in	wampimuk	the			
		wampimuk	1985			



"word2vec Explained..."

Goldberg & Levy, arXiv 2014

- "Negative Sampling"
- SGNS samples k contexts c' at random as negative examples
- "Random" = unigram distribution

$$P(c) = \frac{\#c}{|D|}$$

• **Spoiler:** Changing this distribution has a significant effect



• Take SGNS's embedding matrices (*W* and *C*)





- Take SGNS's embedding matrices (*W* and *C*)
- Multiply them
- What do you get?





- A $V_W \times V_C$ matrix
- Each cell describes the relation between a specific wordcontext pair

$$\vec{w} \cdot \vec{c} = ?$$





- Levy&Goldberg [2014] **proved** that for large enough *d* and enough iterations ...
- ... one obtains the word-context PMI matrix





- Levy&Goldberg [2014] proved that for large enough d and enough iterations ...
- ... one obtains the word-context PMI matrix ...
- shifted by a global constant





- SGNS is doing something very similar to the older approaches
- SGNS factorizes the traditional word-context PMI matrix
- So does SVD!
- GloVe factorizes a similar word-context matrix



But embeddings are still better, right?

- Plenty of evidence that embeddings outperform traditional methods
 - "Don't Count, Predict!" (Baroni et al., ACL 2014)
 - GloVe (Pennington et al., EMNLP 2014)
- How does this fit with our story?



The Big Impact of "Small" Hyperparameters

- word2vec & GloVe are more than just algorithms...
- Introduce **new hyperparameters**
- May seem minor, but make a big difference in practice



New Hyperparameters

•	 Preprocessing Dynamic Context Windows Subsampling Deleting Rare Words 	(word2vec)
•	Postprocessing – Adding Context Vectors	(GloVe)
•	Association Metric – Shifted PMI	(SGNS)

Context Distribution Smoothing

JNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME IM FOCUS DAS LEBEN 43

Dynamic Context Windows

Marco saw a furry little wampimuk hiding in the tree.



Dynamic Context Windows

saw a furry little wampimuk hiding in the tree



IM FOCUS DAS LEBEN 45

saw a furry little wampimuk hiding in the tree

Word2vec:	$\frac{1}{4}$	$\frac{2}{4}$	3 4	$\frac{4}{4}$	$\frac{4}{4}$	3 4	$\frac{2}{4}$	$\frac{1}{4}$
GloVe:	$\frac{1}{4}$	$\frac{1}{3}$	<u>1</u> 2	$\frac{1}{1}$	$\frac{1}{1}$	<u>1</u> 2	$\frac{1}{3}$	$\frac{1}{4}$
Aggressive:	$\frac{1}{8}$	$\frac{1}{4}$	<u>1</u> 2	$\frac{1}{1}$	$\frac{1}{1}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$

The Word-Space Model (Sahlgren, 2006)



- SGNS creates word vectors \vec{w}
- SGNS creates auxiliary context vectors \vec{c}
 - So do GloVe and SVD



- SGNS creates word vectors \vec{w}
- SGNS creates auxiliary context vectors \vec{c}
 - So do GloVe and SVD
- Instead of just \vec{w}
- Represent a word as: $\vec{w} + \vec{c}$
- Introduced by Pennington et al. (2014)
- Only applied to GloVe



Context Distribution Smoothing

- SGNS samples $c' \sim P$ to form **negative** (w, c') examples
- Our analysis assumes P is the unigram distribution

$$P(c) = \frac{\#c}{\sum_{c' \in V_C} \#c'}$$



Context Distribution Smoothing

- SGNS samples $c' \sim P$ to form **negative** (w, c') examples
- Our analysis assumes *P* is the unigram distribution
- In practice, it's a **smoothed** unigram distribution

$$P^{0.75}(c) = \frac{(\#c)^{0.75}}{\sum_{c' \in V_c} (\#c')^{0.75}}$$

• This little change makes a big difference



Context Distribution Smoothing

- We can adapt context distribution smoothing to PMI!
- Replace P(c) with $P^{0.75}(c)$:

$$PMI^{0.75}(w,c) = \log \frac{P(w,c)}{P(w) \cdot P^{0.75}(c)}$$

- Consistently improves **PMI** on **every task**
- Always use Context Distribution Smoothing!



Represent the meaning of **sentence/text**

- Paragraph vector (2014, Quoc Le, Mikolov)
 - Extend word2vec to text level
 - Also two models: add paragraph vector as the input







Don't Count, Predict! [Baroni et al., 2014]

- "word2vec is better than count-based methods"
- Hyperparameter settings account for most of the reported gaps
- Embeddings do **not** really outperform count-based methods
- No unique conclusion available



The Contributions of Word Embeddings

Novel Algorithms

(objective + training method)

- Skip Grams + Negative Sampling
- CBOW + Hierarchical Softmax
- Noise Contrastive Estimation
- GloVe
- ..

New Hyperparameters

(preprocessing, smoothing, etc.)

- Subsampling
- Dynamic Context Windows
- Context Distribution Smoothing
- Adding Context Vectors

What's really improving performance?



Improving Distributional Similarity with Lessons Learned from Word Embeddings, Omer Levy, Yoav Goldberg, Ido Dagan