Web-Mining Agents

Word Semantics and Latent Relational Structures

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

Tanya Braun (Übungen)



IM FOCUS DAS LEBEN

Word-Word Associations in Document Retrieval

Recap

- LSI: Documents as vectors, dimension reduction
- Topic Modeling
 - Topic = Word distribution
 - From LDA-Model: P(Z | w)
 - Assumption: Bag of words model (independence, naïve Bayes, unigram distribution)

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 word semantics



Approaches for Representing Word Semantics

Beyond bags of words

 Distributional Semantics	 Word Embeddings (Predict) Inspired by deep learning word2vec				
(Count) Used since the 90's Sparse word-context	(Mikolov et al., 2013) GloVe				
PMI/PPMI matrix Decomposed with SVD	(Pennington et al., 2014)				
Underlying Theory: The Distributional Hypothesis (<i>Harris, '54; Firth, '57</i>) "Similar words occur in similar contexts"					

https://www.tensorflow.org/tutorials/word2vec https://nlp.stanford.edu/projects/glove/



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)





The Contributions of Word Embeddings

Novel Algorithms New Hyperparameters (preprocessing, smoothing, etc.) (objective + training method) Skip Grams + Negative Sampling **Subsampling** • CBOW + Hierarchical Softmax **Dynamic Context Windows** • Noise Contrastive Estimation **Context Distribution Smoothing** • • Adding Context Vectors GloVe •

What's really improving performance?



Embedding Approaches

- Represent each word with a low-dimensional vector
- Word similarity = vector similarity
- Key idea: Predict surrounding words of every word
- Faster and can easily incorporate a new sentence/document or add a word to the vocabulary



Represent the meaning of **word** – word2vec

- 2 basic network models:
 - Continuous Bag of Word (CBOW): use a window of word 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

























Logistic function

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),
- x_0 = the *x*-value of the sigmoid's midpoint,
- L = the curve's maximum value, and
- *k* = the steepness of the curve.^[1]





[Wikipedia]

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.









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



Some interesting results

Word Analogies

Test for linear relationships, examined by Mikolov et al. (2014)



UNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

Word analogies



20 IM FUCUS DAS LEBEN

What is word2vec?

- word2vec is not a single algorithm
- It is a software package for representing words as vectors, containing:

(SG)

(NS)

- Two distinct models
 - CBoW
 - Skip-Gram
- Various training methods
 - Negative Sampling
 - Hierarchical Softmax
- A rich preprocessing pipeline
 - Dynamic Context Windows
 - Subsampling
 - Deleting Rare Words



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 learns a similar auxiliary matrix C of context vectors
- In fact, each word has two embeddings







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



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

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





• We **proved** that for large enough *d* and enough iterations





- 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

- Context Distribution Smoothing

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

UNIVERSITÄT ZU LÜBECK IINSTITUT FÜR INFORMATIONSSYSTEME

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 41

Dynamic Context Windows

saw a furry little wampimuk hiding in the tree

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

The Word-Space Model (Sahlgren, 2006)



Adding Context Vectors

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



Adding Context Vectors

- 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'~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'~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



Latent Relational Structures

Processing natural language data:

- ✓ Tokenization/Sentence Splitting
- ✓ Part-of-speech (POS) tagging
- Phrase chunking
- Named entity recognition
- Coreference resolution
- Semantic role labeling



Phrase Chunking

JNIVERSITÄT ZU LÜBECK

TUT FÜR INFORMATIONSSYSTEME

• Identifies phrase-level constituents in sentences

[NP Boris] [ADVP regretfully] [VP told] [NP his wife][SBAR that] [NP their child] [VP could not attend] [NP night school] [PP without] [NP permission].

- Useful for filtering: identify e.g. only noun phrases, or only verb phrases
- Used as source of features, e.g. distance, (abstracts away determiners, adjectives, for example), sequence,...
 - More efficient to compute than full syntactic parse
 - Applications in e.g. Information Extraction getting (simple) information about concepts of interest from text documents
- Hand-crafted chunkers (regular expressions/finite automata)
- HMM/CRF-based chunk parsers derived from training data

An Introduction to Machine Learning and Natural Language Processing Tools, V. Srikumar, M. Sammons, N. Rizzolo

Named Entity Recognition

- Identifies and classifies strings of characters representing proper nouns
- [PER Neil A. Armstrong], the 38-year-old civilian commander, radioed to earth and the mission control room here: "[LOC Houston], [ORG Tranquility] Base here; the Eagle has landed."
- Useful for filtering documents
 - "I need to find news articles about organizations in which Bill Gates might be involved..."
- **Disambiguate tokens:** "Chicago" (team) vs. "Chicago" (city)
- Source of abstract features
 - E.g. "Verbs that appear with entities that are Organizations"
 - E.g. "Documents that have a high proportion of Organizations"



Named Entity Recogniton: Definition

- NE involves identification of proper names in texts, and classification into a set of predefined categories of interest
 - Three universally accepted categories: person, location and organisation
 - Other common tasks: recognition of date/time expressions, measures (percent, money, weight etc), email addresses etc.
 - Other domain-specific entities: names of drugs, medical conditions, names of ships, bibliographic references etc
- NER ist not easy



Named Entity Classification

- Category definitions are intuitively quite clear, but there are many grey areas.
- Many of these grey area are caused by **metonymy**.
 Person vs. Artefact: "The **ham sandwich** wants his bill." vs "Bring me a **ham sandwich**."
- Organisation vs. Location : "England won the World Cup" vs. "The World Cup took place in England".
 Company vs. Artefact: "shares in MTV" vs. "watching MTV" Location vs. Organisation: "she met him at Heathrow" vs.

"the **Heathrow** authorities"



Basic Problems in NE

- Variation of NEs e.g. John Smith, Mr Smith, John.
- Ambiguity of NE types
 - John Smith (company vs. person)
 - May (person vs. month)
 - Washington (person vs. location)
 - 1945 (date vs. time)
- Ambiguity with common words, e.g. "may"



More complex problems in NER

- Issues of style, structure, domain, genre etc.
 - Punctuation, spelling, spacing, formatting,all have an impact

Dept. of Computing and Maths Manchester Metropolitan University Manchester United Kingdom

> Tell me more about Leonardo> Da Vinci



List Lookup Approach

- System that recognises only entities stored in its lists (gazetteers).
- Advantages Simple, fast, language independent, easy to retarget
- Disadvantages collection and maintenance of lists, cannot deal with name variants, cannot resolve ambiguity



Shallow Parsing Approach

 Internal evidence – names often have internal structure. These components can be either stored or guessed.

location:

CapWord + {City, Forest, Center} *e.g. Sherwood Forest* Cap Word + {Street, Boulevard, Avenue, Crescent, Road} *e.g. Portobello Street*



Shallow Parsing Approach

 External evidence - names are often used in very predictive local contexts

Location:

"to the" COMPASS "of" CapWord e.g. to the south of Loitokitok
"based in" CapWord e.g. based in Loitokitok
CapWord "is a" (ADJ)? GeoWord e.g. Loitokitok is a friendly city



Difficulties in Shallow Parsing Approach

• **Ambiguously capitalised words** (first word in sentence) [All American Bank] vs. All [State Police]

Semantic ambiguity

"John F. Kennedy" = airport (location) "Philip Morris" = organisation

Structural ambiguity

[Cable and Wireless] vs. [Microsoft] and [Dell]

[Center for Computational Linguistics] vs. message from [City Hospital] for [John Smith].



Coreference

- Identify all phrases that refer to each entity of interest i.e., group mentions of concepts
- [Neil A. Armstrong], [the 38-year-old civilian commander], radioed to [earth]. [He] said the famous words, "[the Eagle] has landed"."
- The Named Entity Recognizer only gets us part-way...
- ...if we ask, "what actions did Neil Armstrong perform?", we will miss many instances (e.g. "He said...")
- Coreference resolver abstracts over different ways of referring to the same person
 - Useful in feature extraction, information extraction



Semantic Role Labeling (SRL)

Input Text:

A car bomb that exploded outside the U.S. military base in Beniji killed 11 Iraqi citizens.

Result: Complete!

General Explanation of Argument Labels



- SRL reveals relations

 and arguments in the
 sentence (where
 relations are expressed
 as verbs)
- Cannot abstract over variability of expressing the relations – e.g. kill vs. murder vs. slay…



An Introduction to Machine Learning and Natural Language Processing Tools, V. Srikumar, M. Sammons, N. Rizzolo

Why is SRL Important – Applications

- Question Answering
 - Q: When was Napoleon defeated?
 - Look for: [PATIENT Napoleon] [PRED defeat-synset] [ARGM-TMP *ANS*]
- Machine Translation

English (SVO) [_{AGENT} The little boy] [_{PRED} kicked] [_{THEME} the red ball] [_{ARGM-MNR} hard] Farsi (SOV)[AGENT PESAR koocholo] boy-little[THEME toop germezi] ball-red[ARGM-MNR moqtam] hard-adverb[PRED zaad-e] hit-past

- Document Summarization
 - Predicates and Heads of Roles summarize content

Information Extraction

SRL can be used to construct useful rules for IE



Automatic Semantic Role Labeling S. Wen-tau Yih, K. Toutanova

Some History

- Minsky 74, Fillmore 1976: Frames describe events or situations
 - Multiple participants, "props", and "conceptual roles"
 - E.g., agent, instrument, target, time, ...
- Levin 1993: verb class defined by sets of frames (meaningpreserving alternations) a verb appears in
 - {break, shatter,..}: Glass X's easily; John Xed the glass, ...
 - *Cut* is different: *The window broke*; **The window cut*.
- FrameNet, late '90s: based on Levin's work: large corpus of sentences annotated with *frames*
- PropBank



Automatic Semantic Role Labeling S. Wen-tau Yih, K. Toutanova



[Agent Kristina] hit [Target Scott] [Instrument with a baseball] [Time yesterday].



IM FOCUS DAS LEBEN

Proposition Bank (PropBank) [Palmer et al. 05]

- Transfer sentences to propositions

 Kristina hit Scott → hit(Kristina,Scott)
- Penn TreeBank → PropBank
 - Add a semantic layer on Penn TreeBank
 - Define a set of semantic roles for each verb
 - Each verb's roles are numbered

...[A0 the company] to ... offer [A1 a 15% to 20% stake] [A2 to the public]

...[A0 Sotheby's] ... offered [A2 the Dorrance heirs] [A1 a money-back guarantee]

- ...[A1 an amendment] *offered* [A0 by Rep. Peter DeFazio] ...
- ...[A2 Subcontractors] will be offered [A1 a settlement] ...



Latent Relational Structures

Input Text:

A car bomb that exploded outside the U.S. military base in Beniji killed 11 Iraqi citizens.

Result: Complete!

General Explanation of Argument Labels



- SRL reveals relations

 and arguments in the
 sentence (where
 relations are expressed
 as verbs)
- Cannot abstract over variability of expressing the relations – e.g. kill vs. murder vs. slay…



An Introduction to Machine Learning and Natural Language Processing Tools, V. Srikumar, M. Sammons, N. Rizzolo

Collective Learning on Multi-Relational Data





Towards Latent Relational Structures

- Modelling binary relations as a tensor: Two modes of a tensor refer to the entities, one mode to the relations.
- The entries of the tensor are 1 when a relation between two entities exists and 0 otherwise
- We use the RDF formalism to model relations as (subject, predicate, object) triples





Maximilian Nickel, Volker Tresp, Hans-Peter Kriegel A Three-Way Model for Collective Learning on Multi-Relational Data In Proc. 28th International Conference on Machine Learning, **2011**

Motivation

Why Tensors?

- Modelling simplicity: Multiple binary relations can be expressed straightforwardly as a three-way tensor
- **No structure learning**: Not necessary to have information about independent variables, knowledge bases, etc. or to infer it from data
- Expected performance: Relational domains are high-dimensional and sparse, a setting where factorization methods have shown very good results



Tensor Factorization with Rescal

 RESCAL takes the inherent structure of dyadic relational data into account, by employing the tensor factorization

$$X_k \approx A R_k A^T$$

- A is a n × r matrix, representing the global entity-latent-component space
- R_k is an asymmetric $r \times r$ matrix that specifies the interaction of the latent components per predicate



Solving Relational Learning Tasks

- Link Prediction: To predict the existence of a relation between two entities, it is sufficient to look at the rank-reduced reconstruction of the appropriate slice AR_kA^T
- Collective Classification: Can be cast as a link prediction problem by including the classes as entities and adding a classOf relation. Alternatively, standard classification algorithms could be applied to the entites' latent-component representation A
- Link-based Clustering: Since the entities'latent-component representation is computed considering all relations, Link-based clustering can be done by clustering the entities in the latent-component space A


$$\|\mathbf{A}\|_F \equiv \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}.$$

• To compute the factorization, we solve the optimization problem

$$\arg\min_{A,R_k} loss(A,R_k) + reg(A,R_k)$$

where loss is the loss function

$$loss(A, R_k) = \frac{1}{2} \sum_k \|\mathcal{X}_k - AR_k A^T\|_F^2$$

and reg is the regularization term

$$reg(A, R_k) = \frac{1}{2}\lambda\left(\|A\|_F^2 + \sum_k \|R_k\|_F^2\right)$$



© M. Nickel https://github.com/mnick/rescal.py

Predict party membership of US (vice) presidents



• Helpful to consider element-wise version of the loss function f

$$f(A, R_k) = \frac{1}{2} \sum_{i,j,k} \left(\mathcal{X}_{ijk} - \boldsymbol{a}_i^T R_k \boldsymbol{a}_j \right)^2$$







• Helpful to consider element-wise version of the loss function f

$$f(A, R_k) = \frac{1}{2} \sum_{i,j,k} \left(\mathcal{X}_{ijk} - \boldsymbol{a}_i^T R_k \boldsymbol{a}_j \right)^2$$



Predict party membership of US (vice) presidents



- Helpful to consider element-wise version of the loss function f

$$f(A, R_k) = rac{1}{2} \sum_{i,j,k} \left(\mathcal{X}_{ijk} - \boldsymbol{a}_i^T R_k \boldsymbol{a}_j \right)^2$$



Predict party membership of US (vice) presidents



 \blacksquare Helpful to consider element-wise version of the loss function f

$$f(A, R_k) = \frac{1}{2} \sum_{i,j,k} \left(\mathcal{X}_{ijk} - \boldsymbol{a}_i^T R_k \boldsymbol{b}_j \right)^2$$



Predict party membership of US (vice) presidents



• Helpful to consider element-wise version of the loss function f

$$f(A, R_k) = \frac{1}{2} \sum_{i,j,k} \left(\mathcal{X}_{ijk} - \boldsymbol{a}_i^T R_k \boldsymbol{b}_j \right)^2$$



- Collective learning is performed via the entities' latent-component representation
- Important aspect of the model: Entities have a unique latent-component representation, regardless of their occurrence as subjects or objects



Using Learned Relational Networks for IR

- Query answering: indirect queries requiring chains of reasoning
- KB Completion: exploits redundancy in the KB + chains to infer missing facts

Freebase 15k benchmar	k	
Methods	Hits@1	0
Unstructured [Bordes et al., 2014]	4.5	baseline method
RESCAL [Nickel et al., 2011]	28.4	tensor factorization
SE [Bordes et al., 2011]	28.8	
SME [Bordes et al., 2014]	31.3	
LFM [Jenatton et al., 2012]	26.0	
TransE [Bordes et al., 2013]	34.9	deen NN
ConvNets [Shi and Zhu, 2015]	37.7	embedding
TransH [Wang et al., 2014b]	45.7	5
TransR [Lin et al., 2015b]	48.4	
PTransE [Lin et al., 2015a]	51.8	



Using Learned Relational Networks for IR

