
Text Mining Seminar

GloVe: Global Vectors for Word Representation

Munich, March 26, 2018

DEPARTMENT OF STATISTICS
Ludwig Maximilian University of Munich



Degree course: M.Sc. Statistics

Student:

Daniel Schalk
(11470019)

Advisor:

Prof. Dr. Christian Heumann

Contents

1	Introduction	1
2	GloVe Model	2
2.1	Terminology	2
2.2	Word-Word Co-Occurrence Matrix	2
2.3	The Model	3
3	Evaluating Word Vectors	7
3.1	Evaluation Metrics	7
3.2	p-Norm in high Dimensions	7
3.3	Semantic and Analogies	7
3.4	Questions Words File	8
3.5	Hyperparameter Tuning	8
4	About the Data	10
4.1	The Language	10
4.2	Common Sources	10
5	Real Word Vectors	12
5.1	Own Word Vectors	12
5.2	Pre Trained Word Vectors	12
6	Outlook and Conclusion	13
6.1	Outlook	13
6.2	Conclusion	13
	List of Figures	15
	List of Tables	16

1 Introduction

An important thing when it comes to text mining is to create word embeddings. An important technique within the context of text mining is to map words to vectors. Those so called word embeddings are used for further analyses or more general as features in statistical models. GloVe is a technique to create word embeddings out of a given corpus. A very interesting thing is that we want to learn from text without having labels. Hence, we have an unsupervised task.

Word vectors should be able to represent the human understanding of text in a very basic way. That means that word vectors should be able to display the analogy of different words as well as general semantic questions. How this is achieved shows image .

- word vectors (image + explanation)

In the following we want to discuss some important topics related to GloVe. After deriving the model with a more theoretical point of view in section 2. In section 3 we also want to take a look at how to evaluating given word vectors which is quite interesting since we handle an unsupervised task. After that we take a short look at the data and common sources for text corpora in section 4.

Then we know how to evaluate word vectors and have an idea about the data. With that knowledge we may ask how different methods to create word embeddings or different data sources influences the quality of the word vectors. This is discussed shortly in chapter 5. After that we take a short look how we can use GloVe embeddings for further text classification techniques followed from a small conclusion in section 6.

2 GloVe Model

2.1 Terminology

Notation	Symbol	Description
Vocabulary	V	Set of unique words of the given corpus
Word vector	w (w_i)	Word vector (of corresponding word i)
Context vector	\tilde{w} (\tilde{w}_i)	Context vector (of corresponding word i)
Dimension	d	Dimension of word vectors w
Word-word co-occurrence matrix	$X \in \mathbb{R}^{ V \times V }$	Matrix of context counts X_{ij}
Loss function		
Empirical risk		

Table 2.1: Overview about the used terminology.

2.2 Word-Word Co-Occurrence Matrix

Base of the model is the word-word co-occurrence matrix $X \in \mathbb{R}^{|V| \times |V|}$, where $|V|$ is the number of different words which occurs within a given corpus. In X every entry X_{ij} describes how often word j occurs in context of word i with a given window size. Therefore we have the following properties:

- $X_i = \sum_k X_{ik}$: Number of words which occur in the context of i .
- $P_{ij} = P(j|i) = X_{ij}/X_i$: "Probability" of word j occurs in context of i .
- The ratio P_{ik}/P_{jk} describes if word k is more related to ...
 - ... word i if $P_{ik}/P_{jk} > 1$
 - ... word j if $P_{ik}/P_{jk} < 1$
 - ... or similar related if $P_{ik}/P_{jk} \approx 1$

Example: word-word co-occurrence Matrix

To understand how an entry of the word-word co-occurrence matrix is computed we take a look at the following sample corpus:

2 GloVe Model

$$A D C E A D E B A C E D \Rightarrow V = \{A B C D E\}, |V| = 5$$

Furthermore, we need to specify a window size which indicate how much words we want to look at around a specific word. We choose a window size of 2 (the 2 words of the either side) for this example. With the given corpus and the window size we now compute the counts of how often word D occurs in context of word A (X_{AD}):

A D C E A D E B A C E D \Rightarrow D occurs 1 time

A D C E A D E B A C E D \Rightarrow D occurs 1 time

A D C E A D E B A C E D \Rightarrow D occurs 0 times

\Rightarrow D occurs 2 times in context of A: $X_{AD} = 2$

Finally, we can do this for every word combination of $(i, j) \in V \times V$ and calculate the ratio P_{AD} :

$$X = \begin{array}{c|ccccc} & A & B & C & D & E \\ \hline A & 0 & 1 & 3 & 2 & 4 \\ B & 1 & 0 & 1 & 1 & 1 \\ C & 3 & 1 & 0 & 2 & 2 \\ D & 2 & 1 & 2 & 0 & 4 \\ E & 4 & 1 & 2 & 4 & 0 \end{array} \Rightarrow P_{AD} = X_{AD}/X_A = 2/10$$

Sidenote:

Some packages (like R's `text2vec`) have the option to give a weight vector with as many weights as the window size. It is then common to weight context words closer to the inspected word higher than words which aren't that close.

2.3 The Model

GloVe was introduced by Jeffrey Pennington, Richard Socher and Christopher D. Manning [PSDM14]:

"... GloVe is an unsupervised learning algorithm for obtaining vector representations for words. Training is performed on aggregated global word-word co-occurrence statistics from a corpus, and the resulting representations showcase interesting linear substructures of the word vector space."

Remember: We want to create word vectors $w_i \in \mathbb{R}^d$ and want to use the ratio P_{ik}/P_{jk} since the ratio is more appropriate to take words into context.

The idea is to take a function F which takes the three words i , j and k , given by using the ratio, and map F to that ratio:

$$F(w_i, w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$

Here w_i, w_j are context vectors while \tilde{w}_k is the context vector. F is unknown at this stage and basically can be any function. But we remember that we want to keep the linear

2 GloVe Model

structure of the vector space \mathbb{R}^d . So it is important to choose a function F which is able to hold the linear structure. For instance we can choose F to be a neural net. But that wouldn't guarantee that F keeps the linear structure.

We now parameterize every word vectors, therefore we have $d \cdot |V|$ parameter to estimate to get word vectors. But how should we estimate those parameters without a specific function F ? There are several steps to come up with a solution which fulfills the desired behaviour:

1. Reduce number of possible outputs of F by taking $w_i - w_j$ instead of two vectors w_i and w_j :

$$F(w_i - w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$

2. We want to keep the linear structure between the word vectors and context vector:

$$F((w_i - w_j)^T \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$

3. Keep exchange symmetry. We want be able to exchange word vectors w with context vectors \tilde{w} without getting different results. To obtain this property we apply to tricks:

- a) We restrict F to be a homomorphism between $(\mathbb{R}, f) = (\mathbb{R}, +)$ and $(\mathbb{R}_+, g) = (\mathbb{R}_+, \cdot)$. That means we expect F to fulfill:

$$F(a + b) = \underbrace{F(f(a, b)) = g(F(a), F(b))}_{\text{Definition of homomorphism}} = F(a)F(b)$$

With $a = w_i^T \tilde{w}_k$ and $b = -w_j^T \tilde{w}_k$.

We notice, that the upper equation implies a functional equation which has just one solution:

$$F = \exp_a$$

We choose $a = \exp(1)$, therefore $F = \exp$.

$$F((w_i - w_j)^T \tilde{w}_k) = \exp((w_i - w_j)^T \tilde{w}_k) = \frac{\exp(w_i^T \tilde{w}_k)}{\exp(w_j^T \tilde{w}_k)} = \frac{P_{ik}}{P_{jk}}$$

$$\begin{aligned} \Rightarrow \exp(w_i^T \tilde{w}_k) &= P_{ik} = \frac{X_{ik}}{X_i} \\ \Leftrightarrow w_i^T \tilde{w}_k &= \log(X_{ik}) - \log(X_i) \\ \Leftrightarrow w_i^T \tilde{w}_k + \log(X_i) &= \log(X_{ik}) \end{aligned}$$

But: How to handle $X_{ik} = 0$? We handle this later.

- b) Since $\log(X_i)$ is independent of k we can put this in a bias term. This bias term can be decomposed into a term b_i from w_i and \tilde{b}_k from \tilde{w}_k :

$$\Rightarrow w_i^T \tilde{w}_k + b_i + \tilde{b}_k = \log(X_{ik})$$

This is now an model equation which we can use for training the word vectors. We also note, that we have transformed the unsupervised task into a supervised task which we know how to handle.

Sidenote:

Point 2. looks a bit like the model equation of a generalized linear model (glm) with linear predictor $(x_i - w_j)^T \tilde{w}_k$, response function F and response P_{ik}/P_{jk} . The glm takes a linear structure and transforms the output on an interval of plausible values. This also holds for GloVe by using a homomorphism as F which preserves the structure between the two used algebraic structures. Therefore, it is possible to keep the linearity of the used vector space \mathbb{R}^d .

Model Equation

A problem appears for $X_{ik} = 0$. This is definitely the case since X is a sparse matrix. We don't want to drop the sparsity due to convenient memory handling. Therefore, we use an additive shift in the logarithm:

$$\log(X_{ik}) \rightarrow \log(X_{ik} + 1)$$

This maintains the sparsity:

$$X_{ik} = 0 \Rightarrow \log(X_{ik} + 1) = \log(1) = 0$$

The final model equation which we use for training is

$$w_i^T \tilde{w}_k + b_i + \tilde{b}_k = \log(X_{ik} + 1)$$

Empirical Risk

Kurz Loss funktion erwähnen!

To estimate the word vectors GloVe uses a weighted least squares approach:

$$J = \sum_{i=1}^{|V|} \sum_{j=1}^{|V|} f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log(X_{ij} + 1))^2$$

[PSDM14] proposed as weight function:

$$f(x) = \begin{cases} (x/x_{\max})^\alpha, & x < x_{\max} \\ 1, & x \geq x_{\max} \end{cases}$$

Plotten wie funktion aussieht!

Using this function also introduces two new hyperparameters α and x_{\max} . The ordinary way to find good values for α and x_{\max} would be to use a tuning method. The problem now is that for all tuning methods we need to evaluate the model. Since we want to handle an unsupervised task evaluating isn't straightforward.

Algorithm used for Fitting

Basically, GloVe uses a gradient descend technique to minimize the objective J . Nevertheless, using ordinary gradient descend would be way too expensive in practice. A more appropriate gradient based method is used here called adaptive gradient descent (AdaGrad). A short description of the most important points of AdaGrad transferred to text mining are:

2 GloVe Model

- Individual learning rate in each iteration.
- Adaption of the learning rate to the parameters, performing larger updates for infrequent and smaller updates for frequent parameter.
- Well-suited for sparse data, improve robustness of (stochastic) gradient descent.
- For GloVe, infrequent words require much larger updates than frequent ones.

For a more detailed explanation see [\[Rud16\]](#).

3 Evaluating Word Vectors

3.1 Evaluation Metrics

We need something to measure the distance between vectors:

- **Euclidean Distance:**

$$d_{\text{euclid}}(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- **Cosine Distance**

$$d_{\text{cosine}}(x, y) = 1 - \cos(\angle(x, y)) = 1 - \underbrace{\frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}}_{\text{Cosine Similarity}}$$

In general we could use every metric to evaluate the model. But the cosine similarity is more robust in terms of curse of dimensionality.

Images of word cloud and description (cosine similarity, which wv etc.

3.2 p-Norm in high Dimensions

[AHK01] have shown, that in high dimensions the ratio of the maximal norm divided by the minimal norm of n points x_1, \dots, x_n which are randomly drawn converges in probability to 1 for increasing dimension d :

$$\text{p} \lim_{d \rightarrow \infty} \frac{\max_k \|x_k\|_2}{\min_k \|x_k\|_2} = 1$$

\Rightarrow Points are concentrated on the surface of a hyper sphere using the euclidean norm. The same holds for every p -Norm.

Image of p norm in high dimension

3.3 Semantic and Analogies

One thing we now can do is to ask for semantic analogies between words. Something like:

paris behaves to france like berlin to ? animal behaves to animals like people to ? i behaves to j like k to l

3 Evaluating Word Vectors

Therefore, we have 3 given word vectors w_i , w_j and w_k . To get the desired fourth word l we use the linearity of the word vector space:

$$w_l \approx w_j - w_i + w_k$$

Furthermore, we obtain \hat{l} from our model and a given metric $d(w_i, w_j)$ (mostly $d = d_{\text{cosine}}$) by computing:

$$\hat{l} = \arg \min_{l \in V} d(w_j - w_i + w_k, w_l)$$

3.4 Questions Words File

To evaluate trained word vectors, [MCCD13] provide a word similarity task. This task is given within a question words file which contains about 19544 semantic analogies:

```
: capital-common-countries
Athens Greece Baghdad Iraq
Athens Greece Bangkok Thailand
Athens Greece Beijing China
Athens Greece Berlin Germany
Athens Greece Bern Switzerland
Athens Greece Cairo Egypt
Athens Greece Canberra Australia
Athens Greece Hanoi Vietnam
Athens Greece Havana Cuba
Athens Greece Helsinki Finland
```

3.5 Hyperparamter Tuning

- Now tuning is “possible” for a given task specified in the questions words file.
- [PSDM14] have tuned the model and came to good values (just empirical without a proof):
 - $\alpha = 0.75$
 - $x_{\text{max}} = 100$ (does just have a weakly influence on performance)
- Note that we are dependent on this file and can just test on this file. If we want to test other properties of the model we need other files.

3 Evaluating Word Vectors

Category	Number of Test Lines	Example
capital-common-countries	506	Athens Greece Baghdad Iraq
capital-world	4524	Abuja Nigeria Accra Ghana
currency	866	Algeria dinar Angola kwanza
city-in-state	2467	Chicago Illinois Houston Texas
family	506	boy girl brother sister
gram1-adjective-to-adverb	992	amazing amazingly apparent apparently
gram2-opposite	812	acceptable unacceptable aware unaware
gram3-comparative	1332	bad worse big bigger
gram4-superlative	1122	bad worst big biggest
gram5-present-participle	1056	code coding dance dancing
gram6-nationality-adjective	1599	Albania Albanian Argentina Argentinean
gram7-past-tense	1560	dancing danced decreasing decreased
gram8-plural	1332	banana bananas bird birds
gram9-plural-verbs	870	decrease decreases describe describes

Table 3.1: Examples for questions per category within the question word file.

4 About the Data

4.1 The Language

Be careful with the language:

German term list	English term list
Wassermolekuel	hydrogen
Wasserstoff	hydrogen-bonding
Wasserstoffatom	
Wasserstoffbindung	
Wasserstoffbrueckenbildung	
Wasserstoffbrueckenbindung	
Wasserstoffhalogenid	
Wasserstoffverbindung	

Table 4.1: Copied from Grammar & Corpora 2009 [Kon11]. Excerpt of the resulting English and German term lists focusing on the term *hydrogen* (German: *Wasserstoff*).

For instance, German has much more rare words and a bigger variety (harder for modelling) than English.

4.2 Common Sources

- We need a lot of words to train the model.
- Often crawled from the web. This is mostly followed by a lot of preprocessing (regular expressions, filtering stop words etc.).
- How big should the corpus and the vocabulary be to get good word vectors?
- Does different corpora imply a different quality of the word vectors?
- **Wikipedia Dump + Gigaword 5:** Wikipedia gives access to all articles collected within one XML file (unzipped about 63 GB). [PSDM14] combines this with Gigaword 5, an archive of newswire text data.

→ 6 billion tokens and 400 thousand word vocabulary
- **Common Crawl:** Published by Amazon Web Services through its Public Data Sets program in 2012. The data was crawled from the whole web and contains 2.9 billion

4 *About the Data*

web pages and over 240 TB of data.

→ 42 billion tokens and 1.9 million word vocabulary or

→ 820 billion tokens and 2.2 million word vocabulary

- **Twitter:** [PSDM14] crawled 2 billion tweets.

→ 27 billion tokens and 1.2 million word vocabulary

5 Real Word Vectors

5.1 Own Word Vectors

5.2 Pre Trained Word Vectors

6 Outlook and Conclusion

6.1 Outlook

6.2 Conclusion

6 *Outlook and Conclusion*

List of Figures

List of Tables

- 2.1 Overview about the used terminology. 2
- 3.1 Examples for questions per category within the question word file. 9
- 4.1 Term list comparison between German and English 10

Bibliography

- [AHK01] Charu C Aggarwal, Alexander Hinneburg, and Daniel A Keim. On the surprising behavior of distance metrics in high dimensional space. In *International conference on database theory*, pages 420–434. Springer, 2001.
- [Kon11] Marek Konopka. *Grammar & Corpora 2009*, volume 1. BoD–Books on Demand, 2011.
- [MCCD13] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [PSDM14] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
- [Rud16] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.