

# Limits of Deep Learning

Simon Jenni  
University of Bern  
Email: simujenni@students.unibe.ch

Simon Brulhart  
University of Fribourg  
Email: simon.brulhart@unifr.ch

**Abstract**—Current natural language models are able to capture semantic and syntactic relationships between words. This allows them to answer questions in the form of analogies: "X is to Y as U is to ?". We study the performance of the word2vec tool on this task. To this end, we analyse the sensitivity to different hyper-parameter choices. We generate new questions and demonstrate the models (in)ability to learn and answer on domain-specific knowledge. Using new questions we also highlight the general limits of the model's capabilities. To offer a more detailed view of the model's performance, we introduce a broader definition of accuracy which is taking into account how far off the model is in a failure case. Finally, we experiment with different distance functions to compute the similarity between words and demonstrate when they can benefit the model's accuracy.

## I. INTRODUCTION

Recent state of the art natural language models represent words by embedding them as vectors in a continuous vector space [1] [2]. These embeddings are learned as weights of a recurrent neural network language model. It has been demonstrated that the distributed representation of words as vectors in a common vector-space captures not just syntactic similarities (e.g. *cat* is close to *cats*) but semantic similarities as well. Concretely, the model's ability to perform "word-arithmetic" has been demonstrated in [3]. To exemplify this, let  $f$  be the word-embedding function. Ideally, we would then have an embedding where distances in vector-space give a measure of semantic-similarity between the words. For two word pairs such as (*king* : *man*) and (*queen* : *woman*) which have the same relation (or similarity), we would obtain  $f(\textit{king}) - f(\textit{man}) = f(\textit{queen}) - f(\textit{woman})$  or equivalently  $f(\textit{king}) - f(\textit{man}) + f(\textit{woman}) = f(\textit{queen})$ . This of course is equivalent to the analogy "*king* is to *man* as *queen* is to *woman*".

The goal of this work is to further evaluate the performance of the popular word2vec tool on these tasks. We evaluate the effect of different choices for model-architectures and hyper-parameters on the performance and test the limits of the models semantic-analysis capabilities on new challenging test data. The ability to learn analogies has only really been demonstrated on rather general and hand-picked test-sets. We test the models ability to learn domain-specific knowledge by generating questions that test relationships such as *Brand*  $\sim$  *Country* and *Movie*  $\sim$  *Director*. Previous work only considers the nearest resulting vector in the embedding space when evaluating the accuracy of this word-arithmetic. We generalise the notion of accuracy by considering not only

the nearest word-vector, but consider the set of the  $k$  nearest vectors and test whether the target word is present therein.

The content of this paper is structured as follows: Section II gives an overview of relevant previous work. In Section III we give a brief overview of the underlying model and in Section IV we outline the experimental setup and test-data used. The results of the experiments are then reported in Section V before we conclude in VI.

## II. RELATED WORK

The work by Tomas Mikolov [1], where the models were introduced, reports results on a test set they themselves created. The set contained five types of semantic questions where relations such as *Country*  $\sim$  *Currency* or *Country*  $\sim$  *Capital* were tested. In [3] the models performance on semantical similarity has been demonstrated on SemEval-2012 Task 2 [4]. The objective in this task is to decide how similar two target word-pairs such as (*glass* : *break*) and (*soldier* : *fight*) are with respect to a relation "an  $X$  typically  $Y$ ". Though the model is not directly trained on this task, it achieved better results than the previous best methods.

In [5] Levy et. al. provide a comprehensive study of several word embedding techniques on the two tasks *Word Similarity* and *Analogy*. They conclude that much of the performance gains of new embedding methods stem from design-choices and hyper-parameter optimisation rather than the underlying algorithms. Performance estimation was based only on the test-data sets introduced in [1] and [3].

In a previous work [6] they show that generalising the Skip-Gram model's negative sampling to arbitrary contexts results in the model learning functional-similarities. They demonstrate this by manually inspecting the top  $k$  most similar words to a given query word. A query e.g. for *Hogwarts* returns a list of famous schools in their modified Skip-Gram implementation while CBOW returns a list of Harry Potter characters.

## III. THE MODELS OF WORD2VEC

word2vec is an implementation of the Continuous Bag of Words (CBOW) and Skip-Gram models introduced in [1]. Both of these models are based on two-layer recurrent neural networks with architectures depicted in Figure 1 and Figure 2. The two architectures reflect the objectives of both models: CBOW is trained to predict the centre-word from its context in a sentence while Skip-Gram tries to predict the context from its centre. Both architectures take one-hot encoded words as inputs. Let  $x_i$  denote the input corresponding to word  $i$  and

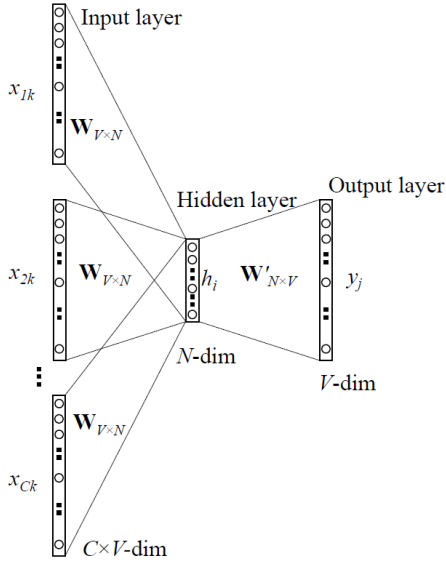


Fig. 1: Architecture of the CBOW model. The three one-hot input vectors represent the context of a centre-word in a sentence and the output vector is a vector with probabilities over possible centre-words.

let the input vocabulary be  $V$ , then the inputs are given by the  $|V|$ -dimensional vectors:

$$w_a = \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, w_{abandon} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \dots, w_{zone} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \quad (1)$$

These one-hot vectors get mapped into a  $N$ -dimensional vector via the  $|V| \times N$  weight matrix  $\mathbf{W}$ . The hidden unit then accumulates the sum of these vectors for all input words in the case of CBOW and consists of only the vector associated with the single input in the Skip-Gram model. The output of the model is generated by multiplying the hidden-units by weights  $\mathbf{W}'$  and applying softmax to the result in order to generate a valid probability-distribution over words in the vocabulary. Concretely, given a sequence of training words  $w_1, w_2, \dots, w_T$  the hidden-states are computed as

$$h(t) = \sigma(\mathbf{W}w(t) + \mathbf{H}h(t-1)), \quad (2)$$

where  $\mathbf{H}$  is the  $N \times N$  matrix of hidden-to-hidden weights and  $\sigma$  the sigmoid function:

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (3)$$

Given the values of the hidden units, the output can be computed as

$$y(t) = P_t(\mathbf{W}'h(t)) \quad (4)$$

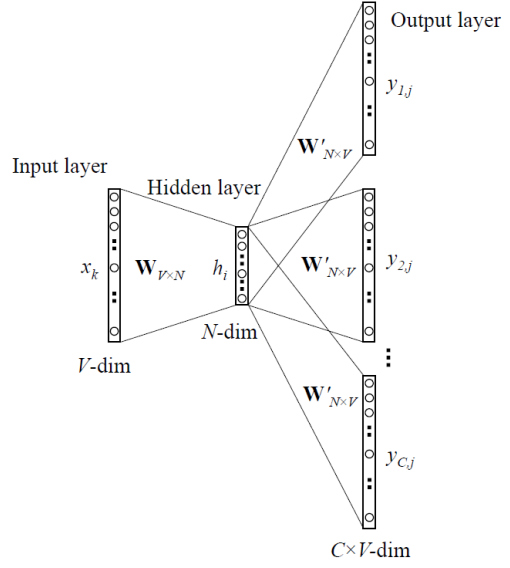


Fig. 2: Architecture of the Skip-Gram model. The output of this model are vectors with probabilities for the words occurring in the context of the input word.

with the softmax function  $P_t$  of a  $N$ -dimensional vector  $x$  given by:

$$P_t(x)_i = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} \quad (5)$$

The learnt word representation are then given by the columns of  $\mathbf{W}$ .

The models are learnt in an unsupervised fashion where either the centre-word of a series of words is predicted from its context (CBOW) or the context is predicted from the centre-word (Skip-Gram). Both models try to maximise the log-likelihood of their output. Training is done via gradient descent using back-propagation [7] to compute derivatives.

At test-time when looking for the word most similar to a given word  $w$  we solve for

$$x = \operatorname{argmin}_{v \in V} d(f(w), f(v)) \quad (6)$$

Where  $f$  is the embedding and  $d$  is a distance function. We therefore search for the word in the vocabulary which has the minimal distance in the embedding space. The model as proposed by [1] chose  $d$  to be the cosine-distance, but we investigate different choices such as  $L^p$ - or Dice-norms.

#### IV. EXPERIMENTS

In this sections we first outline our experimental setup, introduce novel test-data and then report results obtained for different hyper-parameter choices and test scenarios.

##### A. Setup

We used the Python re-implementation of word2vec contained in the Gensim library [8]. This implementation is functionally identical to the original C-implementation provided

by [1], but is more readable and has shown to perform just as well.

To compare the influence of hyper-parameters on the performance we train several models with different values for one of the parameters and keep all other parameters fixed. We identified the following set of hyper-parameters for our experiments:

- 1) Context-window size
- 2) Vector size (i.e dimension of hidden-layer)
- 3) Learning rate  $\alpha$
- 4) Training corpus
- 5) Hierarchical softmax vs. negative sampling for Skip-Gram

### B. Training Text Corpora

We trained models using multiple text corpora to check if hyper-parameters influence performance differently depending on the corpus. These corpora include:

- 1) *Wiki*: The first 900 megabytes of the English Wikipedia dump on 03/03/2006.<sup>1</sup>
- 2) *Text8*: Wiki set truncated to 100 megabytes. word2vec provides facilities for downloading and using this set.<sup>2</sup>
- 3) *Brown*: English corpus of one million words assembled from 500 sources of different genres, created in 1961 at Brown University. It is considerably smaller than our other corpora. We used this set through the NLTK python library.<sup>3</sup>

### C. New Test Data

To test the models performance on domain-specific questions we constructed several questions sets. The questions are all of the form "X is to Y as U is to ?", i.e. they take the form of analogies. We constructed the following sets of questions:

- 1) *Brand-Origin*: A set of questions testing for the relationship between well-known brands and their countries of origin. To construct this questions set we gathered names and country of origin for the top 500 brands according to brandfinance.com. We excluded all the brands from countries with names composed of multiple words (e.g. United States of America, ...) due to the model working on single words. The result is a set of 38'220 questions.
- 2) *Movie-Director*: A set of analogy-questions for the relationship between famous movies and their respective directors. The set was constructed by gathering movie-titles and last-name of the director for some of the top rated movies according to IMDb. Similar to the Brand-Origin questions we excluded movies with titles composed of more than one word with the exception of titles starting with "The ..." or similar. The result is a set of 4'557 questions.

## V. RESULTS

### A. Influence of Hyper-Parameters on Model Performance

We first wanted to see the influence of hyper-parameters when using the generalised question set described in [1]. For

<sup>1</sup>Instructions to download and build the test data on this website: <http://mattmahoney.net/dc/textdata.html>

<sup>2</sup>Ready-to-use corpus: <http://mattmahoney.net/dc/text8.zip>

<sup>3</sup>See <http://www.nltk.org/book/ch02.html>

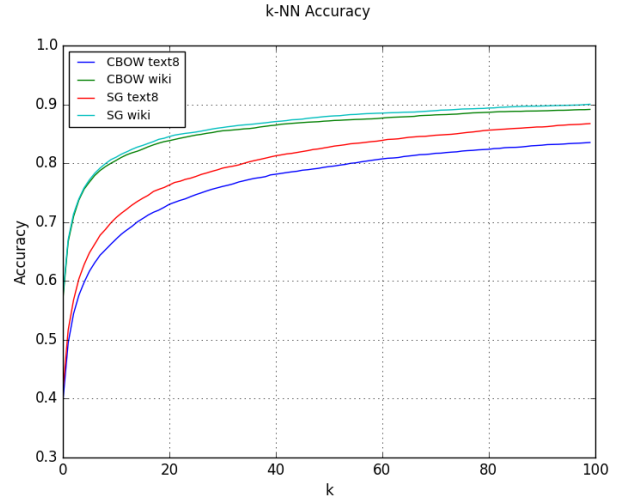


Fig. 3: k-NN accuracy according to k.

this purpose we plotted the accuracy obtained for different values of a chosen hyper-parameter. Here, the accuracy is defined as the ratio of right answers to the total number of questions.

In Figure 6 we see that the window size has only a limited influence on the accuracy. It is also difficult to pinpoint a common best value but window sizes between 5 and 8 seem to work well on all corpora. Generally the model is quite robust with respect to this hyper-parameter. Interestingly, a bigger window size seems beneficial when using CBOW on the wiki corpus, but detrimental when using Skip-Gram.

Conversely, we see in Figure 5 that a large enough vector size is essential for good performance. A larger value generally improves accuracy, although we get diminishing returns above a certain threshold. This is not surprising as the vector size defines the complexity of the model and therefore also matters more for larger corpora. We noticed that the vector size also has some effect on training time, so it is wise to restrain the vector size.

The learning rate ( $\alpha$ ) greatly affects performance as we see in Figure 4. There is a sweet-spot for  $\alpha$  around which accuracy drops. However, a learning rate around 0.031 yields good results for every model we tried.

From a broader perspective, we observe that the corpus size is by far what influences performance the most. We also notice that the Skip-Gram model tends to outperform CBOW for larger corpora.

### B. K-NN Accuracy

When a model gives the wrong answer, our accuracy definition ignores how far the model was from the right answer. We tried a different definition of accuracy where we ask the model for its  $k$  closest answers. We call it a success if this set of answers contains the right one.

We plotted this k-NN accuracy for different values of  $k$  in Figure 3. We see that models are close to the right answer

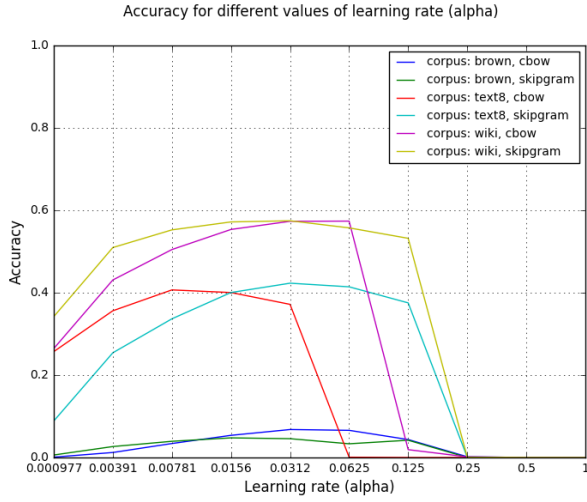


Fig. 4: Accuracy for different values of the learning rate  $\alpha$ .

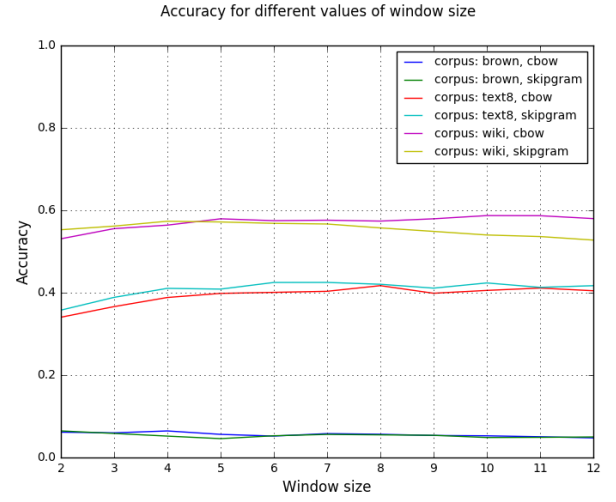


Fig. 6: Accuracy for different values of the context-window size.

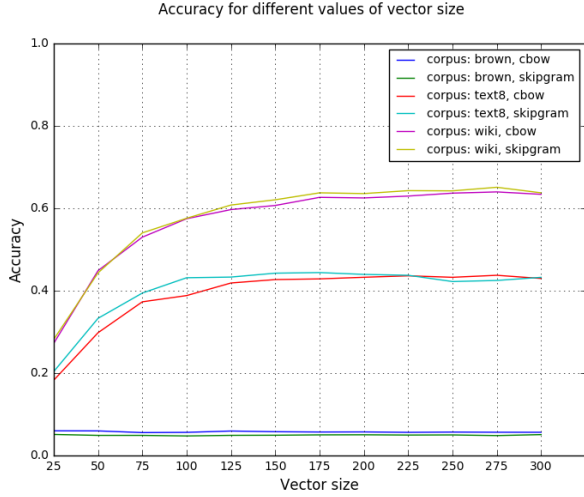


Fig. 5: Accuracy for different values of vector-size (i.e. dimension of embedding space).

most of the time. To highlight this, the accuracy of the *Skip-Gram Text8* model improves by 50% (20pt) with a  $k$  of 3. It would be considered as accurate as the *Skip-Gram Wiki* model with a  $k$  of 1.<sup>4</sup>

### C. Distance Functions

Results for new distance-functions will go here...

## VI. CONCLUSION

The conclusion goes here.

## REFERENCES

- [1] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.

- [2] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *EMNLP*, vol. 14, 2014, pp. 1532–1543.
- [3] T. Mikolov, W.-t. Yih, and G. Zweig, "Linguistic regularities in continuous space word representations," in *HLT-NAACL*, 2013, pp. 746–751.
- [4] D. A. Jurgens, P. D. Turney, S. M. Mohammad, and K. J. Holyoak, "Semeval-2012 task 2: Measuring degrees of relational similarity," in *Proceedings of the First Joint Conference on Lexical and Computational Semantics-Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation*. Association for Computational Linguistics, 2012, pp. 356–364.
- [5] O. Levy, Y. Goldberg, and I. Dagan, "Improving distributional similarity with lessons learned from word embeddings," *Transactions of the Association for Computational Linguistics*, vol. 3, pp. 211–225, 2015.
- [6] O. Levy and Y. Goldberg, "Dependency-based word embeddings," in *ACL* (2), 2014, pp. 302–308.
- [7] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Cognitive modeling*, vol. 5, no. 3, p. 1, 1988.
- [8] R. Řehůřek and P. Sojka, "Software Framework for Topic Modelling with Large Corpora," in *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. Valletta, Malta: ELRA, May 2010, pp. 45–50, <http://is.muni.cz/publication/884893/en>.

<sup>4</sup>1-NN accuracy is the same as the classic accuracy definition