

# 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 try 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 hyperparameter 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 shown 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  $f(\textit{king}) - f(\textit{man}) \approx f(\textit{queen}) - f(\textit{woman})$  or equivalently  $f(\textit{king}) - f(\textit{man}) + f(\textit{woman}) \approx f(\textit{queen})$ . 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 hyperparameters on the performance and test the limits of the models semantic-analysis capabilities on new challenging test data. The ability to learn analogies has 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 focuses on 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 present a brief overview of the underlying model and in Section IV we outline the experimental setup and testdata 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 wordpairs 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 parameter optimisation rather than the underlying algorithms. Performance estimation was based only on the testdata 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

The word2vec tool 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 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 vectors of one-hot encoded words as inputs. Let  $x_i$  denote an input-vector corresponding



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.

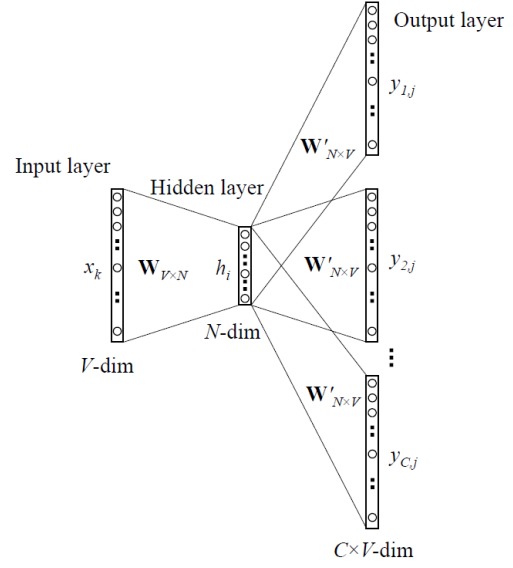


Fig. 2: Architecture of the Skip-Gram model. The output of this model are vectors with probabilities for the words occurring in the context of the input word.

to word  $i$  and let the input vocabulary be  $V$ , then the inputs to the model 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 with a  $N \times |V|$  weight matrix  $\mathbf{W}'$  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

$$h = \mathbf{W}w. \quad (2)$$

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

$$y = P_t(\mathbf{W}'h) \quad (3)$$

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

$$P_t(x)_i = \frac{e^{x_i}}{\sum_{j=1}^{|V|} e^{x_j}}, \quad (4)$$

for  $i \in \{1, \dots, |V|\}$ . The learnt word representations 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 backpropagation [7].

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

$$x = \operatorname{argmax}_{v \in V} d(f(w), f(v)) \quad (5)$$

Where  $f$  is the embedding and  $d$  is a similarity measure. 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 similarity, but we investigate different choices such as  $L^p$ -norm or Dice similarity.

#### A. Word Arithmetic for Solving Analogies

Consider the case of an analogy question "X is to Y as U is to V" where we are interested in solving for the word V. Let  $x = f(X)$ ,  $y = f(Y)$  and  $u = f(U)$  be the embedded words. The straightforward way of computing  $v = f(V)$  is then to use Equation 5 and solving for

$$v = \operatorname{argmax}_{v \in f(V)} d(v, x - y + u) \quad (6)$$

As it turns out however depending on the similarity measure  $d(\cdot)$  it can be a good idea to consider different approaches such as

$$v = \operatorname{argmax}_{v \in f(V)} d(v, x) - d(v, y) + d(v, u) \quad (7)$$

The intuition behind this approach is less one of performing word-vector arithmetic and more one of finding the word  $V$  most similar to  $X$  and  $U$  and dissimilar to  $Y$ . This approach turns out to work considerably better when using the cosine-similarity. Note that Equations 6 and 7 are equivalent for unit vectors when using the cosine similarity.

### B. Considered Similarity Measures

In our experiments we considered the following similarity measures:

$$L^2(x, y) = -\|x - y\| = -\sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (8)$$

$$L^1(x, y) = -\sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (9)$$

$$Dice(x, y) = \frac{x \cdot y}{\|x\| + \|y\|} \quad (10)$$

$$Cos(x, y) = \frac{x \cdot y}{\|x\| \cdot \|y\|} \quad (11)$$

We denote the approach in Equation 7 with cosine similarity by *cosAdd* in the results.

### C. Mean-Shift for Better Discrimination

Solving for the closest word in Equation 5 based on the cosine similarity can be interpreted as projecting on an  $N$ -dimensional unit sphere and looking for the projected vector with smallest angle to the given query vector. This projection can lead to similar vectors being very close on the sphere and therefore hard to distinguish based on angular distance. This observation is also backed up by the plots shown Figure 7 where we observe that the correct word is often very close by. Being able to better distribute out the the points on the sphere would therefore likely improve the predictions.

We therefore propose the following mean-shift procedure:

- 1) Normalise all vectors (*i.e.* project on the unit sphere)
- 2) Compute the mean of the normalised vectors
- 3) Shift by mean
- 4) Normalise again (*i.e.* re-project on unit sphere)
- 5) Compute query vector with respect to the new normalised word vectors
- 6) Compute dot-products with query vector to get cosine-similarities

The intuition behind this procedure is depicted in Figure 3 which demonstrates how points very close by get distributed out. Note that the word arithmetic to compute the query vector is performed after the normalisation (otherwise no change in the results would be expected). We denote this approach with *CosMS* in our experiments.

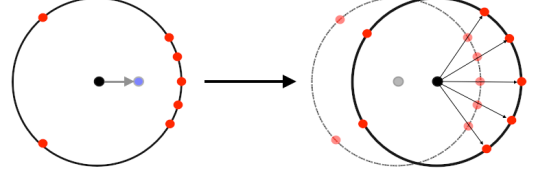


Fig. 3: Sketch of the proposed mean-shift procedure in 2D. *Left*: Unit sphere with projected vectors in red. The points on the right are relatively close together. The mean of the red points is indicated in blue. *Right*: The shifted and reprojected points. The points are now farther apart and better distinguishable.

### D. Combining Analogy Interpretations

We have already touched upon the idea that there are different ways of interpreting analogy exercises with respect to some similarity measure in Subsection III-A. The idea for Equation 7 for example is that we are looking for the vector  $v$  most similar to  $x$  and  $u$  but dissimilar to  $y$ . We could however also pose it as looking for

$$v = \operatorname{argmin}_{v \in f(V)} |d(x, y) - d(u, v)|, \quad (12)$$

that is  $x$  and  $y$  should have the same similarity as  $u$  and  $v$ . Yet another option is to consider

$$v = \operatorname{argmin}_{v \in f(V)} |d(x, u) - d(y, v)|. \quad (13)$$

Note that these objectives all capture different aspects of the analogy, *i.e.* they are not implying each other. We propose to combine these objectives in a weighted sum of these terms. Let  $\phi$  denote the term in Equation 7, *i.e.*  $\phi(v) = d(v, x) - d(v, y) + d(v, u)$ . We then want to maximise the following expression:

$$\phi(v) - \lambda \cdot (|d(x, y) - d(u, v)| + |d(x, u) - d(y, v)|) \quad (14)$$

We set the weight parameter  $\lambda = 0.1$  in our experiments and denote models using this similarity measure (including mean-shift) with *CosMS+*.

## IV. EXPERIMENTS

In this section we first outline our experimental setup, introduce novel test data and then report results obtained for different 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 easier to extend and has shown to perform just as well.

To compare the influence of 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 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 hyperparameters 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*: The first 100 megabytes of the English Wikipedia dump. Facilities for downloading and using this set are provided with word2vec.<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 Parameters on Model Performance

We first evaluate the influence of parameter choices using Google's generalised analogy dataset described in [1]. For this purpose we plotted the accuracy obtained for different values of a chosen parameter, while keeping all other parameters fixed. The baseline model for the comparison (with default

<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>

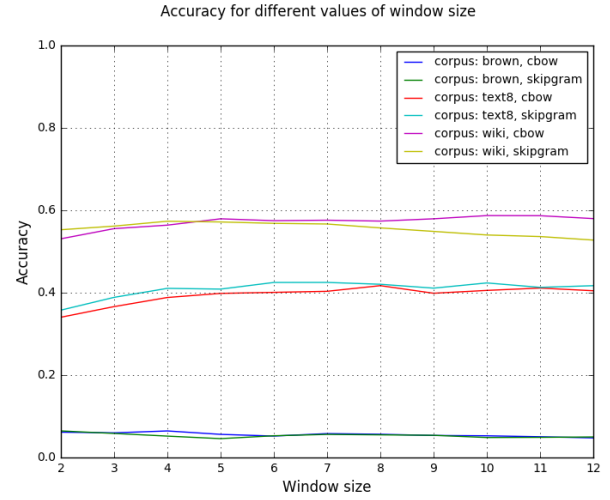


Fig. 4: Accuracy for different values of the context-window size.

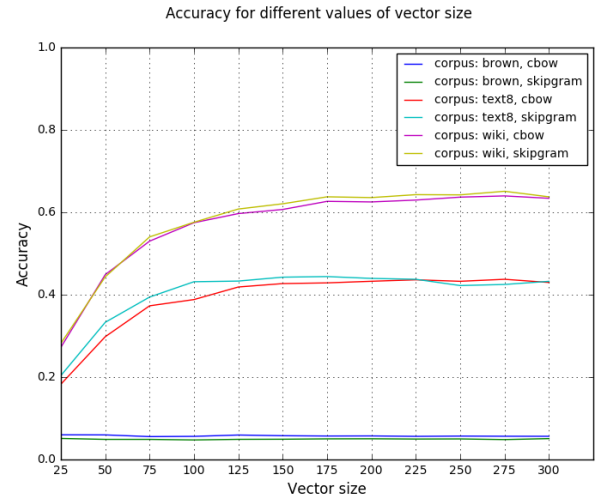


Fig. 5: Accuracy for different values of vector-size (i.e. dimension of embedding space).

parameter values) is shown in Table I. Here, the accuracy is defined as the ratio of the number of right answers over the total number of questions.

Model-Architecture	Learning Rate $\alpha$	Window-Size	Vector-Size
Skip-Gram	0.025	5	100

TABLE I: The parameter values for the baseline model used in our experiments.

In Figure 4 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 parameter. Interestingly, a bigger window size seems beneficial when using CBOW on the wiki corpus,

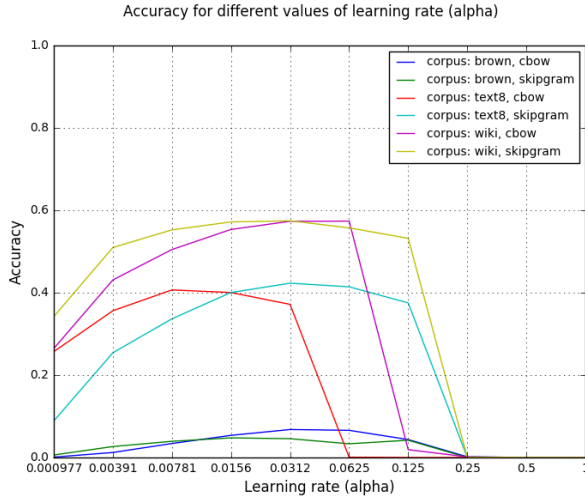


Fig. 6: Accuracy for different values of the learning rate  $\alpha$ .

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 point. This is not surprising as the vector size defines the complexity of the model and therefore also matters more for larger corpora. Obviously 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 can be seen in Figure 6. There appears to be a critical spot for  $\alpha$  around which the accuracy drops dramatically (values around 0.03 - 0.15). However, a learning rate around 0.007 - 0.03 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. This is not surprising given that neural networks are dependent on large amounts of training data. This is especially apparent when looking at the models trained on the *Brown* corpus, which all fail to get reasonable accuracies.

Considering the choice of model architecture we also notice that Skip-Gram slightly outperforms CBOW in most cases. It should be noted however that the Skip-Gram model is slightly more complex to train.

### B. K-NN Accuracy

In cases when the model gives the wrong answer, the common definition of accuracy ignores how far off the model's prediction was from the correct answer. To get an idea of how wrong the model really is, we tried a different definition of accuracy where we consider the  $k$  nearest neighbours of the prediction. We count the answer as a success if this set of neighbours contains the correct answer.

We plotted this k-NN accuracy for different values of  $k$  in Figure 7. It is clearly observable that the accuracy increases dramatically when changing from  $k = 1$  to 2 or 3 and only slowly starts saturating at around  $k > 10$ . This suggests that

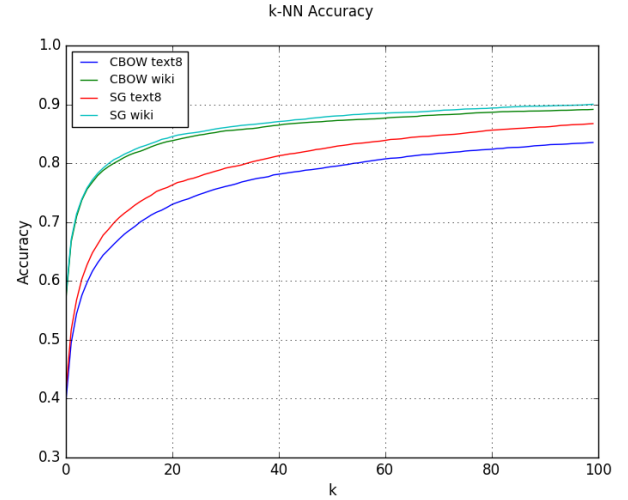


Fig. 7: k-NN accuracy according to  $k$ .

the prediction is in fact not far off most of the time and suggests that small modifications to the similarity computations (such as the mean-shift) might have a positive effect on the performance. We also notice how models trained on larger corpora saturate faster and that the difference in performance between CBOW and Skip-Gram is more apparent in models trained on small corpora.

### C. Distance Functions

To compare the different similarity measures introduced in Subsections III-B and III-C we show results on selected categories of the generalised analogy dataset [1] obtained with the baseline Skip-Gram model with parameters shown in Table I. The resulting accuracies are listed in Table II.

We observe that overall cosine similarities perform significantly better than the other similarity measures (with a few exceptions). Possible explanations for this circumstance are difficult to provide. For one we note that the Euclidean distance between weight-vectors are never explicitly used in the training of the model. In fact, the weight vectors only occur in dot-products with the columns of the hidden-to-out weight matrix  $W'$ . Another explanation is that the norm of a word vector is related to the word-frequency as found by [9] and is therefore not directly helpful in the analogy task.

### D. Mean Shift

Tables III and IV show results comparing models with and without the proposed mean-shift. We observe a slight overall improvement compared to the standard model for all the models we compared.

### E. Domain Specific Questions

Results for the domain specific questions will go here...

## VI. CONCLUSION

The conclusion goes here.

Method	Capital-Country	Currency-Country	City-State	Family	Adjective-Adverb	Opposite	Superlative	Nationality-Adjective	Past-tense	Plural	Total
<i>CosAdd</i>	47.4%	21.3%	22.2%	63.7%	15.1%	15.4%	40.1%	77.8%	32.9%	48.3%	43.0%
<i>Dice</i>	42.8%	20.5%	19.5%	65.0%	13.5%	15.4%	37.7%	<b>79.2%</b>	32.1%	44.8%	41.1%
$L^2$	44.0%	14.9%	20.6%	60.1%	15.2%	11.8%	27.7%	74.9%	30.2%	42.8%	39.9%
$L^1$	39.5%	17.5%	17.4%	56.9%	12.8%	8.8%	25.7%	72.8%	24.8%	36.2%	36.1%
<i>CosMul</i>	43.4%	20.5%	20.6%	56.5%	13.2%	11.8%	39.7%	73.5%	28.6%	39.9%	39.5%
<i>CosMS</i>	47.8%	<b>21.6%</b>	23.1%	63.4%	15.1%	15.0%	<b>42.7%</b>	77.9%	33.0%	<b>49.3%</b>	43.3%
<i>CosMS+</i>	<b>49.2%</b>	21.3%	<b>24.1%</b>	<b>65.4%</b>	<b>15.5%</b>	<b>16.0%</b>	<b>43.9%</b>	78.4%	<b>33.9%</b>	47.7%	<b>44.0%</b>

TABLE II: Comparison of different similarity-measures. Results obtained on the generalised question set of [1] with a Skip-Gram model trained on the *Text8* corpus. We can observe that cosine-similarity overall performs better compared to the other measures. Although not a similarity-measure by its own we list the results of our proposed mean-shift approach for comparison and note a slight increase in accuracy. Note that the total is computed as the micro-average (mean over all questions).

Model	Capital Country	Currency Country	City State	Family	Adjective Adverb	Opposite	Superlative	Nationality Adjective	Past-tense	Plural	Total
CBOW	63.6%	19.1%	42.8%	<b>79.7%</b>	28.4%	36.0%	62.3%	84.3%	45.3%	67.8%	58.7%
CBOW <i>CosMS</i>	63.9%	19.7%	43.9%	<b>79.7%</b>	28.4%	35.7%	<b>63.4%</b>	83.9%	45.3%	67.8%	59.1%
CBOW <i>CosMS+</i>	64.9%	20.2%	<b>44.5%</b>	79.2%	28.8%	34.2%	<b>63.4%</b>	83.7%	44.4%	67.8%	59.1%
Skip-Gram	70.4%	21.3%	34.6%	70.5%	<b>29.2%</b>	<b>38.6%</b>	58.7%	89.4%	49.2%	69.3%	59.2%
Skip-Gram <i>CosMS</i>	70.5%	21.3%	35.9%	70.8%	28.3%	<b>38.6%</b>	58.7%	<b>89.5%</b>	<b>49.6%</b>	70.0%	59.4%
Skip-Gram <i>CosMS+</i>	<b>71.6%</b>	<b>21.9%</b>	37.4%	73.7%	28.0%	37.7%	60.1%	88.9%	48.6%	<b>70.1%</b>	<b>60.0%</b>

TABLE III: Performance comparison of models with and without our proposed mean-shift approach. Results obtained on the generalised question set of [1] with models trained on the *Wiki* corpus.

Model	Capital Country	Currency Country	City State	Family	Adjective Adverb	Opposite	Superlative	Nationality Adjective	Past-tense	Plural	Total
CBOW	33.5%	16.4%	14.7%	<b>79.7%</b>	13.9%	<b>20.3%</b>	42.3%	65.4%	30.3%	44.7%	39.5%
CBOW <i>CosMS</i>	35.3%	17.5%	15.8%	78.4%	14.9%	19.9%	42.3%	64.9%	31.0%	45.3%	39.8%
CBOW <i>CosMS+</i>	35.5%	16.4%	16.3%	79.4%	14.4%	19.0%	44.1%	64.9%	30.3%	46.1%	40.0%
Skip-Gram	47.4%	21.3%	22.2%	63.7%	15.1%	15.4%	40.1%	77.8%	32.9%	48.3%	43.0%
Skip-Gram+RC	47.8%	<b>21.6%</b>	23.1%	63.4%	15.1%	15.0%	42.7%	77.9%	33.0%	<b>49.3%</b>	43.3%
Skip-Gram <i>CosMS+</i>	<b>49.2%</b>	21.3%	<b>24.1%</b>	65.4%	<b>15.5%</b>	16.0%	<b>43.9%</b>	<b>78.4%</b>	<b>33.9%</b>	47.7%	<b>44.0%</b>

TABLE IV: Performance comparison of models with and without our proposed mean-shift approach. Results obtained on the generalised question set of [1] with models trained on the *Text8* corpus.

## 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>.
- [9] A. M. J. Schakel and B. J. Wilson, "Measuring word significance using distributed representations of words," 08 2015. [Online]. Available: <http://arxiv.org/abs/1508.02297>