

# Distributed Representations of Words and Phrases and their Compositionality

Authors of the paper: Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, Jeffrey Dean

Student poster created by: Felix Buchert, Daniela Geisinger, Johann Jahner, Alexander Ladwein, Beike Lu, Alexander Schwer

## 1. Motivation

- Basis: Skip-gram neural network model using vector representations for words
- Goal: Improving the quality of the vectors and the training speed
- Approach: Extending the Skip-gram model with subsampling of frequent words, negative sampling, as well as the learning of idiomatic phrases

## 2. The Skip-gram Model

- Goal: Find **word representations** that are useful for predicting the surrounding words, the context, given an input word.
- Approach: Training objective is to maximize the average log of the probability
 
$$p(w_o | w_I) = \frac{\exp(v'_{w_o} \cdot v_{w_I})}{\sum_{w=1}^W \exp(v'_{w_o} \cdot v_{w_I})}$$
- Output layer: Softmax layer that computes the probabilities for each word in the vocabulary to be inside the context
- Limitations: **cost of computation**  $\propto W$ , the number of words in the vocabulary



Figure 1: Example of the context for a given input word and a window size of two.

## 2.1 Hierarchical Softmax (HS)

- Goal: Increase the computational efficiency
- Approach: Replaces the flat Softmax output layer with a **binary tree** layer. Each leaf represents a word in the vocabulary and the product of all node probabilities on the path from the root to the leaf represents the probability  $p(w_o | w_I)$  which is defined as follows

$$p(w | w_I) = \prod_{j=1}^{L(w)-1} \sigma([n(w, j+1) = \text{ch}(n(w, j))]) * v'_{n(w, j)} \cdot v_{w_I}$$

- Tree Structure: **Binary Huffman tree** is used for fast training; it uses short codes for frequent words
- Increase in efficiency with HS: **cost of computation**  $\propto \log(W)$

## 2.2 Negative Sampling (NEG)

- Goal: Computationally inexpensive objective for learning word vector representations

- Alternative to Hierarchical Softmax: The Negative Sampling objective

$$\log(v'_{w_o} \cdot v_{w_I}) + \sum_{i=1}^k E_{w_i \sim P_n(w)} [\log \sigma(-v'_{w_i} \cdot v_{w_I})]$$

replaces every  $\log P(w_o | w_I)$  term in the Skip-gram model

- Task: Distinguish the target word  $w_o$  from draws from the noise distribution  $P_n(w)$  using logistic regression
- Negative Samples: For every data sample  $k$  negative samples are drawn from the noise distribution
- Noise distribution  $P_n(w)$ : In general the choice of the distribution is free, but the unigram distribution raised to the 3/4rd power has empirically found to significantly outperform other choices

## 2.3 Subsampling of Frequent Words

- Motivation: Vector representations of frequent words do not change significantly after training on large datasets, as they usually provide less information than rare words
- Goal: Approach that efficiently discards frequent words while preserving the accuracy of the learned vectors
- Approach: Each word  $w_i$  is discarded with probability

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$

with frequency  $f(w_i)$  of word  $w_i$  and threshold  $t$  (around  $10^{-5}$ )

## 2.4 Results

- Evaluation: The described methods were evaluated on the basis of the analogical reasoning task considering syntactic and semantic analogies
- Training parameters: Vocabulary size of 692K, Dimension size of 300

Method	Time [min]	Syntactic [%]	Semantic [%]	Total accuracy [%]
NEG-5	38	63	54	59
NEG-15	97	63	58	61
HS-Huffman	41	53	40	47
The following results use $10^{-5}$ subsampling				
NEG-5	14	61	58	60
NEG-15	36	61	61	61
HS-Huffman	21	52	59	55

Table 1: Computation time and accuracy of various Skip-gram models on the analogical reasoning task. NEG- $k$ , with  $k$  denoting the number of negative samples for each positive sample

## 3. Learning Phrases

- Motivation: Meaning of phrases is not given by the mere composition of its words
- Approach: Find a single representation for every phrase by detecting words that appear frequently together and rarely in other context
- Result: Training with different Skip Gram Models (hyperparameter) shows that subsampling improves performance also significantly for phrases
- Size of training data: By evaluating the results of analogy task, it was found that the amount of data is crucial (33B  $\rightarrow$  72%; 6B  $\rightarrow$  66%)

## 4. Additive Compositionality

- Linearity: due to the linear structure of the representation, also vector addition gives reasonable output (distr. are multiplied  $\rightarrow$  AND-function)
- e.g.:  $\text{vec}(\text{"Berlin"}) - \text{vec}(\text{"Germany"}) + \text{vec}(\text{"France"}) = \text{vec}(\text{"Paris"})$

## 5. Conclusion

- the model architecture and the subsampling of frequent words both result in a computationally more efficient training and improve the quality of the word representations, in particular for rare and uncommon words, respectively
- the performance is sensitive towards both the training algorithm and the hyperparameters, which is why they should be selected depending on the specific task