

# REPRESENTATIONS IN NATURAL LANGUAGE PROCESSING

BEN RUPPIK

ABSTRACT. Extended notes for the lecture “Representations in Natural Language Processing”, HHU Düsseldorf, Winter Term 2022/2023.

- **Lecture:** Wednesdays, 12:30 - 14:00, Building 25.22 / 2522.HS 5G
- **Exercise:** Thursdays, 08:30 - 10:00, Building 25.22 / 2522.U1.52
- Lecture course repository: <https://github.com/ben300694/word-embeddings>
- Live stream on WebEx will probably be available for registered participants.
- **Requirements for credit points:**
  - Successful participation in the exercises.
  - Submission of a  $\text{\LaTeX}$ -file and compiled PDF of a final report
    - \* 2000 words
    - \* [ACL Overleaf template](#)
    - \* has to include references
- **Prerequisites:** Some background in deep learning, linear algebra.

## 1. INTRODUCTION

In a *word embedding* we encode language by associating to a word  $w$  a point or vector  $v_w \in \mathbb{R}^n$  in some ambient space, such that words which have similar meaning are close together in the space. When words are modelled as vectors or points in an ambient high-dimensional space, this is called an *embedding*.

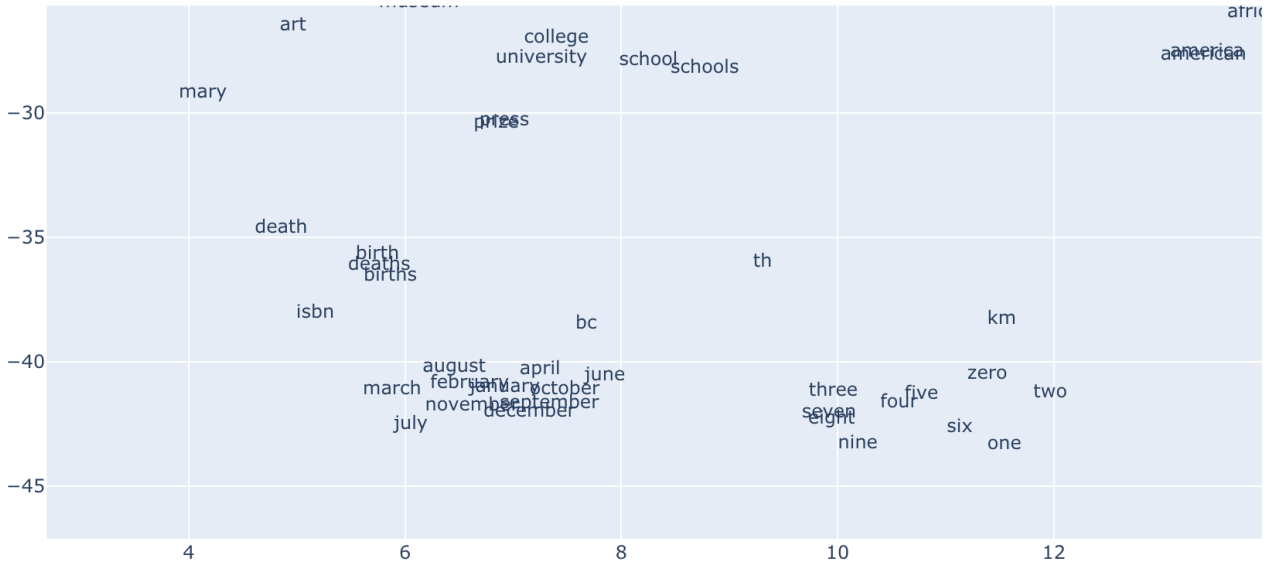


FIGURE 1. Some of the word vectors from a 100 dimensional fastText embedding trained on a Wikipedia corpus; projected to 2 dimensions using t-SNE.

Words which are used in similar contexts usually have similar meanings.<sup>1</sup> The *distributional hypothesis* states that words which frequently appear in similar contexts have similar meanings. There is *first-order co-occurrence*, also called *syntagmatic association*, if the words are typically nearby each other; for example: ‘wrote’ and ‘book’, ‘wrote’ and ‘poem’. We also observe *second order co-occurrence*, called *paradigmatic association*, if the words have similar neighbors; for example: ‘wrote’ and ‘said’, ‘wrote’ and ‘remarked’.

This makes it possible to acquire meaningful representations from unlabeled data – and there is lots of unlabeled data, e.g., in the form of text from Wikipedia, books, or news headlines. Thus, we can make both *distributed* and *distributional representations*: Distributional means that it is based on counts of words in context; distributed means that the representation is vector based.

The task of finding good word embeddings becomes hard because there are many subtleties in the usage of words in natural language. Words (or *lemmas*) like ‘mole’ can have multiple meanings, which are its *word senses*; a word which has multiple meanings is *polysemous*. The embedding should also be able to capture concepts such as *synonyms* (words have the same meaning in many or all contexts) and *antonyms* (the word senses are opposites with respect to one feature of the meaning), *similarity*, *relatedness*, *connotation* or *sentiment*. There is usually a *super-subordinate* relation (also called *hyperonymy*, *hyponymy* or “is-a” relation). On the other hand, *stop words* like ‘a’, ‘the’, ‘is’, ‘are’, ‘and’ are so commonly used that they carry very little useful information.

There is a big difference between *static* and *contextualized* embeddings: In a *static embedding* (Section 2) there is one fixed embedding for each word in the vocabulary. In a *contextual embedding* (Section 3) the vector we associate to a word is different if it is used in a different context.

<sup>1</sup>This is not always true, as the words ‘good’ and ‘bad’ demonstrate, more on this below.

We will also talk about how one can try to evaluate the quality of a word embedding. In *extrinsic evaluation* in a downstream task, we use a word embedding for a specific task and measure whether it increases performance; *intrinsic evaluation* is a harder challenge.

### 1.1. General References and Links.

- For a general introduction see the book [JM09], but take a look at the [draft of the 3rd edition](#), in particular Chapter 6 on “Vector Semantics and Embeddings”. There are corresponding [slides](#) as well.
- [Eis19, Chp. 14] on “Distributional and Distributed Semantics”; draft available [here](#).
- Michael Heck’s slides on the Space of Word Embeddings.
- Lena Voita’s NLP Course For You, in particular the lesson on Word Embeddings and the lesson on Sequence to Sequence (seq2seq) and Attention.
- Rasa Algorithm Whiteboard YouTube playlist.
- AI Coffee Break with Letitia YouTube channel.

1.1.1. *Acknowledgments.* The content of this lecture is based on the seminar “Word embedding spaces” run at HHU Düsseldorf in the Summer Term 2022. I am grateful for the suggestions and feedback collected during the sessions.

Thank you to Milica Gašić, Chris Geishauser, Michael Heck for additional input.

Sections of this write-up are based on the final reports written by the students after their talks:

- Nick Rucks: fastText and GloVe

## 2. BACKGROUND AND STATIC WORD EMBEDDINGS

2.1. **Frequency based methods.** Co-occurrence counts; term-document matrix: each document is represented by a vector of word counts; **term frequency** – **inverse document frequency** (tf-idf): sparse vectors, words are represented as a simple function of the counts of neighbors; **pointwise positive mutual information** (PPMI):  $\text{PMI}(w_1, w_2) = \log \frac{p(w_1, w_2)}{p(w_1) \cdot p(w_2)}$ ; dense versus sparse embeddings; Latent semantic analysis: singular value decomposition applied to term-document matrix (weighted by log frequency and normalized by entropy); cosine similarity.

### References:

- General references in Section 1.1;
- Using “tricks” from word2vec to improve count-based models [LGD15];
- Skip-gram with negative-sampling (SGNS, see word2vec below) implicitly factorizes the shifted pointwise mutual information (PMI) matrix [LG14].

2.2. **word2vec.** static embedding; skip-gram with **negative sampling** (SGNS); Idea: train a classifier to a prediction task ‘the word  $w$  is likely to show up near the word  $c$ ’; self-supervision: can use the next word in a corpus of running text as the supervision signal

- get positive examples from taking a target word  $w$  and its neighboring context words (have to choose a context window size  $L$ , these notes [Gol15] discuss the effect of the window size).
- randomly sample other words to get negative examples (based on the empirical distribution of words, but usually modified to sample less frequent words more often).
- train a classifier (for example by logistic regression, i.e. the probability score is given as the sigmoid of the dot product) to distinguish the two cases; skip-gram makes simplifying assumption that all context words are independent.
- the learned weights of the classifier give the embedding of the word  $w$ : target embedding in matrix  $W$  and context embedding in matrix  $C$ .

Variations of word2vec training methods:

- Skip-Gram: predict context words given the central word.
- Continuous **Bag-of-Words** (CBOW): predict the central word from the sum of context words (this sum is called the “bag of words”).

- Simple, but powerful approach: Represent a sentence with the average of its word embeddings. Hybrid: word embeddings are computed distributionally, and the sentence embedding computed by composition.

Evaluating word embeddings: *Intrinsic evaluation* measures how well the word vectors capture meaning, for example by evaluating on word similarity and word analogy tasks; *extrinsic evaluation* is performed by looking at the performance of a downstream task. Linear structure in the embedding space can be observed from word analogies, such as that the vector of

$$v_{\text{king}} - v_{\text{man}} + v_{\text{woman}} \text{ is close to } v_{\text{queen}}.$$

#### References:

- The original sources are [Mik+13a], [Mik+13b].
- The survey [Ron14] has a detailed explanation of the parameter update equations in word2vec.
- [AH19] try to find a theoretical explanation for these linear relations.
- Similarities across languages can be detected in the embeddings [MLS13], which allows constructions of mappings between semantic spaces even without having parallel data [Con+17].
- Such an analysis can lead to ideas for improving the embeddings, as in [MBV17] (they subtract the mean from the word vectors and eliminate top PCA components).
- [General Language Understanding Evaluation \(GLUE\) benchmark](#), [SuperGLUE](#)

**2.3. GloVe.** GloVe [PSM14] stands for “global vectors”, which refers to the fact that the model captures global corpus statistics directly. This is similar to most matrix factorization methods but GloVe additionally provides dense vector representations and has a linear substructure in the vector space. Therefore, it combines the advantages of global factorization methods and local context window methods. On the other hand, it cannot handle OOV words. However, by efficiently training on huge corpora, it can at least decrease the occurrence probability for OOV words.

Starting point for learning the GloVe embeddings are the ratios of co-occurrence probabilities

$$(2.1) \quad \frac{P_{ik}}{P_{jk}},$$

where  $P_{ik}$  is the probability that word  $k$  appears in the context of word  $i$ . Further the probability is defined by

$$(2.2) \quad P_{ij} = \frac{X_{ij}}{\sum_k X_{ik}},$$

where  $X_{ij}$  is the frequency that word  $j$  occurs in the context of word  $i$  while  $\sum_k X_{ik}$  is the sum of co-occurrence frequencies of all words  $k$  in the context of word  $i$ . The ratios from Equation (2) are able to discriminate relevant from irrelevant words in the corpus. Hence from these ratios, the authors derive a weighted least squares objective to learn the GloVe embeddings, i.e., they optimize

$$(2.3) \quad \mathcal{J} = \sum_{i,j=1}^V f(X_{ij})(w_i^\top \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

where  $w_i$  and  $\tilde{w}_j$  are the randomly initialized vector representations for the target word  $i$  and context word  $j$  while  $b_i$  and  $\tilde{b}_j$  are bias terms that are learned together with the embeddings. The loss imposes that the inner product of the word vectors  $w_i^\top \tilde{w}_j$  should be approximately equal to their log co-occurrence counts  $\log X_{ij}$ . Hence, GloVe can be viewed as a matrix factorization method. The weighting function  $f$  penalizes irrelevant word pairs, either if they are too rare or too frequent. During optimization of  $\mathcal{J}$ , the authors sample non-zero entries  $X_{ij}$  from the co-occurrence matrix  $X$ . This matrix tabulates how often words co-occur with each other and can be obtained by a single pass through the whole corpus.

The GloVe embeddings are successfully evaluated on the word analogy task, validating that the vector space contains linear sub-structures. Moreover, it is demonstrated that GloVe embeddings achieve

better results on the word similarity task than embeddings from Word2Vec. In general, GloVe is able to consistently outperform Word2Vec in a similar training setting. Moreover, since it can leverage bigger amounts of data it further provides a better representation for rare words than Word2Vec. Finally, during the extrinsic evaluation the performance on a Named entity recognition (NER) task could be improved using GloVe embeddings as features compared to a CBOW embeddings from Word2Vec.

static embedding; captures global corpus statistics (so it considers global context), based on ratios of probabilities from word-word co-occurrence matrix  $X_{ij}$ ; these counts are used to define the loss function, which leads to minimization of a sum of squares

$$J = \sum_{i,j} \frac{X_{ij}}{X_{\max}} (w_i \tilde{w}_j^T - \log X_{ij})^2$$

$w_i$  is the word vector and  $\tilde{w}_j$  the context vector; can add in a weighting function to penalize rare events and not over-weight frequent events, and add in bias terms for each weight vector; does not handle out of vocabulary (OOV) words;

#### References:

- Original source [PSM14], GloVe project page.

**Acknowledgments:** Nick Rucks, HHU, 2022

**2.4. fastText.** FastText [Boj+16] is an extension of the Skip-gram model introduced by [Mik+13b] and therefore belongs to the family of local context window methods. In contrast to Skip-gram and Continuous Bag-of-Words (CBOW) [Mik+13b] it is able to embed out-of-vocabulary (OOV) words. Furthermore, by considering character level information it improves the representation for rare and compound words which leads to an increased performance, especially for morphologically rich languages like German or Finnish.

FastText computes a bag of constituent  $n$ -grams  $\mathcal{G}_w$  for each word  $w$  in the corpus, containing the word itself plus all  $n$ -grams for a predefined range of  $n$ . Then, the word representation  $\mathbf{u}_w$  is the sum of  $n$ -gram embeddings  $\mathbf{z}_g$  obtained by a pre-trained Skip-gram model, i.e.,

$$(2.4) \quad \mathbf{u}_w = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g.$$

Besides the different representation  $\mathbf{u}_w$ , the training of the objective is adopted from the Word2Vec framework.

The authors show in an intrinsic evaluation that the inclusion of character level information enables an improved performance on the word similarity task compared to Word2Vec. This also holds for the syntactic questions of the word analogy task. Moreover, the representations that FastText computes for OOV words always improves the performance on the word similarity task which validates that the representations are reasonable. During further evaluation, the authors showed that the most important  $n$ -grams in the representation actually correspond to valid morphemes, e.g., prefixes and suffixes. For that, the authors omitted a single  $n$ -gram from the representation  $\mathbf{u}_w$  to get a restricted representation  $\mathbf{u}_{w \setminus g}$ . Then, they ranked all  $n$ -grams by their cosine distance between  $\mathbf{u}_w$  and  $\mathbf{u}_{w \setminus g}$ , to determine the importance of each  $n$ -gram for  $\mathbf{u}_w$ . Furthermore, the better usage of subword information enables FastText to obtain better embeddings from smaller datasets than Word2Vec. Additionally, the word representation trained with subword information outperforms plain Skip-gram embeddings on a language modeling task.

static embedding; deals better with unknown words and word sparsity; add special boundary symbols  $<$  and  $>$  to the word, example for  $n = 3$  and the word ‘where’:  $<\text{where}>$  and  $<\text{wh, whe, her, ere, re}>$ ; learn skip-gram embedding for each constituent  $n$ -gram, word is represented by the sum of all the embeddings of the constituent  $n$ -grams;

#### References:

- Original source [Boj+16]

**Acknowledgments:** Nick Rucks, HHU, 2022

### 3. CONTEXTUAL WORD EMBEDDINGS

#### 3.1. Sequence to sequence tasks and Recurrent networks.

**3.1.1. Motivation.** Modeling language using static embeddings like word2vec or Glove, there are a few problems. Different meanings of the same word cannot be accurately captured, as each word has only a single vector embedding. Thus, depending on the training corpus, the embedding will either only capture a single meaning, or the embedding will be the average of different meanings. Furthermore, these embeddings allow for simple arithmetic operation that capture some meaning, but these lack the flexibility to capture the interaction between words that are so essential to language. To tackle more complex tasks like translation we need to rethink language modelling to integrate a source input. One way to do this is by using an Encoder-Decoder architecture [Voi].

**3.1.2. Encoder and Decoder.** The idea of the Encoder-Decoder architecture, is that the Encoder takes an input and produces a latent representation. The Decoder uses this representation to (re-)create information. More specifically, using translation as an example (throughout this abstract) it looks like this: The Decoder uses the given representation of a source sentence to create a translation. For translations recurrent neural networks (RNNs) have proven to be effective implementation for this.

**3.1.3. Recurrent Neural Networks.** A RNN is a neural network that is getting repeated throughout time, and at each time step  $t$  there can be a new input and output. The crux because RNNs are eligible for NLP task, is that the hidden layer at time step  $t - 1$  influences the hidden layer at time step  $t$ . Thus, the problem of lacking the ability to capture language dependencies between words is being tackled. A RNN can then be used both as the Encoder and as the Decoder. RNNs in their basic form, can run into problems during model training because of vanishing or exploding gradients [SVL14]. Furthermore, the information propagation is not ideal for longer sequences and performance suffers. To address this a commonly used modification is using a Long Short-Term Memory (LSTM) RNN.

**3.1.4. Long Short-Term Memory.** A LSTM is a different kind of RNN, where the internal structure of each reoccurring unit is different. Besides the hidden state, there is an additional cell state. This cell state gets passed through each unit, with little interference. It helps to keep the gradient in appropriate ranges as well as transfer information through many units [Kar]. With these LSTMs Peters et al. designed a model named Embeddings from Language Models (ELMo), to tackle the problem of finding embeddings for polysemic words.

**3.1.5. Embeddings from Language Models (ELMo).** ELMo produces a different embedding for each word dependent on the sentence it is used in. This is implemented as a multi-layered bidirectional LSTM. Using an LSTM for each direction ensures to capture information of the entire input sequence. The model then combines the hidden state representation of both LSTMs by concatenating them, and then adds up the concatenated states of each layer up with a task-specific weighting [Pet+18a]. The embeddings created by ELMo have the potential to improve the language understanding of many NLP models. Still, using a single vector representation for an entire source sentence is bound to fail at capturing all information. To reduce this dependency Bahdanau et al. added the concept of *attention*.

**3.1.6. Attention.** The idea behind attention is that the model learns to take in more specific information from the single source inputs. This can be implemented with a context vector. A context vector is the result of calculating attentions scores between a single hidden state of the decoder and all encoder states. Using a softmax function, these attention scores can be turned into probabilities or weights. Summing up each encoder state weighted by its softmax value, we calculate the context vector [BCB15]. The decoder at each time step uses this context and is not solely dependent on the single representation from the Encoder. If implemented with an appropriate attention score, the whole process is differentiable and can be learned task-specific by the model.



Embeddings from **Language Models (ELMo)**: The ELMo architecture is based on earlier ideas using Recurrent Neural Nets (RNNs) to model sequences, such as sentences in natural language. ELMo produces deep contextualized word representations, which are learned functions of internal states of a deep bidirectional language model with Long short-term memory (LSTM) units. ELMo embeddings are context-sensitive, i.e. each token representation is a function of the entire input sentence. Here lower layers appear to capture syntactics, while higher layers capture semantics. ELMo is unsupervised, can handle OOV, but is not truly bidirectional.

#### References:

- **Original ELMo research article:** ‘Deep Contextualized Word Representations’ [Pet+18b].
- The concept of attention, which lets a model focus on different parts of the input while processing a sequence, was introduced in [Gar+19].
- Lecture notes: Sequence to Sequence (seq2seq) and Attention

**Acknowledgments:** Lukas Rücker, HHU, 2022

#### (1) Transformers:

Encoder and decoder block; self-attention mechanism; query, key and value matrices; multi-head attention (combining multiple self-attention mechanisms); positional embeddings; autoregressive models: predict the next token in a sequence, for this we have to mask the self-attention in the decoder block so that the model can only attend to tokens in the preceding sequence. The original transformer model was trained on machine translation, which is a sequence to sequence task.

**Original research article:** ‘Attention is all you need’ [Vas+17].

**Further references:** It is probably a good idea to look at alternative expositions of the Transformer architecture concurrently while going through the original paper, here are some suggestions:

- Blog post: ‘Visualizing A Neural Machine Translation Model (Mechanics of Seq2seq Models With Attention)’
- Blog post: ‘The Illustrated Transformer’
- Blog post: ‘Transformers from scratch’
- PyTorch implementation ‘The Annotated Transformer’
- Vision transformers [Dos+21]
- Visualizations of attention maps: BertViz

#### (2) Bi-directional Encoder Representations from Transformers (BERT) and its variants:

The objective of the BERT model is a *Masked language modeling task*, in which we randomly mask words in the training text and try to predict the masked words from their context. Additionally, the model is trained on a *next sentence prediction* task. For the input embedding: WordPiece tokenization; special <cls> token for the representation of the entire input sequence.

**Original research article:** [Dev+18].

#### **Additional references:**

- (Introduction to) Transfer Learning in the NLP for you course.
- [Cla+19] analyze the role of the different attention heads in the trained BERT model.
- A Robustly Optimized BERT Pretraining Approach (RoBERTa): The authors observe that BERT was significantly under-trained, and update the model [Liu+19].
- GPT for language generation.

#### (3) Sentence embeddings:

Sentence-BERT: [RG19] and SentenceTransformers find an embedding on the sentence level so that semantically similar sentences are close in the embeddings space; adds a pooling operation to the output of BERT or RoBERTa to derive a fixed sized sentence embedding; siamese and triplet network structures.

*Contrastive learning* (pull semantically close neighbors together and push apart non-neighbors) for sentence embeddings:

TransEncoder: For pairwise comparison of sentences, [Liu+21] alternate between a bi-encoder (which produces fixed-dimensional sentence embeddings) and a cross-encoder (which has attention heads that span the complete sequence of both sentences), so that they learn from each other to train an unsupervised sentence representation model.

SimCSE [GYC21]: Uses contrastive learning to improve sentence embeddings, where the positive pairs are created by adding dropout in the encoder as the noise.

**Further references:**

- Dialog utterance embeddings via contrastive learning [Zho+22]
- Video: [Transformer Neural Networks - an Overview!](#)

### 3.2. Multimodal Embeddings. Joint text and image space: [Zha+21]

TODO: Add more references

## 4. ADDITIONAL TOPICS

### (4) Geometry of the word embedding space:

[MT17] find that in skip-gram with negative sampling the word vectors all point in roughly the same direction, i.e. lie in a narrow cone (this is sometimes called *anisotropy*, the opposite of this is *isotropy* where the vectors are directionally uniform); the context vectors point in the other direction (have negative dot product with the mean of the word vectors), in contrast to this for GloVe the word and context vectors point in the same direction;

Anisotropy can also be studied in contextualized word embeddings, see [Eth19], or for sentence embeddings [Li+20].

[NF18] observe that mappings between semantic spaces usually are only locally linear, and not globally linear;

Method for finding words which are used differently in two text corpora: train embedding on each corpus, find a (e.g. linearly via orthogonal Procrustes) mapping between the embedding spaces, then words which do not match well under this mapping appear to have different meanings. Using historical texts, one can see how the meaning changes through time [HLJ16];

Different method: again train embeddings on two corpora separately, then for each word find the closest vectors in each embedding space, if these differ a lot the embeddings are different [Gon+20].

### (5) Bias in word embeddings and ethics:

A desirable association could be “a man to a woman is as a king to a queen”, whereas an undesirable association is “a man to a woman is as a physician to a nurse”. Unfortunately, it appears that word embeddings trained on a corpus of text extracted from books and websites tend to incorporate and amplify the biases contained in our language. When such biased embeddings are used for downstream tasks, this can lead to bias amplification. To investigate gender bias systematically, [Bol+16] take a seed pair ( $v_{he}, v_{she}$ ) and find other pairs of word vectors which have a difference similar to the seed vector. They also provide methods for debiasing the embedding, but they argue that bias ultimately appears to come from the training data.

**Further references:**

- Iterative nullspace projection to debias word embeddings [Rav+20];
- Visualizations of bias mitigation techniques [Rat+21], their [code](#) for creating interactive demos is available online;
- Gender bias in the contextualized embeddings from ElMo [Zha+19];
- Bias in BERT embeddings [Kur+19];
- Debiasing contextualized embeddings [KB21].

### (6) Sentiment specific word embeddings:

*Sentiment* is a positive or negative evaluation, and we can enhance our *semantic* specific word embeddings with *sentiment* [Yu+17]. Words with similar vector representations can have an



opposite sentiment polarity (e.g., ‘good’ vs ‘bad’ or ‘happy’ vs ‘sad’). To capture this in the embedding one can add a post-processing step to adapt the pre-trained vectors to sentiment applications, where we move the word closer to a set of both semantically and sentimentally similar nearest neighbors (i.e., those with the same polarity) and further away from sentimentally dissimilar neighbors (i.e., those with an opposite polarity).

**Further references:**

- Using sentiment aware word embeddings for emotion classification [Mao+19];
- TweetEval dataset, which includes emotion recognition and sentiment analysis tasks [Bar+20], the authors released their [pre-trained RoBERTa models](#);
- German sentiment classification with using fastText and BERT embeddings [Guh+20];
- EmoWOZ [Fen+21] is a task-oriented dialogue dataset annotated with emotions.

(7) Multilingual word embeddings:

We want to use a shared embedding space for words across two or more languages, so that words from different languages with different meanings are close in the cross-lingual embedding space; the hope is that this can improve performance because the model can acquire and transfer knowledge from other languages; also look at the references from the ‘Geometry of the word embedding space’-topic, in particular [Con+17] where they aligned embeddings without having parallel data; MUSE is an implementation of both supervised and unsupervised multilingual word embeddings with fastText; [CC18] train unsupervised MWEs; [Fen+20] build language-agnostic BERT Sentence Embedding.

**Further references:**

- Video: [Language Agnostic BERT](#)
- Please also take a look at the document in the appendix with further references on multilingual transformer architectures.

#### 4.1. Word Embeddings in Dialog Systems. TODO

### 5. TOPOLOGICAL DATA ANALYSIS AND WORD EMBEDDINGS

These last two topics are more advanced, in that they require more background knowledge, and it might be more challenging to prepare an accessible presentation.

(8) (\*) Hyperbolic embeddings:

Instead of embedding data into an Euclidean vector space, we can associate it to points in a more general ambient manifold. For example, we can use so-called *negatively curved spaces*, which leads to *hyperbolic embeddings*. The geometry of these spaces allows a more efficient representation of tree-like structures and hierarchies, as was observed in [NK17]; see also this [blog post](#) discussing the paper, and [this blog post](#) with implementation details. In our case, where we associate words to points in a hyperbolic space, these are called *hyperbolic word embeddings* (which can be constructed e.g. with Poincaré GloVe [TBG18]). Further references include an earlier paper on hyperbolic text embeddings [Dhi+18], these in turn were inspired by hyperbolic image embeddings [Khr+19]. There is also a proposed hyperbolic version of fastText [Zhu+20]. One can also ‘probe’ the contextualized embeddings of a BERT model by mapping into a hyperbolic space [Che+21], and analyzing the structure of the projected data.

(9) (\*) Topological Polysemy and singularities in word embeddings:

The paper [JGZ20] studies singular points in a fixed word embedding trained via fastText. They find that the dimension of the zeroth local homology group of a punctured neighborhood of a word correlates with the number of different meanings of a polysemous word.

There have been two main developments in Topological Data Analysis (TDA), which might lead to further interesting application of TDA to word embeddings in the future:

- *Persistent homology* is a description of the topological features of a space at various scales. Usually the data comes with a *filtration*, and we compute the homology of the pieces which assemble into a *persistence module*. These can be encoded into a *bar-code*, which then is used in

further machine learning pipelines. For example [RCB21] develop a Transformer model which takes persistence diagrams as input.

- The MAPPER algorithm is a projection technique which also uses clustering to capture more of the global features of the dataset. An interesting projection choice is  $t$ -Stochastic Neighborhood Embedding ( $t$ -SNE, [MH08]).

Here are some general references for background reading on Topological Data Analysis:

- [MR19]: Survey paper with an introduction to “Topological Data Analysis for Physicists”.
- [DW21]: Book (draft) on Topological Data Analysis.
- Vidit Nanda’s [Computational Algebraic Topology Notes](#).

## REFERENCES

- [AH19] C. ALLEN and T. M. HOSPEDALES. Analogies Explained: Towards Understanding Word Embeddings. In: *CoRR* abs/1901.09813 (2019). arXiv: [1901.09813](#) (↑ 4).
- [BCB15] D. BAHDANAU, K. CHO, and Y. BENGIO. Neural Machine Translation by Jointly Learning to Align and Translate. In: *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*. Ed. by Y. BENGIO and Y. LECUN. 2015 (↑ 6).
- [Bar+20] F. BARBIERI, J. CAMACHO-COLLADOS, L. ESPINOSA ANKE, and L. NEVES. TweetEval: Unified Benchmark and Comparative Evaluation for Tweet Classification. In: *Findings of the Association for Computational Linguistics: EMNLP 2020*. Online: Association for Computational Linguistics, Nov. 2020, pp. 1644–1650. DOI: [10.18653/v1/2020.findings-emnlp.148](#) (↑ 9).
- [Boj+16] P. BOJANOWSKI, E. GRAVE, A. JOULIN, and T. MIKOLOV. Enriching Word Vectors with Subword Information. In: *CoRR* abs/1607.04606 (2016). arXiv: [1607.04606](#) (↑ 5, 6).
- [Bol+16] T. BOLUKBASI, K. CHANG, J. Y. ZOU, V. SALIGRAMA, and A. KALAI. Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. In: *CoRR* abs/1607.06520 (2016). arXiv: [1607.06520](#) (↑ 8).
- [Che+21] B. CHEN, Y. FU, G. XU, P. XIE, C. TAN, M. CHEN, and L. JING. Probing BERT in Hyperbolic Spaces. In: *CoRR* abs/2104.03869 (2021). arXiv: [2104.03869](#) (↑ 9).
- [CC18] X. CHEN and C. CARDIE. Unsupervised Multilingual Word Embeddings. In: *CoRR* abs/1808.08933 (2018). arXiv: [1808.08933](#) (↑ 9).
- [Cla+19] K. CLARK, U. KHANDLWAL, O. LEVY, and C. D. MANNING. What Does BERT Look at? An Analysis of BERT’s Attention. In: *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Florence, Italy: Association for Computational Linguistics, Aug. 2019, pp. 276–286. DOI: [10.18653/v1/W19-4828](#) (↑ 7).
- [Con+17] A. CONNEAU, G. LAMPLE, M. RANZATO, L. DENOYER, and H. JÉGOU. Word Translation Without Parallel Data. In: *CoRR* abs/1710.04087 (2017). arXiv: [1710.04087](#) (↑ 4, 9).
- [Dev+18] J. DEVLIN, M. CHANG, K. LEE, and K. TOUTANOVA. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In: *CoRR* abs/1810.04805 (2018). arXiv: [1810.04805](#) (↑ 7).
- [DW21] T. K. DEY and Y. WANG. *Computational Topology for Data Analysis*. 2021 (↑ 10).
- [Dhi+18] B. DHINGRA, C. J. SHALLUE, M. NOROUZI, A. M. DAI, and G. E. DAHL. Embedding Text in Hyperbolic Spaces. In: *CoRR* abs/1806.04313 (2018). arXiv: [1806.04313](#) (↑ 9).
- [Dos+21] A. DOSOVITSKIY et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In: *International Conference on Learning Representations*. 2021 (↑ 7).
- [Eis19] J. EISENSTEIN. *Introduction to Natural Language Processing*. Adaptive Computation and Machine Learning series. MIT Press, 2019. ISBN: 9780262042840 (↑ 3).
- [Eth19] K. ETHAYARAJH. How Contextual are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 55–65. DOI: [10.18653/v1/D19-1006](#) (↑ 8).
- [Fen+20] F. FENG, Y. YANG, D. CER, N. ARIVAZHAGAN, and W. WANG. Language-agnostic BERT Sentence Embedding. In: *CoRR* abs/2007.01852 (2020). arXiv: [2007.01852](#) (↑ 9).
- [Fen+21] S. FENG, N. LUBIS, C. GEISHAUSER, H. LIN, M. HECK, C. VAN NIEKERK, and M. GASIC. EmoWOZ: A Large-Scale Corpus and Labelling Scheme for Emotion in Task-Oriented Dialogue Systems. In: *CoRR* abs/2109.04919 (2021). arXiv: [2109.04919](#) (↑ 9).
- [GYC21] T. GAO, X. YAO, and D. CHEN. SimCSE: Simple Contrastive Learning of Sentence Embeddings. In: *CoRR* abs/2104.08821 (2021). arXiv: [2104.08821](#) (↑ 8).
- [Gar+19] S. GARG, S. PEITZ, U. NALLASAMY, and M. PAULIK. Jointly Learning to Align and Translate with Transformer Models. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 4453–4462. DOI: [10.18653/v1/D19-1453](#) (↑ 7).
- [Gol15] Y. GOLDBERG. A Primer on Neural Network Models for Natural Language Processing. In: *CoRR* abs/1510.00726 (2015). arXiv: [1510.00726](#) (↑ 3).

- [Gon+20] H. GONEN, G. JAWAHAR, D. SEDDAH, and Y. GOLDBERG. Simple, Interpretable and Stable Method for Detecting Words with Usage Change across Corpora. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, July 2020, pp. 538–555. DOI: [10.18653/v1/2020.acl-main.51](https://doi.org/10.18653/v1/2020.acl-main.51) (↑ 8).
- [Guh+20] O. GUHR, A.-K. SCHUMANN, F. BAHRMANN, and H. J. BÖHME. Training a Broad-Coverage German Sentiment Classification Model for Dialog Systems. English. In: *Proceedings of the 12th Language Resources and Evaluation Conference*. Marseille, France: European Language Resources Association, May 2020, pp. 1627–1632. ISBN: 979-10-95546-34-4 (↑ 9).
- [HLJ16] W. L. HAMILTON, J. LESKOVEC, and D. JURAFSKY. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics, Aug. 2016, pp. 1489–1501. DOI: [10.18653/v1/P16-1141](https://doi.org/10.18653/v1/P16-1141) (↑ 8).
- [JGZ20] A. JAKUBOWSKI, M. GASIC, and M. ZIBROWIUS. Topology of word embeddings: Singularities reflect polysemy. In: *Proceedings of the Ninth Joint Conference on Lexical and Computational Semantics*. 2020, pp. 103–113. arXiv: [2011.09413](https://arxiv.org/abs/2011.09413) [cs.CL] (↑ 9).
- [JM09] D. JURAFSKY and J. H. MARTIN. *Speech and Language Processing*. MIT Press, 2009. ISBN: 978-0-13-187321-6 (↑ 3).
- [KB21] M. KANEKO and D. BOLLEGALA. Debiasing Pre-trained Contextualised Embeddings. In: *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Online: Association for Computational Linguistics, Apr. 2021, pp. 1256–1266. DOI: [10.18653/v1/2021.eacl-main.107](https://doi.org/10.18653/v1/2021.eacl-main.107) (↑ 8).
- [Kar] KARPATY. *The Unreasonable Effectiveness of Recurrent Neural Networks*. [http://karpathy.github.io/2015/05/21/rnn-effectiveness/](https://github.com/karpathy/2015-05-21-rnn-effectiveness/). Accessed: 2022-05-05 (↑ 6).
- [Khr+19] V. KHRULKOV, L. MIRVAKHABOVA, E. USTINOVA, I. V. OSELEDETS, and V. S. LEMPITSKY. Hyperbolic Image Embeddings. In: *CoRR* abs/1904.02239 (2019). arXiv: [1904.02239](https://arxiv.org/abs/1904.02239) (↑ 9).
- [Kur+19] K. KURITA, N. VYAS, A. PAREEK, A. W. BLACK, and Y. TSIVETKOV. Measuring Bias in Contextualized Word Representations. In: *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*. Florence, Italy: Association for Computational Linguistics, Aug. 2019, pp. 166–172. DOI: [10.18653/v1/W19-3823](https://doi.org/10.18653/v1/W19-3823) (↑ 8).
- [LG14] O. LEVY and Y. GOLDBERG. Neural Word Embedding as Implicit Matrix Factorization. In: *Advances in Neural Information Processing Systems*. Ed. by Z. GHAHRAMANI, M. WELLING, C. CORTES, N. LAWRENCE, and K. Q. WEINBERGER. Vol. 27. Curran Associates, Inc., 2014 (↑ 3).
- [LGD15] O. LEVY, Y. GOLDBERG, and I. DAGAN. Improving Distributional Similarity with Lessons Learned from Word Embeddings. In: *Transactions of the Association for Computational Linguistics* 3 (2015), pp. 211–225. DOI: [10.1162/tac1\\_a\\_00134](https://doi.org/10.1162/tac1_a_00134) (↑ 3).
- [Li+20] B. LI, H. ZHOU, J. HE, M. WANG, Y. YANG, and L. LI. On the Sentence Embeddings from Pre-trained Language Models. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Online: Association for Computational Linguistics, Nov. 2020, pp. 9119–9130. DOI: [10.18653/v1/2020.emnlp-main.733](https://doi.org/10.18653/v1/2020.emnlp-main.733) (↑ 8).
- [Liu+21] F. LIU, Y. JIAO, J. MASSIAH, E. YILMAZ, and S. HAVRYLOV. Trans-Encoder: Unsupervised sentence-pair modelling through self- and mutual-distillations. In: *CoRR* abs/2109.13059 (2021). arXiv: [2109.13059](https://arxiv.org/abs/2109.13059) (↑ 8).
- [Liu+19] Y. LIU, M. OTT, N. GOYAL, J. DU, M. JOSHI, D. CHEN, O. LEVY, M. LEWIS, L. ZETTLEMOYER, and V. STOYANOV. RoBERTa: A Robustly Optimized BERT Pretraining Approach. In: *CoRR* abs/1907.11692 (2019). arXiv: [1907.11692](https://arxiv.org/abs/1907.11692) (↑ 7).
- [MH08] L. VAN DER MAATEN and G. HINTON. Visualizing Data using t-SNE. In: *Journal of Machine Learning Research* 9.86 (2008), pp. 2579–2605 (↑ 10).
- [Mao+19] X. MAO, S. CHANG, J. SHI, F. LI, and R. SHI. Sentiment-Aware Word Embedding for Emotion Classification. In: *Applied Sciences* 9.7 (2019), p. 1334. ISSN: 2076-3417. DOI: [10.3390/app9071334](https://doi.org/10.3390/app9071334) (↑ 9).
- [Mik+13a] T. MIKOLOV, K. CHEN, G. CORRADO, and J. DEAN. Efficient Estimation of Word Representations in Vector Space. In: *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*. Ed. by Y. BENGIO and Y. LECUN. 2013 (↑ 4).

- [MLS13] T. MIKOLOV, Q. V. LE, and I. SUTSKEVER. Exploiting Similarities among Languages for Machine Translation. In: *CoRR* abs/1309.4168 (2013). arXiv: [1309.4168](#) (↑ 4).
- [Mik+13b] T. MIKOLOV, I. SUTSKEVER, K. CHEN, G. CORRADO, and J. DEAN. Distributed Representations of Words and Phrases and their Compositionality. In: *CoRR* abs/1310.4546 (2013). arXiv: [1310.4546](#) (↑ 4, 5).
- [MT17] D. MIMNO and L. THOMPSON. The strange geometry of skip-gram with negative sampling. In: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Copenhagen, Denmark: Association for Computational Linguistics, Sept. 2017, pp. 2873–2878. DOI: [10.18653/v1/D17-1308](#) (↑ 8).
- [MBV17] J. MU, S. BHAT, and P. VISWANATH. All-but-the-Top: Simple and Effective Postprocessing for Word Representations. In: *CoRR* abs/1702.01417 (2017). arXiv: [1702.01417](#) (↑ 4).
- [MR19] J. MURUGAN and D. ROBERTSON. *An Introduction to Topological Data Analysis for Physicists: From LGM to FRBs*. 2019. arXiv: [1904.11044 \[astro-ph.IM\]](#) (↑ 10).
- [NF18] N. NAKASHOLE and R. FLAUGER. Characterizing Departures from Linearity in Word Translation. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Melbourne, Australia: Association for Computational Linguistics, July 2018, pp. 221–227. DOI: [10.18653/v1/P18-2036](#) (↑ 8).
- [NK17] M. NICKEL and D. KIELA. Poincaré Embeddings for Learning Hierarchical Representations. In: *CoRR* abs/1705.08039 (2017). arXiv: [1705.08039](#) (↑ 9).
- [PSM14] J. PENNINGTON, R. SOCHER, and C. MANNING. GloVe: Global Vectors for Word Representation. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Doha, Qatar: Association for Computational Linguistics, Oct. 2014, pp. 1532–1543. DOI: [10.3115/v1/D14-1162](#) (↑ 4, 5).
- [Pet+18a] M. E. PETERS, M. NEUMANN, M. IYYER, M. GARDNER, C. CLARK, K. LEE, and L. ZETTLEMOYER. Deep Contextualized Word Representations. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*. Ed. by M. A. WALKER, H. JI, and A. STENT. Association for Computational Linguistics, 2018, pp. 2227–2237. DOI: [10.18653/v1/n18-1202](#) (↑ 6).
- [Pet+18b] M. E. PETERS, M. NEUMANN, M. IYYER, M. GARDNER, C. CLARK, K. LEE, and L. ZETTLEMOYER. Deep contextualized word representations. In: *CoRR* abs/1802.05365 (2018). arXiv: [1802.05365](#) (↑ 7).
- [Rat+21] A. RATHORE, S. DEV, J. M. PHILLIPS, V. SRIKUMAR, Y. ZHENG, C. M. YEH, J. WANG, W. ZHANG, and B. WANG. VERB: Visualizing and Interpreting Bias Mitigation Techniques for Word Representations. In: *CoRR* abs/2104.02797 (2021). arXiv: [2104.02797](#) (↑ 8).
- [Rav+20] S. RAVFOGEL, Y. ELAZAR, H. GONEN, M. TWITON, and Y. GOLDBERG. Null It Out: Guarding Protected Attributes by Iterative Nullspace Projection. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, July 2020, pp. 7237–7256. DOI: [10.18653/v1/2020.acl-main.647](#) (↑ 8).
- [RG19] N. REIMERS and I. GUREVYCH. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In: *CoRR* abs/1908.10084 (2019). arXiv: [1908.10084](#) (↑ 7).
- [RCB21] R. REINAUER, M. CAORSI, and N. BERKOUK. *Persformer: A Transformer Architecture for Topological Machine Learning*. 2021. arXiv: [2112.15210 \[cs.LG\]](#) (↑ 10).
- [Ron14] X. RONG. word2vec Parameter Learning Explained. In: *CoRR* abs/1411.2738 (2014). arXiv: [1411.2738](#) (↑ 4).
- [SVL14] I. SUTSKEVER, O. VINYALS, and Q. V. LE. Sequence to sequence learning with neural networks. In: *Advances in neural information processing systems* 27 (2014) (↑ 6).
- [TBG18] A. TIFREA, G. BÉCIGNEUL, and O. GANEA. Poincaré GloVe: Hyperbolic Word Embeddings. In: *CoRR* abs/1810.06546 (2018). arXiv: [1810.06546](#) (↑ 9).
- [Vas+17] A. VASWANI, N. SHAZEER, N. PARMAR, J. USZKOREIT, L. JONES, A. N. GOMEZ, L. KAISER, and I. POLOSUKHIN. Attention Is All You Need. In: *CoRR* abs/1706.03762 (2017). arXiv: [1706.03762](#) (↑ 7).
- [Voi] VOITA. *Sequence to Sequence (seq2seq) and Attention*. [https://lena-voita.github.io/nlp\\_course/seq2seq\\_and\\_attention.html](https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html). Accessed: 2022-05-05 (↑ 6).



- [Yu+17] L.-C. YU, J. WANG, K. R. LAI, and X. ZHANG. Refining Word Embeddings for Sentiment Analysis. In: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Copenhagen, Denmark: Association for Computational Linguistics, Sept. 2017, pp. 534–539. DOI: [10.18653/v1/D17-1056](#) (↑ 8).
- [Zha+21] Y. ZHANG, M. CHOI, K. HAN, and Z. LIU. Explainable Semantic Space by Grounding Language to Vision with Cross-Modal Contrastive Learning. In: *CoRR* abs/2111.07180 (2021). arXiv: [2111.07180](#) (↑ 8).
- [Zha+19] J. ZHAO, T. WANG, M. YATSKAR, R. COTTERELL, V. ORDONEZ, and K.-W. CHANG. Gender Bias in Contextualized Word Embeddings. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, June 2019, pp. 629–634. DOI: [10.18653/v1/N19-1064](#) (↑ 8).
- [Zho+22] Z. ZHOU, D. ZHANG, W. XIAO, N. DINGWALL, X. MA, A. O. ARNOLD, and B. XIANG. *Learning Dialogue Representations from Consecutive Utterances*. 2022. DOI: [10.48550/ARXIV.2205.13568](#) (↑ 8).
- [Zhu+20] Y. ZHU, D. ZHOU, J. XIAO, X. JIANG, X. CHEN, and Q. LIU. HyperText: Endowing FastText with Hyperbolic Geometry. In: *Findings of the Association for Computational Linguistics: EMNLP 2020*. Online: Association for Computational Linguistics, Nov. 2020, pp. 1166–1171. DOI: [10.18653/v1/2020.findings-emnlp.104](#) (↑ 9).

DIALOG SYSTEMS AND MACHINE LEARNING GROUP, HHU DÜSSELDORF, GERMANY

Email address: [benjamin.ruppik@hhu.de](mailto:benjamin.ruppik@hhu.de)

URL: <http://ben300694.github.io>