

Understanding and exploiting subspace organization in contextual word embeddings

David Yenicelik
ETH Zürich

*A dissertation submitted to ETH Zürich
in partial fulfilment of the requirements for the degree of
Master of Science ETH in Computer Science*

Under the supervision of
Florian Schmidt
Yannic Kilcher
Prof. Dr. Thomas Hofmann

Eidgenössische Technische Hochschule Zürich
Data Analytics Lab
Institute of Machine Learning
CAB F 42.1, Universitätsstrasse 6, 8006, Zürich
Zürich 8006
SWITZERLAND

Email: yedavid@ethz.ch

April 29, 2020

Declaration

I David Yenicecik of ETH Zürich, being a candidate for the M.Sc. ETH in Computer Science, hereby declare that this report and the work described in it are my own work, unaided except as may be specified below, and that the report does not contain material that has already been used to any substantial extent for a comparable purpose.

Total word count: 14,235

Signed:

Date:

This dissertation is copyright ©2020 David Yenicecik.

All trademarks used in this dissertation are hereby acknowledged.

Abstract

Write a summary of the whole thing. Make sure it fits in one page.

How can we arrive at meaning embeddings? Compared to context-embeddings, these should encode meaning better and exclusively.

Better represent what a computer means, can be used for interpretability, downstream tasks, better unbiased embeddings which do not entail sentiment but meaning only with the words.

Contents

1	Introduction	3
1.0.1	Main Contributions	4
2	Background	5
2.1	Linguistic Features	6
2.2	Word Embeddings	7
2.2.1	Static Word Embeddings	13
2.2.2	Context Embeddings	19
2.2.3	Other methods	25
2.3	Resources and Datasets	28
3	Related Work	37
3.1	Structure inside BERT	37
3.1.1	Attention mechanism	39
3.1.2	From token-vectors to word-vectors	39
3.1.3	Bias in BERT vectors	39
3.1.4	Change of meanings over time	39
3.1.5	Clinical concept extration	40
3.1.6	Discovering for semantics	40
3.1.7	Embeddings for translation	40
3.2	Metric Learning and Disentanglement	41
3.3	Clustering Algorithms	47
3.4	Applications of word vector	48
4	Analysing the current state of the art	53
4.1	On the Linear Separability of meaning within sampled BERT vectors	53
4.1.1	Motivation	53
4.1.2	Experiment setup	53
4.1.3	Results	55

4.2	On the Clusterability of meaning within sampled BERT vectors	57
4.2.1	Motivation	57
4.2.2	Experiment setup	57
4.2.3	Results	60
4.3	Correlation between Part of Speech and Context within BERT	71
4.3.1	Motivation	71
4.3.2	Experiment setup	71
4.3.3	Results	71
5	Exploiting subspace organization of semantics of BERT embeddings	73
5.0.1	BERnie PoS	74
5.0.2	BERnie Meaning	75
5.0.3	BERnie Meaning with additional pre-training	80
5.1	Compressing the non-lexical out	81
6	Conclusion	83

List of Figures

2.1	Figure taken from [28]. The CBOW architecture predicts the current word based on the context. The Skip-gram predicts surrounding words given the current word.	11
2.2	Figure taken from [29]. A 2-dimensional PCA projection of the 1000-dimensional skip-gram vectors of countries and their capital cities. The proposed model is able to automatically organize concepts and learn implicit relationships between them. No supervised information was provided about what a capital city means.	14
2.3	Figure taken from [52]. The transformer module architecture. The transformer encapsulates multiple attention layers. . . .	22
2.4	Figure taken from [10]. BERT takes as input multiple tokens, including a position embedding, the token embedding and the segment embedding. This allows BERT to distinguish between the location of the word within a sentence, and which word token was provided and which sentence the word token is a part of.	23
2.5	Figure taken from [30]. Word-forms F_1 , and F_2 are synonyms of each other, as they share one word meaning M_1 . Word-form F_2 , as it entails more than one meaning, namely M_1 and M_2 . .	29
2.6	Example output for WordNet 3.1 noun propositions for the word "bank". In total, 18 different concepts are recorded. . . .	30
2.7	Example output for WordNet 3.1 noun propositions for the word "bank". In total, 18 different concepts are recorded. . . .	31

2.8	Shows that the SemCor data is biased. Words with a low WordNet sense index (i.e. close to 0) occur much more often than words that have a high WordNet sense index (i.e. above 5). The x-axis shows the WordNet sense index for a chosen word, while the y-axis shows the log-frequency within SemCor. This is a cumulative plot over all words with WordNet senses within SemCor 3.0. The skew could be a natural effect of how word lower WordNet indecies are assigned to more commonly used words.	32
2.9	Hello	33
3.1	From [21], visualizing clustering of the encoder representations of all languages, based on ther SVCCA similarity.	41
3.2	Taken from [50]. The individual tiles show the kNN prediction regions by color for every point in the image. Using the unmodified euclidean distancee, this would result in classification regions on the left. The reader can see, learning an appropriate distance, the classification is much more effective (middle). Finally, the dimensionality is also reducable with this dimension while still matching the classification accuracy (right).	44
3.3	Taken from [20]. An illustration of deep metric learning. The space is transformed in such a way, that similar object are closer to each other, and dissimilar objects are moved away from each other.	45
3.4	Taken from [9]. Each	48
4.1	A famous equation	62
4.2	A famous equation	63
4.3	A famous equation	63
4.4	A famous equation	64
4.5	A famous equation	64
4.6	A famous equation	65
4.7	A famous equation	65
4.8	A famous equation	66
4.9	A famous equation	66
4.10	A famous equation	67
4.11	A famous equation	67
4.12	A famous equation	68
4.13	A famous equation	68
4.14	A famous equation	69

4.15	A famous equation	69
4.16	A famous equation	70
4.17	A famous equation	70
4.18	plots of....	72
4.19	plots of....	72
4.20	plots of....	72
5.1	The BERT model takes as input a sentence s . The sentence s is converted to a sequence of BERT tokens t_1, \dots, t_m as defined in a given vocabulary V . Each item in the vocabulary V has a corresponding embedding vector inside the embedding layer of the transformer. This embedding vector is used by the intermediate layers of the transformer, and thus affects the downstream pipeline of the transformer for any subsequent layers of the transformer.	74
5.2	The modified pipeline. The BERNie model takes as input a sentence s . The sentence s is converted to a sequence of BERT tokens t_1, \dots, t_m as defined in a given vocabulary V . For each target token t_{target} , we make the token more specific by converting the token to a more specialized token-representation, which specifies the part-of-speech information as part of the token. In this case, <i>run</i> becomes <i>run-VERB</i> . Again, each item in the vocabulary V has a corresponding embedding vector inside the embedding layer of the transformer. This embedding vector is used by the intermediate layers of the transformer, and thus affects the downstream pipeline of the transformer for any subsequent layers of the transformer.	75
5.3	Inside the embedding layer of BERT, we introduce more specific embeddings <i>run- VERB</i> and <i>run- NOUN</i> . The BERT model should intuitively now capture more expressiveness, as the model size increased. The original <i>run</i> embedding is removed.	76
5.4	plots of....	77
5.5	77
5.6	plots of....	78

5.7	The BERT model takes as input a sentence s . The sentence s is converted to a sequence of BERT tokens t_1, \dots, t_m as defined in a given vocabulary V . Each item in the vocabulary V has a corresponding embedding vector inside the embedding layer of the transformer. This embedding vector is used by the intermediate layers of the transformer, and thus affects the downstream pipeline of the transformer for any subsequent layers of the transformer.	78
-----	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	----

List of Tables

2.1	Taken from [55] Table to test captions and labels	7
4.1	Mean and standard deviation of the accuracy of a linear classifier trained on the 2 most common classes of WordNet meanings for the word <i>was</i>	55
4.2	Mean and standard deviation of the accuracy of a linear classifier trained on the 2 most common classes of WordNet meanings for the word <i>is</i>	56
4.3	Mean and standard deviation of the accuracy of a linear classifier trained on the 2 most common classes of WordNet meanings for the word <i>one</i>	56
4.4	Mean and standard deviation of the accuracy of a linear classifier trained on the the 4 most common classes of WordNet meanings for the word <i>was</i>	56
5.1	asj	79
5.2	Mean and standard deviation of the accuracy of a linear classifier trained on the the 4 most common classes of WordNet meanings for the word <i>was</i>	80
5.3	Mean of mulitple experiments for BERNie with addiitonally trained embeddings and weights.	81

Chapter 1

Introduction

One of the goals of Natural Language Understanding *NLU* is to formalize written text in a way, such that linguistic features can be captured through computational models, which can then again be used for downstream machine learning tasks. The most popular methodology in NLU is the use of vectors that represent written text, including word-vectors, sentence-vectors and more. These vectors are often referred to as *embedding* vectors, as they embed more complex concepts into a vector representation. Optimally, the algebra of the space the we operate in, including the various operations such as multiplication and addition should have meaningful relations to the concepts captured by these embeddings.

Language models *LM* are a generalization of NLU models that can both generate word-embeddings, sentence embeddings, but also intermediate representations which can be used for downstream tasks. Most commonly, language models take into consideration the context that a word token appears and, and includes this in the calculation when creating the embedding.

Because embeddings - due to their strong usage in downstream tasks - can be considered as the fundamental unit of machine learning in the domain

of natural language, any shortcomings will propagate to the performance of the target task. Thus, improving the way these embeddings have strong implications for any subtasks in NLU, including but not limited to named entity recognition *NER*, sentiment analysis and translation between languages. However, most of the modern LMs are so complex that these are considered black-box models.

Although this work does not try to imminently push the performance of these language models, we believe that providing a better understanding of the underlying principles will pave the way for future research to better track these issues. For the sake of focus, and because we believe that meaning is the most important linguistic feature for communication, we will focus on investigating how modern LMs capture semantics.

1.0.1 Main Contributions

Our main contribution lies in conducting experiments which test how semantics is captured in one of the most popular language models. Because discriminative tasks are simpler than generative ones, first we analyse the linear separability of semantics within embedding vectors produced by BERT. Then we check how natural this separability is entailed, by trying to match an implicit generative distribution onto it through clustering. We then investigate to what effect introducing additional embedding vectors inside modern language models has. We aim to use simple and easy to understand models for best interpretability of the results.

In conclusion, we show that semantics in modern language models is not interpretable using simple linear models, implying that model complex language models are still required for downstream tasks. We also show that linguistic features such as sentiment are much more strongly captured in modern language models, allowing for language stereotypes to easily manifest itself, leading to stronger cases of bias in language models.

Chapter 2

Background

TODO: Take some more from here

- Some of the main points behind [14] are that there is inherent structure in language, and that this structure can be formalized. - [14] mentions that the relation between the linguistic representation in terms of sounds and tokens are related to the meaning that the representation entail. - However, despite the obvious relationship, the distinction between distributional structure of the language and meaning is not always clear. and that there is a "parallel" meaning structure, and argues that there is not a one-to-one relation between vocabulary and classification of meaning. - Also argues that the meaning of a word is entailed by it's environment, and the words that it occurs with. - Other viewpoints (CITE HOFMANNs "teachers") In the end, language captures the evolution of human thought.

- cite paper "meaning is classified by its context"

Chronologically, the following data structures are manifestations of the above idea of defining a word by it's neighbourhood.

2.1 Linguistic Features

Polysemy: a word may have multiple meanings. in a cross-lingual embedding space, this feeling is amplified. there's some work in multi-sense embedding. this should enable to capture more fine-grained embeddings embeddings.

Over the past decades, linguists have tried to inject structure into language. Similar to features in machine learning, there are properties in language. There is a vast number of linguistic features, going from phonological features, morphological and syntactic features to semantic features. To keep this analysis concise, we will give a short introduction of what features that are relevant for the topic of finding suitable word embeddings.

Lexical features include word classes, such as nouns, verbs and adjectives, which follow certain grammatical rules in western languages. Whenever we are looking to parse for such lexical features, we use the python spaCy package [16]. As defined in [1], some of the linguistic features that are available to us include but are not limited to

1. adjectives *ADJ*
2. adposition *ADP*
3. adverbs *ADV*
4. auxiliary *AUX*
5. conjunction *CONJ*
6. coordinating conjunction *CCONJ*
7. determiner *DET*
8. interjection noun *INTJ*
9. noun *NOU*
10. numeral *NUM*
11. particle *PART*

12. pronoun *PRON*
13. proper noun *PROP*
14. punctuation *PUNCT*
15. subordinating conjunction *SCONJ*
16. symbol *SYM*
17. verb *VERB*

These features can alter as we look at different languages

Other linguistic phenomena are well described in [55], which is one of the benchmarking resources we are using (described below). [55] records the following linguistic phenomena.

Coarse-Grained Categories	Fine-Grained Categories
Lexical Semantics	Lexical Entailment, Morphological Negation, Factivity, Symmetry/Collectivity, Redundancy, Named Entities, Quantifiers, Core Predicate Argument Structure & Core Arguments, I
Phrases, Ellipsis Implicits,	Anaphora/Coreference Active/Passive, Nominalization, Genitives/Partitives, Datives, Relative Clauses, Coordination Scope, Intersectivity, Restrictivity
Logic	Negation, Double Negation, Intervals/Numbers, Conjunction, Disjunction, Conditionals, Universal, Existential Temporal
Knowledge	Downward Monotone, Non-Monotone Common Sense, World Knowledge

Table 2.1: Taken from [55] Table to test captions and labels

2.2 Word Embeddings

Word-Embeddings In general, we want to find a mapping

$$w^t \mapsto (x_w, b_w) \in \mathcal{X}^{d+1} \quad (2.1)$$

where w^t is a token representation from a vocabulary $w^t \in V$ and where this token representation is transformed into a d -dimensional vector representation x_w , and a bias-term b_w that is specific to the token w^t . Whenever we will talk about *word-vectors*, *(word)-embeddings*, or *feature-vectors*, we will refer to the image of the above map. For convenience and unless stated otherwise, we will assume that x_w absorbs the bias term b_w as an additional vector-element. Also, for simplicity and unless otherwise stated, we will use the euclidean real-valued vector-space \mathbb{R}^{d+1} to describe the resulting embedding vectors. Please note, however, that the choice of the embedding space \mathcal{X} is not fixed in general, and as such, work in other spaces have also been conducted.

Distance Now our goal is to build an embedding space where the relationship between words are meaningful. Specifically, we want to incorporate the notion of *meaning* into these word-embeddings, which should follow following properties. Formally, we introduce the concept of *distance* $d : \mathcal{X} \times \mathcal{X} \mapsto [0, \infty)$, which should capture the relation between different elements. For a set of elements $x, y, z \in \mathcal{X}$, following properties must hold for our distance metric to be a valid one:

1. $d(x, y) \geq 0$ (non-negativity)
2. $d(x, y) = 0 \iff x = y$ (identity, i.e. if two embedding vectors are identical, they must capture the same underlying word instance)
3. $d(x, y) \leq d(x, z) + d(z, y)$ (triangle inequality)

One consequence of the above rules is that for words a, b, c , each having a word embedding x, y, z respectively, $d(x, y) < d(x, z) \iff$ word instance b is conceptually closer to word a than word c . For convenience, whenever I input a, b, c into the distance function d , the word-embeddings for the corre-

sponding word-tokens shall be used. Also, please notice that I left out the notion of symmetry for distance measures.

Learning a distance Often in machine learning, one wants to maximize a certain probability distribution given some data \mathbf{X} . In the context of word vectors, we want to maximize the probability that w occurs in the context window of w' through some parameters θ . Implicitly, this corresponds to minimizing the distance between w and w' , while keeping the distance between w and all other words constant. We call w' a *context word* for w .

$$p(w|w') \tag{2.2}$$

and

$$\forall w, w' : \quad \max p(w|w') \iff \min d(w|w') \tag{2.3}$$

Going from the distributional structure of sentences to learning distances between words. One of the early formalizations of representing the distributional structure of words through was expressed in [4], which argues that a sequential statistical model can be constructed to estimate this true posterior. In it's most naive form, this would imply that we can estimate the probability of a word $w^{(t)}$ after words $w^{(t-1)}, \dots, w^{(1)}$ as

$$p(\mathbf{w}) = p(w^{(t)}, \dots, w^{(1)}) \tag{2.4}$$

$$= p(w^{(t)} | w^{(t-1)}, \dots, w^{(1)}) \tag{2.5}$$

where $\mathbf{w} = w^{(1)}, \dots, w^{(T)}$ is the sequence of words, $p(\mathbf{w})$ the probability of this sentence occurring.

Because inference for the above equation would have been computationally

infeasible, the notion of *n-grams* was introduced to estimate this conditional probability, following a markov assumption.

$$p(w^{(t)}|w^{(t-1)}, \dots, w^{(1)}) \approx p(w^{(t)}|w^{(t-1)}, \dots, w^{(t-n+1)}) \quad (2.6)$$

Often, the above equation will be simplified even further using the independence of words assumption

However, because we are only interested in the distance between two words at a time (and not the full sequence) [4] simplify the above function using the assumption of marginalized independence of words amongst each other w.r.t. a target word (t) (we can do this because we marginalize over all possible combinations of target words and context words).

$$p(\mathbf{w}) = \prod_{t=1}^T \prod_{\Delta < t} p(w^{(t)}|w^{(t-\Delta)}) \quad (2.7)$$

whose probability we wish to estimate (and maximize if this is in the given training dataset), $\mathcal{I} = \{-R, \dots, -1, 1, \dots, R\}$ is the so-called *context window* which includes all the words that R index units to the left, and R index units to the right of the word of interest $w^{(t)}$.

However, we do not need to only look at only the previous words. One can consider the full neighbourhood of a word. If we describe this by a loss-measure we would like to minimize, this would result in

$$\mathcal{L}(\theta; \mathbf{w}) = \prod_{t=1}^T \prod_{\Delta \in \mathcal{I}} p_{\theta}(w^{(t)}|w^{(t+\Delta)}) \quad (2.8)$$

where $\mathbf{w} = w^{(1)}, \dots, w^{(T)}$ is the sequence of words whose probability we wish to estimate (and maximize if this is in the given training dataset),

$\mathcal{I} = \{-R, \dots, -1, 1, \dots, R\}$ is the so-called *context window* which includes all the words that R index units to the left, and R index units to the right of the word of interest $w^{(t)}$.

If p_θ is a parametric function, one can then optimize for the most optimal $\hat{\theta} = \operatorname{argmax}_\theta \mathcal{L}(\theta; \mathbf{w})$ using a maximum likelihood estimation approach.

In general, one does not necessarily need to interpret \mathbf{w} as a sequence of word-tokens, but can also interpret $w^{(1)}, \dots, w^{(T)}$ as *n-grams*, which are triplets of n-characters forming a token-unit. Generally, the definition of a token is open for interpretation. We will generally assume that a single word is a token unless otherwise stated.

Intuitively the above models follow the idea expressed by [14] very well, speaking that the meaning of a word is captured by its neighborhood. However, the above equations follow a *continuous bag of words* (CBOW) approach. A continuous-bag-of-words approach is an approach where we want to predict a target word w^t given some context words w' . However, it is also possible to follow a *skip-gram* approach. Here, we want to predict context words w' given a target word w^t .

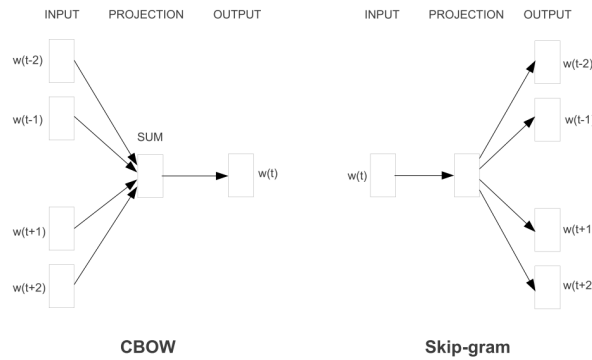


Figure 2.1: Figure taken from [28]. The CBOW architecture predicts the current word based on the context. The Skip-gram predicts surrounding words given the current word.

Using a skip-gram approach, (2.2) for example would be reformulated into

$$\mathcal{L}(\theta; \mathbf{w}) = \prod_{t=1}^T \prod_{\Delta \in \mathcal{I}} p_{\theta}(w^{(t+\Delta)} | w^{(t)}) \quad (2.9)$$

The question now is, how do we parametrize the distribution $p_{\theta}(w, w')$, as well as the word feature vectors x_w and b_w . The following will show a few methods of how this can be achieved. In the following methods shown, we will aim to provide a loss function that we try to minimize, and interesting properties for each method. However, we will not go into too much details, as to how the loss function is optimized, as almost all of these methods can be solved using gradient-based methods.

2.2.1 Static Word Embeddings

Here we will talk about word-embeddings where each word-token only has a single x_w . Specifically, the mapping (2.2) is not a probabilistic function with an implicit random factor, but rather a deterministic one.

Basic Model

The first model we are going to look at is a most basic model which fulfills the properties of the distance metrics shown in (2.2).

Here, we can introduce a *log-bilinear* model where the log-probability is defined as

$$\log p_\theta(w|w') = \langle \mathbf{x}_w, \mathbf{x}_{w'} \rangle + b_w + \text{const.} \quad (2.10)$$

To arrive at the actual probability, we can exponentiate the log-probability as such

$$p_\theta(w|w') = \frac{\exp[\langle \mathbf{x}_w, \mathbf{x}_{w'} \rangle + b_w]}{Z_\theta(w')}$$

where $Z_\theta(w') := \sum_{v \in \mathcal{V}} \exp[\langle \mathbf{x}_v, \mathbf{x}_{w'} \rangle + b_v]$ is a normalization constant such that the probability mass sums to 1, and the model parameters entail the word-embeddings $\theta = (x_w, b_w) \in \mathbb{R}^{d+1}$

Word2Vec

During training, the above basic model comes with certain drawbacks. The distance would be minimized if all word-embeddings would collapse onto a single point. Also, there is no term that forces unrelated words to move away from each other, a property that we are interested in as unrelated words should form embedding vectors that are far from each other.

One of the most prominent example of word vectors manifests itself in the work of [28] and [29]. Here, a neural network with a single embedding layer

can be trained to transform one-hot-vectors $\in \{0, 1\}^{|\mathcal{V}|}$ which represents a word w in vocabulary \mathcal{V} into a latent vector representation $w \in \mathbf{R}^{d+1}$ using both a continuous bag of word and also a continuous skip-gram approach. The skip-gram approach is preferred in practice.

Specifically, the loss-function that is optimized looks as follows

$$\mathcal{L}(\theta; \mathbf{w}) = \sum_{t=1}^T \sum_{\Delta \in \mathcal{I}} [\quad \quad \quad] \quad (2.11)$$

$$b_{w^{t+\Delta}} + \langle \mathbf{x}_{w^{t+\Delta}}, \mathbf{x}_{w^{(t)}} \rangle \quad (2.12)$$

$$- \log \sum_{v \in \mathcal{V}} \exp [\langle \mathbf{x}_v, \mathbf{x}_{w^{(t)}} \rangle + b_v] \quad (2.13)$$

As one can see, the loss function takes in the bilinear loss from the basic model, and complements this by adding a term, such that random samples are not put next to each other. This principle is often referred to as *negative sampling*.

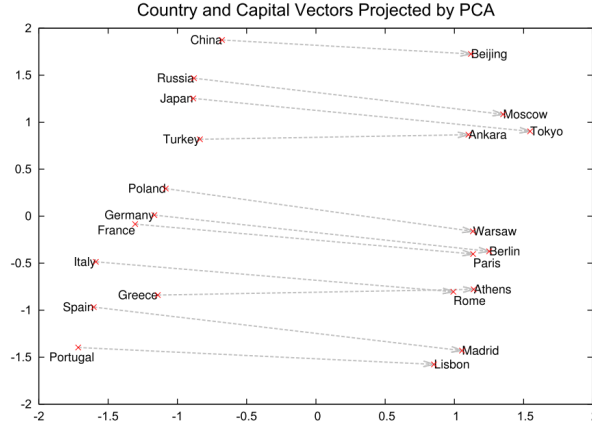


Figure 2.2: Figure taken from [29]. A 2-dimensional PCA projection of the 1000-dimensional skip-gram vectors of countries and their capital cities. The proposed model is able to automatically organize concepts and learn implicit relationships between them. No supervised information was provided about what a capital city means.

GloVe

For the global vectors for word representation (GloVe) [35], the authors follow a more straight-forward matrix factorization approach.

First, a global co-occurrence matrix is created $\mathbf{N} = (n_{ij}) \in \mathbb{N}^{|\mathcal{V}||\mathcal{C}|}$ where each entry n_{ij} is determined by the number of occurrences of word $w_i \in \mathcal{V}$ in context $w_j \in \mathcal{C}$. Given that the vocabulary size can exceed multiple thousand items, this practically results in a sparse matrix.

$$\mathcal{H}(\theta; \mathbf{N}) = \sum_{i,j} f(n_{ij}) \underbrace{(\log n_{ij})}_{\text{target}} - \underbrace{\log \tilde{p}_\theta(w_i|w_j)}_{\text{model}})^2 \quad (2.14)$$

$$= \sum_{i,j} f(n_{ij}) (\log n_{ij} - \langle x_i, y_j \rangle)^2 \quad (2.15)$$

$$(2.16)$$

where $f(n_{ij})$ is the weighting function which assigns a weight for each entry in the co-occurrence matrix. In the second line, we also again use a bilinear probability density model $\tilde{p}_\theta(w_i|w_j) = \exp[\langle \mathbf{x}_i, \mathbf{y}_j \rangle + b_i + c_j]$. The constants b_i, c_j are left out and are assumed to be absorbed in the embedding vectors.

A popular choice for the weighting function is

$$f(n) = \min \left\{ 1, \left(\frac{n}{n_{\max}} \right)^\alpha \right\}$$

with $\alpha \in (0, 1]$. The motivation behind this is that frequent words do not receive a weighting which is "too high" (there is a cutoff at some point) and small counts are considered noise and slowly cancelled out.

Gaussian Embeddings

Gaussian Embeddings have been first proposed in the context of words, although prior work has been adapted in embedding matrix-rows into a mean

and standard deviation [53]

It is a continuous probabilistic relaxation of the otherwise common discrete point vectors. Each word is represented by a Gaussian distribution in high-dimensional space, with the aim to better capture uncertainty and a representation and it's relationships. It can also express asymmetric relations more naturally than dot-products or cosine similarity, and enables for better-parametrized rule between decision boundaries.

Training is done using an energy function which is incorporated within a loss function that we try to minimize. The energy function describes similarity between two items. The authors propose the following energy functions to derive a Gaussian word-embedding.

$$L_m(w, c_p, c_n) = \max(0, m - E_\theta(w, c_p) + E_\theta(w, c_n)) \quad (2.17)$$

Here, w is the word we want to sample, c_p is a "positive" context word, i.e. a word that co-occurs with the word w , and c_n is a "negative" context word, i.e. a word that does not co-occur with the word w . Usually the negative context words is sampled randomly from the corpus. The loss function reminds of a hinge-loss in logistic regression.

The authors propose two possible ways to learn the mean and variance of the Gaussian embeddings. They argue that the empirical covariance is not the most effective method of deriving the words as Gaussian embeddings. This does not allow for inclusion between ellipsoids

Symmetric similarity: expected likelihood or probability product kernel We can use any kernel (which is symmetric by definition) to derive at an energy function. For two Gaussians $f(x)$, $g(x)$, the inner product is defined as:

$$E_\theta(w, c) = \int_{x \in \mathcal{R}^d} f(x)g(x)dx \quad (2.18)$$

$$= \int_{x \in \mathcal{R}^d} \mathcal{N}(x; \mu_w, \Sigma_w) \mathcal{N}(x; \mu_c, \Sigma_c) dx \quad (2.19)$$

$$= \mathcal{N}(0; \mu_w - \mu_c, \Sigma_w + \Sigma_c) \quad (2.20)$$

For numerical feasibility and easy of differentiation, we usually maximize the $\log E_\theta(w, c)$ for a given dataset with $w \in \mathcal{W}, c \in \mathcal{C}$.

Asymmetric divergence: KL-Divergence We can use more directional supervision to exploit directional supervision, such as a knowledge graph.

Following energy-function is optimized:

$$-E(w_i, c_j) = D_{KL}(c_j || w_i) \quad (2.21)$$

$$= \int_{x \in \mathcal{R}^d} \mathcal{N}(x; \mu_{w_i}, \Sigma_{w_i}) \log \frac{\mathcal{N}(x; \mu_{c_j}, \Sigma_{c_j})}{\mathcal{N}(x; \mu_{w_i}, \Sigma_{w_i})} dx \quad (2.22)$$

$$= \frac{1}{2} \left(\text{tr}(\Sigma_i^{-1} \Sigma_j) + (\mu_i - \mu_j)^\top \Sigma_i^{-1} (\mu_i - \mu_j) - d - \log \frac{\det(\Sigma_j)}{\det(\Sigma_i)} \right) \quad (2.23)$$

Because of the loss function, this can entail information such as "y entails x" as a soft form of inclusion between two datasets (if KL divergence is used). If a symmetric loss function is used, then this would most likely lead to overlap (IS THIS TRUE...???)

Uncertainty calculation: In contrast to the empirical standard deviation as an uncertainty measure, we can now calculate the uncertainty of the inner product (i.e. the distribution $P(z = x^T y)$ using the following formula

$$\mu_z = \mu_x^T \mu_y \Sigma_z = \mu_x^T \Sigma_x \mu_x + \mu_y^T \Sigma_y \mu_y + \text{tr}(\Sigma_x \Sigma_y) \quad (2.24)$$

We then get an uncertainty bound, where c denotes the number of standard deviations away from the mean.

$$\mu_x^\top \mu_y \pm c \sqrt{\mu_x^\top \Sigma_x \mu_x + \mu_y^\top \Sigma_y \mu_y + \text{tr}(\Sigma_x \Sigma_y)} \quad (2.25)$$

We can learn the parameters Σ and μ for each of these embeddings using a simple gradient-based approach, where we set hard constraints on

$$\|\mu_i\|_2 \leq C, \forall i \quad (2.26)$$

$$mI < \Sigma_i < MI \quad (2.27)$$

The method shows competitive scores to the Skip-Gram model, although usually only with minor improvements depending on the benchmark-dataset.

2.2.2 Context Embeddings

These language models also include capturing the more probability of $p(\text{output}|\text{input})$. The MQAN demonstrates that a single language task is able to infer and perform many different tasks on examples with this type of format.

Now that we have seen a few methods that aim to capture one embedding x_w for each word-token, we now turn our attention to context embeddings.

Context embeddings are more modern applications of embeddings which

A single transformer takes as input a sequence x and outputs softmax probabilities. In the context of sentences, it is able to take as input a full sequence of words. This has implications for the complexity of the model. Instead of calculating the raw probability of (w, w') , and even using a Markovian assumption because computation power is expensive, the transformer architecture implicitly calculates the joint probability of an entire sequence of words without any explicit conditioning

$$p(w^{(t-d)}) = p(w^{(t)}, \dots, w^{(t-d+1)}, w^{(t-d-1)}, \dots, w^{(1)}) \quad (2.28)$$

where $p(w^{(t-d)})$ is the probability of the target word $w^{(t-d)}$ which we want to calculate.

This modifies our initial task of calculating $p(w, w')$ to not including only a single context word w' , but rather a wider context (i.e. $w^{(t)}, \dots, w^{(t-d+1)}, w^{(t-d-1)}, \dots, w^{(1)}$). Theoretically, and because we want to arrive at some distance measure, $p(w, w')$ can still be calculated by marginalizing out all the context-variables except the one of interest. Specifically, if we are interested in word w , we can marginalize out all input words except w' . Practically this is infeasible due to the high space, and as such, we make use of some other mechanisms that we will talk about in the "BERT" section below.

Because the following models not merely static word embeddings, but language models, we shall provide a short comparison between for downstream tasks between the individual language models after a short introduction on the details of each language models.

ELMo

The simplest idea of context embeddings, which take into account more than just a single word, but all tokens from a predefined context is the *Embeddings from Language Models* (ELMo) model proposed in [37].

In its underlying architecture, ELMo is using LSTMs. LSTMs are sequential recurrent neural networks .

Specifically, LSTMs [15] are good at estimating the probability of word $w^{(t)}$ occurring given a previous history of words $w^{(t-1)}, \dots w^{(1)}$. In practice, the LSTM is a popular choice for sequence modeling tasks, as it is able to capture very well the following probability distribution.

$$p(w_1, w_2, \dots, w_N) = \prod_{k=1}^N p(w_k | w_{k-1}, \dots, w_2, w_1) \quad (2.29)$$

The ELMo language model now makes use of a backward LSTM, which does not evaluate the probability of a certain token occuring with its *past* words (i.e. the context with context tokens $w^{(k)}$ with $k < t$ for a word $w^{(t)}$ whose probability we want to estimate), but rather considers only the *future context*. Specifically, this results in

$$p(w_1, w_2, \dots, w_N) = \prod_{k=1}^N p(w_k | w_{k+1}, w_{k+2}, \dots, w_N) \quad (2.30)$$

A bidirectional language model (biLM) combines the above two LSTM models, thus jointly maximizing the log likelihood of the forwards, as well as backward directions.

$$\sum_{k=1}^N \left(\log p \left(w_k | w_{k-1}, \dots, w_1; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s \right) \right) \quad (2.31)$$

$$+ \log p \left(w_k | w_{k+1}, \dots, w_N; \Theta_x, \overleftarrow{\Theta}_{LSTM}, \Theta_s \right) \quad (2.32)$$

Because the output of a sequential neural network - which the LSTM is part of - includes $2L + 1$ output representations, we receive the following set of vectors, one for each input token, which we now consider the *context embeddings*. Each context embeddings outputs a vector representing x_k for the word token $w^{(k)}$, given its context $w_1, \dots, w_{k-1} w_{k+1}, \dots, w_N$.

$$R_k = \left\{ \mathbf{x}_k^{LM}, \vec{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} | j = 1, \dots, L \right\} \quad (2.33)$$

$$= \left\{ \mathbf{h}_{k,j}^{LM} | j = 0, \dots, L \right\} \quad (2.34)$$

Notice that one can stack the above biLM such that the output of one module is taken as input to another module. The biLM module is stacked twice ELMo. If a downstream task is trained end-to-end, the all output vectors can be stacked together and provided as input to a final downstream-specific neural network module.

The Transformer Architecture

All of the below presented models use the transformer architecture which initially was presented in [52], which shows how far simple *attention* modules can go.

The following architectures are all based on the transformer module.

BERT

As introduced in [10], the Bidirectional Encoder Representations from Transformers model (BERT) combines the above concepts of forward and back-

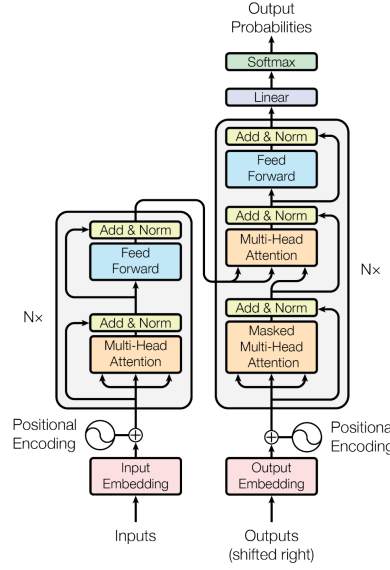


Figure 2.3: Figure taken from [52]. The transformer module architecture. The transformer encapsulates multiple attention layers.

wards sequential models with attention, as marked in the transformer architecture above.

[10] argue that the main limitations is that standard language models are unidirectional, and that .

BERT introduces itself as a pre-trained language model. Specifically, it shall be used to fine-tune on specific downstream tasks.

BERT is pre-trained on multiple gigabytes of text .

Multiple versions of BERT are provided, including $BERT_{LARGE}$ which includes , and $BERT_{BASE}$ which includes

Ever since BERT came out, a wide variety of similar models came out that mimic the main logic of BERT, and improve upon it, such as DistillBERT [44].

Numerous works have been published on making BERT representations more compact, for it to be used for on-device text-processing applications, such as more compact sentence-representations [46].

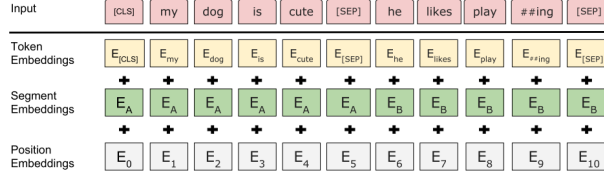


Figure 2.4: Figure taken from [10]. BERT takes as input multiple tokens, including a position embedding, the token embedding and the segment embedding. This allows BERT to distinguish between the location of the word within a sentence, and which word token was provided and which sentence the word token is a part of.

GPT and GPT-2

First proposed in [39] and further extended in [40].

[39] introduces the concept of *generative pre-training* and discriminative fine-tuning *generative pre-training*.

Unsupervised pre-training consists of an unsupervised corpus of tokens $\mathcal{U} = \{u_1, \dots, u_n\}$ and use the common objective of

$$L(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \theta) \quad (2.35)$$

where k describes a context window, and the conditional probability P is modeled using a neural network with parameters θ . The transformer decoder [24] is used which is a variant of the transformer [52], which introduces a multi-headed self-attention operation over the input context tokens followed by position-wise feedforward layers to produce an output distribution over target tokens as such:

$$h_0 = UW_e + W_p \quad (2.36)$$

$$h_l = \text{transformer_block}(h_{l-1}) \forall i \in [1, n] \quad (2.37)$$

$$P(u) = \text{softmax}(h_n W_e^T) \quad (2.38)$$

where $U = (u_{-k}, \dots, u_{-1})$ is the context vector of tokens, n is the number of layers, W_e is the token embedding matrix, and W_p is the position embedding matrix.

Supervised fine-tuning consists of leveraging a labeled dataset \mathcal{C} , where each instance consists of a sequence of input tokens, x^1, \dots, x^m along with a label y . The inputs are passed through the pre-trained model to obtain the final transformer block’s activation h_l^m which is then fed into an added linear output layer with parameters W_y to predict y .

$$P(y|x^1, \dots, x^m) = \text{softmax}(h_l^m W_y) \quad (2.39)$$

which then gives the following objective to maximize:

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y|x^1, \dots, x^m) \quad (2.40)$$

As suggested in [42] and [36], an auxiliary language modeling loss (unsupervised training) is added which accelerates convergence and improves generalization

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C}) \quad (2.41)$$

with a regularization parameter λ .

GPT-2 improves on the first version of GPT. First, the training dataset uses a crawled text of all outbound texts from Reddit resulting in the WebText corpus. The crawling also prioritizes documents by quality.

The second major modification is in input representations. GPT-2 uses a Byte Pair Encoding (BPE) scheme, [45], which operates on unicode code points as the most fundamental unit of language. The BPE tokenization is modified by avoiding merging across character categories for any byte

sequence (i.e. for words that occur in many variations, such as `dog.`, `dog!`, and `dog?`) This is similar to the character level rnns (cite jason).

The base model of GPT is modified by moving the layer normalization [3] layer to the input of each sub-block and an additional layer normalization after the final self-attention block. Also, initialization was modified to account for the accumulation of the residual path with model depth. Also, vocabulary is expanded to 50257 tokens.

2.2.3 Other methods

Although the above presented methods are the prevalent methods for word-embeddings, there are also other methods which do not clearly fit into one of the above categories.

Generating "static" word-embeddings through contextual embeddings

Some work has been done in extracting word-embeddings from contextual language models like BERT or ELMo.

CITE (BERT WEARS GLOVES: DISTILLING STATIC EMBEDDINGS FROM PRETRAINED CONTEXTUAL REPRESENTATIONS)

(1) Uses *pooling* between BERT tokens to arrive at a single representation between words.

Here, sentences are split by space (tokenized). Words are tokenized further into a subword as defined by WordPiece (Wu et al. 2016).

The defined pooling operations looks as follows to arrive at the word from the individual subwords:

$$\mathbf{w}_c = f(\mathbf{w}_c^1, \dots, \mathbf{w}_c^k); f \in \{\text{min, max, mean, last}\}$$

where we have subwords w^1, \dots, w^k such that $\text{cat}(w^1, \dots, w^k) = w$

Why would any of these pooling operations result in a meaningful source-word? This is just squishing tokens together!

-¿ This is a major limitation for which we may need to use ELMo -¿ However this may be needed for "unseen concepts" (which are unseen words...) -¿ Perhaps check what fasttext does...?

(2) Uses *context combination* to map from different contexts c_1, \dots, c_n to a single static embedding w that is agnostic of context.

Proposed are two ways to represent context.

Decontextualization For a single word-context, we simply feed-in the word by itself to the model.

Aggregated combine w in multiple contexts. n sentences are sampled from the dictionary \mathcal{D} . From the multiple sampled words, we then apply pooling to arrive at a single representation that aggregates the different tokens into one.

$$\mathbf{w} = g(\mathbf{w}_{c_1}, \dots, \mathbf{w}_{c_n}); g \in \{\text{min, max, mean}\}$$

This post extracts (token?) word-embeddings: (<https://towardsdatascience.com/nlp-extract-contextualized-word-embeddings-from-bert-keras-tf-67ef29f60a7b>)

This seems to be a way to extract embeddings for tokens from BERT (https://github.com/immortalc/word_embeddings)

(-¿How can we create a (parametric) probability density from a point-cloud distribution?)

perhaps not necessarily interpretable in standard euclidean space (<https://www.kdnuggets.com/2019/04/bert-features-interbertible.html>) original (<https://medium.com/thelocalminima/are-bert-features-interbertible-250a91eb9dc>)

Perhaps we can mask all but the target token to arrive at one vector per token (and then combine them somehow...). But how do they extract the

singular word-embeddings...?

(-; you could be like "acquiring bedeutung" is a big problem in many tasks. especially useful when we try to map one concept to another. we look at the NLP task for concreteness)

Generally, really good critique on this paper:

(<https://openreview.net/forum?id=SJg3T2EFvr>)

usually, we have sentence-embeddings, and do not look at word-embeddings.

(-; we don't want to add more and more context. we want a model which contains the polysemy of different contexts, which could allow for probability maximization..., otherwise we have to look at bigger and bigger documents to build more accurate language models, which becomes infeasible at some point. (although this would be the way humans work, because they live in context as well))

This blog aims to generate word-embeddings (and sentence-embeddings) from the BERT model. (<https://mccormickml.com/2019/05/14/BERT-word-embeddings-tutorial/>)

create word-vectors by taking out what BERT predicts at the nth token.
create word-vectors by concatenating or summing multiple layer's outputs.

the cosine similarity between these vectors seem pretty well-done!

(-; Does it make sense to use BERT and then on calculate word-embeddings through an extended fully-connected model)

-; ELMo may provide a better tokenizer, maybe better ot use this? What about GPT? ELMo uses moSES tokenizer which seems word-level enough

-; How to solve this tokenization problem....

-; Can also analyze only words that exist.

2.3 Resources and Datasets

WordNet

The online lexical database WordNet was originally introduced in [30]. WordNet is a reference system whose design is inspired by current psycholinguistic theories of human lexical memory. English nouns, verbs and adjectives are organized into synonym sets, each representing one underlying lexical concept. Different relations link the synonym sets. WordNet 1 contains of 95,600 different word forms (51,500 simple words and 44,100 collocations), and include a total of 70,100 word meanings, or sets of synonyms. Compared to a standard dictionary, WordNet divides the lexicon into five categories, namely nouns, verbs, adjectives, adverbs and function words. The authors argue that the rest of language is probably stored separately as part of the syntactic component of language.

The most ambitious feature of WordNet, however, is to attempt to organize lexical information in terms of word meanings, rather than word forms.

The problem with alphabetical thesaurus is redundant entries. If word W_a and word W_b are synonyms, the pair should be entered twice. The problem with a topical thesaurus is that two look-ups are required, first on an alphabetical list and again in the thesaurus proper.

Lexical matrix: Words are often referred to both the utterance and to its associated concept. As such, [30] specify the difference between the two concepts as "word form", which refers to the physical utterance or inscription, and "word meaning", which refers to the lexicalized concept that a form can be used to express.

The following table captures the concept of a lexical matrix, which encodes word-forms (columns), word-meanings (rows), and the existence of the expression of the word-meanings through the corresponding word-form (cell). If multiple entries exist for a single column, then the single word-form encodes multiple meanings, implying that the word-form is polysemous (it encodes multiple word-forms). If multiple entries exist for a single row, the two words

express the same underlying concept, and thus the two words-forms are synonyms.

Word Meanings	Word Forms				
	F_1	F_2	F_3	\dots	F_n
M_1	$E_{1,1}$	$E_{1,2}$			
M_2		$E_{2,2}$			
M_3			$E_{3,3}$		
\vdots				\ddots	
M_m					$E_{m,n}$

Figure 2.5: Figure taken from [30]. Word-forms F_1 , and F_2 are synonyms of each other, as they share one word meaning M_1 . Word-form F_2 , as it entails more than one meaning, namely M_1 and M_2 .

WordNet also includes synonym sets (synsets), which are a collection of word-forms that together determine a single meaning. Thus, a meaning is represented as a synset. WordNet is organized by semantic relations. Thus, WordNet includes a set of semantic relations. WordNet defines synonyms as a set of words, where one word is interchangeable for another. This implies that in terms of substitutability, WordNet clusters words into nouns, verbs, adjectives and adverbs, as the part of speech must be valid at the position of the sentence.

Although the authors mention that synonyms should be best thought of a continuum along which similarity of meaning can be graded, the authors decide to determine similarity in terms of a binary event, something which is either present, or not present.

This is mainly handcrafted by linguistics over the last 2 decades. Meanings have different levels of detail.

Part of Speech	Definition
noun	(sloping land (especially the slope beside a body of water)) "they pulled the canoe up on the bank"; "he sat on the bank of the river and watched the currents"
noun	depository financial institution, bank, banking concern, banking company (a financial institution that accepts deposits and channels the money into lending activities) "he cashed a check at the bank"; "that bank holds the mortgage on my home"
noun	(an arrangement of similar objects in a row or in tiers) "he operated a bank of switches"
verb	(tip laterally) "the pilot had to bank the aircraft"
verb	(do business with a bank or keep an account at a bank) "Where do you bank in this town?"

Figure 2.6: Example output for WordNet 3.1 noun propositions for the word "bank". In total, 18 different concepts are recorded.

Now we go to "was"

Part of Speech	Definition
noun	Washington, Evergreen State, WA, Wash. (a state in north-western United States on the Pacific)
verb	(have the quality of being; (copula, used with an adjective or a predicate noun)) "John is rich"; "This is not a good answer" bank"; "that bank holds the mortgage on my home"
verb	(an arrangement of similar objects in a row or in tiers) "he operated a bank of switches"
verb	(form or compose) "This money is my only income"; "The stone wall was the backdrop for the performance"; "These constitute my entire belonging"; "The children made up the chorus"; "This sum represents my entire income for a year"; "These few men comprise his entire army"
verb	(work in a specific place, with a specific subject, or in a specific function) "He is a herpetologist"; "She is our resident philosopher"

Figure 2.7: Example output for WordNet 3.1 noun propositions for the word "bank". In total, 18 different concepts are recorded.

SemCor dataset

The SemCor 3.0 corpus is a collection of the Brown corpora [11], where each noun, adjective and verb is tagged with the WordNet 3.0 senses. It was part of the early WordNet project, initially introduced in [31].

Some example snippets look as follows, where SemCor includes sentences, and each sentence is annotated by WordNet senses.

```
A Texas halfback who does n't even know the team 's plays,
Eldon_Moritz, ranks fourth in Southwest_conference scoring after three games.
```

The corresponding part-of-speech labels and wordnet sense ids are

```
DT, NN, NN, WP, VBZ, RB, RB, VB, DT, NN,
POS, NN, NNP, VB, JJ, IN, NNP, VP, IN, JJ, NN
```

None, 1, 1, None, 0, 1, 7, None, 1,
None, 3, 1, 1, 1, None, 1, 1, 1, 1

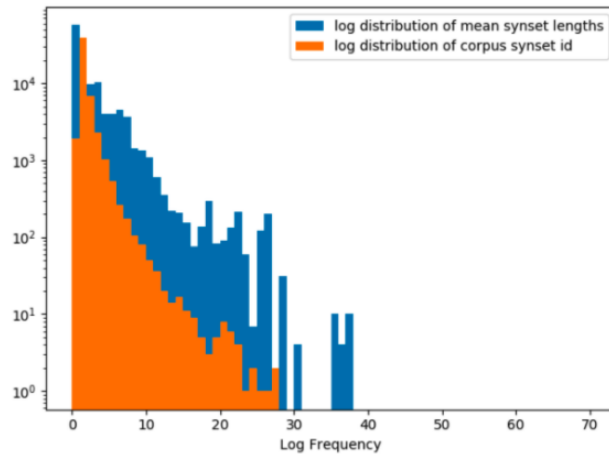


Figure 2.8: Shows that the SemCor data is biased. Words with a low WordNet sense index (i.e. close to 0) occur much more often than words that have a high WordNet sense index (i.e. above 5). The x-axis shows the WordNet sense index for a chosen word, while the y-axis shows the log-frequency within SemCor. This is a cumulative plot over all words with WordNet senses within SemCor 3.0. The skew could be a natural effect of how word lower WordNet indecies are assigned to more commonly used words.

News dataset

GLUE benchmark dataset

The GLUE benchmark dataset was first introduced by [55]

Additional work has been done through the introduction of SuperGLUE [54] to train for an even bigger variety of language tasks, after most NLU models started surpassing non-expert human performance on the initially-proposed GLUE-tasks .

Corpus	Train	Test	Task	Metrics	Domain
CoLA	8.5k	1k	acceptability	Matthews corr.	misc
SST-2	67k	1.8k	sentiment	accuracy	movie reviews
MRPC	3.7k	1.7k	paraphrase	accuracy / F1	news
STS-B	7k	1.4k	sentence similarity	Pearson/Spearman corr.	misc.
QQP	364k	391k	paraphrase	accuracy / F1	social QA questions
MNLI	393k	20k	NLI	mis-/matched accuracy	misc
QNLI	105k	5.4k	QA/NLI	accuracy	Wikipedia
RTE	2.5k	3k	NLI	accuracy	news, Wikipedia
WNLI	634	146	coreference/NLI	accuracy	fiction books

Figure 2.9: Hello

Single Sentence Tasks

The Corpus of Linguistic Acceptability **CoLA** consists of english sentences, where the task of the model is to predict whether or not the sentence is a valid english sentence. The Matthews correlation evaluates the model from a score of -1 to 1, where 0 corresponds to uniformly random guessing. The test set includes out of domain sentences.

label:1

The professor talked us into a stupor.

label:0

The professor talked us.

The Stanford Sentiment Treebank **SST-2** consists of sentences from movie reviews and human annotations of their sentiment. This is a binary classification (positive / negative) task.

label:1

comes from the brave , uninhibited performances

label:0

excruciatingly unfunny and pitifully unromantic

Similarity and paraphrase tasks

The Microsoft Research Paraphrase Corpus **MRPC** is a corpus of sentence pairs automatically extracted from online news sources, with human annotations whether the sentences in the pair are semantically equivalent. Here, an F1-score is taken, because the class-sets are imbalanced.

quality:1

Prosecutors filed a motion informing Lee they intend to seek the death penalty.

He added that prosecutors will seek the death penalty.

quality:0

A former teammate , Carlton Dotson , has been charged with the murder.

His body was found July 25 , and former teammate Carlton Dotson has been charged.

The Quora Question Pairs **QQP** dataset is a collection of question and answer pairs from the website Quora. The task is to check whether a pair of questions are semantically equivalent. Again, the F1-score is taken as a measure because the class-sets are unbalanced.

My Galaxy ace is hang?

Why are the people on Staten Island are racist?

0

Where can I learn to invest in stocks?

How can I learn more about stocks?

1

The Semantic Textual Similarity Benchmark **STS-B** is a corpus of pairs of news headlines, video and image captions. Each pair is annotated with a similarity score from 1 to 5, and the task is to predict these scores. The evaluation is done using Pearson and Spearman correlation, as the determining factor is the relative ranking amongst scores.

The man is riding a horse.

A man is riding on a horse.

5.000

A man is playing a guitar.

A girl is playing a guitar.

2.800

Inference Tasks

The Multi-Genre Natural Language Inference Corpus **MNLI** is a crowd-sourced collection of sentence pairs with textual entailment annotations. For each premise and hypothesis sentence, the task is to predict whether the premise entails the hypothesis, contradicts the hypothesis, or neither. Thus, this is a 3-class classification problem. The matched MNLI version tests on data which is in the same domain as the training set, while the mis-matched MNLI tests on a training set which is cross-domain. The MNLI contains parse-trees, which is the reason we do not display these here.

The Stanford Question Answering Dataset is a question-answering dataset. The authors of GLUE use this dataset, to create the Question answering NLI **QNLI**, by removing the requirements that the model selects the exact answer. On top of that, they also remove the simplifying assumption that the answer is always present in the input and that lexical overlap is a reliable cue. The task is to determine whether the context sentence contains the answer to the question.

Who did the children work beside?

In many cases, men worked from home.

not_entailment

How many alumni does Olin Business School have worldwide?

Olin has a network of more than 16,000 alumni worldwide.

entailment

The Recognizing Textual Entailment **RTE** dataset comes from a series of annual textual entailment challenges. Data is combined from RTE1, RTE2,

RTE3 and RTE5. The classes to be predicted by the model include *neutral*, *contradiction* and *no entailment*.

Oil prices fall back as Yukos oil threat lifted
Oil prices rise.
not_entailment

Money raised from the sale will go into a trust for Hepburn's family.
Proceeds go to Hepburn's family.
entailment

The Winograd NLI **WNLI** dataset uses the original Winograd Schema Challenge dataset [23], which is a reading comprehension task where the reader must read a sentence with a pronoun and select the referent of that pronoun from a list of choices. Sentence pairs are constructed by replacing the ambiguous pronoun with each possible referent. The task is to predict if the sentence with the pronoun substituted is entailed by the original sentence.

Some example sentences include

Bob was playing cards with Adam and was way ahead.
If Adam hadn't had a sudden run of good luck, he would have won.
Adam would have won.
label:0

Mark told Pete many lies about himself, which Pete included in his book.
He should have been more truthful.
Mark should have been more truthful.
label:1

Chapter 3

Related Work

3.1 Structure inside BERT

(How Contextual are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings)

going down the drain of "geometry" of BERT and ELMo.

could also go down the drain of bias (we would prefer to have uniform space over gender etc.)

-¿ does projection into some subspace which has same metric properties perhaps not make it isotropic?

pretty ok summary of what kind of properties we want from word-embeddings... (<https://devopedia.org/word-embedding>)

Especially in Named Entity Recognition (NER), there is a lot of use for static word-embeddings. I guess this is because we need static embeddings which represent the individual clusters?

-¿ Using pooling for some

-¿ Character level operation

- ¿ Perhaps make good sense to work towards a word-embeddings where different vectors are close to each other?
- ¿ perhaps find a metric space warping the vectors, s.t. an isotropic representation is achieved?
- ¿ Perhaps tokenization is a big problem, but perhaps other architecture..? but retraining is too difficult.. probably best to just stick to BERT? one way or the other, we need good word-embeddings derived from good language models to form a probabilistic prediction of the concept
- ¿ Could perhaps also try to make an adversarial autoencoder after the BERT layer (or continue training, s.t. a second loss is minimized as a downstream task?)
- ¿ Perhaps distilling with "correct" tokens? i.e. another network which copies BERT, but instead of outputting `##end`, it outputs one of most frequent 20k words
- ¿ thesaurus using a (set of) words. a little like sentence-generation, but generating most-probable examples
- ¿ Everyone just averages token-embeddings..
- ¿ perhaps fitting a GMM to the contextualized representations of BERT may give a good probability space..?
- ¿ Perhaps make sense to apply MUSE to this?
- ¿ Artetxe 2019 uses language models to generate embeddings. we also do this, but do it using 1) better language models, and 2) better
- ¿ QUESTION: Which factors (mapping algo, embedding) is delimiting in automated embedding matching
- ¿ perhaps create a GMM for each concept, based on how many modals we identify? how to estimate the number of clusters? by graph-clustering perhaps! (this could be very consuming)
- ¿ Adversarial autoencoder on BERTs last models to enforce it to some better

distribution

3.1.1 Attention mechanism

Some analysis has been done with looking at the similarity between the produced context-vectors for biLM [38].

3.1.2 From token-vectors to word-vectors

[27] analyse the social biases that are part of the context embeddings and show that depending on the language model (ELMo, BERT), the amount of social bias deviates. Because context embeddings capture more than just semantics, bias is a natural implication of context vectors. A by-product of their research includes aggregating token-embeddings (i.e. aggregating the tokens *hav* and *###ing* to result at the word *having*). Although not very extensive on the topic of aggregation, the authors use techniques of mean-pooling, max-pooling, and last-pooling to determine arrive at a single context-vector, if the given token is intrinsically split-up by the language-model tokenizer (incl. BERT, ELMo, GPT).

Finally, evaluation on the corpora also yield inclusion of human-like biases [19] by names, races, male-vs-female.

3.1.3 Bias in BERT vectors

3.1.4 Change of meanings over time

[17] analyse how meanings quantitatively using corpora from different historical epochs, including the 1890s, 1940s, 1960s, 1970s, 1990s and 2000s. Specifically, they measure the similarity scores for given context-vectors based on differently trained corpora. They demonstrate through examples of the word *gay* and *alien* that the similarity between context-vectors change.

3.1.5 Clinical concept extration

The task of concept-extration is heavily applied in the field in the clinical domain.

[58] trains a LSTM-Conditional-Random-Field to extract concepts. They used annotated corpora from doctors notes to train the model and model this as a named-entity-recognition task.

3.1.6 Discovering for semantics

[7] apply a co-clustering framework that discovery multiple semantic and visual senses of a given noun phrase. They use an structure-EM-like algorithm to achive this.

3.1.7 Embeddings for translation

A lot of work is done in the field of analysing embeddings for translation tasks, to further mitigate the black-box behvaior of neural network.

This includes [21] who use singular value canonical correlation analysis to compare hidden representations of language models between different languages. To arrive at a sentence-representation, they do mean-pooling over the outputted embedding-vectors of a sentence. They find interesting behavior in the arrangement of context-embedding in the Transfoerm-Big architecture proposed in the [52] paper, which is also inherently used in BERT and GPT. Interestingly, they show that different languages do not project into similar spaces, thus making zero-shot learning tasks between languages difficult.

[22] uses a cyclic loss to apply unsupervised machine translation suign monolingual corpora only. Here, sentence-embeddings are learned over time, which follow the cyclic loss constraint, and minimize the sunsupervised translation loss.

Other work tries to disentangle the different linguisitc, visual and audio-based features by supplying multi-modal input at once [25]. A translation model

As is also mentioned in [32], the proximity relation can also be captured through *relation triplets*

$$R = \{(x_i, x_j, j_l) : x_i \text{ is more similar to } x_j \text{ than } x_l\} \quad (3.3)$$

$$(3.4)$$

The general notion of a distance in the euclidean space can be defined as

$$d_E(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j)} \quad (3.5)$$

or equivalent as the l_2 -norm

$$d_E(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_2 \quad (3.6)$$

Most distance learning techniques propose different ways of learning a Mahalanobis distance [26]

$$d_M(\mathbf{x}_i - \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{A}^{-1} (\mathbf{x}_i - \mathbf{x}_j)} \quad (3.7)$$

Because calculating d^2 is computationally simpler, and d^2 maintains most mathematical properties as calculating d , often d^2 is used for downstream calculations (i.e. optimization, distance measure), rather than d .

In the case of the Mahalanobis distance, A is a linear operator such as a matrix. However, A can also be generalized to a nonlinear operator by using the following formulation of the distance metric

$$d_M^2(x_i - x_j) = (x_i - x_j)^\top A^{-1} (x_i - x_j) \quad (3.8)$$

$$d_M^2(x_i - x_j) = (L(x_i - x_j))^\top (L(x_i - x_j)) \quad (3.9)$$

$$d_M(x_i - x_j) = \|L(x_i - x_j)\|_2 \quad (3.10)$$

$$d_M(x_i - x_j) = \|Lx_i - Lx_j\|_2 \quad (3.11)$$

where $A = L^\top L$ and L can now be any function $L \in \mathbf{R}^n \rightarrow \mathbf{R}^d$. Whenever, $d < n$, the data is projected to a lower dimensional space d , which results in dimensionality reduction, which however breaks the convexity property if L is not invertible. In comparison to a dimensionality reduction such as PCA, it has been observed that learning a similarity measure can sometimes be more effective [6].

Specifically, the data is projected into another. The authors in [32] argue that due to this property of measuring similarity between pairs, the metric learning models are able to generalize better across classes.

Non-Deep Metric Learning

Although we focus on deep metric learning during our experiments, due to the additional flexibility that deep neural networks provide, a good understanding of the underlying mechanisms is useful.

[50] provides a good introduction into the literature of metric learning algorithms.

We first define by a proper definition of a distance

Definition 1 *Let X be a non empty set. A distance or metric over X is a map $d : X \times X \rightarrow \mathbb{R}$ that satisfies the following properties:*

1. *Coincidence:* $d(x, y) = 0 \iff x = y$, for all $x, y \in X$
2. *Symmetry:* $d(x, y) = d(y, x)$ for all $x, y \in X$
3. *Triangle inequality:* $d(x, z) \leq d(x, y) + d(y, z)$ for all $x, y \in X$

The ordered pair (X, d) is called a metric space. We will be assuming a euclidean metric space unless otherwise specified.

Because the coincidence property is not always important to us, but because word-embeddings are more interested in measuring relative distances, we will also consider mappings konwn as *pseudoinstances* which relax the coincidence property to merely fulfill $d(x, x) = 0$. In the d dimensional euclidean space, the set of *Mahalanobis distances*, which is parametrized by positive semidefinite matrices are useful.

Using above similarity, dissimilarity and relation triplets, an optimal distance d from a set of distances \mathcal{D} can be found

$$\min_{d \in \mathcal{D}} l(d, S, D, R) \quad (3.12)$$

where l is one of the loss metrics shown below.

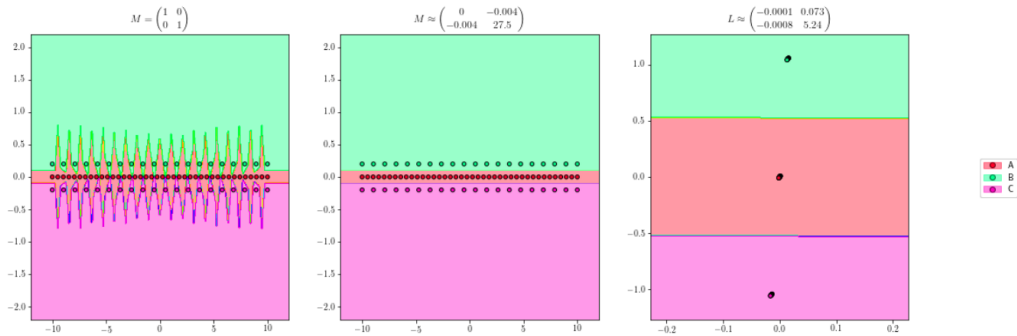


Figure 3.2: Taken from [50]. The individual tiles show the kNN prediction regions by color for every point in the image. Using the unmodified euclidean distancee, this would result in classification regions on the left. The reader can see, learning an appropriate distance, the classification is much more effective (middle). Finally, the dimensionality is also reducable with this dimension while still matching the classification accuracy (right).

Although we will not go into too much detail, we will outline some algorithms that [50] mentions in their work.

We will first start with dimensionality reduction techniques, then move to

nearest neighbors classifiers, follow with nearest centroids classifiers and finally go over some information theory based algorithms. Notice that for almost each one of these methods, kernelized versions are available.

Dimensionality Reduction techniques include principal component analysis *PCA*, LDA ANMM

Deep Metric Learning

Besides the non-deep proposed in the work above models , in[20] also introduces the deep generalization to these frameworks.

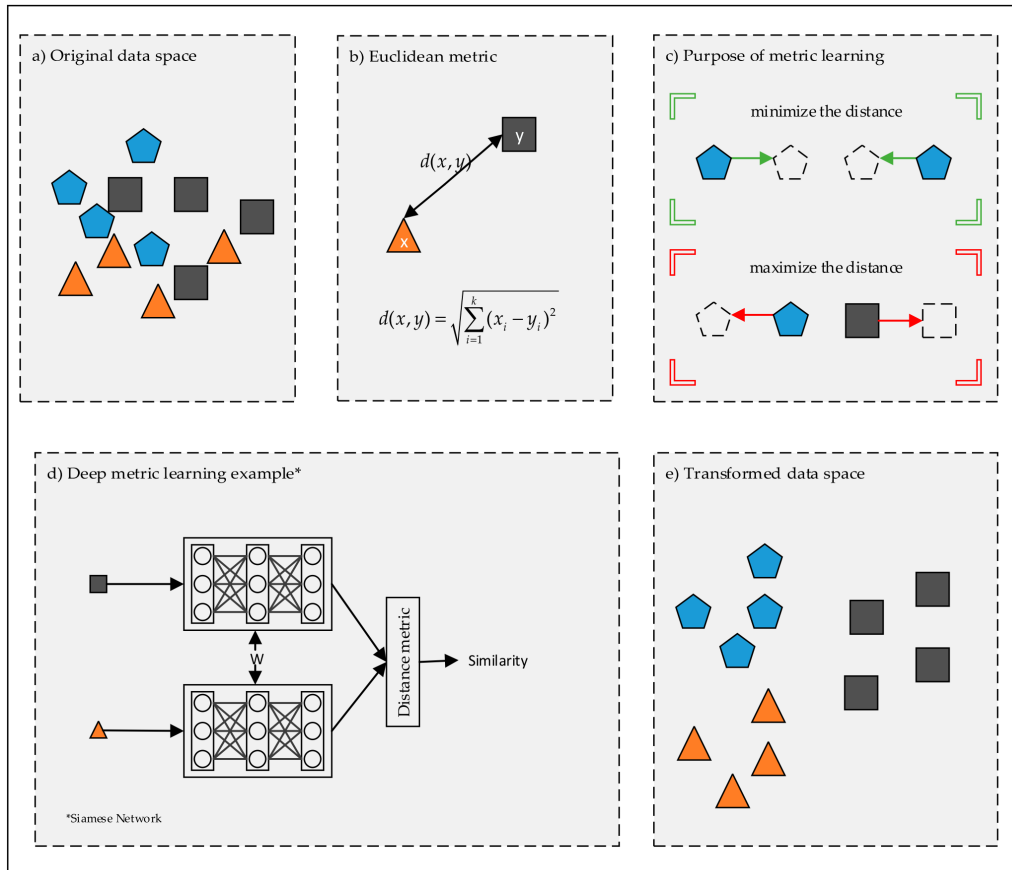


Figure 3.3: Taken from [20]. An illustration of deep metric learning. The space is transformed in such a way, that similar object are closer to each other, and dissimilar objects are moved away from each other.

The most prominent model is the siamese network [5], [8], [13], where the network consists of two identical sub-networks joined at their outputs. We provide an example the architecture looks that looks as follows:

```
class Siamese:

    def __init__(input_dim, latent_dim):
        self.fully_connected = nn.Linear(input_dim, latent_dim)

    def forward(input1, input2):
        latent_1 = self.fully_connected(input1)
        latent_2 = self.fully_connected(input2)

        distance = torch.sqrt(
            torch.sum((latent_1 - latent_2) * (latent_1 - latent_2))
        )

        return distance, latent_1, latent_2
```

Two samples are ingested by the neural network simultaneously. Both samples are passed through identical weights. Training occurs by minimizing the contrastive loss.

$$L_{\text{Contrastive}} = (Y)\frac{1}{2}(D_W)^2 + (1 - Y)\frac{1}{2}\{\max(0, m - D_W)\}^2 \quad (3.13)$$

where we follow the convention that $Y = 1$ whenever the given sample-pairs are of the same class, and $Y = 0$ whenever the given sample-pairs are from different classes, and D_W is the output of the above neural network. Training can be highly unstable, and as such, the learning rate must be properly set, batches must have a balanced amount of both similar, and dissimilar pairs. Finally, training can fail if a single batch includes multiple domains of data.

For the triplet networks [?] which makes use of relation triplets, including an

input X , a point similar to the input X^p and a point dissimilar to the input X^n , the following loss can be used for minimization.

$$L_{\text{Triplet}} = \max(0, \|G_W(X) - G_W(X^p)\|_2 - \|G_W(X) - G_W(X^n)\|_2 + \alpha) \quad (3.14)$$

where α forms a distance margin, similar to the length of the support vector in SVMs.

As the authors suggest, using a triplet logic results in more stable training and bigger margins between similarity classes. Analogous logic goes for quadrupel loss etc., but it is unsure to what extend these extensions increase performance. This general idea can also be extended to angular distance (using cosine distances), and non-euclidean spaces, however, we will not go into further detail on this as this is not the focus of this work.

Although we will not go into further detail, other loss functions include histogram loss [51], the structured loss [49] n-pair loss [47], magnet loss [43], angular loss [56], quadruple loss [33], clustering loss [48], hierarchical triplet loss [12] and the multi-similarity loss [57].

3.3 Clustering Algorithms

We will later on use different clustering algorithms to evaluate how well semantics is interpretable inside BERT vectors. We will give a short introduction on the following clustering algorithms.

- Affinity Propagation
- Chinese Whispers
- DBScan
- HDBScan
- MeanShift

- Optics

3.4 Applications of word vector

Although word-vectors can be used for a variety of tasks, we will focus on some varieties of how these vectors are used for machine translation (of human languages) tasks.

Amongst the most prominent method is the MUSE model introduced by [9]. Here, a mapping W is found using either an unsupervised methodology by minimizing the batch-size cosine-distance between sampled word-vectors between two languages A and B , or a supervised methodology is devised using the procrustes algorithm.

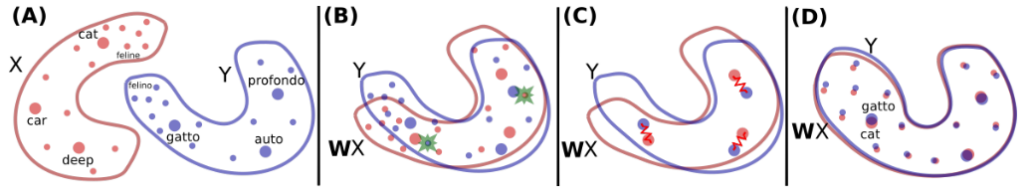


Figure 3.4: Taken from [9]. Each

Specifically, the following loss-function optimizes over the space of possible mapping matrices W^* :

$$W^* = \operatorname{argmin}_{W \in O_d(\mathbb{R})} \|WX - Y\|_F = UV^T, \text{ with } U\Sigma V^T = \operatorname{SVD}(YX^T) \quad (3.15)$$

Unsupervised training is conducted using a discriminator and a mapping function, where the discriminator's task is to identify the origin-distribution, and the mapping function is supposed to fool the discriminator, resulting in a mapping function which maps all word-embeddings into a common global embedding space.

A generalization of these word vectors to point-vector probability-distributions

are captured by [2]. Here, the point-vectors are interpreted as point-delta-distributions in the $d + 1$ dimensional space. Common optimal transport methods are coupled with the Gromov-Wasserstein loss to find an orthogonal mapping which maps from vector-space X to vector space Y , and thus from the word-token of one language to another. One way to achieve this is to use the supervised loss-metric and solve the procrustes problem, finding correspondences between the columns of X and columns of Y through:

$$\min_{\mathbf{P} \in O(n)} \|\mathbf{X} - \mathbf{P}\mathbf{Y}\|_F^2 \quad (3.16)$$

with $O(n) = \{\mathbf{P} \in \mathbb{R}^{n \times n} | \mathbf{P}^\top \mathbf{P} = \mathbf{I}\}$. The resulting algorithm achieves similar performance to MUSE while requiring only a fraction of the computational cost, being able to get competitive results on CPU.

The unsupervised approach deals with finding a transportation map T that fulfills

$$\inf_T \left\{ \int_{\mathcal{X}} c(\mathbf{x}, T(\mathbf{x})) d\mu(\mathbf{x}) | T_{\#}\mu = \nu \right\} \quad (3.17)$$

,

where μ and ν are the point-delta distributions describing the word-embeddings in the probability space. The relaxed Kantorovich's formulation is then finally used to reduce this problem to finding a set of transportation plans in a polytopes.

$$\Pi(\mathbf{p}, \mathbf{q}) = \{\Gamma \in \mathbb{R}_+^{n \times m} | \Gamma \mathbb{1}_n = \mathbf{p}, \Gamma^\top \mathbf{1}_m = \mathbf{q}\}$$

Finally, the cost function

$$\min_{\Gamma \in \Pi(\mathbf{p}, \mathbf{q})} \langle \Gamma, \mathbf{C} \rangle$$

is minimized to find an optimal transportation plan.

Another extension is the work by

Word2Vec

BERT conditions on the rest of the input sentence.

BERT uses words, subwords and individual characters (in total 30'000) that are then used as tokens.

Idea is to do the following: Concepts (and thus words), are represented across multiple contexts. We can create probabilistic word-embeddings by sampling from a corpus the context of a specific word. From multiple samples of the context-embedding-vector, we take the mean and stddev, or create a GMM if there are multimodal logic (we can check this multimodality by running some sort of rejection sampling algorithm). Then we have a probability density function (one for each language), and can map one into another.

Perhaps we could split up too high-variance embeddings to multimodal embeddings, depending on their use-cases.

This allows for interpretability in polysemy sentences.

Not using more complex flows implies that the flow itself is not the bottleneck (they probably tried out other flows as well).

Are the individual word-embeddings going to be one big blob with high variance, or is it going to be a multi-modal distribution...?

Another task we may look at is, from a given word-embedding, sample it's context. Not entirely sure how to do this with BERT and co.

At what point do we overfit, and at what point do we generalize?

Artetxe bilingual token matching through unsupervised machine translation

- Input is cross-lingual word-embeddings - Build an unsupervised phrase-based statistical machine translation system - Derive phrase-based embeddings from input-word-embeddings by taking top 400'000 bigrams and 400'000 trigrams - take arithmetic mean of word-embedding - score top 100 close phrases using softmax cosine similarity - generation of synthetic parallel cor-

pus using this approach - Then use FastAlign to use parallel corpus to align words -

Chapter 4

Analysing the current state of the art

Upadhyay et al. argue that the choice of data is more important than the actual algorithm.

Definitely also look into this, [Analyzing the Limitations of Cross-lingual Word Embedding Mappings] seems to be an analysis of the difficulties etc.

4.1 On the Linear Separability of meaning within sampled BERT vectors

4.1.1 Motivation

To see if there is any structure within BERT vectors w.r.t. the different meaning of one word, we ask ourselves whether or not different part of the meaning are at different locations of the embedding space produced by BERT.

4.1.2 Experiment setup

Let us regard a word w in sentence s is indexed by i . When passed through the BERT model, the word w produces an embedding x_w . When we sample

BERT embeddings for a single word w for $n = 500$ sentences.

Let us restrict the choice of words w on polysemous and ambiguous words which carry multiple meanings. Specifically, w is chosen in such a way that each meaning has more than 30 samples in each class when chosen from within the SemCor dataset.

Words used to test BERT-sampled vectors for linear separability of lexical features		
was	is	be
are	more	one
first	only	time

When we sample Specifically, the experiment setup looks as follows.

Algorithm 1: Checks sampled BERT vectors for linear interpretability by meaning

Input: A target word w_{target} ; The latent dimensionality k for PCA;

Result: Accuracy of a logistic regression classifier

$\mathbf{D}, \mathbf{y} \leftarrow$ sample up to 500 sentences (as much as available) from the SemCor corpus which include the word w_{target} , along with the corresponding WordNet meaning-id;

$\mathbf{X} \leftarrow \text{BERT}(\mathbf{D})$ i.e. pass each sentence through BERT and retrieve the resulting word embedding x_w as defined in the above section;

$\mathbf{X}, \mathbf{y} \leftarrow \text{oversample}(\mathbf{X}, \mathbf{y})$ such that we don't have dominating classes;

$\mathbf{X}_{\text{train}}, \mathbf{X}_{\text{test}}, \mathbf{y}_{\text{train}}, \mathbf{y}_{\text{test}} \leftarrow \text{trainTestSplit}(\mathbf{X}, \mathbf{y}, \text{testProportion} = 0.4)$;

$\mathbf{X}_{\text{train}} \leftarrow \text{StandardScaler}(\mathbf{X}_{\text{train}})$ such that all the data is normalized;

$\mathbf{X}_{\text{train}} \leftarrow \text{PCA}(\mathbf{X}_{\text{train}}, k)$ such that all the data is projected to a lower latent dimensionality k ;

$\text{model} \leftarrow \text{LogisticRegression}(\mathbf{X}_{\text{train}}, \mathbf{y}_{\text{train}})$;

$\hat{\mathbf{y}}_{\text{test}} \leftarrow \text{model.predict}(\mathbf{X}_{\text{test}})$;

$\text{accuracy}, \text{confusionMatrix} \leftarrow \text{loss}(\mathbf{y}_{\text{test}}, \hat{\mathbf{y}}_{\text{test}})$;

return $\text{accuracy}, \text{confusionMatrix}$;

We use a binary

First of all, we sample $n = 500$ sentences from the news corpus for a target

word \hat{w} , which has index i in sentence s .

Each sentence, we pass through the

The BERT model takes a sequence of size (up to) 512 elements as input.

We sample $n = 500$ vectors from BERT. This is the vector at the output embedding layer of BERT, which has the same location as the input to the BERT model.

We apply t-fold cross validation, and measure the mean accuracy as well as the standard deviation of the accuracy. We also note down the variance kept after projecting PCA on the lower dimensionality k . This is the sum of the first k largest normalized eigenvalues.

4.1.3 Results

The number of words we can use to reliably estimate a planar separation of semantics within the sampled BERT vectors is small.

We run the above experiment for a set of different k to build up an intuition of how well different dimensionalities still capture the meaning in different locations of the vector space.

Table 4.1: Mean and standard deviation of the accuracy of a linear classifier trained on the 2 most common classes of WordNet meanings for the word *was*.

dimensionality	variance kept	accuracy (mean / \pm stddev)
10	0.27	0.82/ \pm 0.03
20	0.41	0.81/ \pm 0.04
30	0.50	0.85/ \pm 0.03
50	0.63	0.92/ \pm 0.03
75	0.73	0.94/ \pm 0.02
100	0.81	0.95/ \pm 0.02

Table 4.2: Mean and standard deviation of the accuracy of a linear classifier trained on the 2 most common classes of WordNet meanings for the word *is*.

dimensionality	variance kept	accuracy (mean / \pm stddev)
2	0.09	0.57/ \pm 0.02
10	0.29	0.82/ \pm 0.03
20	0.42	0.82/ \pm 0.04
30	0.51	0.83/ \pm 0.03
50	0.72	0.85/ \pm 0.04
75	0.78	0.84/ \pm 0.04
100	0.79	0.85/ \pm 0.03

Table 4.3: Mean and standard deviation of the accuracy of a linear classifier trained on the 2 most common classes of WordNet meanings for the word *one*.

dimensionality	variance kept	accuracy (mean / \pm stddev)
2	0.10	0.55/ \pm 0.10
3	0.14	0.51/ \pm 0.05
10	0.34	0.59/ \pm 0.08
20	0.50	0.76/ \pm 0.03
30	0.62	0.77/ \pm 0.02
50	0.76	0.83/ \pm 0.06
75	0.87	0.87/ \pm 0.05
100	0.94	0.87/ \pm 0.05

There are many permutations We now also employ multi-class classification.

Table 4.4: Mean and standard deviation of the accuracy of a linear classifier trained on the the 4 most common classes of WordNet meanings for the word *was*.

dimensionality	variance kept	accuracy (mean / \pm stddev)
2	0.08	0.38/ \pm 0.03
3	0.11	0.38/ \pm 0.04
10	0.28	0.65/ \pm 0.03
20	0.43	0.76/ \pm 0.04
30	0.53	0.83/ \pm 0.03
50	0.67	0.93/ \pm 0.01
75	0.77	0.95/ \pm 0.01
100	0.83	0.95/ \pm 0.01

We get very similar accuracies for "time", "made", "thought". However, these tables are left out as we do not deem these to be statistically significant.

4.2 On the Clusterability of meaning within sampled BERT vectors

4.2.1 Motivation

To see if there is any structure within BERT vectors w.r.t. the different meaning of one word, we ask ourselves whether or not different part of the meaning are at different locations of the embedding space produced by BERT.

This is an extension to linear separability. We want to see how obvious these clusterings and differences are.

4.2.2 Experiment setup

The general algorithm to cluster the dataset is shown below.

Algorithm 2: Checks sampled BERT vectors for clusters by meaning

Input: w_{target} : The target word whose BERT vectors we want to cluster;

$DimRed$: The dimensionality reduction method;

k : The latent dimensionality for the dimensionality reduction method;

$ClusterAlgorithm$: The clustering algorithm to use;

Result: The associated cluster-id with each sampled sentence, and the adjusted random index.

$\mathbf{D}, \mathbf{y} \leftarrow$ sample up to 500 sentences from the SemCor corpus which include the word w_{target} , along with the corresponding WordNet meaning-id, and compensate more sentences by sampling from the news corpus if less than 500 sentences are available in the SemCor sentence. We set $y = -1$ whenever no labeling information is available, which is the case if we don't sample data from the SemCor corpus.;

$\mathbf{X} \leftarrow BERT(\mathbf{D})$ i.e. pass each sentence through BERT and retrieve the resulting word embedding x_w as defined in the above section;

$\mathbf{X}, \mathbf{y} \leftarrow oversample(\mathbf{X}, \mathbf{y})$ such that we don't have dominating classes (all except for $y = -1$).;

$\mathbf{X} \leftarrow StandardScaler(\mathbf{X})$ such that all the data is normalized;

$\mathbf{X} \leftarrow DimRed(\mathbf{X}, k)$ such that all the data is projected to a lower latent dimensionality k ;

$model \leftarrow ClusterAlgorithm(\mathbf{X})$;

$\hat{\mathbf{y}} \leftarrow model.predict(\mathbf{X})$;

$score \leftarrow AdjustedRandomIndex(\hat{\mathbf{y}}, \mathbf{y})$;

return $score, \hat{\mathbf{y}}$;

However, because some simple algor.

We had a few constraints. When clustering, we are not given the number of cluster to be found. Thus, the algorithm must solve the multi-modal detection problem intrinsically, which is considered a hard problem in machine learning, as it falls in the same category as global probability density estimation. Also, because we use the lexical distinction of semantics as defined in

WordNet, our algorithm needs to adapt to the granularity that was defined by linguists.

Because we want to evaluate how well our clustering method can mimic the semantic definitions in WordNet, but also generalize to unseen words, we train on an unsupervised dataset. We then test our resulting clustering on a labeled dataset using the random adjusted index [41], [18].

The adjusted random index calculates the overlap and as such the similarity between two clustering assignments. Specifically,

$$ARI = \frac{\sum_{ij} (n_{ij}) - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}} \quad (4.1)$$

where $n_{i,j}$ is the number of items that are present in both clustering assignment, $a_i = \sum_j n_{i,j}$, $b_j = \sum_i n_{i,j}$ for two clustering a and b .

where the resulting score is between $[-1.0, 1.0]$, where a score of 0 implies completely assignment of cluster-labels into random buckets, -1.0 implies a negative correlation, and 1.0 implies perfect similarity. The adjusted random index is adjusted for chance, by taking all possible permutations of the labels that the clustering algorithm can have.

Although we will not go into too much detail with the algorithm descriptions here, we will

The general algorithm to cluster the dataset is shown below.

1. s

model	ARI)
Affinity Propagation	0.001
DBScan	0.000
HDBScan	0.328
MeanShift	0.004
Optics	0.000

4.2.3 Results

Repeating the experiment with 1000 datapoint does not change the results by much.

model	ARI)
Affinity Propagation	0.000
DBScan	0.139
HDBScan	0.271
MeanShift	0.005
Optics	0.070

Repeating the experiment with 1000 items and projecting PCA to 100 results in

model	ARI)
Affinity Propagation	0.000
DBScan	0.215
HDBScan	0.359
MeanShift	0.003
Optics	0.000

We see that including more dimensions make the clustering better w.r.t. the adjusted random score.

Here, we also introduce the chinese whispers algorithm. The chinese whispers algorithm was used to cluster for word-senses in the context of static word embeddings, such as in [34].

We now apply forceful clustering using Bayesian Optimization.

100 latent dimensions, 500 samples

model	ARI)
Affinity Propagation	0.168
Chinese Whispers	0.298
DBScan	0.201
HDBScan	0.242
MeanShift	0.167
Optics	0.167

100 components, 1000 samples

model	ARI)
Affinity Propagation	0.170
Chinese Whispers	0.249
DBScan	0.260
HDBScan	0.234
MeanShift	0.167
Optics	0.197

Projecting to lower dimensions increases accuracy strongly. 20 components, 1000 samples

model	ARI)
Affinity Propagation	0.165
Chinese Whispers	0.349
DBScan	0.167
HDBScan	0.273
MeanShift	0.226
Optics	0.167

Projecting to lower dimensions increases accuracy strongly. and now we add a SEP token 20 components, 1000 samples

We now analyse the embeddings manually

Qualitative evaluation

Given that the above numbers likely don't mean anything, here are some clusters found for the best model configuration (chinese whispers), applied

- - - -

She swooped him up into her arms and kissed him madly

- - - -

I place my son in her arms and I pray that it somehow comforts her
perfect babies ... into the loving arms of middle class+ Americans
which he falls back into her arms like a baby
Sometimes I took it into my arms and felt its surprising heft

- - - -

Figure 4.2: A famous equation

- - - -

and shuttle robotic arms of a solar array and truss

- - - -

a contingent of young arms that will allow us to win now

By and large, those arms remained as fictional as those in "The War ..."
and extensive use of robotic arms operating at their limits
staff of strong young arms that might have tamed the National League East

- - - -

(4.3)

Figure 4.3: A famous equation

$$\begin{array}{c}
\text{this country is so polarized that people spring to arms against any proposal} \\
\text{At least they are carrying arms to protect themselves} \\
\text{His organization issued a call to arms} \\
\text{he shoves it in the faces of his comrades in arms} \\
\text{people who had taken up arms against the United States} \\
\text{Mostly non-Arab rebels took up arms in early 2003}
\end{array}
\tag{4.4}$$

Figure 4.4: A famous equation

$$\begin{array}{c}
\text{The classic years of the arms race, the 1950s and '60s before} \\
\text{that concerns over nuclear arms proliferation in the Middle East} \\
\text{Russian adherence to another arms control treaty} \\
\text{he will press for peace and an eventual arms cut for the states} \\
\text{payments to the companies that supplied arms to Iraq were often delayed} \\
\text{Mr. Safar denies any wrongdoing, including any arms dealings}
\end{array}
\tag{4.5}$$

Figure 4.5: A famous equation

- - - -

leaned back in his chair and, with arms crossed,

- - - -

God, fully in name, is at the bottom with his arms out wide.

If you feel yourself falling, spread your arms

Agamemnon, arms raised ... barely contained violence

Mr. James sat with his arms folded, his head lowered

she felt tears in her eyes and held her arms out in simple joy

- - - -

(4.6)

Figure 4.6: A famous equation

NOW MOVING OVER TO THE BANK EXAMPLE

- - - -

heavy withdrawals from the British bank Northern Rock reignited concern

- - - -

credit card and other consumer loans, forcing the bank to set aside \$3.4 billion

by a mainland Chinese commercial bank in a U.S. bank

Investors fear the bank will be forced to write down

investors hoped that the bank had disclosed the

- - - -

(4.7)

Figure 4.7: A famous equation

$$\begin{array}{c}
\text{-----} \\
\text{I would expect the bank by the trail to the left of the road to} \\
\text{-----} \\
\text{The current slowed and swirled alongside a mud bank where cows had trodden to the water} \\
\text{But the bank had positioned itself well} \\
\text{provide plant scientists and farmers with a bank of genes} \\
\text{Upstairs at the large bank of cashiers} \\
\text{-----}
\end{array}
\tag{4.8}$$

Figure 4.8: A famous equation

$$\begin{array}{c}
\text{-----} \\
\text{of their family's naturalization — bank deposit by bank deposit} \\
\text{-----} \\
12 \text{ that the company might suffer a run on the bank because of mortgage concerns} \\
\text{Third, per several bank managers of major national banks} \\
\text{Local party bosses gained broad powers over state bank lending, taxes} \\
\text{government contracts, and a web of bank accounts} \\
\text{Prince Bandar's Washington bank accounts} \\
\text{-----}
\end{array}
\tag{4.9}$$

Figure 4.9: A famous equation

- - - -
 to lift some of the mystery surrounding the central bank and improve communications with Wall Street
 - - - -
 many economists had predicted that the bank would not cut its rate
 expectations that the central bank will raise interest
 a hint that the central bank plans to hold rates
 China's central bank has stepped up its already huge purchases
 fertilizer prices in African countries, but that the bank itself had often failed to recognize
 - - - -
 (4.10)

Figure 4.10: A famous equation

- - - -
 citing challenges for the investment bank and the potential for an above-average credit burden
 - - - -
 93 percent drop in profits at its investment bank last week
 UBS said it did not expect its investment bank to return to profitability
 defrauded by the investment bank in 1998 when
 said in an interview that the investment bank approached him last month
 said the leaders of Citigroup's investment bank and alternative
 - - - -
 (4.11)

Figure 4.11: A famous equation

EXAMPLES FOR KEY CLUSTERING

$$\begin{array}{c}
\text{-----} \\
\text{Many of the key Arab states} \\
\text{-----} \\
\text{in two months and Australia's key S\&P ASX 200 shed 1.9 per cent} \\
\text{Wall Street rebounded Wednesday after key earnings reports from JPMorgan Chase} \\
\text{The Democratic candidate hires a key strategist} \\
\text{[CLS] Mr. Jones "quickly established a good rapport with key donors"} \\
\text{able to meet the two key officials in the government} \\
\text{-----}
\end{array}
\tag{4.12}$$

Figure 4.12: A famous equation

$$\begin{array}{c}
\text{-----} \\
\text{former president of Trinity College, who played a key role in designing the test} \\
\text{-----} \\
\text{seen in the West as a key yardstick of the fairness of an election} \\
\text{treat ... as a key factor in its decisions about regulatory} \\
\text{which policy makers have called the key test for deciding whether to lower interest rates} \\
\text{9. 11. as a key element in pitch meetings} \\
\text{-----}
\end{array}
\tag{4.13}$$

Figure 4.13: A famous equation

$$\begin{array}{c}
\text{-----} \\
\text{Interstate 5 is a key route connecting Southern and Northern California} \\
\text{-----} \\
\text{A key piece of new functionality for Ops Center} \\
\text{Youssef Squali at Jefferies \& Co. says two key factors are driving the stock up} \\
\text{transforming connection with believers is a key element of evangelical Christianity} \\
\text{What would you say was the key element of your management style that allowed you} \\
\text{is a key indicator of retailer performance} \\
\text{in the West the key players were not a small group of intellectuals reading} \\
\text{-----} \\
(4.14)
\end{array}$$

Figure 4.14: A famous equation

$$\begin{array}{c}
\text{-----} \\
\text{times change and technology advances, the key to the city symbolizes} \\
\text{-----} \\
\text{And an official, five-and-three-quarters-inch-long gold-plated pewter key to prove it} \\
\text{but it is small enough to fit onto a key chain} \\
\text{In it lay three keys on a key chain in the shape of a red} \\
\text{-----} \\
(4.15)
\end{array}$$

Figure 4.15: A famous equation

$$\begin{array}{c}
\text{The Red Sox will now have all their key players from their 2007 championship t} \\
\text{Mike Green scored 12 points and had a key assist in overtime as No. 22 Butler bea} \\
\text{or taking the chance of losing a key player to injury} \\
\text{But they never led, could not get a key basket at crucial times and played like} \\
\text{but Cam Long stole the ball near the top of the key and ran out the clock}
\end{array}
\tag{4.16}$$

Figure 4.16: A famous equation

$$\begin{array}{c}
\text{Three of their key players played more than 4} \\
\text{A key for the Giants on Sunday} \\
\text{chemical reactions on solid surfaces, which are key to understanding questions} \\
\text{but is she part of the conspiracy or the key to Sim's salvation?} \\
\text{whose ability to play on a sprained ankle against the Eagles key to that matchup} \\
\text{Horses have been the lifelong key to satisfying the real femin} \\
\text{Connecticut cornerback said the key to defeating Louisville woul}
\end{array}
\tag{4.17}$$

Figure 4.17: A famous equation

4.3 Correlation between Part of Speech and Context within BERT

4.3.1 Motivation

Because there are more resources in NLP with Part of Speech (PoS) tags, as compared to semantics, we want to analyse to what extent BERT sees similarities between PoS and, and because we assume a strong correlation between PoS and semantics, we analyse to what extent this is visible within BERT vectors.

4.3.2 Experiment setup

We test the hypothesis "semantics implies PoS" by conducting the following experiment. For a chosen target word w_t , we fixate one of the wordnet meanings. We then sample n sentences for the target word w_t where w_t has semantic meaning m in the occurring sentence. After we have sampled all the sentences, we determine the PoS for the target word w_t . We then calculate the percentage occurrence of the majority PoS class and record this as a percentage. If all of the sampled target words w_t for all the sentences have the same assigned PoS tag, then the score results in a value of 1.0. If the dominant PoS tag occurs only half the time, this number decreases to 0.5. Please notice that in this experiment, we only view simple PoS tags (i.e. "noun", "verb", "adjective", "pronoun"), and not the more complex ones listed above.

4.3.3 Results

It is apparent that there is a strong relation between PoS and meaning. Especially "erstarre" Verben are a strong part of this

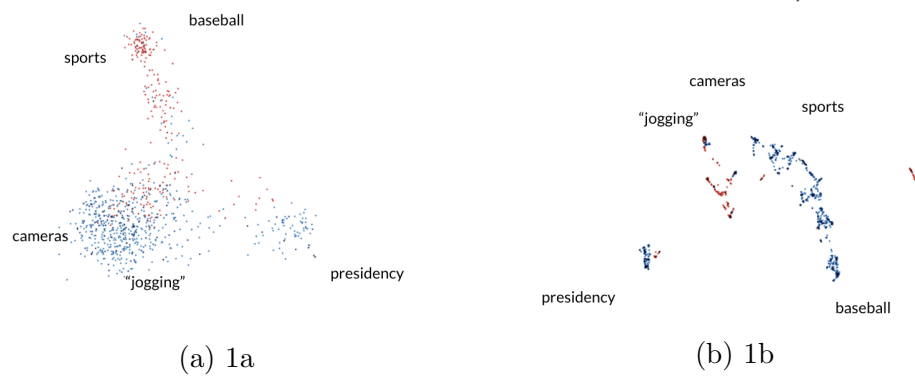


Figure 4.18: plots of....

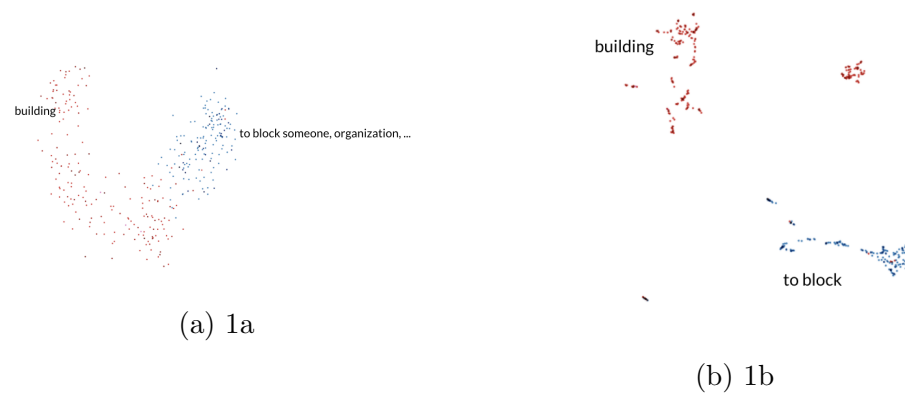


Figure 4.19: plots of....

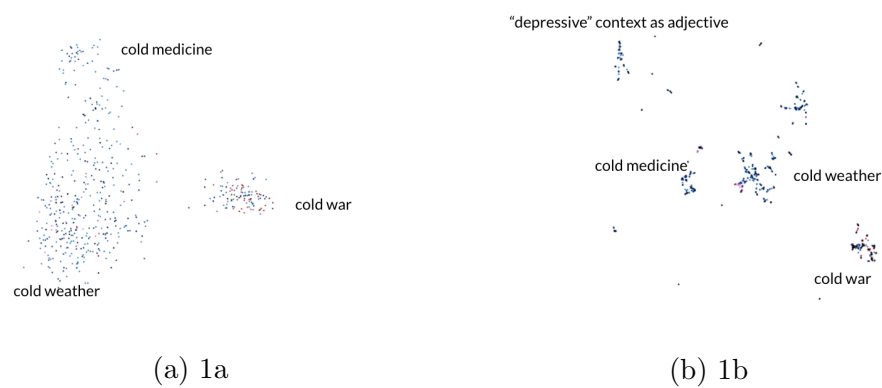


Figure 4.20: plots of....

Chapter 5

Exploiting subspace organization of semantics of BERT embeddings

The general pipeline of how BERT processes sentences is as follows. We have a sentence s which includes a set of words w_1, \dots, w_n . This sequence of words is now given as input to the BERT model. First, this sequence of words is split into a set of tokens t_1, \dots, t_m , where $m \geq n$. Now these tokens, which are in the vocabulary of the BERT tokenizer, are converted to indices, which correspond to index of each individual embedding inside BERT.

We introduce *split-words*, for which we will be generating more specific embeddings. In particular, The idea behind this is that introducing more specialized embeddings for certain tokens will allow to model more complex distributions.

Because the non-modified (vanilla) BERT model uses a certain workflow, we will shortly introduce BERTs pipeline.

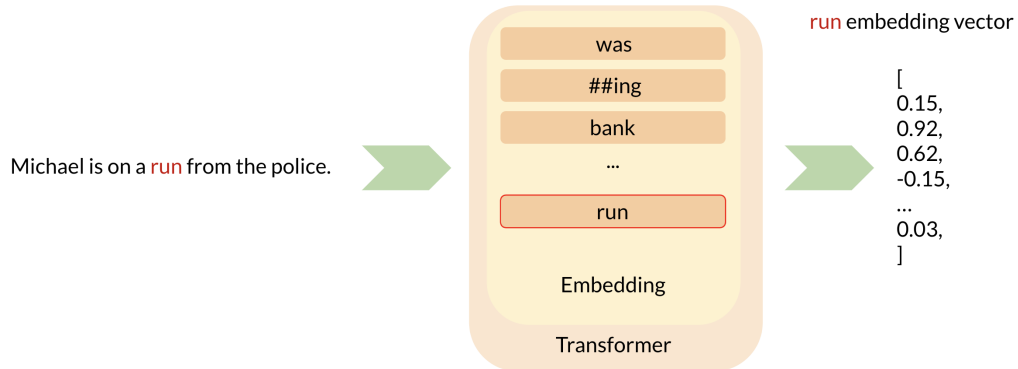


Figure 5.1: The BERT model takes as input a sentence s . The sentence s is converted to a sequence of BERT tokens t_1, \dots, t_m as defined in a given vocabulary V . Each item in the vocabulary V has a corresponding embedding vector inside the embedding layer of the transformer. This embedding vector is used by the intermediate layers of the transformer, and thus affects the downstream pipeline of the transformer for any subsequent layers of the transformer.

5.0.1 BERNie PoS

Motivation

We want to start with a simple model first. Because we have seen in the above section, that there is a strong correlation between semantics and part-of-speech, the initial idea is to introduce additional embedding vectors which are able to capture semantics (rather than just purely word tokens) better. In this case, part-of-speech is used as a proxy for semantics.

Experiment setup

BERnie PoS introduces one new embedding vector for each possible PoS configuration of the split-words.

As an example, instead of having a single embedding vector for the word *run*, we introduce two new embeddings *run_ VERB* and *run_ NOUN*, which both replace the initial *run*-embedding. We will refer to *run_ VERB* and *run_ NOUN* as *sub-embeddings*, as these exploit the subspace structure of the context embeddings sampled for the word *run*. As for the tokenizer, we

also introduce a mechanism which turns any occurrence of *run* into one of the sub-embeddings.

Specifically, our desired pipeline would not look as follows.

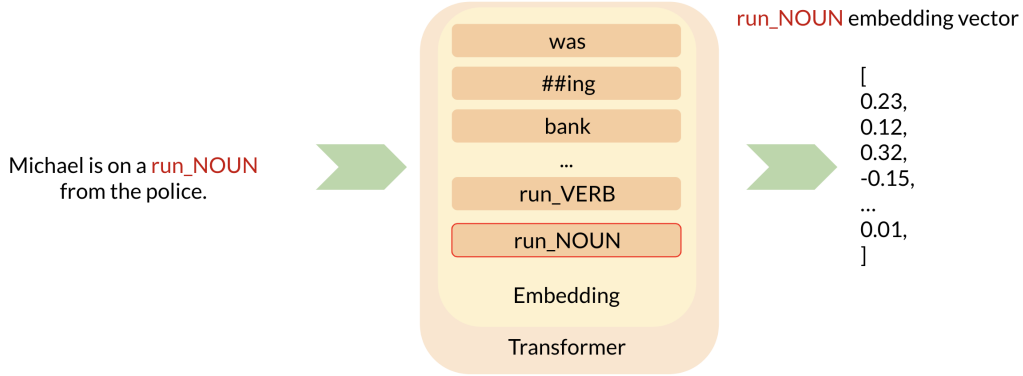


Figure 5.2: The modified pipeline. The BERNie model takes as input a sentence s . The sentence s is converted to a sequence of BERT tokens t_1, \dots, t_m as defined in a given vocabulary V . For each target token t_{target} , we make the token more specific by converting the token to a more specialized token-representation, which specifies the part-of-speech information as part of the token. In this case, *run* becomes *run_VERB*. Again, each item in the vocabulary V has a corresponding embedding vector inside the embedding layer of the transformer. This embedding vector is used by the intermediate layers of the transformer, and thus affects the downstream pipeline of the transformer for any subsequent layers of the transformer.

To identify whether the occurring *run* in an example sentence is a verb or a noun, we use the spaCy part-of-speech tagger [1], which claims 92.6 % accuracy for this task.

5.0.2 BERNie Meaning

Motivation

Similar to the above model, we are introducing new BERT embeddings. This time however, we do not introduce new tokens by part-of-speech, as we did in the previous section, but rather by semantics. For this, we use the a similar methodology as the model from the section 5.0.1. However, instead

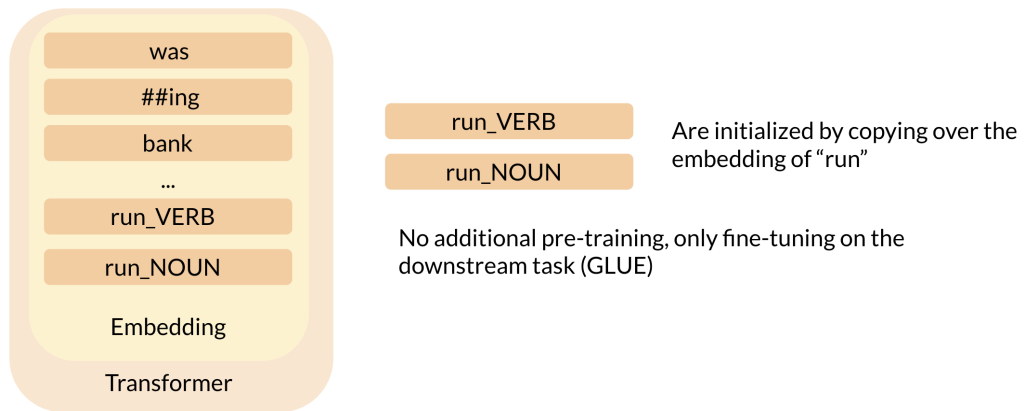


Figure 5.3: Inside the embedding layer of BERT, we introduce more specific embeddings *run_ VERB* and *run_ NOUN*. The BERT model should intuitively now capture more expressiveness, as the model size increased. The original *run* embedding is removed.

of replacing word by their part-of-speech specialization, we replace words by some semantic specialization. The idea here is to adhere more strongly to the subspace structure as seen in section 4.3.

Thus, we introduce the clustering methodology introduced in 4.2.

Experiment setup

Again, we start with the standard BERT model, whose pipeline looks as follows.

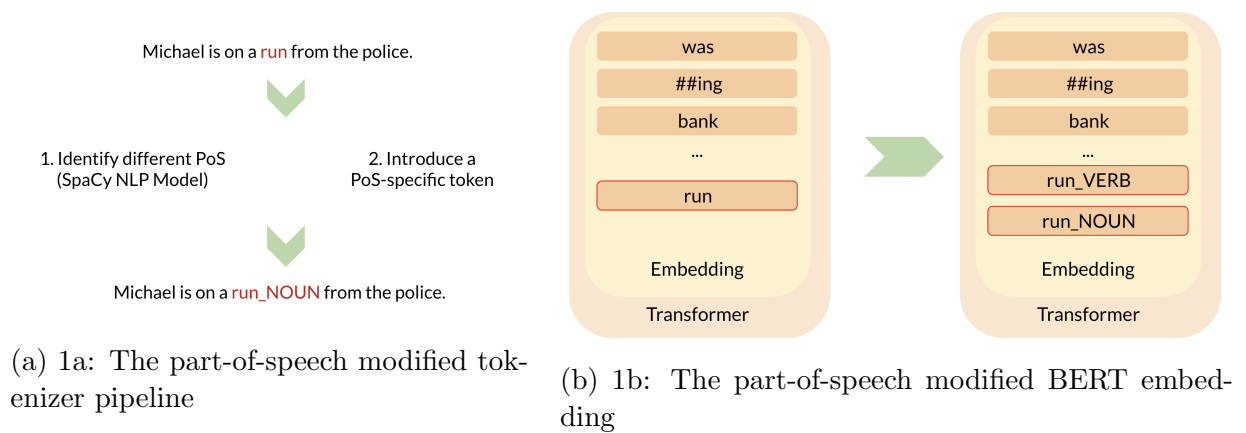


Figure 5.4: plots of...

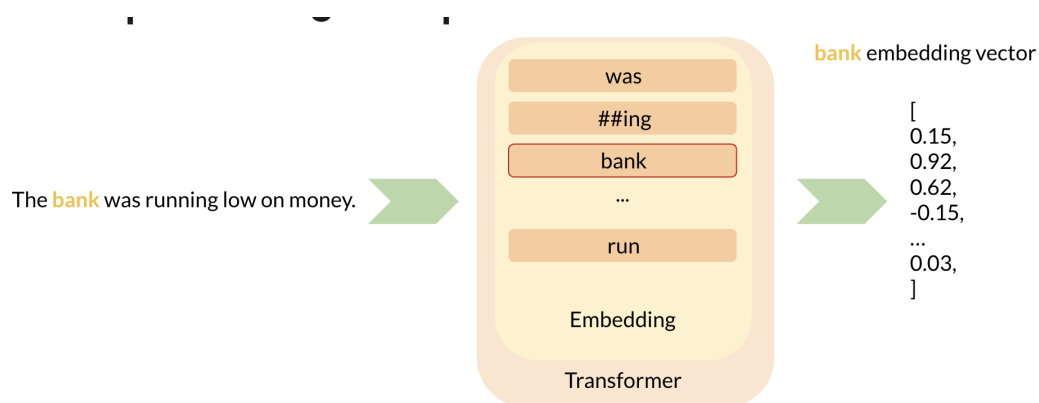


Figure 5.5

We now modify the tokenizer in such a way, that ambiguous and polysemous words are replaced with a more specific token. As an example, we want to replace the word *bank* with a token, which captures whether or not the *bank* that we refer to implies a 1) financial institution, or 2) a river bank.

We use the our clustering approach from 4.2 as the intermediate model to distinguish on which bank the sentence refers to. We also introduce new embedding vectors that the new tokens correspond to.

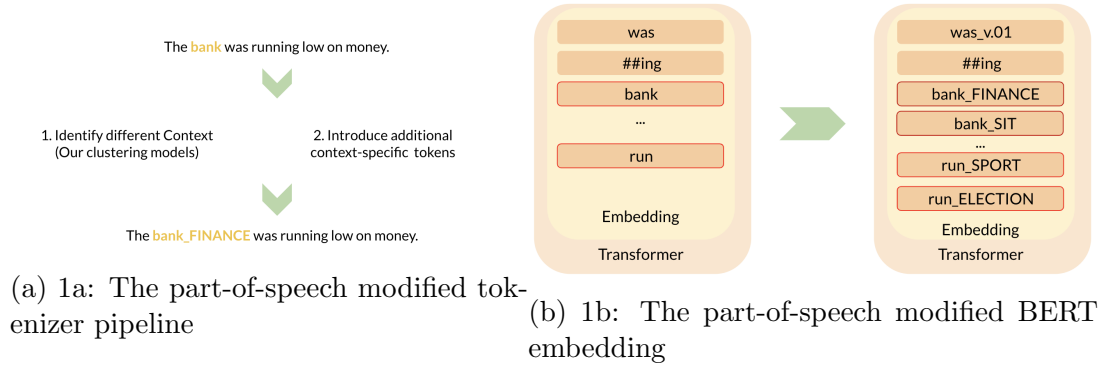


Figure 5.6: plots of....

After we have concluded these changes, we arrive at the following model, which has a modified tokenizer, and a modified BERT model. We call the corresponding model **BERnie**.

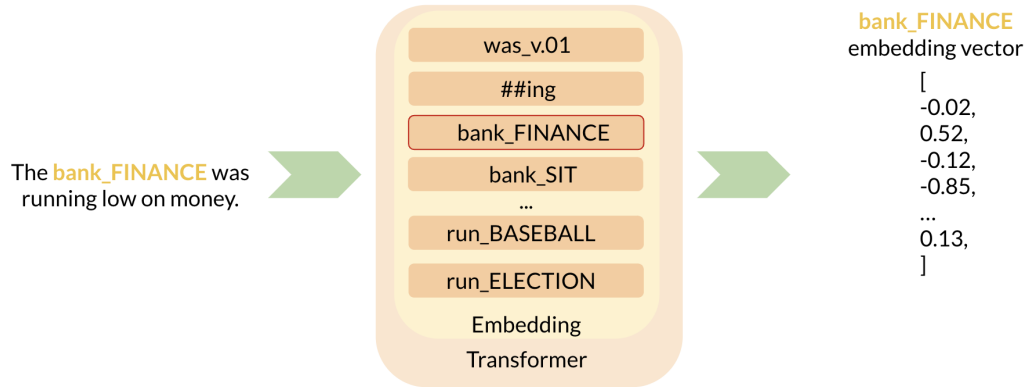


Figure 5.7: The BERT model takes as input a sentence s . The sentence s is converted to a sequence of BERT tokens t_1, \dots, t_m as defined in a given vocabulary V . Each item in the vocabulary V has a corresponding embedding vector inside the embedding layer of the transformer. This embedding vector is used by the intermediate layers of the transformer, and thus affects the downstream pipeline of the transformer for any subsequent layers of the transformer.

Results

We now do a short evaluation of the introduced BERNie PoS and BERNie model. To understand on a wide scale, we use the GLUE benchmarks to

		BERT	BERnie PoS	BERnie
CoLA	Accuracy	0.5739	0.5263	0.5457
MRPC	Accuracy	0.8223	0.8199	0.8064
	F1	0.8778	0.8614	0.8684
	Mixed	0.8501	0.8579	0.8374
SST-2	Accuracy	0.9214	0.9203	0.9266
STS-B	Correlation	0.8841	0.8615	0.8574
	Pearson	0.8860	0.8621	0.8587
	Spearman	0.8822	0.8601	0.8567
QNLI	Accuracy	0.9126	0.9090	0.9020
RTE	Accuracy	0.6462	0.6083	<i>0.5722</i>
WNLI	Accuracy	<i>0.35915</i>	0.3947	0.4649
MNLI	Accuracy	0.8453		0.8334
SNLI	Accuracy	0.8461		0.8367
QQP	Accuracy	0.9113		0.9042
	F1	0.8809		0.8732
	Mixed	0.8961		0.8887

Table 5.1: asj

understand what effect the extension of the embedding layer on points of high variance has.

One can see that in general, the BERnie models performs slightly worse.

All values are the average of two runs. Weights which are instantiated specifically for the GLUE tasks are once instantiated with a random seed of 42, and once with a random seed of 101.

Although most experimental results are similar between standard BERT and the modified BERnie PoS and BERnie models, the WNLI and RTE experiments are considerably deviating in performance through the addition of the additional embedding vectors. For the BERnie model, WNLI performs **10.56%** better than standard BERT, and RTE performs **-7.40%** worse than the standard BERT implementation. Similar results are observable for the difference between the unmodified BERT and the BERnie PoS model with **3.52%** improvement for WLNI and **-3.79%** decline in performance for the RTE task.

Also, compared to BERnie PoS, the standard BERnie model amplifies the

performance-difference to BERT in both the negative and positive direction.

5.0.3 BERNie Meaning with additional pre-training

Motivation

In section 5.0.2, we can see that extending the BERT model generally leads to worse performance.

We assume that the BERNie require additional pre-training, as we are adding more specific embedding vectors, without fine-tuning them. We assume that this is necessary, as this is how the weights of BERT are also trained in the original paper [10] through the masked language model approach.

Experiment setup

We evaluate the models of the previous experiments.

Table 5.2: Mean and standard deviation of the accuracy of a linear classifier trained on the the 4 most common classes of WordNet meanings for the word *was*.

dimensionality	variance kept	accuracy (mean / \pm stddev)
2	0.08	0.38/ \pm 0.03
3	0.11	0.38/ \pm 0.04
10	0.28	0.65/ \pm 0.03
20	0.43	0.76/ \pm 0.04
30	0.53	0.83/ \pm 0.03
50	0.67	0.93/ \pm 0.01
75	0.77	0.95/ \pm 0.01
100	0.83	0.95/ \pm 0.01

WNLI is a comprehension task where a model must read a sentence with a pronoun, and select the referent of that pronoun from a list of choices.

Each one is contingent on contextual information provided by a single word or phrase in the sentence. To convert the problem into sentence pair classification, we construct sentence pairs by replacing the ambiguous pronoun with each possible referent.

To train an ablation study, whether or not additional pre-training improves the instantiated word-vectors, we conduct the following experiment. We instantiate the additional embeddings as described in the previous section. We then run one full epoch of training with the same training parameters as used for GLUE on the news.corpus.2007 dataset using a masked language model methodology . We save this model and load it for the GLUE tasks. BERNie full Pre-training trains the entire embeddings, whereas BERNie partial Pre-Training fixates all embedding vectors except the ones for the ones that were newly added.

Table 5.3: Mean of multiple experiments for BERNie with additionally trained embeddings and weights.

		BERnie	BERnie full Pre	BERnie partial Pre
CoLA	Accuracy	0.5457	0.5418	0.5731
MRPC	Accuracy	0.8064	0.3824	0.3162
	F1	0.8684	0.1720	0.0000
	Mixed	0.8374	0.2775	0.1581
SST-2	Accuracy	0.9266	0.9169	0.9266
QNLI	Accuracy	0.9020	0.5374	0.6700
RTE	Accuracy	0.5722	0.4784	0.4729
WNLI	Accuracy	0.4649	0.4225	0.4507

5.1 Compressing the non-lexical out

Motivation

Experiment setup

We conduct three experiments.

First, we use a simple algorithm to what extent we can disentangle the semantic space within a single sampled word.

The, we use a simple algorithm to what extent we can disentangle the semantic space within multiple sampled words, where we assume that the underlying semantics should again be mutually exclusive (i.e. the sampled words are not independent towards each other).

Finally, we use a simple algorithm to what extent we can disentangle the semantic space within multiple sampled words, where we don't inject any assumptions upon the underlying semantics of the words, i.e. the words are independent towards each other.

Results

Chapter 6

Conclusion

Bibliography

- [1] spacy: Part-of-speech tagging. <https://spacy.io/api/annotation#pos-tagging>. Accessed: 2020-04-13.
- [2] David Alvarez-Melis and Tommi S. Jaakkola. Gromov-wasserstein alignment of word embedding spaces. 2018.
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *arXiv*, 2016.
- [4] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic language model. *Journal of Machine Learning Research*, 3:1137–1155, 2003.
- [5] Jane Bromley, Isabelle Guyon, and Yann LeCun. Signature verification using a "siamese" time delay neural network. 1994.
- [6] Gal Chechik, Varun Sharma, Uri Shalit, and Samy Bengio. Large scale online learning of image similarity through ranking. *JMLR*, 2010.
- [7] Xinlei Chen, Alan Ritter, Abhinav Gupta, and Tom Mitchell. Sense discovery via co-clustering on images and text. 2019.
- [8] Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively, with application to face verification. 2005.
- [9] Alexis Conneau, Guillaume Lample, Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. Word translation without parallel data. 2017.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *Arxiv*, 2018.
- [11] W. N. Francis and H. Kucera. A standard corpus of present-day edited american english, for use with digital computers. 1964.
- [12] Weifeng Ge, Weilin Huang, Dengke Dong, and Matthew R. Scott. Deep metric learning with hierarchical triplet loss. 2018.
- [13] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. 2006.
- [14] Zellig Harris. Distributional structure. *Word*, 10, 1954.
- [15] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Arxiv*, 1997.

- [16] Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017.
- [17] Renfen Hu, Shen Li, and Shichen Liang. Diachronic sense modeling with deep contextualized word embeddings: An ecological view. 2019.
- [18] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of Classification*, pages 193–218, 1985.
- [19] Sophie Jentzsch, Constantin Rothkopf, and Patrick Schramowski Kristian Kersting. On measuring social biases in sentence encoders. *AIES 2019 - Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 37–44, 2019.
- [20] Mahmut Kaya and Hasan Sakir Bilge. Deep metric learning : A survey. *MDPI Symmetry*, 2019.
- [21] Sneha Kudugunta, Ankur Bapna, Isaac Caswell, Naveen Arivazhagan, and Orhan Firat. Investigating multilingual nmt representations at scale. 2018.
- [22] Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. Unsupervised machine translation using monolingual corpora only. *ICLR*, 2018.
- [23] Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In *Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning*, 2012.
- [24] Peter J. Liu, Mohammad Saleh, Etienne Pot†, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences. 2018.
- [25] Shuang Ma, Daniel McDuff, and Yale Song. M3 d-gan: Multi-modal multi-domain translation with universal attention. 2019.
- [26] P. C. Mahalanobis. On the generalized distance in statistics. *Proc. Nat. Inst. Sci.*, 1936.
- [27] Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. On measuring social biases in sentence encoders. 2019.
- [28] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, 2013.
- [29] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. *CoRR*, abs/1310.4546, 2013.
- [30] George A. Miller, Richard Beckwith, Christiane Fellbaum, Derek Gross, and Katherine J. Miller. Introduction to wordnet: An on-line lexical database. *International Journal of Lexicography*, pages 235–244, 1990.
- [31] George A. Miller, Martin Chodorow, Shari Landes, Claudia Leacock, and Robert G. Thomas. Using a semantic concordance for sense identi-

- fication. pages 240–243, 1994.
- [32] Panagiotis Moutafis, Mengjun Leng, and Ioannis A. Kakadiaris. An overview and empirical comparison of distance metric learning methods. *IEEE Transactions on Cybernetics*, 2017.
 - [33] Jiazhi Ni, Jie Liu, Chenxin Zhang, Dan Ye, and Zhirou Ma. Fine-grained patient similarity measuring using deep metric learning. 2017.
 - [34] Maria Pelevina, Nikolay Arefyev, Chris Biemann, and Alexander Panchenko. Making sense of word embeddings. *Journal of Classification*, 2016.
 - [35] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. *EMNLP*, page 1532–1543, 2014.
 - [36] Matthew E. Peters, Waleed Ammar, Chandra Bhagavatula, and Russell Po. Semi-supervised sequence tagging with bidirectional language models. *ACL*, 2017.
 - [37] Matthew E Peters, Mark Neumann, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. *Arxiv*, 2018.
 - [38] Matthew E Peters, Mark Neumann, Luke Zettlemoyer, and Wen-Tau Yih. Dissecting contextual word embeddings: Architecture and representation. 2018.
 - [39] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. *Arxiv*, 2018.
 - [40] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *Arxiv*, 2019.
 - [41] William M. Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical Association*, pages 846–850, 1971.
 - [42] Marek Rei. Semi-supervised multitask learning for sequence labeling. 2017.
 - [43] Oren Rippel, Manohar Paluri, Piotr Dollar, and Lubomir Bourdev. Metric learning with adaptive density discrimination. 2016.
 - [44] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. 2019.
 - [45] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. *arXiv*, 2015.
 - [46] Dinghan Shen, Pengyu Cheng, Dhanasekar Sundararaman, Xinyuan Zhang, Qian Yang, Meng Tang, Asli Celikyilmaz, and Lawrence Carin. Learning compressed sentence representations for on-device text pro-

- cessing. 2019.
- [47] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. 2016.
 - [48] Hyun Oh Song, Stefanie Jegelka, Vivek Rathod, and Kevin Murphy. Deep metric learning via facility location. 2017.
 - [49] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. *CVPR*, 2016.
 - [50] Juan Luis Suarez, Salvador Garcia, and Francisco Herrera. A tutorial on distance metric learning: Mathematical foundations, algorithms and experiments. 2019.
 - [51] Evgeniya Ustinova and Victor Lempitsky. Learning deep embeddings with histogram loss. *NIPS*, 2016.
 - [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Neural Information Processing Systems*, 2017.
 - [53] Luke Vilnis and Andrew McCallum. Word representations via gaussian embedding. *CoRR*, abs/1412.6623, 2014.
 - [54] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. *NeurIPS*, 2019.
 - [55] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *ICLR*, 2019.
 - [56] Jian Wang, Feng Zhou, Shilei Wen, Xiao Liu, and Yuanqing Lin. Deep metric learning with angular loss. 2017.
 - [57] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R. Scott. Multi-similarity loss with general pair weighting for deep metric learning. 2019.
 - [58] Henghui Zhu, Ioannis Ch Paschalidis, and Amir Tahmasebi. Clinical concept extraction with contextual word embedding. 2018.