

Text Analysis

With Embeddings



Text Features and Embeddings

Doc2Vec

BERT



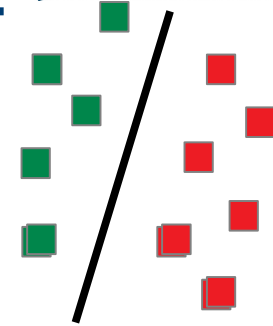
Image by: [displayr.com](https://www.displayr.com)

Motivation

Recall the importance of the feature extraction

In our supervised coffee cup classification:

Function f operates in
feature space



(X, y)

Feature extraction

$\begin{pmatrix} 2.06 \\ 1.76 \\ 0.06 \\ -1.7 \\ \vdots \\ 1.6 \end{pmatrix}$

Feature space

$$X \in \mathbb{R}^n$$

ML

$y': \{-1, 1\}$

Motivation

Recall the importance of the feature extraction

In our supervised coffee cup classification:

- Image as tensor
- Feature extraction is function on tensor



(X, y)

Feature extraction

$\begin{pmatrix} 2.06 \\ 1.76 \\ 0.06 \\ -1.7 \\ \vdots \\ 1.6 \end{pmatrix}$

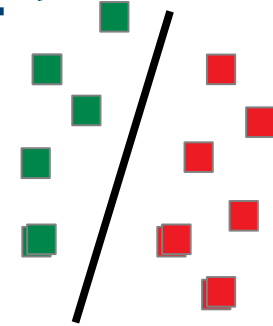
Feature space

$$X \in \mathbb{R}^n$$

ML

$y': \{-1, 1\}$

Function f operates in feature space

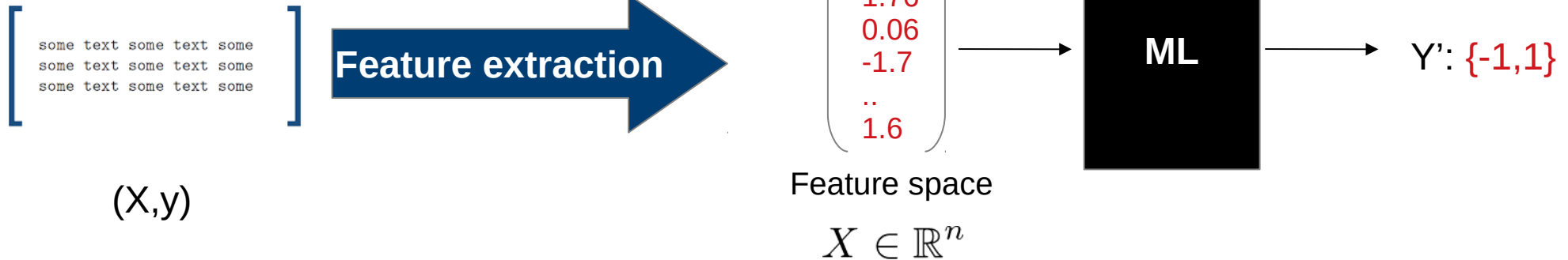
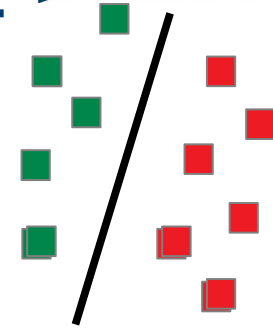


Motivation

How can we deal this non-numeric inputs?

What if we have text?

Function f operates in
feature space



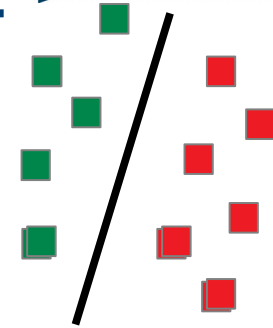
Motivation

How can we deal this non-numeric inputs?

What if we have text?

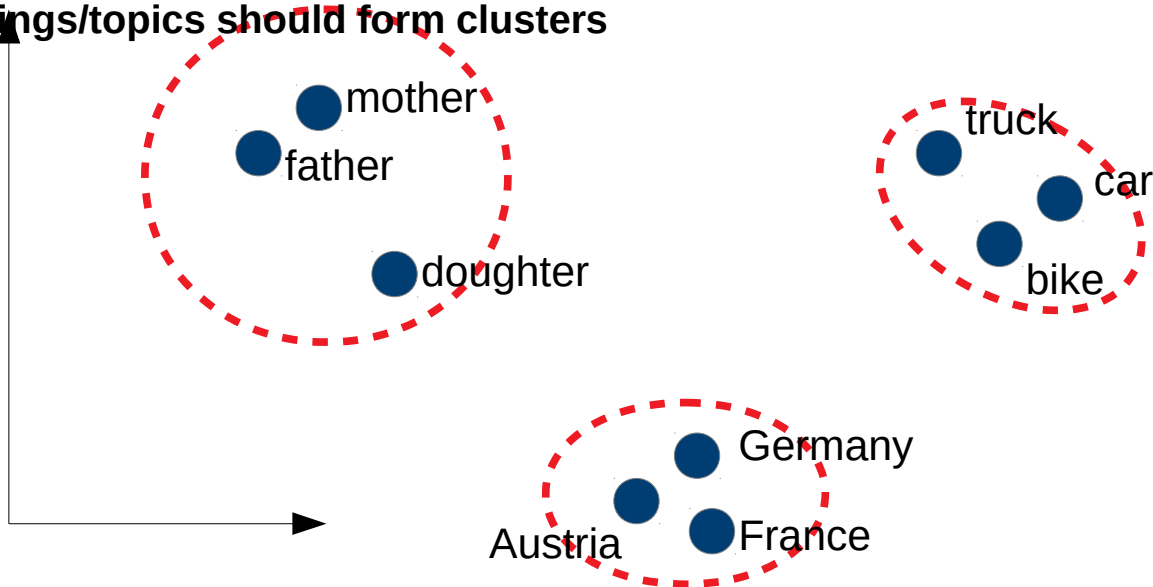
- How can we even define a (feature) function on non-numerical inputs?

Function f operates in feature space



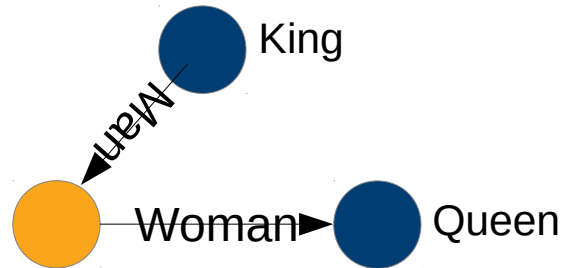
Desired Properties for a “Text Space” :

- **Unsupervised learning of feature extraction**
- Compact, e.g. low dimension
- Similar words should be close
- Categories of things/topics should form clusters



Desired Properties for a “Text Space” :

- Unsupervised learning of feature extraction
- Compact, e.g. low dimension
- Similar words should be close
- Categories of things/topics should form clusters
- Dimensions should encode certain properties
- Basic “Text Calculus” would be very nice, e.g.:

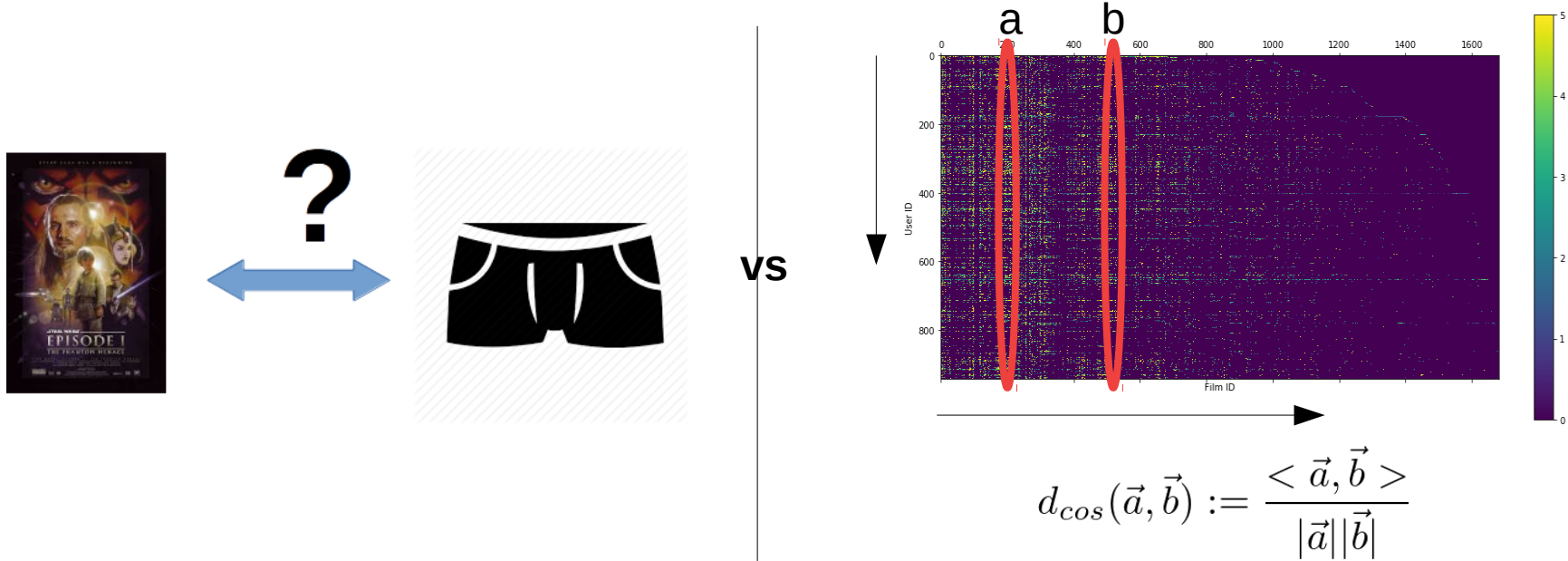


$$\textit{King} - \textit{man} + \textit{woman} = \textit{queen}$$

How to measure complex similarities?

Recall Collaborative Filters

Instead of using complex item-to-item measures, we used the **context** of very simple features, like user ratings:



How to measure complex similarities?

Similarity by co-occurrence

Basic idea: similar things occur in similar context

How to measure complex similarities?

Similarity by co-occurrence

Basic idea: similar things occur in similar context

Text example: which words can we fill in the blank?

“The _____ is climbing on the tree...”

How to measure complex similarities?

Similarity by co-occurrence

Basic idea: similar things occur in similar context

Text example: which words can we fill in the blank?

“The _____ is climbing on the tree...”

More context:

“...His sister is 10 years old”

How to measure complex similarities?

Similarity by co-occurrence

Basic idea: similar things occur in similar context

For Text: similar words have a similar context

youngster
lad
guy

... ▼

*“The **boy** is climbing on the tree...”*

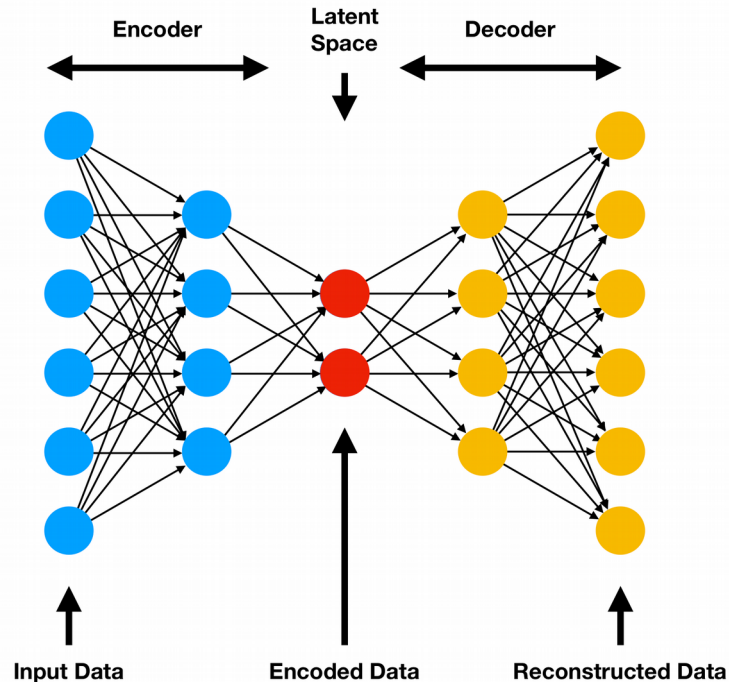
Other words that also fit are **similar**

More context:

“...His sister is 10 years old”

Recall Auto Encoder: Latent Space

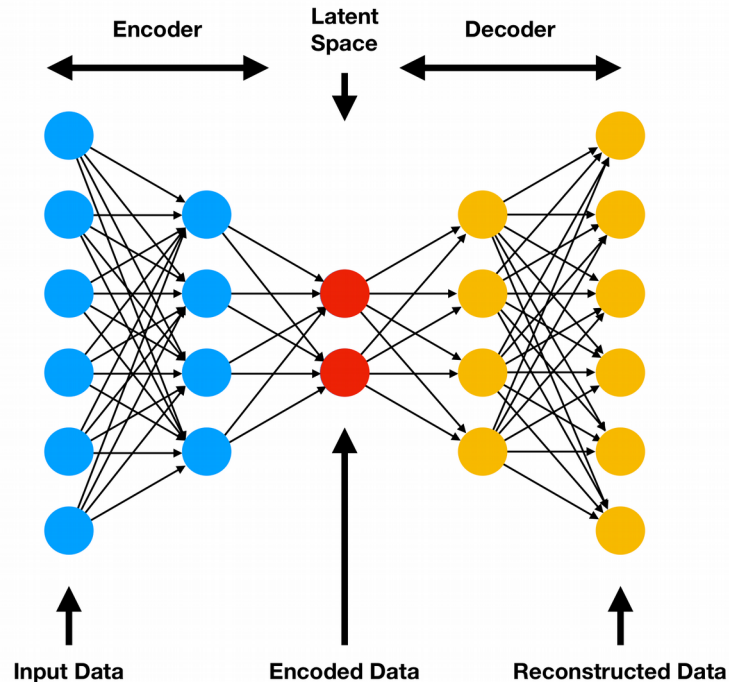
Use context to learn a space which in „embedding“ the words



- Neural Network
- learning a compact latent space representation

Recall Auto Encoder: Latent Space

Use context to learn a space which in „embedding“ the words



- Neural Network
- learning a compact latent space representation
- Main question:

I) What should be the optimization criteria ?

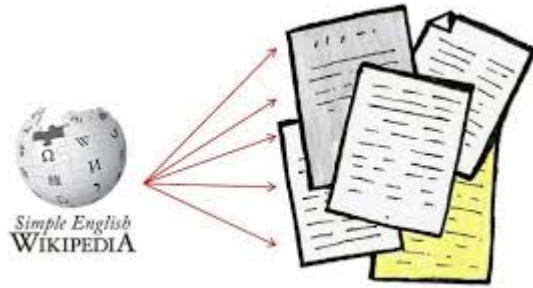
II) How to handle Text input ?

III) Where to get the context information ?

Q3. Where do we get the context from?

Answer: analyze a lot of text!

Open text sources:



Full Wikipedia (available for many languages)

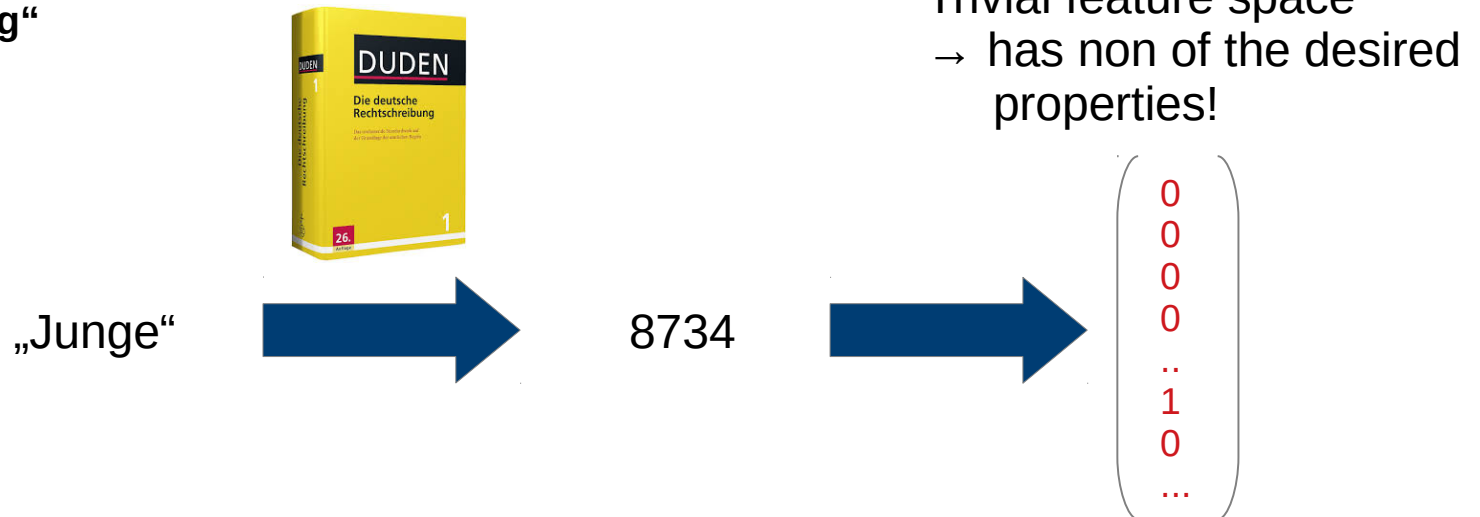
Q2: How to present Words as numbers?

Answer: (initially) just enumerate them → compute embedding

A - Use a dictionary and count all words

B - Assign unique number to all words

C- use „one-hot-encoding“



The Word2Vec Algorithm

Q1: How to meet our context objective?

The famous Word2Vec Algorithm comes in two variants, *CBOW* and *skip-gram*.

The objective in both cases, to use the context of a word to learn the embedding indirectly by constructing a prediction Task:

Efficient Estimation of Word Representations in Vector Space

Tomas Mikolov

Google Inc., Mountain View, CA
tmikolov@google.com

Kai Chen

Google Inc., Mountain View, CA
kaichen@google.com

Greg Corrado

Google Inc., Mountain View, CA
gcorrado@google.com

Jeffrey Dean

Google Inc., Mountain View, CA
jeff@google.com

Abstract

We propose two novel model architectures for computing continuous vector representations of words from very large data sets. The quality of these representations is measured in a word similarity task, and the results are compared to the previously best performing techniques based on different types of neural networks. We observe large improvements in accuracy at much lower computational cost, i.e. it takes less than a day to learn high quality word vectors from a 1.6 billion words data set. Furthermore, we show that these vectors provide state-of-the-art performance on our test set for measuring syntactic and semantic word similarities.

The Word2Vec Algorithm

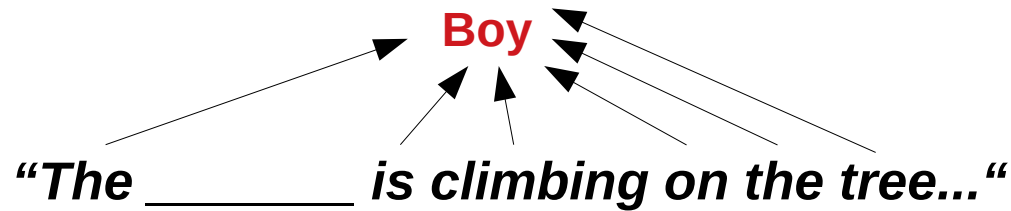
Q1: How to meet our context objective?

The famous Word2Vec Algorithm comes in two variants,
CBOW and *skip-gram*.

The objective in in both cases, to use the context of a word
to learn the embedding indirectly by constructing a prediction

Task: **Continuous bag of words (CBOW)**

→ predict a word from it's context



The Word2Vec Algorithm

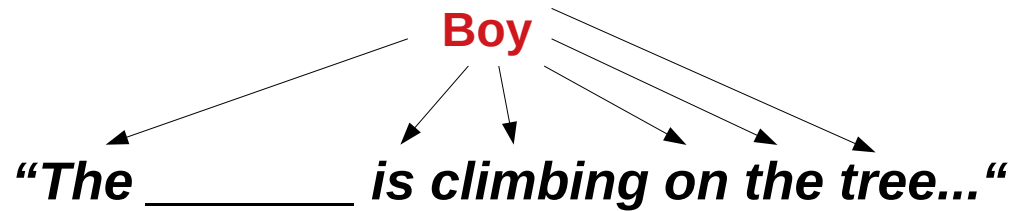
Q1: How to meet our context objective?

The famous Word2Vec Algorithm comes in two variants, *CBOW* and *skip-gram*.

The objective in in both cases, to use the context of a word to learn the embedding indirectly by constructing a prediction

Task: **Skip-gram**

→ predict the context from a word

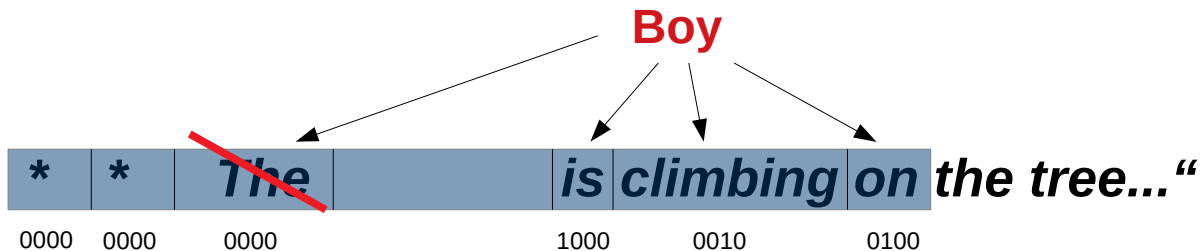


The Word2Vec Algorithm

Context modeling:

In Both cases:

- Words are represented by „one-hot“ encodings of their dictionary numbers
 - Pre-processing: remove fill-words like „a“, „the“, ...
- The context is defined by a window of a fixed size
 - Fill window at beginning and end of text with some placeholder

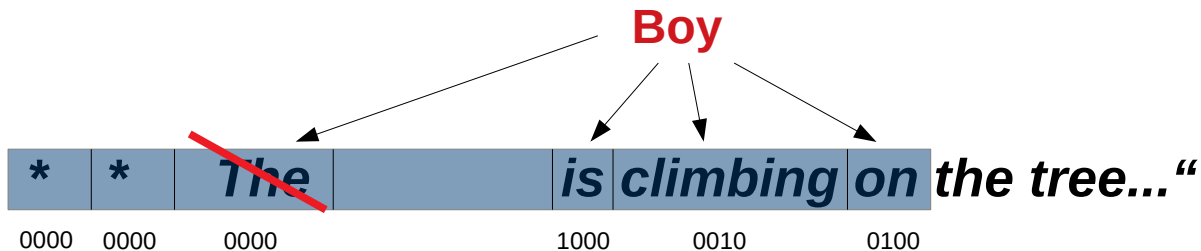


The Word2Vec Algorithm

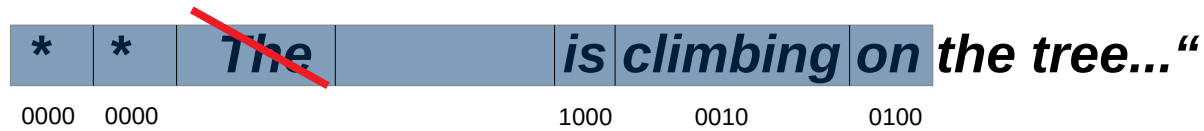
Context modeling:

In Both cases:

- Words are represented by „one-hot“ encodings of their dictionary numbers
 - Pre-processing: remove fill-words like „a“, „the“, ...
- The context is defined by a window of a fixed size
 - Fill window at beginning and end of text with some placeholder
- Use a simple NN classification model for the prediction (next slide)



Neural Network Model for the CBOW case



Input window

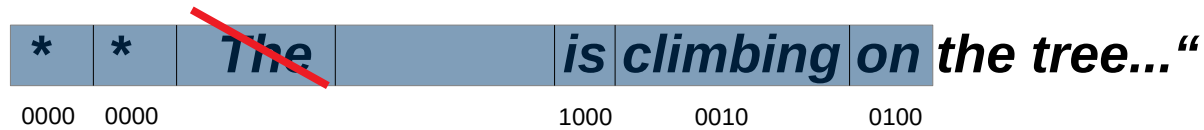


As one-hot matrix

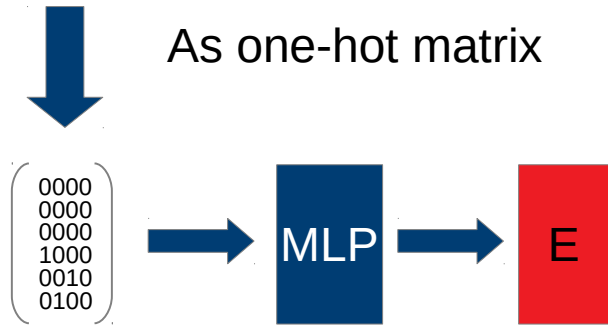
$$\begin{pmatrix} 0000 \\ 0000 \\ 0000 \\ 1000 \\ 0010 \\ 0100 \end{pmatrix}$$

Window size x number of words in vocabulary
→ Up to 100k dimensions

Neural Network Model for the *Skip-Gram* case



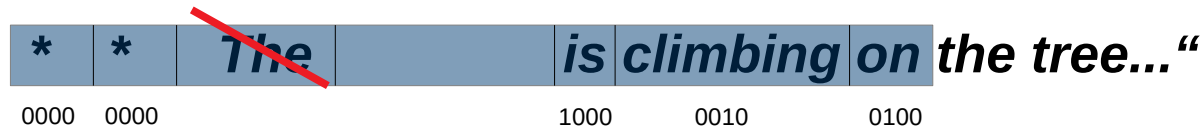
Input window



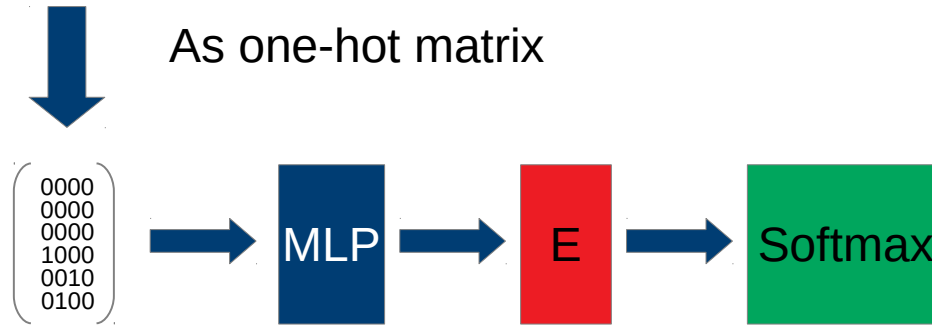
One or more fully connected layers

Output: low dimension **embedding space** → 100 – 300 dims

Neural Network Model for the *Skip-Gram* case



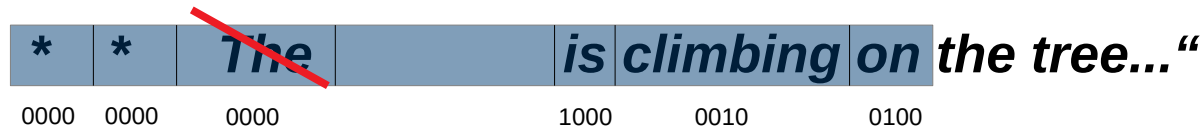
Input window



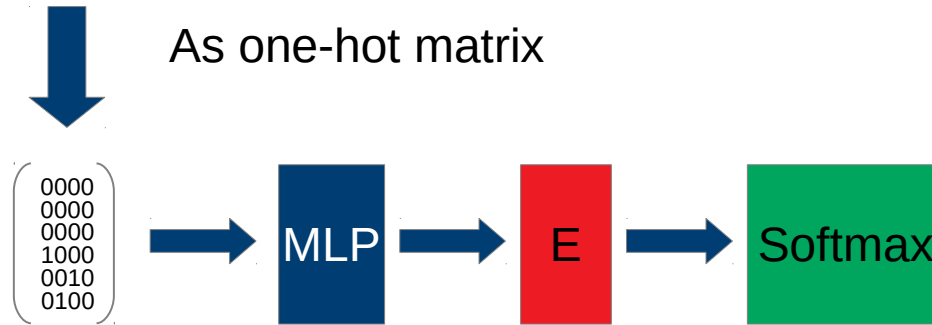
As one-hot matrix

$$L(y^i, y^{i'}) := \frac{e^{y^i y^{i'}}}{\sum_j^k e^{y^i y^{j'}}}$$

Neural Network Model for the *Skip-Gram* case

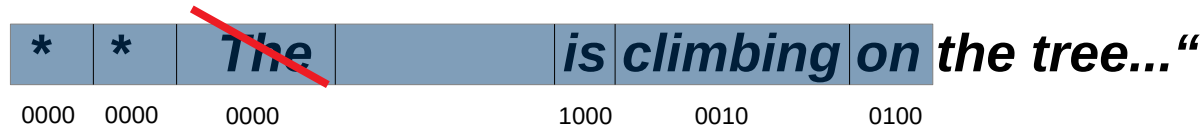


Input window

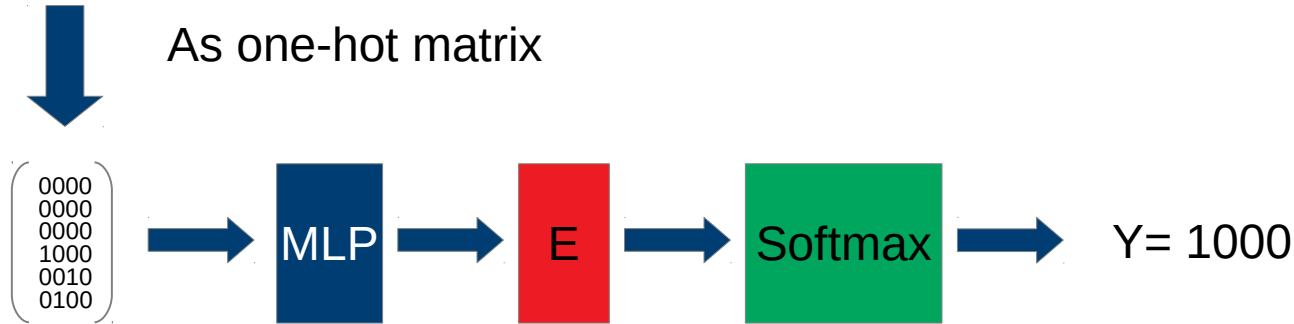


$$L(y^i, y^{i'}) := \frac{e^{y^{i'}}}{\sum_j^k e^{y^{j'}}}$$

Neural Network Model for the *Skip-Gram* case

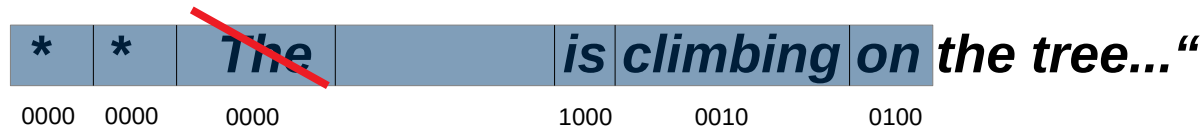


Input window

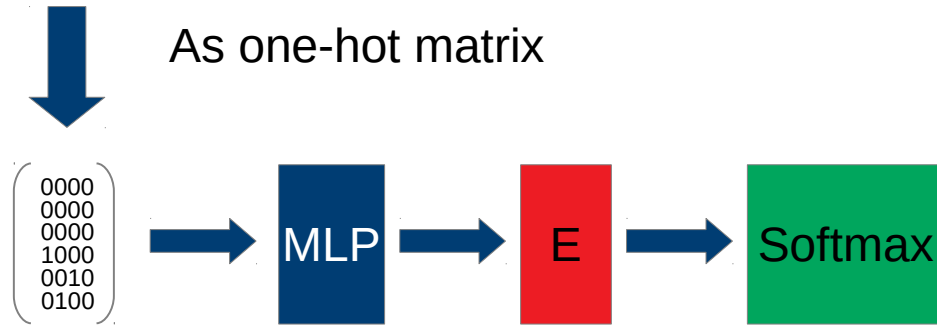


Standard supervised NN training

Neural Network Model for the *Skip-Gram* case



Input window

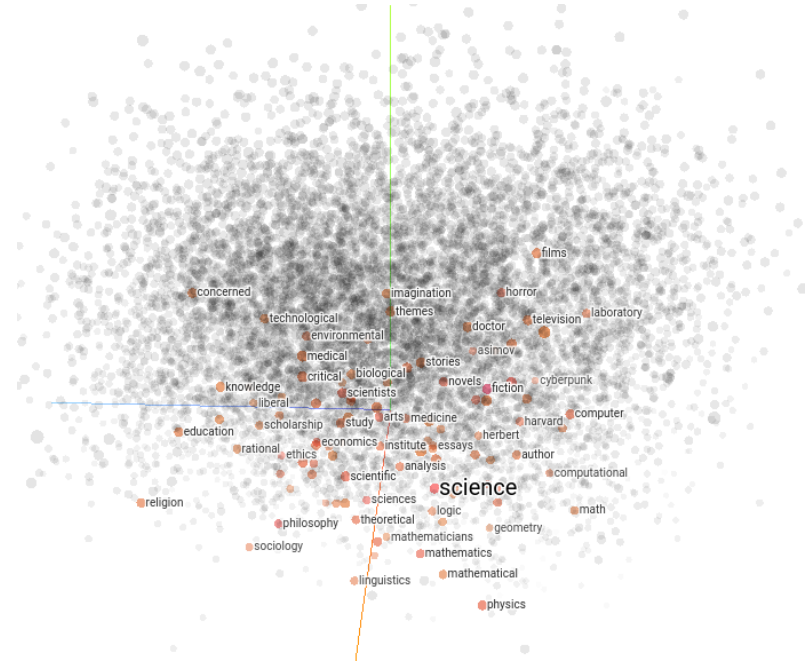


As one-hot matrix

$$L(y^i, y^{i'}) := \frac{e^{y^i y^{i'}}}{\sum_j^k e^{y^i y^{j'}}}$$

Problem: **number** of classes is very large!
→ **solution**: use random sub-set

Show similarity



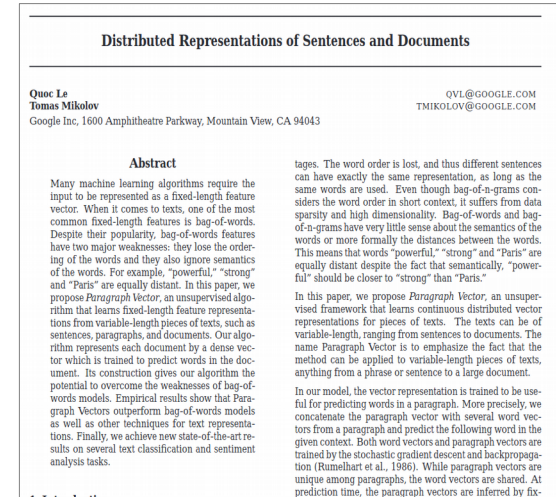
<https://projector.tensorflow.org/>

How to embed text documents – not only words?

The Doc2Vec algorithm is a simple extension of Word2Vec:

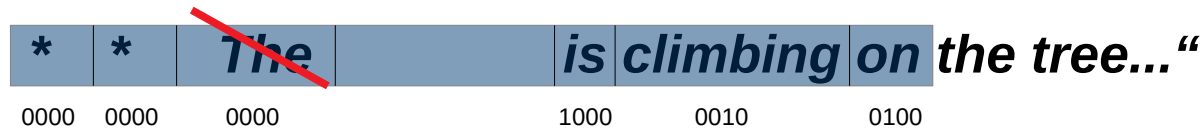


embedding

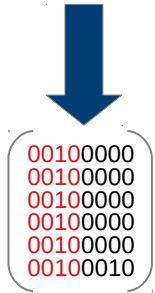


How to embed text documents – not only words?

The Doc2Vec algorithm is a simple extension of Word2Vec:



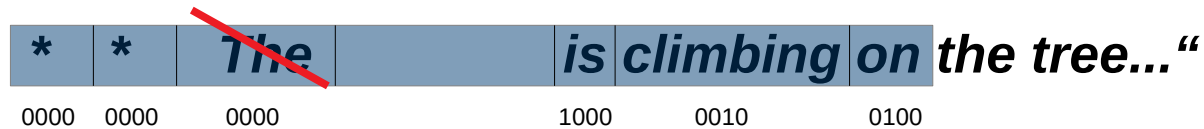
Input window



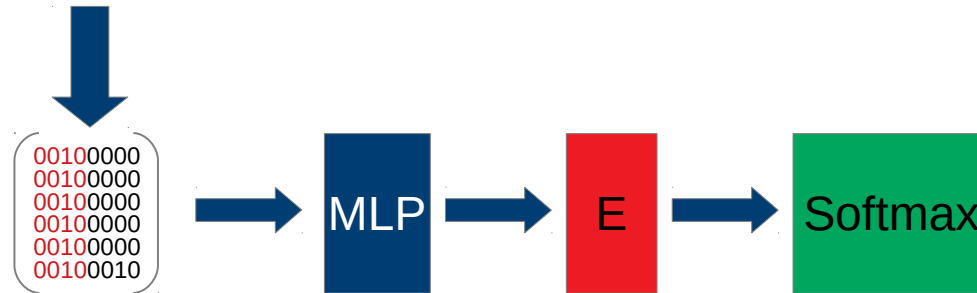
Add document ID

How to embed text documents – not only words?

The Doc2Vec algorithm is a simple extension of Word2Vec:



Input window



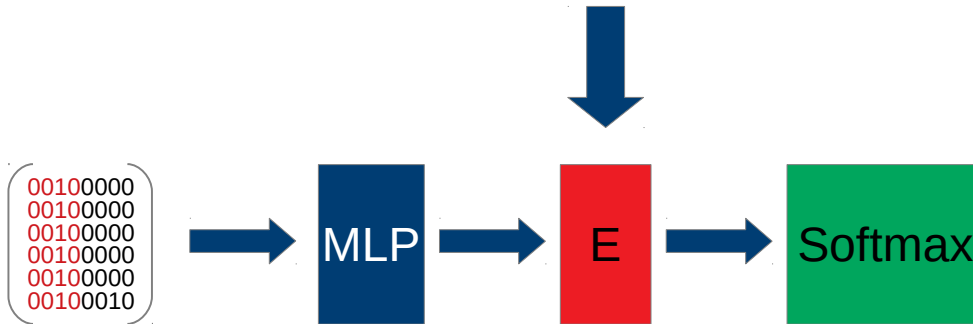
The rest stays the same!

Add document ID

How to embed text documents – not only words?

Step 2: Clustering by document ID

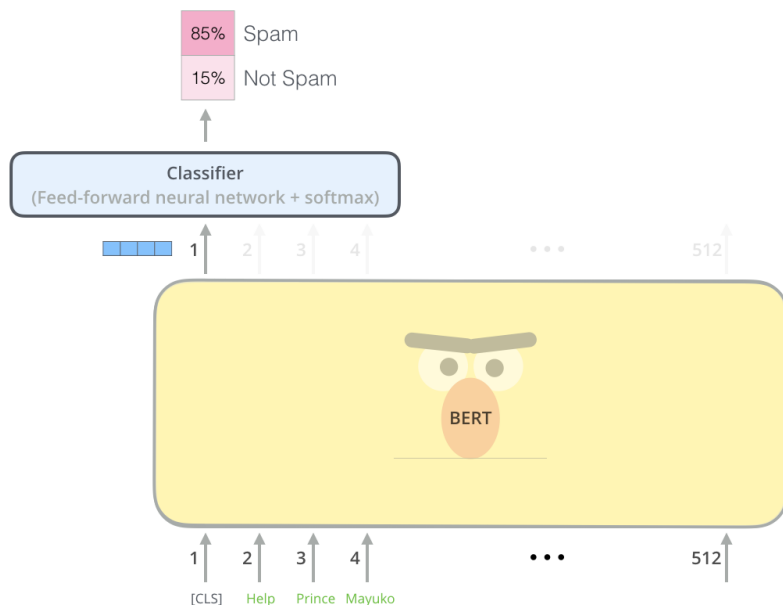
Compute center of samples per document in embedding space



Add document ID

State of the Art Deep Learning Approach

BERT: Bidirectional Encoder Representations from Transformers



BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
Google AI Language
{jacobdevlin, mingweichang, kentonl, kristout}@google.com

Abstract

We introduce a new language representation model called **BERT**, which stands for **Bidirectional Encoder Representations from Transformers**. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications.

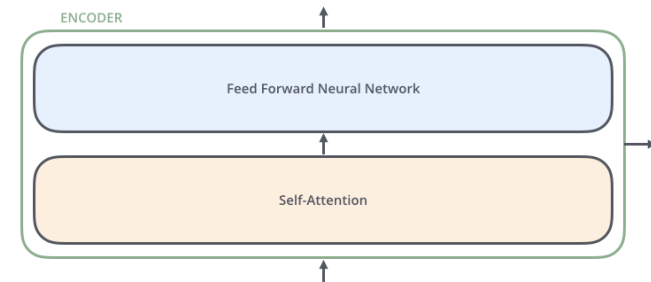
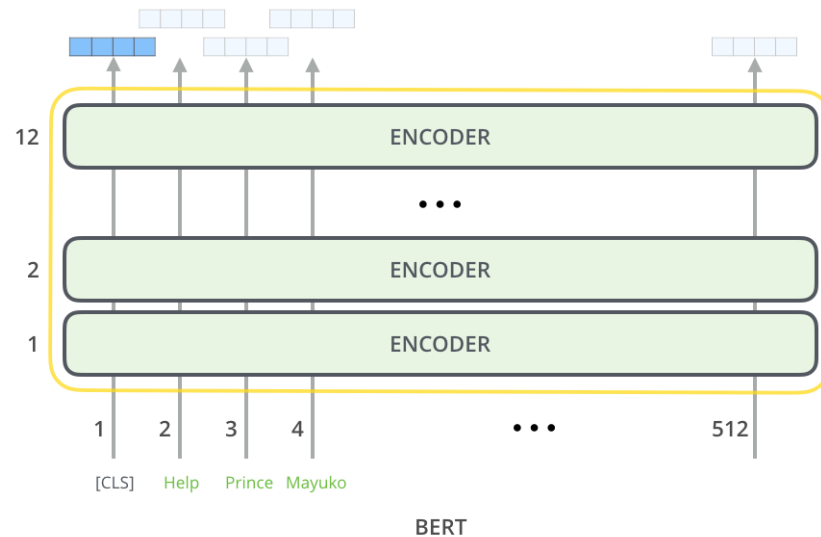
BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing the GLUE score to 80.5% (7.7% point absolute improvement), MultiNLI accuracy to 86.7% (4.6% absolute improvement), SQuAD v1.1 question answering Test F1 to 93.2 (1.5 point absolute improvement) and SQuAD v2.0 Test F1 to 92.1

There are two existing strategies for applying pre-trained language representations to downstream tasks: *feature-based* and *fine-tuning*. The feature-based approach, such as ELMo (Peters et al., 2018a), uses task-specific architectures that include the pre-trained representations as additional features. The fine-tuning approach, such as the Generative Pre-trained Transformer (OpenAI GPT) (Radford et al., 2018), introduces minimal task-specific parameters, and is trained on the downstream tasks by simply fine-tuning *all* pre-trained parameters. The two approaches share the same objective function during pre-training, where they use unidirectional language models to learn general language representations.

We argue that current techniques restrict the power of the pre-trained representations, especially for the fine-tuning approaches. The major limitation is that standard language models are unidirectional, and this limits the choice of architectures that can be used during pre-training. For example, in OpenAI GPT, the authors use a left-to-right architecture, where every token can only at

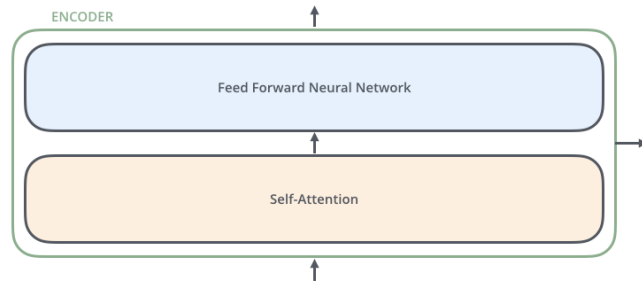
Figures by: <http://jalammar.github.io/illustrated-bert/>

Basic Architecture

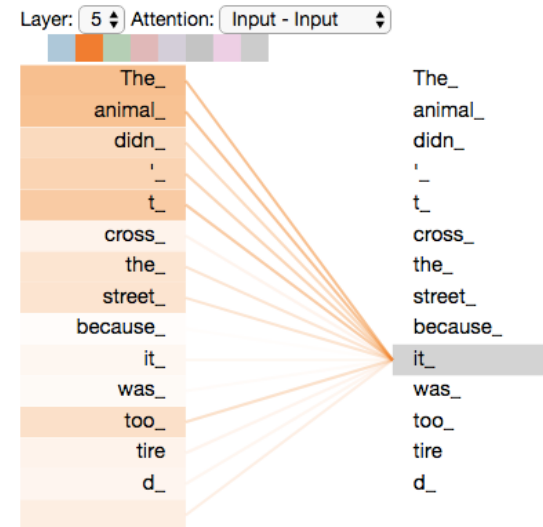


Figures by: <http://jalammar.github.io/illustrated-bert/>

Attention Layers



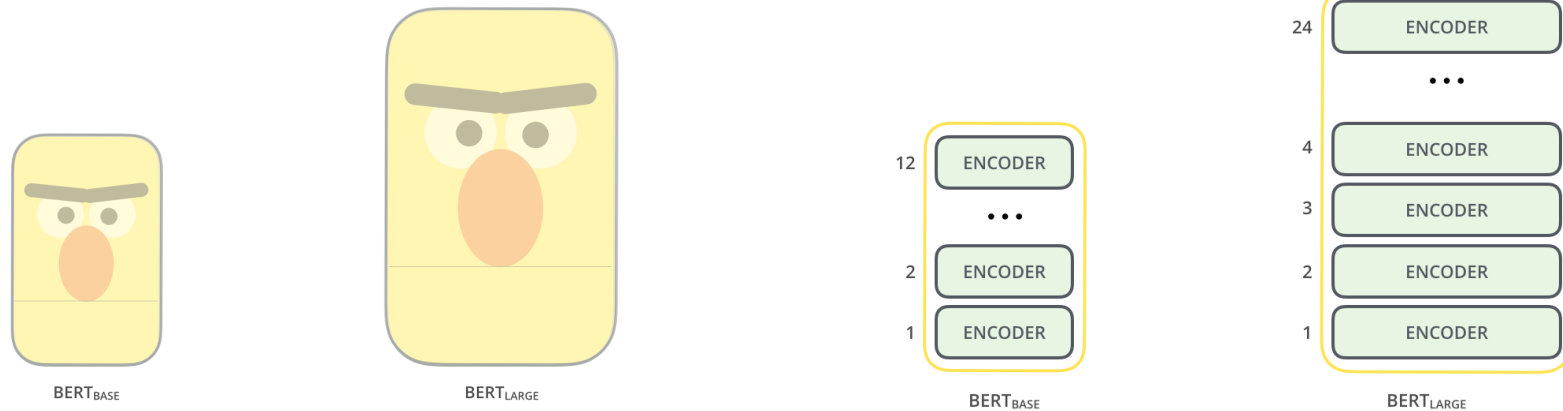
Encoding: to which word does „it“ refer?



Details: <https://jalammar.github.io/illustrated-transformer/>

State of the Art Deep Learning Approach

Then is more than one BERT:



Figures by: <http://jalammar.github.io/illustrated-bert/>

For the project work

BERT on Keras: <https://github.com/CyberZHG/keras-bert>

A Pre-trained Model for German BERT: <https://deepset.ai/german-bert>

