# Building a Search Engine with BERT and TensorFlow

In this experiment, we use a pre-trained BERT model checkpoint to build a general-purpose text feature extractor, which we apply to the task of nearest neighbour search.

Denis Antyukhov    Follow

Jun 28 · 7 min read ★



T-SNE decomposition of BERT text representations (Reuters-21578 benchmark, 6 classes)

Feature extractors based on deep Neural Probabilistic Language Models such as BERT may extract features that are relevant for a wide array of downstream NLP tasks. For that reason, they are sometimes referred to as Natural Language Understanding (NLU) modules.

These features may also be used for computing the similarity between text samples. To demonstrate this, we will build a nearest neighbour search engine for text, using BERT for feature extraction.

The plan for this experiment is:

1. getting the pre-trained BERT model checkpoint

2.  extracting a sub-graph optimized for inference

3.  creating a feature extractor with tf.Estimator

4.  exploring vector space with T-SNE and Embedding Projector

5.  implementing a nearest neighbour search engine

6.  accelerating nearest neighbour queries with math

7.  example: building a movie recommendation system

# Questions and Answers

### What is in this guide?

This guide contains implementations of two things: a BERT text feature extractor and a nearest neighbour search engine.

### Who is this guide for?

This guide should be useful for researchers interested in using BERT for natural language understanding tasks. It may also serve as a worked example of interfacing with tf.Estimator API.

### What does it take?

For a reader familiar with TensorFlow it should take around 30 minutes to finish this guide.

### Show me the code.

The code for this experiment is available in Colab here. Also, check out the repository I set up for my BERT experiments: it contains bonus stuff.

Now, let's start.

# Step 1: getting the pre-trained model

We start with a pre-trained BERT checkpoint. For demonstration purposes, I will be using the uncased English model pre-trained by Google engineers.

For configuring and optimizing the graph for inference we will make use of the awesome bert-as-a-service repository. This repository allows for serving BERT models for remote clients over TCP.

Having a remote BERT-server is beneficial in multi-host environments. However, in this part of the experiment we will focus on creating a local
(in-process) feature extractor. This is useful if one wishes to avoid additional latency and potential failure modes introduced by a client-server architecture.

Now, let us download the model and install the package.

```
!wget
https://storage.googleapis.com/bert_models/2019_05_30/wwm_un
cased_L-24_H-1024_A-16.zip
!unzip wwm_uncased_L-24_H-1024_A-16.zip
!pip install bert-serving-server --no-deps
```

## Step 2: optimizing the inference graph

Normally, to modify the model graph we would have to do some low-level TensorFlow programming. However, thanks to bert-as-a-service, we can configure the inference graph using a simple CLI interface.

```
1    import os
2    import tensorflow as tf
3
4    from bert_serving.server.graph import optimize_graph
5    from bert_serving.server.helper import get_args_parser
6

7
8    MODEL_DIR = '/content/wwm_uncased_L-24_H-1024_A-16/' #
9    GRAPH_DIR = '/content/graph/' #@param {type:"string"}
10   GRAPH_OUT = 'extractor.pbtxt' #@param {type:"string"}
11   GPU_MFRAC = 0.2 #@param {type:"string"}
12
13   POOL_STRAT = 'REDUCE_MEAN' #@param {type:"string"}
14   POOL_LAYER = "-2" #@param {type:"string"}
15   SEQ_LEN = "64" #@param {type:"string"}
16
17   tf.gfile.MkDir(GRAPH_DIR)
18
19   parser = get_args_parser()
20   carg = parser.parse_args(args=['-model_dir', MODEL_DIR
21                                 "-graph_tmp_dir", GRAPH
22                                 '-max_seq_len', str(SEQ
23                                 '-pooling_layer', str(P
```

There are a couple of parameters there to look out for.

For each text sample, BERT encoding layers output a tensor of shape [*sequence_len, encoder_dim*] with one vector per token. We need to apply some sort of pooling if we are to obtain a fixed representation.

**POOL_STRAT** parameter defines the pooling strategy applied to the encoding layer number **POOL_LAYER**. The default value '*REDUCE_MEAN*' averages the vectors for all tokens in a sequence. This strategy works best for most sentence-level tasks when the model is not fine-tuned. Another option is *NONE*, in which case no pooling is applied at all. This is useful for word-level tasks such as Named Entity Recognition or POS tagging. For a detailed discussion of these options check out the Han Xiao's blog post.

**SEQ_LEN** affects the maximum length of sequences processed by the model. Smaller values will increase the model inference speed almost linearly.

Running the above command will put the model graph and weights into a <u>GraphDef</u> object which will be serialized to a *pbtxt* file at **GRAPH_OUT**. The file will often be smaller than the pre-trained model because the nodes and variables required for training will be removed. This results in a very portable solution: for example, the english model only takes 380 MB after serializing.

## Step 3: creating a feature extractor

Now, we will use the serialized graph to build a feature extractor using the tf.Estimator API. We will need to define two things: **input_fn** and **model_fn**

<mark>input_fn</mark> manages getting the data into the model. That includes executing the whole text preprocessing pipeline and preparing a **feed_dict** for BERT.

First, each text sample is converted into a **tf.Example** instance containing the necessary features listed in INPUT_NAMES. <mark>The bert_tokenizer object contains the WordPiece vocabulary and performs the text preprocessing.</mark> After that, the examples are re-grouped by feature name in a **feed_dict**.

```
1   INPUT_NAMES = ['input_ids', 'input_mask', 'input_type_
2   bert_tokenizer = FullTokenizer(VOCAB_PATH)
3
4   def build_feed_dict(texts):
5
6       text_features = list(convert_lst_to_features(
7           texts, SEQ_LEN, SEQ_LEN,
8           bert_tokenizer, log, False, False))
9
10      target_shape = (len(texts), -1)
11
12      feed_dict = {}
13      for iname in INPUT_NAMES:
```

tf.Estimators have a fun feature which makes them rebuild and reinitialize the whole computational graph at each call to the predict function. So, in order to avoid the overhead, we will pass a generator to the predict function, and the generator will yield the features to the model in a never-ending loop. Ha-ha.

```
1   def build_input_fn(container):
2
3       def gen():
4           while True:
5             try:
6                 yield build_feed_dict(container.get())
7             except StopIteration:
8                 yield build_feed_dict(container.get())
9
10      def input_fn():
11          return tf.data.Dataset.from_generator(
12              gen,
13              output_types={iname: tf.int32 for iname in
14              output_shapes={iname: (None, None) for ina
15      return input_fn
16
17  class DataContainer:
18    def __init__(self):
```

. . .

**model_fn** contains the specification of the model. In our case, it is
loaded from the *pbtxt* file we saved in the previous step. The features
are mapped explicitly to the corresponding input nodes via **input_map**.

```
1   def model_fn(features, mode):
2       with tf.gfile.GFile(GRAPH_PATH, 'rb') as f:
3           graph_def = tf.GraphDef()
4           graph_def.ParseFromString(f.read())
5
6       output = tf.import_graph_def(graph_def,
7                                   input_map={k + ':0':
8                                                 for k in I
9                                   return_elements=['fin
```

Now we have almost everything we need to perform inference. Let's do
this!

```
1    def batch(iterable, n=1):
2        l = len(iterable)
3        for ndx in range(0, l, n):
4            yield iterable[ndx:min(ndx + n, l)]
5
6    def build_vectorizer(_estimator, _input_fn_builder, ba
7      container = DataContainer()
8      predict_fn = _estimator.predict(_input_fn_builder(cc
9
10     def vectorize(text, verbose=False):
11       x = []
12       bar = Progbar(len(text))
13       for text_batch in batch(text, batch_size):
14         container.set(text_batch)
15         x.append(next(predict_fn)['output'])
```

A standalone version of the feature extractor described above can be found in the underline{repository}.

```
>>> bert_vectorizer = build_vectorizer(estimator,
build_input_fn)
>>> bert_vectorizer(64*['sample text']).shape
(64, 768)
```

# Step 4: exploring vector space with Projector

Now it's time for a demonstration!

Using the vectorizer, we will generate embeddings for articles from the Reuters-21578 benchmark corpus. To visualize and explore the embedding vector space in 3D, we will use a dimensionality reduction technique called T-SNE.

Let's get the article embeddings first.

```python
1     from nltk.corpus import reuters
2
3     nltk.download("reuters")
4     nltk.download("punkt")
5
6     max_samples = 256
7     categories = ['wheat', 'tea', 'strategic-metal',
8                   'housing', 'money-supply', 'fuel']
9
10    S, X, Y = [], [], []
11
12    for category in categories:
13      print(category)
14
15      sents = reuters.sents(categories=category)
```

The interactive visualization of generated embeddings is available on the Embedding Projector.

From the link you can run T-SNE yourself or load a checkpoint using the bookmark in lower-right corner (loading works only on Chrome).



Reuters-21578 BERT T-SNE

## Step 5: building a search engine

Now, let's say we have a knowledge base of 50k text samples and we need to answer queries based on this data, fast. How do we retrieve the

sample most similar to a query from a text database? The answer is
<u>nearest neighbour search</u>.

Formally, the search problem we will be solving is defined as follows:
given a set of points **S** in vector space **M**, and a query point $Q \in M$, find
the closest point in **S** to **Q**. There are multiple ways to define 'closest' in
vector space, we will use <u>Euclidean distance</u>.

<mark>So, to build an Search Engine for text we will follow these steps:</mark>

1. Vectorize all samples from the knowledge base—that gives **S**

2. Vectorize the query—that gives **Q**

3. Compute euclidean distances **D** between **Q** and **S**

4. Sort **D** in ascending order—providing indices of the most similar
   samples

5. Retrieve labels for said samples from the knowledge base

To make the simple matter of implementing this a bit more exciting, we
will do it in pure TensorFlow.

First, we create the placeholders for **Q** and **S**

```
1    dim = 1024
2    graph = tf.Graph()
3    sess = tf.InteractiveSession(graph=graph)
4
5    Q = tf.placeholder("float", [dim])
```

<mark>Define euclidean distance computation</mark>

```
1    squared_distance = tf.reduce_sum(tf.pow(Q - S, 2), redu
2    distance = tf.sqrt(squared_distance)
```

Finally, get the most similar sample indices

```
1    top_k = 3
2
3    top_neg_dists, top_indices = tf.math.top_k(tf.negative(
4    top dists = tf.negative(top neg dists)
```

## Step 6: accelerating search with math

Now that we have a basic retrieval algorithm set up, the question is: can we make it run faster? With a tiny bit of math, we can.

For a pair of vectors $\mathbf{p}$ and $\mathbf{q}$, the euclidean distance is defined as follows:

$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \cdots + (q_n - p_n)^2}$$

$$= \sqrt{\sum_{i=1}^{n}(q_i - p_i)^2}.$$

Which is exactly how we did compute it in Step 4.

However, since $\mathbf{p}$ and $\mathbf{q}$ are vectors, we can expand and rewrite this:

$$d(\mathbf{p}, \mathbf{q})^2 = \langle \mathbf{p} - \mathbf{q}, \mathbf{p} - \mathbf{q} \rangle = \langle \mathbf{pp} \rangle - 2\langle \mathbf{pq} \rangle + \langle \mathbf{qq} \rangle$$

where $\langle \ldots \rangle$ denotes inner product.

In TensorFlow this can be written as follows:

```
1    Q = tf.placeholder("float", [dim])
2    S = tf.placeholder("float", [None, dim])
3
4    Qr = tf.reshape(Q, (1, -1))
5
6    PP = tf.keras.backend.batch_dot(S, S, axes=1)
7    QQ = tf.matmul(Qr, tf.transpose(Qr))
8    PQ = tf.matmul(S, tf.transpose(Qr))
9
```

Due to the fact that matrix multiplication op is highly optimized, this implementation works slightly faster than the previous one.

By the way, in the formula above **PP** and **QQ** are actually squared L2 norms of the respective vectors. If both vectors are L2 normalized, then **PP = QQ** = 1. That gives an interesting relation between inner product and euclidean distance:

$$d(\mathbf{p}, \mathbf{q})^2 = 2 * (1 - \langle \mathbf{pq} \rangle)$$

However, doing L2 normalization discards the information about vector magnitude which in many cases is undesirable .

Instead, we may notice that as long as the knowledge base does not change, **PP,** its squared vector norm, also remains constant. So, instead of recomputing it every time we can just do it once and then use the precomputed result further accelerating the distance computation.
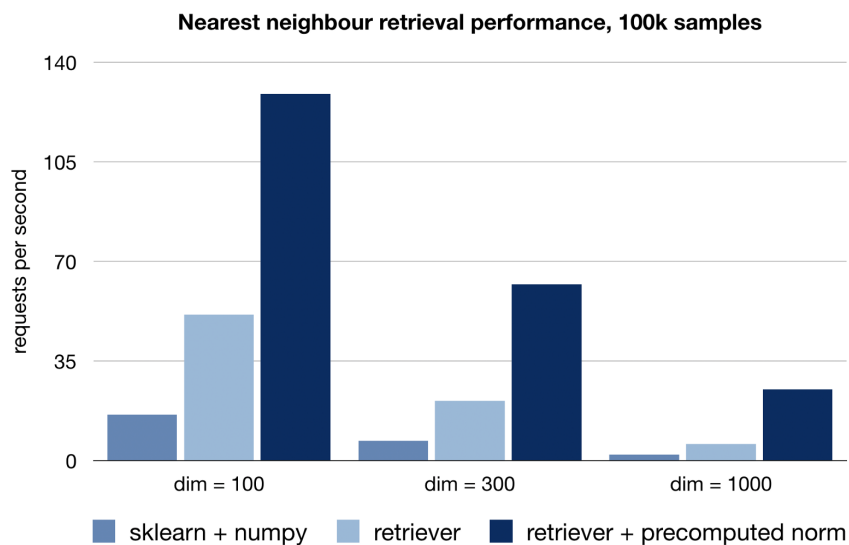
Now let us put it all together.

```python
1    class L2Retriever:
2        def __init__(self, dim, top_k=3, use_norm=False, u
3            self.dim = dim
4            self.top_k = top_k
5            self.use_norm = use_norm
6            config = tf.ConfigProto(
7                device_count={'GPU': (1 if use_gpu else 0)
8            )
9            config.gpu_options.allow_growth = True
10           self.session = tf.Session(config=config)
11
12           self.norm = None
13           self.query = tf.placeholder("float", [self.dim
14           self.kbase = tf.placeholder("float", [None, se
15
16           self.build_graph()
17
18       def build_graph(self):
19
20           if self.use_norm:
21               self.norm = tf.placeholder("float", [None,
22
23           distance = dot_l2_distances(self.kbase, self.q
24           top_neg_dists, top_indices = tf.math.top_k(tf.
25           top_dists = tf.negative(top_neg_dists)
26
27           self.top_distances = top_dists
28           self.top_indices = top_indices
29
30       def predict(self, kbase, query, norm=None):
31           query = np.squeeze(query)
32           feed_dict = {self.query: query, self.kbase: kb
33           if self.use_norm:
34               feed_dict[self.norm] = norm
35
```

Naturally, you could use this implementation with any model, not just BERT. It is quite effective at nearest neighbour retrieval, being able to process dozens of requests per second on CPU.

**Nearest neighbour retrieval performance, 100k samples**



## Example: movie recommendation system

For this example we will use a dataset of movie summaries from IMDB. Using the NLU and Retriever modules, we will build a movie recommendation system that suggests movies with similar plot features.

First, let's download and prepare the IMDB dataset.

```
1   import pandas as pd
2   import json
3
4   !wget http://www.cs.cmu.edu/~ark/personas/data/MovieSu
5   !tar –xvzf MovieSummaries.tar.gz
6
7   plots_df = pd.read_csv('MovieSummaries/plot_summaries.
8   meta_df = pd.read_csv('MovieSummaries/movie.metadata.t
9
10  plot = {}
11  metadata = {}
12  movie_data = {}
13
14  for movie_id, movie_plot in plots_df.values:
15    plot[movie_id] = movie_plot
16
17  for movie_id, movie_name, movie_genre in meta_df[[0,2,
18    genre = list(json.loads(movie_genre).values())
19    if len(genre):
20      metadata[movie_id] = {"name": movie_name,
21                            "genre": genre}
```

Vectorize movie plots with the BERT NLU module:

```
1   X_vect = bert_vectorizer(X, verbose=True)
```

Finally, using the L2Retriever, find movies with plot vectors most similar to the query movie, and return it to user.

```
1   def buildMovieRecommender(movie_names, vectorized_plot
2     retriever = L2Retriever(vectorized_plots.shape[1], u
3     vectorized_norm = np.sum(vectorized_plots**2, axis=1
4
5     def recommend(query):
6       try:
7         idx = retriever.predict(vectorized_plots,
8                                 vectorized_plots[movie_n
9                                 vectorized_norm)[0][1:]
10        for i in idx:
11          print(names[i])
12      except ValueError:
```

Let's check it out!

```
>>> recommend = buildMovieRecommender(names, X_vect)
>>> recommend("The Matrix")
Impostor
Immortel
Saturn 3
Terminator Salvation
The Terminator
Logan's Run
Genesis II
Tron: Legacy
Blade Runner
```

==Even without supervision, the model performs adequately on several classification and retrieval tasks.== While the performance can be further improved with supervised data, the described approach to text feature extraction provides a solid baseline for downstream NLP solutions.

This concludes the guide to building a search engine with BERT and TensorFlow.