

[Get started](#)[Open in app](#)

# Aniruddha Choudhury

[Follow](#)

531 Followers

[About](#)

GOOGLE SOURCES

## Part 2: BERT Fine-Tuning Tutorial with PyTorch for Text Classification on The Corpus of Linguistic Acceptability (COLA) Dataset.

[Get started](#)[Open in app](#)

In this tutorial I'll show you how to use BERT with the hugging face PyTorch library to quickly and efficiently fine-tune a model to get near state of the art performance in sentence classification. More broadly, I describe the practical application of transfer learning in NLP to create high performance models with minimal effort on a range of NLP tasks.



## Introduction

## History

2018 was a breakthrough year in NLP. Transfer learning, particularly models like Allen AI's ELMO, OpenAI's Open-GPT, and Google's BERT allowed researchers to smash multiple benchmarks with minimal task-specific fine-tuning and provided the rest of the NLP community with pretrained models that could easily (with less data and less compute time) be fine-tuned and implemented to produce state of the art results. Unfortunately, for many starting out in NLP and even for some experienced practitioners, the theory and practical application of these powerful models is still not well understood.

## What is BERT?

[Get started](#)[Open in app](#)

smentation models , 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 re- sult, 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 a method of pretraining language representations that was used to create models that NLP practitioners can then download and use for free. You can either use these models to extract high quality language features from your text data, or you can fine-tune these models on a specific task (classification, entity recognition, question answering, etc.) with your own data to produce state of the art predictions.

This post will explain how you can modify and fine-tune BERT to create a powerful NLP model that quickly gives you state of the art results.

## Advantages of Fine-Tuning

In this tutorial, we will use BERT to train a text classifier. Specifically, we will take the pre-trained BERT model, add an untrained layer of neurons on the end, and train the new model for our classification task. Why do this rather than train a train a specific deep learning model (a CNN, BiLSTM, etc.) that is well suited for the specific NLP task you need?

**Quicker Development:** First, the pre-trained BERT model weights already encode a lot of information about our language. As a result, it takes much less time to train our fine-tuned model — it is as if we have already trained the bottom layers of our network extensively and only need to gently tune them while using their output as features for our classification task. In fact, the authors recommend only 2–4 epochs of training for fine-tuning BERT on a specific NLP task (compared to the hundreds of GPU hours needed to train the original BERT model or a LSTM from scratch!).

**Less Data:** In addition and perhaps just as important, because of the pre-trained weights this method allows us to fine-tune our task on a much smaller dataset than would be required in a model that is built from scratch. A major drawback of NLP models built

[Get started](#)[Open in app](#)

dataset creation. By fine-tuning BERT, we are now able to get away with training a model to good performance on a much smaller amount of training data.

**Better Results:** Finally, this simple fine-tuning procedure (typically adding one fully-connected layer on top of BERT and training for a few epochs) was shown to achieve state of the art results with minimal task-specific adjustments for a wide variety of tasks: classification, language inference, semantic similarity, question answering, etc. Rather than implementing custom and sometimes-obscure architectures shown to work well on a specific task, simply fine-tuning BERT is shown to be a better (or at least equal) alternative.



## 1. Setup

### 1.1. Using Colab GPU for Training

Google Colab offers free GPUs and TPUs! Since we'll be training a large neural network it's best to take advantage of this (in this case we'll attach a GPU), otherwise training will take a very long time.

A GPU can be added by going to the menu and selecting:

Edit --> Notebook Settings --> Hardware accelerator --> (GPU)

Then run the following cell to confirm that the GPU is detected.

```
import tensorflow as tf
```

[Get started](#)[Open in app](#)

```
# The device name should look like the following:  
if device_name == '/device:GPU:0':  
    print('Found GPU at: {}'.format(device_name))  
else:  
    raise SystemError('GPU device not found')
```

The default version of TensorFlow in Colab will soon switch to TensorFlow 2.x.

☞ The default version of TensorFlow in Colab will soon switch to TensorFlow 2.x.  
We recommend you [upgrade](#) now or ensure your notebook will continue to use TensorFlow 1.x via the `@tensorflow_version 1.x` magic: [more info](#).

```
Found GPU at: /device:GPU:0
```

In order for torch to use the GPU, we need to identify and specify the GPU as the device. Later, in our training loop, we will load data onto the device.

```
import torch  
  
# If there's a GPU available...  
if torch.cuda.is_available():  
  
    # Tell PyTorch to use the GPU.  
    device = torch.device("cuda")  
  
    print('There are %d GPU(s) available.' %  
          torch.cuda.device_count())  
  
    print('We will use the GPU:', torch.cuda.get_device_name(0))  
  
# If not...  
else:  
    print('No GPU available, using the CPU instead.')  
    device = torch.device("cpu")
```

☞ There are 1 GPU(s) available.  
We will use the GPU: Tesla T4

## 1.2. Installing the Hugging Face Library

[Get started](#)[Open in app](#)

pretrained language models like OpenAI's GPT and GPT-2.) We've selected the pytorch interface because it strikes a nice balance between the high-level APIs (which are easy to use but don't provide insight into how things work) and tensorflow code (which contains lots of details but often sidetracks us into lessons about tensorflow, when the purpose here is BERT!).

At the moment, the Hugging Face library seems to be the most widely accepted and powerful pytorch interface for working with BERT. In addition to supporting a variety of different pre-trained transformer models, the library also includes pre-built modifications of these models suited to your specific task. For example, in this tutorial we will use BertForSequenceClassification.

The library also includes task-specific classes for token classification, question answering, next sentence predicton, etc. Using these pre-built classes simplifies the process of modifying BERT for your purposes.

```
!pip install transformers
```

```
Collecting transformers
  Downloading https://files.pythonhosted.org/packages/50/10/aeefcfed99c8a59d828a92cc11d213e2743212d3641c87c82d61b1
    |██████████| 450kB 3.5MB/s
Requirement already satisfied: requests in /usr/local/lib/python3.6/dist-packages (from transformers) (2.21.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from transformers) (1.17.5)
```

The code in this notebook is actually a simplified version of the `run_glue.py` example script from `huggingface`.

`run_glue.py` is a helpful utility which allows you to pick which GLUE benchmark task you want to run on, and which pre-trained model you want to use (you can see the list of possible models [here](#)). It also supports using either the CPU, a single GPU, or multiple GPUs. It even supports using 16-bit precision if you want further speed up.

## 2. Loading CoLA Dataset

We'll use [The Corpus of Linguistic Acceptability \(CoLA\)](#) dataset for single sentence classification. It's a set of sentences labeled as grammatically correct or incorrect. It was

[Get started](#)[Open in app](#)

## 2.1. Download & Extract

We'll use the wget package to download the dataset to the Colab instance's file system.

```
!pip install wget
```

```
C  Collecting wget
    Downloading https://files.pythonhosted.org/packages/47/6a/62e288da7bcda82b935ff0c6cfe542970f04e29c756b0e1472511
Building wheels for collected packages: wget
```

The dataset is hosted on GitHub in this repo: <https://nyu-mll.github.io/CoLA/>

```
import wget
import os

print('Downloading dataset')

# The URL for the dataset zip file.
url = 'https://nyu-mll.github.io/CoLA/cola_public_1.1.zip'

# Download the file (if we haven't already)
if not os.path.exists('./cola_public_1.1.zip'):
    wget.download(url, './cola_public_1.1.zip')
```

Unzip the dataset to the file system. You can browse the file system of the Colab instance in the sidebar on the left.

```
# Unzip the dataset (if we haven't already)
if not os.path.exists('./cola_public/'):
    !unzip cola_public_1.1.zip
```

```
Archive:  cola_public_1.1.zip
creating:  cola_public/
inflating:  cola_public/README
creating:  cola_public/tokenized/
inflating:  cola_public/tokenized/in_domain_dev.tsv
```

[Get started](#)[Open in app](#)

available.

We can't use the pre-tokenized version because, in order to apply the pre-trained BERT, we *must* use the tokenizer provided by the model. This is because

- (1) the model has a specific, fixed vocabulary and (2) the BERT tokenizer has a particular way of handling out-of-vocabulary words.

We'll use pandas to parse the "in-domain" training set and look at a few of its properties and data points.

```
import pandas as pd

# Load the dataset into a pandas dataframe.
df = pd.read_csv("./cola_public/raw/in_domain_train.tsv",
delimiter='\t', header=None, names=['sentence_source', 'label',
'label_notes', 'sentence'])

# Report the number of sentences.
print('Number of training sentences: {:.}\n'.format(df.shape[0]))

# Display 10 random rows from the data.
df.sample(10)
```

Number of training sentences: 8,551

<b>sentence_source</b>	<b>label</b>	<b>label_notes</b>	<b>sentence</b>	
<b>5980</b>	c_13	1	NaN	Frank has drunk too much.
<b>2694</b>	I-93	0	*	The package carried to New York.
<b>1884</b>	r-67	1	NaN	Jack will have a hole in his pocket.
<b>6628</b>	m_02	0	*	Dogs chase.
<b>5711</b>	c_13	1	NaN	The officer carefully inspected the license.

The two properties we actually care about are the sentence and its label, which is referred to as the "acceptability judgment" (0=unacceptable, 1=acceptable).

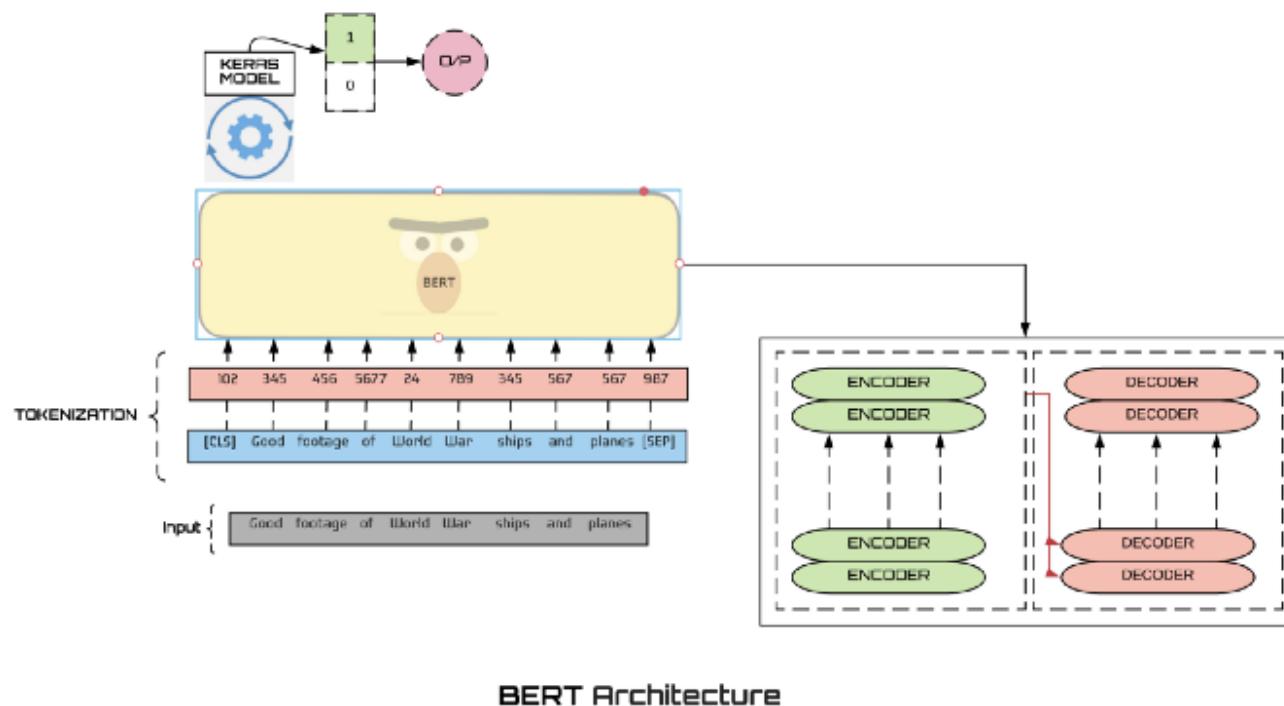
[Get started](#)[Open in app](#)

Let's extract the sentences and labels of our training set as numpy ndarrays.

```
# Get the lists of sentences and their labels.
sentences = df.sentence.values
labels = df.label.values
```

### 3. Tokenization & Input Formatting

We'll transform our dataset into the format that BERT can be trained on.



#### 3.1. BERT Tokenizer

To feed our text to BERT, it must be split into tokens, and then these tokens must be mapped to their index in the tokenizer vocabulary.

The tokenization must be performed by the tokenizer included with BERT—the below cell will download this for us. We'll be using the “uncased” version here. For more details please find my previous [Article](#).

[Get started](#)[Open in app](#)

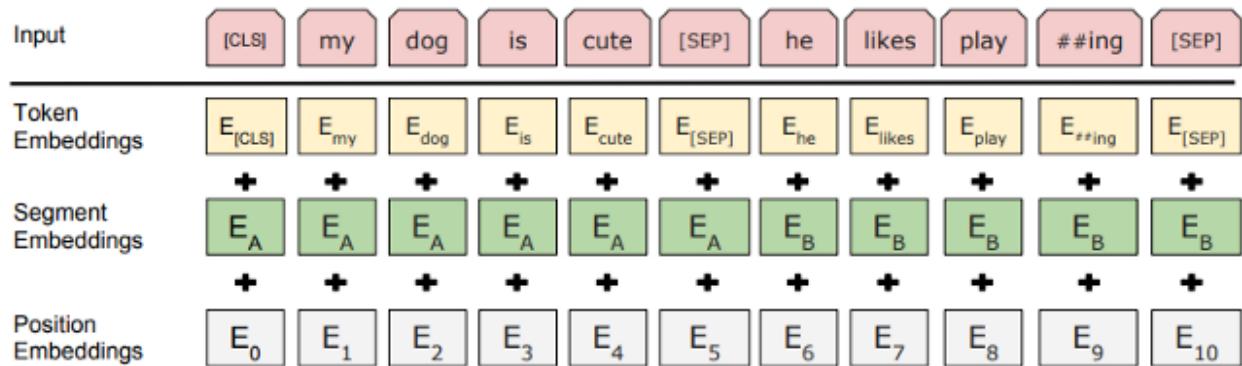
```
print('Loading BERT tokenizer...')
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased',
do_lower_case=True)
```

Original: Our friends won't buy this analysis, let alone the next one we propose.  
Token IDs: [101, 2256, 2814, 2180, 1005, 1056, 4965, 2023, 4106, 1010, 2292, 2894,

Let's apply the tokenizer to one sentence just to see the output. When we actually convert all of our sentences, we'll use the tokenize.encode function to handle both steps, rather than calling tokenize and convert\_tokens\_to\_ids separately. Before we can do that, though, we need to talk about some of BERT's formatting requirements.

### 3.2. Required Formatting

The above code left out a few required formatting steps that we'll look at here.



BERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings and the position embeddings.

1. The first token of every sequence is always a special classification token ([CLS]).
2. The final hidden state corresponding to this token is used as the aggregate sequence representation for classification tasks. Sentence pairs are packed together into a single sequence.

[Get started](#)[Open in app](#)

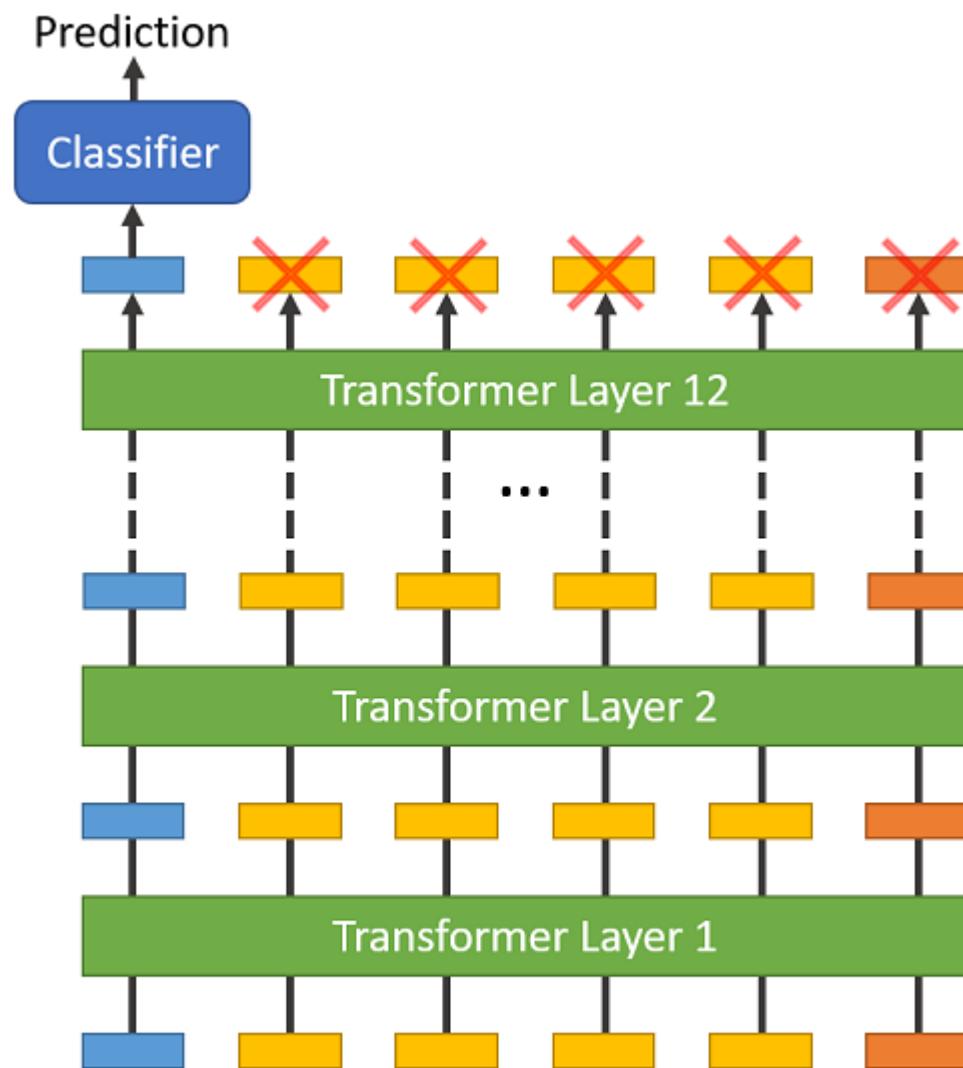
whether it belongs to sentence A or sentence B.

4. A positional embedding is also added to each token to indicate its position in the sequence.

We are required to:

1. Add special tokens to the start and end of each sentence.
2. Pad & truncate all sentences to a single constant length.
3. Explicitly differentiate real tokens from padding tokens with the “attention mask”.

## Special Tokens



[Get started](#)[Open in app](#)

## [SEP]

At the end of every sentence, we need to append the special [SEP] token.

This token is an artifact of two-sentence tasks, where BERT is given two separate sentences and asked to determine something (e.g., can the answer to the question in sentence A be found in sentence B?).

Here we are not certain yet why the token is still required when we have only single-sentence input, but it is!

## [CLS]

For classification tasks, we must prepend the special [CLS] token to the beginning of every sentence.

This token has special significance. BERT consists of 12 Transformer layers. Each transformer takes in a list of token embeddings, and produces the same number of embeddings on the output (but with the feature values changed, of course!).

On the output of the final (12th) transformer, *only the first embedding (corresponding to the [CLS] token) is used by the classifier.*

*“The first token of every sequence is always a special classification token ([CLS]). The final hidden state corresponding to this token is used as the aggregate sequence representation for classification tasks.”*

Also, because BERT is trained to only use this [CLS] token for classification, we know that the model has been motivated to encode everything it needs for the classification step into that single 768-value embedding vector.

## Sentence Length & Attention Mask

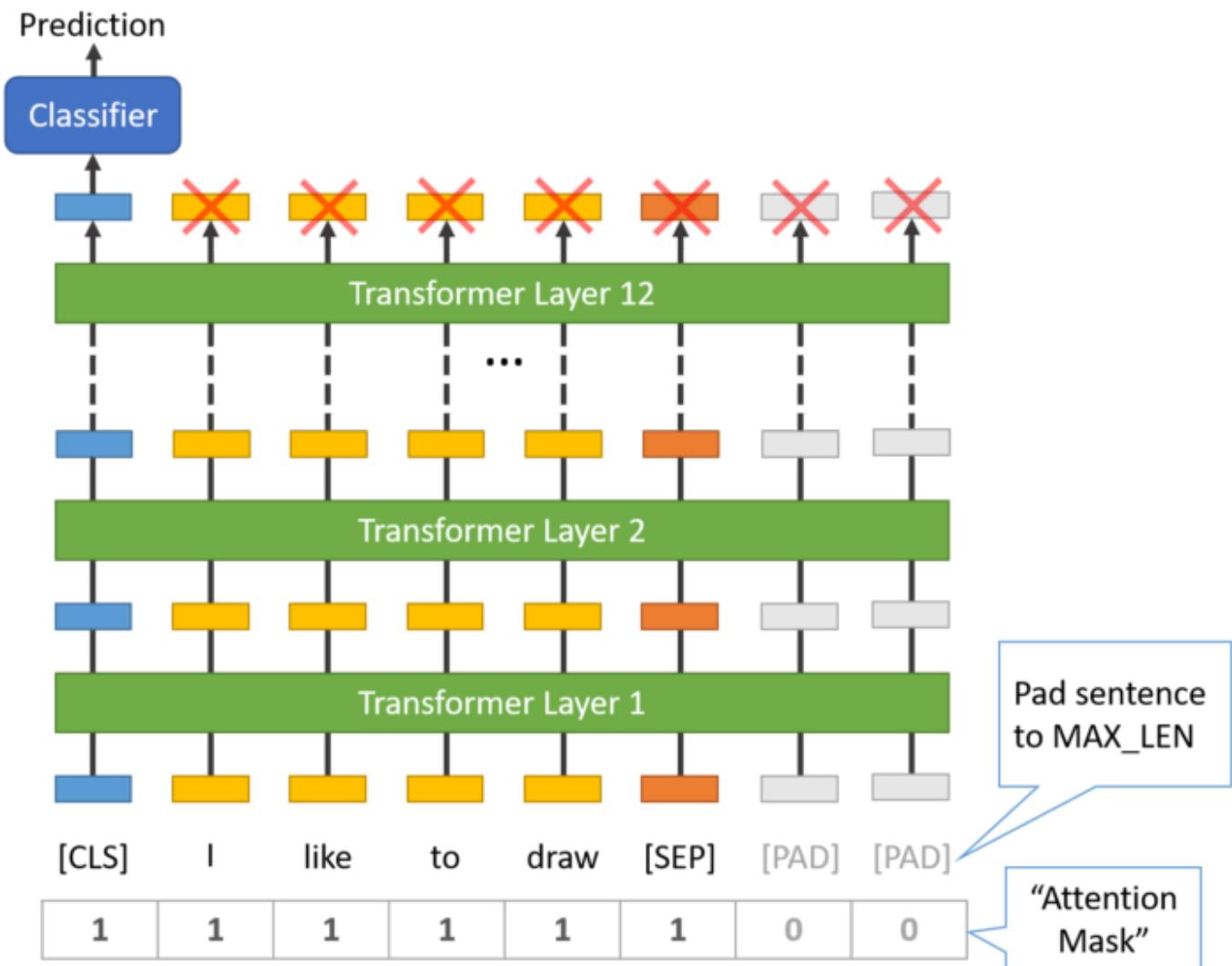
The sentences in our dataset obviously have varying lengths, so how does BERT handle this?

[Get started](#)[Open in app](#)

1. All sentences must be padded or truncated to a single, fixed length.

2. The maximum sentence length is 512 tokens.

Padding is done with a special [PAD] token, which is at index 0 in the BERT vocabulary. The below illustration demonstrates padding out to a “MAX\_LEN” of 8 tokens.



The “Attention Mask” is simply an array of 1s and 0s indicating which tokens are padding and which aren’t (seems kind of redundant, doesn’t it?! Again, I don’t currently know why).

I’ve experimented with running this notebook with two different values of MAX\_LEN, and it impacted both the training speed and the test set accuracy.

[Get started](#)[Open in app](#)

- ~~MAX\_LEN = 120 → Training epochs take ~5.20 each, score is 0.555~~

- MAX\_LEN = 64 → Training epochs take ~2:57 each, score is 0.566

These results suggest that the padding tokens aren't simply skipped over—that they are in fact fed through the model and incorporated in the results (thereby impacting both model speed and accuracy).

## 3.2. Sentences to IDs

The tokenizer.encode function combines multiple steps for us:

1. Split the sentence into tokens.
2. Add the special [CLS] and [SEP] tokens.
3. Map the tokens to their IDs.

Oddly, this function can perform truncating for us, but doesn't handle padding.

```
# Tokenize all of the sentences and map the tokens to thier word IDs.
input_ids = []

# For every sentence...
for sent in sentences:
    # `encode` will:
    #   (1) Tokenize the sentence.
    #   (2) Prepend the `[CLS]` token to the start.
    #   (3) Append the `[SEP]` token to the end.
    #   (4) Map tokens to their IDs.
    encoded_sent = tokenizer.encode(
        sent,                                     # Sentence to
        encode,                                     # Sentence to
        add_special_tokens = True, # Add '[CLS]' and
        '[SEP]'

        # This function also supports truncation and
        conversion
        # to pytorch tensors, but we need to do
        padding, so we
        # can't use these features :( .
        #max_length = 128,                         # Truncate all
        sentences.
        #return_tensors = 'pt',                   # Return pytorch
```

[Get started](#)[Open in app](#)

```
# Add the encoded sentence to the list.
input_ids.append(encoded_sent)
```

```
# Print sentence 0, now as a list of IDs.
print('Original: ', sentences[0])
print('Token IDs:', input_ids[0])
```

↳ Original: Our friends won't buy this analysis, let alone the next one we propose.  
 Tokenized: ['our', 'friends', 'won', "'", 't', 'buy', 'this', 'analysis', ',', 'let', 'alone',  
 Token IDs: [2256, 2814, 2180, 1005, 1056, 4965, 2023, 4106, 1010, 2292, 2894, 1996, 2279, 2028,

### 3.3. Padding & Truncating

Pad and truncate our sequences so that they all have the same length, MAX\_LEN. First, what's the maximum sentence length in our dataset?

```
print('Max sentence length: ', max([len(sen) for sen in input_ids]))
```

↳ Max sentence length: 47

Given that, let's choose MAX\_LEN = 64 and apply the padding.

```
# We'll borrow the `pad_sequences` utility function to do this.
from keras.preprocessing.sequence import pad_sequences

# Set the maximum sequence length.
# I've chosen 64 somewhat arbitrarily. It's slightly larger than the
# maximum training sentence length of 47...
MAX_LEN = 64

print('\nPadding/truncating all sentences to %d values...' % MAX_LEN)

print('\nPadding token: "{}", ID: {}'.format(tokenizer.pad_token,
                                             tokenizer.pad_token_id))

# Pad our input tokens with value 0.
# "post" indicates that we want to pad and truncate at the end of the
# sequence,
```

[Get started](#)[Open in app](#)

```
print('\nDone.\n')
```

```
↳ Padding/truncating all sentences to 64 values.  
Padding token: "[PAD]", ID: 0  
  
Done.  
Using TensorFlow backend.
```

## 3.4. Attention Masks

The attention mask simply makes it explicit which tokens are actual words versus which are padding.

The BERT vocabulary does not use the ID 0, so if a token ID is 0, then it's padding, and otherwise it's a real token.

```
# Create attention masks  
attention_masks = []  
  
# For each sentence...  
for sent in input_ids:  
  
    # Create the attention mask.  
    #   - If a token ID is 0, then it's padding, set the mask to 0.  
    #   - If a token ID is > 0, then it's a real token, set the mask  
    to 1.  
    att_mask = [int(token_id > 0) for token_id in sent]  
  
    # Store the attention mask for this sentence.  
    attention_masks.append(att_mask)
```

## 3.5. Training & Validation Split

Divide up our training set to use 90% for training and 10% for validation.



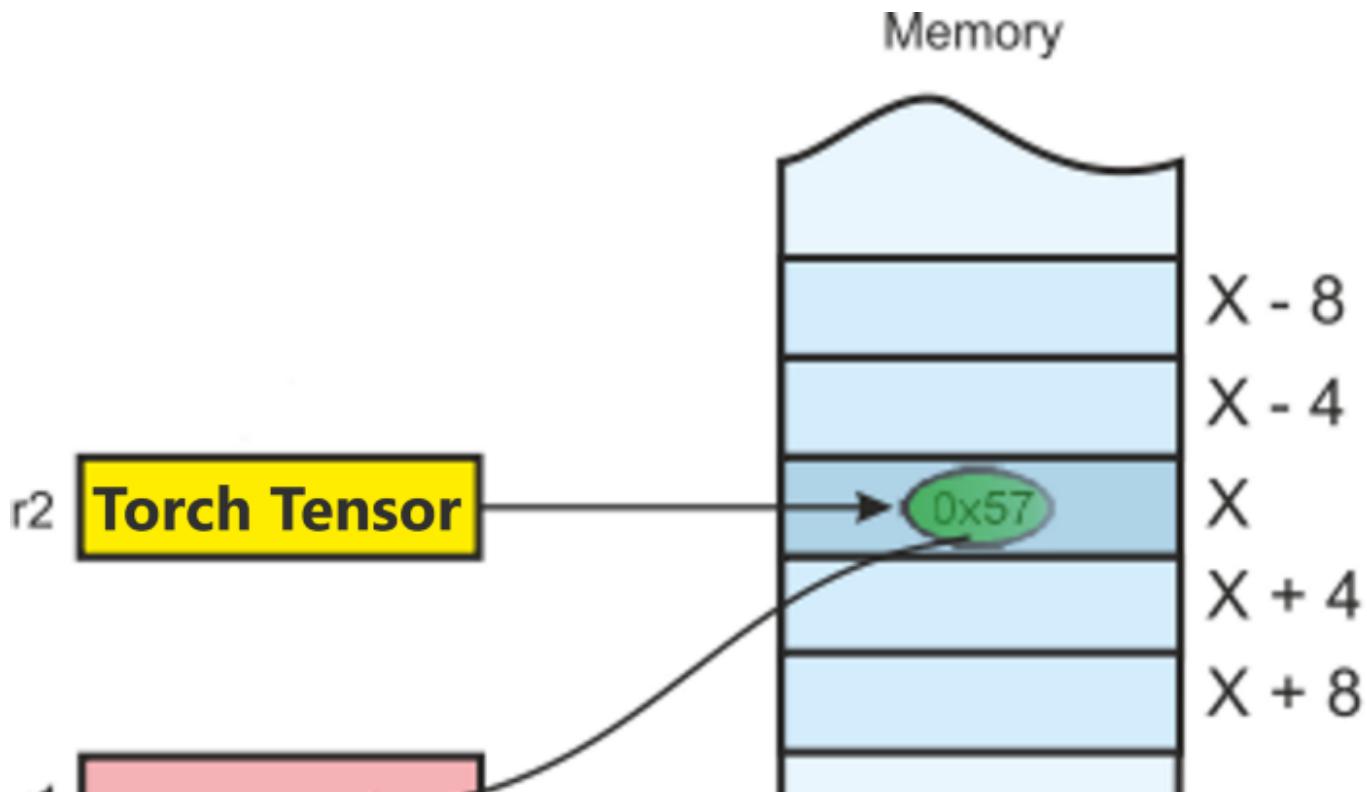
[Get started](#)[Open in app](#)

```
# Use train_test_split to split our data into train and validation sets for
# training
from sklearn.model_selection import train_test_split

# Use 90% for training and 10% for validation.
train_inputs, validation_inputs, train_labels, validation_labels =
train_test_split(input_ids, labels,
random_state=2018, test_size=0.1)
# Do the same for the masks.
train_masks, validation_masks, _, _ =
train_test_split(attention_masks, labels,
random_state=2018,
test_size=0.1)
```

## 3.6. Converting to PyTorch Data Types

Our model expects PyTorch tensors rather than numpy.ndarrays, so convert all of our dataset variables.



[Get started](#)[Open in app](#)

```
# Convert all inputs and labels into torch tensors, the required
datatype
# for our model.
train_inputs = torch.tensor(train_inputs)
validation_inputs = torch.tensor(validation_inputs)

train_labels = torch.tensor(train_labels)
validation_labels = torch.tensor(validation_labels)

train_masks = torch.tensor(train_masks)
validation_masks = torch.tensor(validation_masks)
```

We'll also create an iterator for our dataset using the torch DataLoader class. This helps save on memory during training because, unlike a for loop, with an iterator the entire dataset does not need to be loaded into memory.

```
from torch.utils.data import TensorDataset, DataLoader,
RandomSampler, SequentialSampler

# The DataLoader needs to know our batch size for training, so we
specify it
# here.
# For fine-tuning BERT on a specific task, the authors recommend a
batch size of
# 16 or 32.

batch_size = 32

# Create the DataLoader for our training set.
train_data = TensorDataset(train_inputs, train_masks, train_labels)
train_sampler = RandomSampler(train_data)
train_dataloader = DataLoader(train_data, sampler=train_sampler,
batch_size=batch_size)

# Create the DataLoader for our validation set.
validation_data = TensorDataset(validation_inputs, validation_masks,
validation_labels)
validation_sampler = SequentialSampler(validation_data)
```

[Get started](#)[Open in app](#)

## 4. Train Our Classification Model

Now that our input data is properly formatted, it's time to fine tune the BERT model.

### 4.1. BertForSequenceClassification

For this task, we first want to modify the pre-trained BERT model to give outputs for classification, and then we want to continue training the model on our dataset until that the entire model, end-to-end, is well-suited for our task.

Thankfully, the huggingface pytorch implementation includes a set of interfaces designed for a variety of NLP tasks. Though these interfaces are all built on top of a trained BERT model, each has different top layers and output types designed to accomodate their specific NLP task.



[Get started](#)[Open in app](#)

We'll be using [Bert Classification Model](#). This is the normal BERT model with an added single linear layer on top for classification that we will use as a sentence classifier. As we feed input data, the entire pre-trained BERT model and the additional untrained classification layer is trained on our specific task.

OK, let's load BERT! There are a few different pre-trained BERT models available. "bert-base-uncased" means the version that has only lowercase letters ("uncased") and is the smaller version of the two ("base" vs "large").

```
from transformers import BertForSequenceClassification, AdamW,
BertConfig

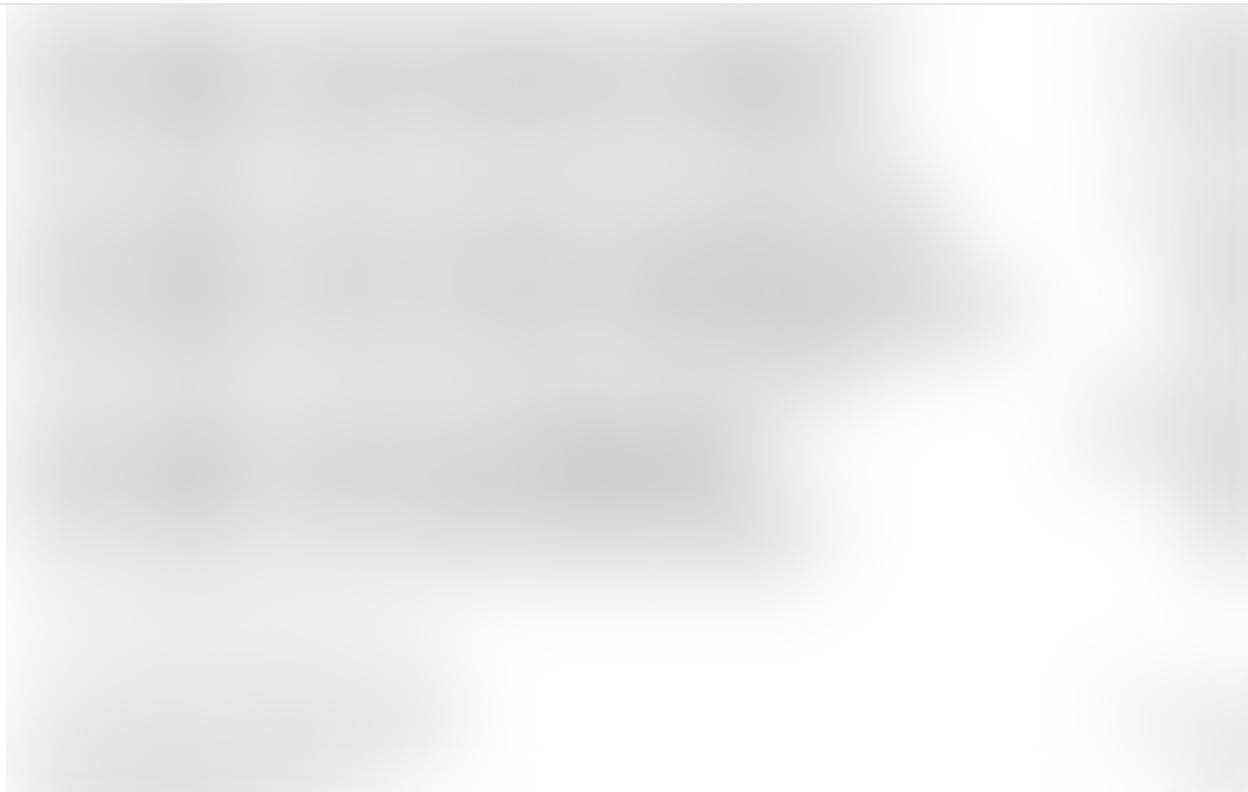
# Load BertForSequenceClassification, the pretrained BERT model with
# a single
# linear classification layer on top.
model = BertForSequenceClassification.from_pretrained(
    "bert-base-uncased", # Use the 12-layer BERT model, with an
    uncased vocab.
    num_labels = 2, # The number of output labels--2 for binary
    classification.
                # You can increase this for multi-class tasks.
    output_attentions = False, # Whether the model returns attentions
    weights.
    output_hidden_states = False, # Whether the model returns all
    hidden-states.
)

# Tell pytorch to run this model on the GPU.
model.cuda()
```

[Get started](#)[Open in app](#)

In the below cell we can check the names and dimensions of the weights for: The embedding layer, The first of the twelve transformers & The output layer.

```
# Get all of the model's parameters as a list of tuples.  
params = list(model.named_parameters())  
  
print('The BERT model has {} different named  
parameters.\n'.format(len(params)))  
  
print('==== Embedding Layer ====\n')  
  
for p in params[0:5]:  
    print("{}:{:<55} {:>12}".format(p[0], str(tuple(p[1].size()))))  
  
print('\n==== First Transformer ====\n')  
  
for p in params[5:21]:  
    print("{}:{:<55} {:>12}".format(p[0], str(tuple(p[1].size()))))  
  
print('\n==== Output Layer ====\n')  
  
for p in params[-4:]:  
    print("{}:{:<55} {:>12}".format(p[0], str(tuple(p[1].size()))))
```

[Get started](#)[Open in app](#)

So we can see the weight and bias of the Layers respectively.

## 4.2. Optimizer & Learning Rate Scheduler

Now that we have our model loaded we need to grab the training hyperparameters from within the stored model.

For the purposes of fine-tuning, the authors recommend choosing from the following values:

- Batch size: 16, 32 (We chose 32 when creating our DataLoaders).
- Learning rate (Adam): 5e-5, 3e-5, 2e-5 (We'll use 2e-5).
- Number of epochs: 2, 3, 4 (We'll use 4).

The epsilon parameter  $\text{eps} = 1\text{e-}8$  is “a very small number to prevent any division by zero in the implementation”

[Get started](#)[Open in app](#)

You can find the creation of the AdamW optimizer in `run_glue.py` [Click here.](#)

```
# Note: AdamW is a class from the huggingface library (as opposed to pytorch)
# I believe the 'W' stands for 'Weight Decay fix"
optimizer = AdamW(model.parameters(),
                  lr = 2e-5, # args.learning_rate - default is 5e-5,
# our notebook had 2e-5
                  eps = 1e-8 # args.adam_epsilon - default is 1e-8.
                  )

from transformers import get_linear_schedule_with_warmup

# Number of training epochs (authors recommend between 2 and 4)
epochs = 4

# Total number of training steps is number of batches * number of epochs.
total_steps = len(train_dataloader) * epochs

# Create the learning rate scheduler.
scheduler = get_linear_schedule_with_warmup(optimizer,
                                              num_warmup_steps = 0, #
# Default value in run_glue.py
                                              num_training_steps =
total_steps)
```

[Get started](#)[Open in app](#)

Below is our training loop. There's a lot going on, but fundamentally for each pass in our loop we have a trianing phase and a validation phase. At each pass we need to:

Training loop:

- Unpack our data inputs and labels
- Load data onto the GPU for acceleration
- Clear out the gradients calculated in the previous pass.
- In pytorch the gradients accumulate by default (useful for things like RNNs) unless you explicitly clear them out.
- Forward pass (feed input data through the network)
- Backward pass (backpropagation)
- Tell the network to update parameters with `optimizer.step()`
- Track variables for monitoring progress

Evalution loop:

- Unpack our data inputs and labels
- Load data onto the GPU for acceleration
- Forward pass (feed input data through the network)
- Compute loss on our validation data and track variables for monitoring progress

Define a helper function for calculating accuracy.

```
import numpy as np

# Function to calculate the accuracy of our predictions vs labels
def flat_accuracy(preds, labels):
    pred_flat = np.argmax(preds, axis=1).flatten()
```

[Get started](#)[Open in app](#)

Helper function for formatting elapsed times.

```
import time
import datetime

def format_time(elapsed):
    """
    Takes a time in seconds and returns a string hh:mm:ss
    """
    # Round to the nearest second.
    elapsed_rounded = int(round((elapsed)))

    # Format as hh:mm:ss
    return str(datetime.timedelta(seconds=elapsed_rounded))
```

We're ready to kick off the training!

```
import random

# This training code is based on the `run_glue.py` script here:
#
https://github.com/huggingface/transformers/blob/5bfcd0485ece086ebcbe\_d2d008813037968a9e58/examples/run\_glue.py#L128

# Set the seed value all over the place to make this reproducible.
seed_val = 42

random.seed(seed_val)
np.random.seed(seed_val)
torch.manual_seed(seed_val)
torch.cuda.manual_seed_all(seed_val)

# Store the average loss after each epoch so we can plot them.
loss_values = []

# For each epoch...
for epoch_i in range(0, epochs):

    # =====#
    #       Training
    # =====#
```

[Get started](#)[Open in app](#)

```
print("")  
print('===== Epoch {} / {} ====='.format(epoch_i + 1,  
epochs))  
print('Training...')  
  
# Measure how long the training epoch takes.  
t0 = time.time()  
  
# Reset the total loss for this epoch.  
total_loss = 0  
  
# Put the model into training mode. Don't be mislead--the call to  
# `train` just changes the *mode*, it doesn't *perform* the  
training.  
# `dropout` and `batchnorm` layers behave differently during  
training  
# vs. test (source:  
https://stackoverflow.com/questions/51433378/what-does-model-train-do-in-pytorch)  
model.train()  
  
# For each batch of training data...  
for step, batch in enumerate(train_dataloader):  
  
    # Progress update every 40 batches.  
    if step % 40 == 0 and not step == 0:  
        # Calculate elapsed time in minutes.  
        elapsed = format_time(time.time() - t0)  
  
        # Report progress.  
        print(' Batch {:>5,} of {:>5,}. Elapsed:  
{}.'.format(step, len(train_dataloader), elapsed))  
  
    # Unpack this training batch from our dataloader.  
    #  
    # As we unpack the batch, we'll also copy each tensor to the  
GPU using the  
    # `to` method.  
    #  
    # `batch` contains three pytorch tensors:  
    # [0]: input ids  
    # [1]: attention masks  
    # [2]: labels  
    b_input_ids = batch[0].to(device)  
    b_input_mask = batch[1].to(device)  
    b_labels = batch[2].to(device)  
  
    # Always clear any previously calculated gradients before  
performing a  
    # backward pass. PyTorch doesn't do this automatically
```

[Get started](#)[Open in app](#)

```

# (source: https://stackoverflow.com/questions/40001590/why-
do-we-need-to-call-zero-grad-in-pytorch)
model.zero_grad()

# Perform a forward pass (evaluate the model on this training
batch).
# This will return the loss (rather than the model output)
because we
# have provided the `labels`.
# The documentation for this `model` function is here:
#
https://huggingface.co/transformers/v2.2.0/model\_doc/bert.html#transformers.BertForSequenceClassification
outputs = model(b_input_ids,
                 token_type_ids=None,
                 attention_mask=b_input_mask,
                 labels=b_labels)

# The call to `model` always returns a tuple, so we need to
pull the
# loss value out of the tuple.
loss = outputs[0]

# Accumulate the training loss over all of the batches so
that we can
# calculate the average loss at the end. `loss` is a Tensor
containing a
# single value; the ` `.item()` ` function just returns the
Python value
# from the tensor.
total_loss += loss.item()

# Perform a backward pass to calculate the gradients.
loss.backward()

# Clip the norm of the gradients to 1.0.
# This is to help prevent the "exploding gradients" problem.
torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0)

# Update parameters and take a step using the computed
gradient.
# The optimizer dictates the "update rule"--how the
parameters are
# modified based on their gradients, the learning rate, etc.
optimizer.step()

# Update the learning rate.
scheduler.step()

```

[Get started](#)[Open in app](#)

```
# Score the loss value for plotting the learning curve.
loss_values.append(avg_train_loss)

print("")
print(" Average training loss: {:.2f}".format(avg_train_loss))
print(" Training epoch took: {}".format(format_time(time.time() - t0)))

# =====
# Validation
# =====
# After the completion of each training epoch, measure our
# performance on
# our validation set.

print("")
print("Running Validation...")

t0 = time.time()

# Put the model in evaluation mode--the dropout layers behave
# differently
# during evaluation.
model.eval()

# Tracking variables
eval_loss, eval_accuracy = 0, 0
nb_eval_steps, nb_eval_examples = 0, 0

# Evaluate data for one epoch
for batch in validation_dataloader:

    # Add batch to GPU
    batch = tuple(t.to(device) for t in batch)

    # Unpack the inputs from our dataloader
    b_input_ids, b_input_mask, b_labels = batch

    # Telling the model not to compute or store gradients, saving
    # memory and
    # speeding up validation
    with torch.no_grad():

        # Forward pass, calculate logit predictions.
        # This will return the logits rather than the loss
because we have
        # not provided labels.
        # token_type_ids is the same as the "segment ids", which
        # differentiates sentence 1 and 2 in 2-sentence tasks.
        # The documentation for this `model` function is here:
```

[Get started](#)[Open in app](#)

```
outputs = model(b_input_ids,
                 token_type_ids=None,
                 attention_mask=b_input_mask)

# Get the "logits" output by the model. The "logits" are the
output
# values prior to applying an activation function like the
softmax.
logits = outputs[0]

# Move logits and labels to CPU
logits = logits.detach().cpu().numpy()
label_ids = b_labels.to('cpu').numpy()

# Calculate the accuracy for this batch of test sentences.
tmp_eval_accuracy = flat_accuracy(logits, label_ids)

# Accumulate the total accuracy.
eval_accuracy += tmp_eval_accuracy

# Track the number of batches
nb_eval_steps += 1

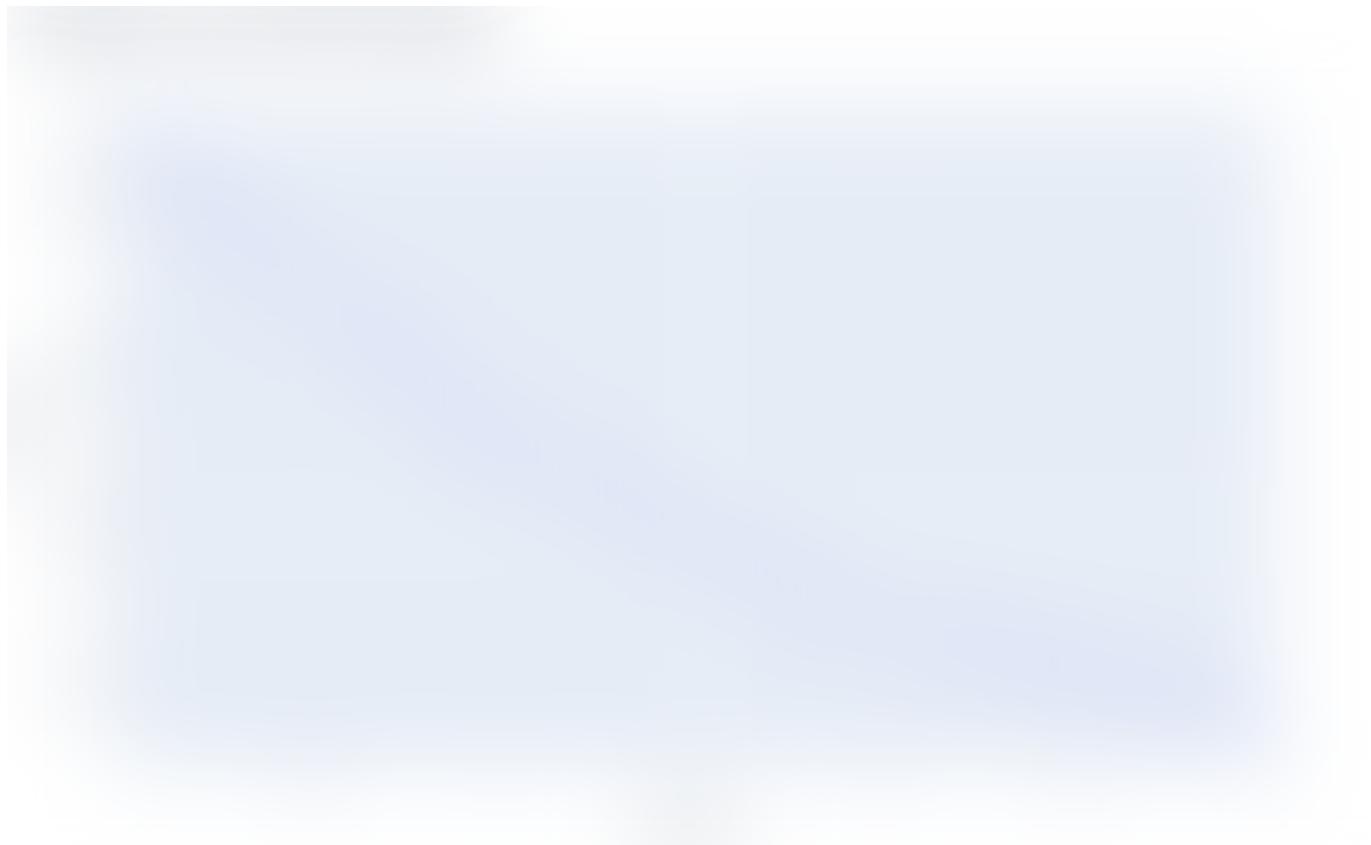
# Report the final accuracy for this validation run.
print(" Accuracy: {:.2f}%".format(eval_accuracy/nb_eval_steps))
print(" Validation took: {}".format(format_time(time.time() - t0)))

print("")
print("Training complete!")
```

[Get started](#)[Open in app](#)

```
import plotly.express as px

f = pd.DataFrame(loss_values)
f.columns=['Loss']
fig = px.line(f, x=f.index, y=f.Loss)
fig.update_layout(title='Training loss of the Model',
                  xaxis_title='Epoch',
                  yaxis_title='Loss')
fig.show()
```



## 5. Performance On Test Set

Now we'll load the holdout dataset and prepare inputs just as we did with the training set. Then we'll evaluate predictions using [Matthews correlation coefficient \(MCC wiki\)](#) because this is the metric used by the wider NLP community to evaluate performance on CoLA. With this metric, +1 is the best score, and -1 is the worst score. This way, we can see how well we perform against the state of the art models for this specific task.

[Get started](#)[Open in app](#)

## 5.1. Data Preparation

We'll need to apply all of the same steps that we did for the training data to prepare our test data set.

```
import pandas as pd

# Load the dataset into a pandas dataframe.
df = pd.read_csv("./cola_public/raw/out_of_domain_dev.tsv",
delimiter='\t', header=None, names=['sentence_source', 'label',
'label_notes', 'sentence'])

# Report the number of sentences.
print('Number of test sentences: {} \n'.format(df.shape[0]))

# Create sentence and label lists
sentences = df.sentence.values
labels = df.label.values

# Tokenize all of the sentences and map the tokens to thier word IDs.
input_ids = []

# For every sentence...
for sent in sentences:
    # `encode` will:
    #   (1) Tokenize the sentence.
    #   (2) Prepend the `[CLS]` token to the start.
    #   (3) Append the `[SEP]` token to the end.
    #   (4) Map tokens to their IDs.
    encoded_sent = tokenizer.encode(
```

[Get started](#)[Open in app](#)

```
[ ०५८ ]
)

input_ids.append(encoded_sent)

# Pad our input tokens
input_ids = pad_sequences(input_ids, maxlen=MAX_LEN,
                           dtype="long", truncating="post",
                           padding="post")

# Create attention masks
attention_masks = []

# Create a mask of 1s for each token followed by 0s for padding
for seq in input_ids:
    seq_mask = [float(i>0) for i in seq]
    attention_masks.append(seq_mask)

# Convert to tensors.
prediction_inputs = torch.tensor(input_ids)
prediction_masks = torch.tensor(attention_masks)
prediction_labels = torch.tensor(labels)

# Set the batch size.
batch_size = 32

# Create the DataLoader.
prediction_data = TensorDataset(prediction_inputs, prediction_masks,
                                prediction_labels)
prediction_sampler = SequentialSampler(prediction_data)
prediction_dataloader = DataLoader(prediction_data,
                                   sampler=prediction_sampler, batch_size=batch_size)
```

## 5.2. Evaluate on Test Set

With the test set prepared, we can apply our fine-tuned model to generate predictions on the test set.

```
# Prediction on test set
```

[Get started](#)[Open in app](#)

```
# Put model in evaluation mode
model.eval()

# Tracking variables
predictions, true_labels = [], []

# Predict
for batch in prediction_dataloader:
    # Add batch to GPU
    batch = tuple(t.to(device) for t in batch)

    # Unpack the inputs from our dataloader
    b_input_ids, b_input_mask, b_labels = batch

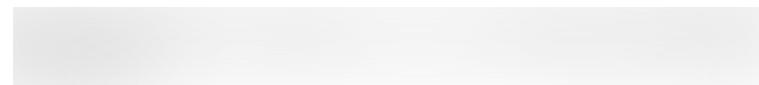
    # Telling the model not to compute or store gradients, saving
    memory and
    # speeding up prediction
    with torch.no_grad():
        # Forward pass, calculate logit predictions
        outputs = model(b_input_ids, token_type_ids=None,
                        attention_mask=b_input_mask)

    logits = outputs[0]

    # Move logits and labels to CPU
    logits = logits.detach().cpu().numpy()
    label_ids = b_labels.to('cpu').numpy()

    # Store predictions and true labels
    predictions.append(logits)
    true_labels.append(label_ids)

print('DONE.')
```



Accuracy on the CoLA benchmark is measured using the Matthews correlation coefficient, We use MCC here because the classes are imbalanced:

```
print('Positive samples: %d of %d (%.2f%%)' % (df.label.sum(),
len(df.label), (df.label.sum() / len(df.label) * 100.0)))
```

[Get started](#)[Open in app](#)

```
from sklearn.metrics import matthews_corrcoef

matthews_set = []

# Evaluate each test batch using Matthew's correlation coefficient
print('Calculating Matthews Corr. Coef. for each batch...')

# For each input batch...
for i in range(len(true_labels)):

    # The predictions for this batch are a 2-column ndarray (one column
    for "0"
    # and one column for "1"). Pick the label with the highest value
    and turn this
    # in to a list of 0s and 1s.
    pred_labels_i = np.argmax(predictions[i], axis=1).flatten()

    # Calculate and store the coef for this batch.
    matthews = matthews_corrcoef(true_labels[i], pred_labels_i)
    matthews_set.append(matthews)
```

The final score will be based on the entire test set, but let's take a look at the scores on the individual batches to get a sense of the variability in the metric between batches. Each batch has 32 sentences in it, except the last batch which has only  $(516 \% 32) = 4$  test sentences in it.

```
matthews_set
```

[Get started](#)[Open in app](#)

```
# Combine the predictions for each batch into a single list of 0s and
# 1s.
flat_predictions = [item for sublist in predictions for item in
                     sublist]
flat_predictions = np.argmax(flat_predictions, axis=1).flatten()

# Combine the correct labels for each batch into a single list.
flat_true_labels = [item for sublist in true_labels for item in
                     sublist]

# Calculate the MCC
mcc = matthews_corrcoef(flat_true_labels, flat_predictions)

print('MCC: %.3f' % mcc)
```

So without doing any hyperparameter tuning (adjusting the learning rate, epochs, batch size, ADAM properties, etc.) we are able to get a good score. we didn't train on the entire training dataset, but set aside a portion of it as our validation set for legibility of code.

Note that (due to the small dataset size?) the accuracy can vary significantly with different random seeds.

## Conclusion

This post demonstrates that with a pre-trained BERT model you can quickly and effectively create a high quality model with minimal effort and training time using the pytorch interface, regardless of the specific NLP task you are interested in.

[Get started](#)[Open in app](#)

## *Sentiment Analysis on Aero Industry Customer Datasets on Twitter using BERT & XLNET.*

### References:

- [Attention Is All You Need](#); Vaswani et al., 2017.
- [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#); Devlin et al., 2018.
- [The Annotated Transformer](#); Alexander Rush, Vincent Nguyen and Guillaume Klein.
- [Universal Language Model Fine-tuning for Text Classification](#); Howard et al., 2018.
- [Improving Language Understanding by Generative Pre-Training](#); Radford et al., 2018.
- [Encoder-Decoder Architecture & Bert Paper](#) on the full research.
- Source of cover picture: Google sources.

[NLP](#)[Machine Learning](#)[Deep Learning](#)[Python Programming](#)[Artificial Intelligence](#)[About](#) [Help](#) [Legal](#)

Get the Medium app



