

In this notebook, You will do amazon review classification with BERT. It contains 5 parts as below. Detailed instructions are given in the each cell. please read every comment we have written.

1. Preprocessing
2. Creating a BERT model from the Tensorflow HUB.
3. Tokenization
4. getting the pretrained embedding Vector for a given review from the BERT.
5. Using the embedding data apply NN and classify the reviews.

instructions:

1. Don't change any Grader Functions. Don't manipulate any Grader functions. If you manipulate any, it will be considered as plagiarised.
2. Please read the instructions on the code cells and markdown cells. We will explain what to write.
3. please return outputs in the same format what we asked. Eg. Don't return List if we are asking for a numpy array.
4. Please read the external links that we are given so that you will learn the concept behind the code that you are writing.
5. We are giving instructions at each section if necessary, please follow them.

Every Grader function has to return True.

```
In [ ]: 1 #all imports
        2 import numpy as np
        3 import pandas as pd
        4 import tensorflow as tf
        5 import tensorflow_hub as hub
        6 from tensorflow.keras.models import Model
```

```
In [ ]: 1 tf.test.gpu_device_name()
```

Grader function 1

```
In [ ]: 1 def grader_tf_version():
        2     assert((tf.__version__)>'2')
        3     return True
        4 grader_tf_version()
```

Part-1: Preprocessing



```
In [ ]: 1 #Read the dataset - Amazon fine food reviews
        2 reviews = pd.read_csv(r"D:\ML\Internal DL\NLP\amazon-fine-food-reviews\Reviews.csv")
        3 #check the info of the dataset
        4 reviews.info()
```

```
In [ ]: 1 #get only 2 columns - Text, Score
        2 #drop the NAN values
```

```
In [ ]: 1 #if score > 3, set score = 1
        2 #if score <= 2, set score = 0
        3 #if score == 3, remove the rows.
```

Grader function 2

```
In [ ]: 1 def grader_reviews():
        2     temp_shape = (reviews.shape == (525814, 2)) and (reviews.Score.value_counts().max() < 4)
        3     assert(temp_shape == True)
        4     return True
        5 grader_reviews()
```

```
In [ ]: 1 def get_wordlen(x):
        2     return len(x.split())
        3 reviews['len'] = reviews.Text.apply(get_wordlen)
        4 reviews = reviews[reviews.len < 50]
        5 reviews = reviews.sample(n=100000, random_state=30)
```

```
In [ ]: 1 #remove HTML from the Text column and save in the Text column only
```

```
In [ ]: 1 #print head 5
```

```
In [ ]: 1 #split the data into train and validation data(20%) with Stratify sampling,
```

```
In [ ]: 1 #plot bar graphs of y_train and y_test
```

```
In [ ]: 1 #saving to disk. if we need, we can load preprocessed data directly.
        2 reviews.to_csv('preprocessed.csv', index=False)
```

Part-2: Creating BERT Model

If you want to know more about BERT, You can watch live sessions on Transformers and BERT. we will strongly recommend you to read [Transformers](https://jalammar.github.io/illustrated-transformer/) (<https://jalammar.github.io/illustrated-transformer/>), [BERT Paper](https://arxiv.org/abs/1810.04805) (<https://arxiv.org/abs/1810.04805>) and, [This blog](https://jalammar.github.io/illustrated-transformer/) (<https://jalammar.github.io/illustrated-transformer/>)

[o/a-visual-guide-to-using-bert-for-the-first-time/](#)).

For this assignment, we are using [BERT uncased Base model \(https://tfhub.dev/tensorflow/bert_en_uncased_L-12_H-768_A-12/1\)](https://tfhub.dev/tensorflow/bert_en_uncased_L-12_H-768_A-12/1). It uses L=12 hidden layers (i.e., Transformer blocks), a hidden size of H=768, and A=12 attention heads.

```
In [ ]: 1  ## Loading the Pretrained Model from tensorflow HUB
2  tf.keras.backend.clear_session()
3
4  # maximum length of a seq in the data we have, for now i am making it as 300
5  max_seq_length = 55
6
7  #BERT takes 3 inputs
8
9  ##this is input words. Sequence of words represented as integers
10 input_word_ids = tf.keras.layers.Input(shape=(max_seq_length,), dtype=tf.int
11
12 #mask vector if you are padding anything
13 input_mask = tf.keras.layers.Input(shape=(max_seq_length,), dtype=tf.int32,
14
15 #segment vectors. If you are giving only one sentence for the classification
16 #If you are giving two sentences with [sep] token separated, first seq segme
17 #second seq segment vector are 1's
18 segment_ids = tf.keras.layers.Input(shape=(max_seq_length,), dtype=tf.int32,
19
20 #bert layer
21 bert_layer = hub.KerasLayer("https://tfhub.dev/tensorflow/bert_en_uncased_L-
22 pooled_output, sequence_output = bert_layer([input_word_ids, input_mask, seg
23
24 #Bert model
25 #We are using only pooled output not sequence out.
26 #If you want to know about those, please read https://www.kaggle.com/questio
27 bert_model = Model(inputs=[input_word_ids, input_mask, segment_ids], outputs
28
```

```
In [ ]: 1  bert_model.summary()
```

```
In [ ]: 1  bert_model.output
```

Part-3: Tokenization



```
In [ ]: 1  #getting Vocab file
2  vocab_file = bert_layer.resolved_object.vocab_file.asset_path.numpy()
3  do_lower_case = bert_layer.resolved_object.do_lower_case.numpy()
```

```
In [ ]: 1  #import tokenization - We have given tokenization.py file
```

```
In [ ]: 1 # Create tokenizer " Instantiate FullTokenizer"
2 # name must be "tokenizer"
3 # the FullTokenizer takes two parameters 1. vocab_file and 2. do_lower_case
4 # we have created these in the above cell ex: FullTokenizer(vocab_file, do_L
5 # please check the "tokenization.py" file the complete implementation
```

Grader function 3

```
In [ ]: 1 #it has to give no error
2 def grader_tokenize(tokenizer):
3     out = False
4     try:
5         out=('[CLS]' in tokenizer.vocab) and ('[SEP]' in tokenizer.vocab)
6     except:
7         out = False
8     assert(out==True)
9     return out
10 grader_tokenize(tokenizer)
```

```
In [ ]: 1 # Create train and test tokens (X_train_tokens, X_test_tokens) from (X_train
2
3 # add '[CLS]' at start of the Tokens and '[SEP]' at the end of the tokens.
4
5 # maximum number of tokens is 55(We already given this to BERT layer above)
6
7 # if it is less than 55, add '[PAD]' token else truncate the tokens length.(
8
9 # Based on padding, create the mask for Train and Test ( 1 for real token, 0
10 # it will also same shape as input tokens (None, 55) save those in X_train_m
11
12 # Create a segment input for train and test. We are using only one sentence
13
14 # type of all the above arrays should be numpy arrays
15
16 # after execution of this cell, you have to get
17 # X_train_tokens, X_train_mask, X_train_segment
18 # X_test_tokens, X_test_mask, X_test_segment
```

Example


```

In [ ]: 1 def grader_alltokens_train():
2         out = False
3
4         if type(X_train_tokens) == np.ndarray:
5
6             temp_shapes = (X_train_tokens.shape[1]==max_seq_length) and (X_train
7             (X_train_segment.shape[1]==max_seq_length)
8
9             segment_temp = not np.any(X_train_segment)
10
11             mask_temp = np.sum(X_train_mask==0) == np.sum(X_train_tokens==0)
12
13             no_cls = np.sum(X_train_tokens==tokenizer.vocab['[CLS]'])==X_train_t
14
15             no_sep = np.sum(X_train_tokens==tokenizer.vocab['[SEP]'])==X_train_t
16
17             out = temp_shapes and segment_temp and mask_temp and no_cls and no_s
18
19         else:
20             print('Type of all above token arrays should be list not numpy array
21             out = False
22             assert(out==True)
23             return out
24
25 grader_alltokens_train()

```

Grader function 5

```

In [ ]: 1 def grader_alltokens_test():
2         out = False
3         if type(X_test_tokens) == np.ndarray:
4
5             temp_shapes = (X_test_tokens.shape[1]==max_seq_length) and (X_test_m
6             (X_test_segment.shape[1]==max_seq_length)
7
8             segment_temp = not np.any(X_test_segment)
9
10            mask_temp = np.sum(X_test_mask==0) == np.sum(X_test_tokens==0)
11
12            no_cls = np.sum(X_test_tokens==tokenizer.vocab['[CLS]'])==X_test_tok
13
14            no_sep = np.sum(X_test_tokens==tokenizer.vocab['[SEP]'])==X_test_tok
15
16            out = temp_shapes and segment_temp and mask_temp and no_cls and no_s
17
18        else:
19            print('Type of all above token arrays should be list not numpy array
20            out = False
21            assert(out==True)
22            return out
23 grader_alltokens_test()

```

Part-4: Getting Embeddings from BERT M

odel

We already created the BERT model in the part-2 and input data in the part-3. We will utilize those two and will get the embeddings for each sentence in the Train and Validation data.

```
In [ ]: 1 bert_model.input
```

```
In [ ]: 1 bert_model.output
```

```
In [ ]: 1 # get the train output, BERT model will give one output so save in  
2 # X_train_pooled_output
```

```
In [ ]: 1 # get the test output, BERT model will give one output so save in  
2 # X_test_pooled_output
```

```
In [ ]: 1 ##save all your results to disk so that, no need to run all again.  
2 pickle.dump((X_train_pooled_output, X_test_pooled_output),open('final_output
```

```
In [ ]: 1 #X_train_pooled_output, X_test_pooled_output= pickle.load(open('final_output
```

Grader function 6

```
In [ ]: 1 #now we have X_train_pooled_output, y_train  
2 #X_test_pooled_output, y_test  
3  
4 #please use this grader to evaluate  
5 def greader_output():  
6     assert(X_train_pooled_output.shape[1]==768)  
7     assert(len(y_train)==len(X_train_pooled_output))  
8     assert(X_test_pooled_output.shape[1]==768)  
9     assert(len(y_test)==len(X_test_pooled_output))  
10    assert(len(y_train.shape)==1)  
11    assert(len(X_train_pooled_output.shape)==2)  
12    assert(len(y_test.shape)==1)  
13    assert(len(X_test_pooled_output.shape)==2)  
14    return True  
15 greader_output()
```

Part-5: Training a NN with 786 features

Create a NN and train the NN.

1. You have to use AUC as metric.
2. You can use any architecture you want.
3. You have to use tensorboard to log all your metrics and Losses. You have to send those logs.

4. Print the loss and metric at every epoch.
5. You have to submit without overfitting and underfitting.

```
In [ ]: 1 ##imports
        2 from tensorflow.keras.layers import Input, Dense, Activation, Dropout
        3 from tensorflow.keras.models import Model
```

```
In [ ]: 1 ##create an NN and
        2
```

```
In [ ]: 1
```