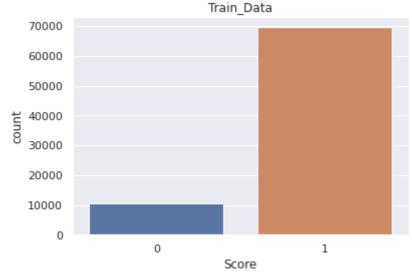
Dataset: Link(https://www.kaggle.com/snap/amazon-fine-food-reviews/data)

```
import keras
import pickle
import pandas as pd
import numpy as np
import tensorflow as tf
import tensorflow_hub as hub
from keras import layers, Model
import torch
from keras.preprocessing.sequence import pad_sequences
import warnings
warnings.filterwarnings('ignore')
print(tf.test.gpu_device_name())
     /device:GPU:0
def grader_tf_version():
    assert((tf.__version__)>'2')
    return True
grader_tf_version()
     True
!wget --header="Host: storage.googleapis.com" --header="User-Agent: Mozilla/5.0 (Windows N
!unzip /content/archive.zip
#Read the dataset - Amazon fine food reviews
reviews = pd.read csv("/content/Reviews.csv")
#check the info of the dataset
reviews.info()
     <class 'pandas.core.frame.DataFrame'>
     RangeIndex: 568454 entries, 0 to 568453
     Data columns (total 10 columns):
         Column
      #
                                  Non-Null Count
                                                   Dtype
         _____
                                  -----
                                                   ----
                                  568454 non-null int64
      0
          Ιd
      1
         ProductId
                                  568454 non-null object
      2
         UserId
                                  568454 non-null object
                                  568438 non-null object
      3
          ProfileName
         HelpfulnessNumerator
                                  568454 non-null int64
      5
         HelpfulnessDenominator 568454 non-null int64
      6
         Score
                                  568454 non-null int64
      7
                                  568454 non-null int64
          Time
      8
          Summary
                                  568427 non-null object
          Text
                                  568454 non-null object
     dtypes: int64(5), object(5)
     memory usage: 43.4+ MB
```

```
review = reviews.drop(['Id','ProductId','UserId','ProfileName','HelpfulnessNumerator','Hel
review.info()
     <class 'pandas.core.frame.DataFrame'>
     RangeIndex: 568454 entries, 0 to 568453
     Data columns (total 2 columns):
         Column Non-Null Count
         -----
          Score 568454 non-null int64
      1
          Text
                 568454 non-null object
     dtypes: int64(1), object(1)
     memory usage: 8.7+ MB
reviews = review[review['Score']!=3]
reviews.head(2)
         Score
                                                     Text
      0
             5
                  I have bought several of the Vitality canned d...
             1 Product arrived labeled as Jumbo Salted Peanut...
      1
scores=[]
for score in reviews['Score']:
  if score>3:
    scores.append(1)
  if score<=2:
    scores.append(0)
reviews['Score']=scores
reviews.info()
     <class 'pandas.core.frame.DataFrame'>
     Int64Index: 525814 entries, 0 to 568453
     Data columns (total 2 columns):
         Column Non-Null Count Dtype
                 -----
          Score 525814 non-null int64
                  525814 non-null object
          Text
     dtypes: int64(1), object(1)
     memory usage: 12.0+ MB
def grader_reviews():
    temp shape = (reviews.shape == (525814, 2)) and (reviews.Score.value counts()[1]==4437
    assert(temp shape == True)
    return True
grader reviews()
     True
def get wordlen(x):
    return len(x.split())
reviews['len'] = reviews.Text.apply(get_wordlen)
reviews = reviews[reviews.len<50]</pre>
reviews = reviews.sample(n=100000, random state=30)
```

```
def remove_html(text):
  for i in range(len(text)):
    html = re.compile('<.*?>')
    text[i] = re.sub(html,' ',text[i])
  return text
reviews['Text']=remove_html(reviews['Text'].values)
reviews.to_csv('/content/drive/MyDrive/Assignment/NLP with Transfer Learning/preprocessed.
y = reviews['Score']
X = reviews['Text']
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(X,y, test_size=0.2 , stratify = y, random
X_train.head(2)
     33523
              I had never tried this brand before, so I was ...
              I love these for a snack. I get a nice taste o...
     Name: Text, dtype: object
y_train.value_counts()
     1
          69603
          10397
     Name: Score, dtype: int64
import seaborn as sns
sns.set_theme(style="darkgrid")
sns.countplot(y_train).set_title('Train_Data')
     Text(0.5, 1.0, 'Train Data')
```



```
sns.set_theme(style="darkgrid")
sns.countplot(y_test).set_title('Test_Data')
```

Text(0.5, 1.0, 'Test\_Data')



```
length=[]
for t in X_train:
    length.append(len(t.split()))
max_seq_length = max(length)
print(max_seq_length)
```

## Loading the Pretrained Model from tensorflow HUB
tf.keras.backend.clear\_session()

# maximum length of a seq in the data we have, for now i am making it as 55. You can chang max seq length = 49

#BERT takes 3 inputs

49

#this is input words. Sequence of words represented as integers
input\_word\_ids = tf.keras.layers.Input(shape=(max\_seq\_length,), dtype=tf.int32, name="inpu

#mask vector if you are padding anything

input\_mask = tf.keras.layers.Input(shape=(max\_seq\_length,), dtype=tf.int32, name="input\_ma

#segment vectors. If you are giving only one sentence for the classification, total seg ve #If you are giving two sentenced with [sep] token separated, first seq segment vectors are #second seq segment vector are 1's

segment\_ids = tf.keras.layers.Input(shape=(max\_seq\_length,), dtype=tf.int32, name="segment")

#bert layer

bert\_layer = hub.KerasLayer("https://tfhub.dev/tensorflow/bert\_en\_uncased\_L-12\_H-768\_A-12/
pooled\_output, sequence\_output = bert\_layer([input\_word\_ids, input\_mask, segment\_ids])

#Bert model

#We are using only pooled output not sequence out.

#If you want to know about those, please read <a href="https://www.kaggle.com/questions-and-answers">https://www.kaggle.com/questions-and-answers</a> bert\_model = Model(inputs=[input\_word\_ids, input\_mask, segment\_ids], outputs=pooled\_output

```
bert model.summary()
     Model: "model"
     Layer (type)
                                     Output Shape
                                                           Param #
                                                                       Connected to
     Total params: 109,482,241
     Trainable params: 0
     Non-trainable params: 109,482,241
bert_model.output
     <KerasTensor: shape=(None, 768) dtype=float32 (created by layer 'keras_layer')>
vocab_file = bert_layer.resolved_object.vocab_file.asset_path.numpy()
do_lower_case = bert_layer.resolved_object.do_lower_case.numpy()
!pip install bert-for-tf2
from bert import bert_tokenization
     Installing collected packages: py-params, params-flow, bert-for-tf2
     Successfully installed bert-for-tf2-0.14.9 params-flow-0.8.2 py-params-0.10.2
tokenizer = bert_tokenization.FullTokenizer(vocab_file,do_lower_case)
def grader tokenize(tokenizer):
    out = False
    try:
        out=('[CLS]' in tokenizer.vocab) and ('[SEP]' in tokenizer.vocab)
    except:
        out = False
    assert(out==True)
    return out
grader tokenize(tokenizer)
     True
X_train_tokens,X_train_mask, X_train_segment=[],[],[]
for i in range(len(X_train)):
  tokens = tokenizer.tokenize(X train.values[i])
  if len(tokens)<(max_seq_length-2):</pre>
    tokens=tokens
  else:
    tokens = tokens[0:(max seq length-2)]
  tokens =['[CLS]',*tokens,'[SEP]']
  X_train_tokens.append(np.array(tokenizer.convert_tokens_to_ids(tokens)))
  X_train_mask.append(np.array([1]*len(tokens)+[0]*(max_seq_length-len(tokens))))
  X_train_segment.append([0]*max_seq_length)
X_train_token=pad_sequences(X_train_tokens , maxlen=max_seq_length, padding='post')
X train tokens = np.array(X train token)
```

```
X_train_mask=np.array(X_train_mask)
X_train_segment = np.array(X_train_segment)
def grader_alltokens_train():
    out = False
    if type(X_train_tokens) == np.ndarray:
        temp_shapes = (X_train_tokens.shape[1]==max_seq_length) and (X_train_mask.shape[1]
        (X_train_segment.shape[1]==max_seq_length)
        segment_temp = not np.any(X_train_segment)
        mask_temp = np.sum(X_train_mask==0) == np.sum(X_train_tokens==0)
        no_cls = np.sum(X_train_tokens==tokenizer.vocab['[CLS]'])==X_train_tokens.shape[0]
        no_sep = np.sum(X_train_tokens==tokenizer.vocab['[SEP]'])==X_train_tokens.shape[0]
        out = temp_shapes and segment_temp and mask_temp and no_cls and no_sep
    else:
        print('Type of all above token arrays should be numpy array not list')
        out = False
    assert(out==True)
    return out
grader_alltokens_train()
     True
X_test_tokens,X_test_mask, X_test_segment=[],[],[]
for i in range(len(X_test)):
  tokens = tokenizer.tokenize(X test.values[i])
  if len(tokens)<(max_seq_length-2):</pre>
    tokens=tokens
  else:
    tokens = tokens[0:(max_seq_length-2)]
  tokens =['[CLS]',*tokens,'[SEP]']
  X test tokens.append(np.array(tokenizer.convert tokens to ids(tokens)))
  X_test_mask.append(np.array([1]*len(tokens)+[0]*(max_seq_length-len(tokens))))
  X_test_segment.append([0]*max_seq_length)
X_test_token=pad_sequences(X_test_tokens , maxlen=max_seq_length, padding='post')
X test tokens = np.array(X test token)
X_test_mask=np.array(X_test_mask)
X_test_segment = np.array(X_test_segment)
def grader alltokens test():
    out = False
    if type(X_test_tokens) == np.ndarray:
        temp_shapes = (X_test_tokens.shape[1]==max_seq_length) and (X_test_mask.shape[1]==
        (X_test_segment.shape[1]==max_seq_length)
```

```
segment_temp = not np.any(X_test_segment)
        mask temp = np.sum(X test mask==0) == np.sum(X test tokens==0)
        no_cls = np.sum(X_test_tokens==tokenizer.vocab['[CLS]'])==X_test_tokens.shape[0]
        no_sep = np.sum(X_test_tokens==tokenizer.vocab['[SEP]'])==X_test_tokens.shape[0]
        out = temp shapes and segment temp and mask temp and no cls and no sep
    else:
        print('Type of all above token arrays should be numpy array not list')
        out = False
    assert(out==True)
    return out
grader_alltokens_test()
 ☐→ True
bert model.input
     [<KerasTensor: shape=(None, 49) dtype=int32 (created by layer 'input_word_ids')>,
      <KerasTensor: shape=(None, 49) dtype=int32 (created by layer 'input_mask')>,
      <KerasTensor: shape=(None, 49) dtype=int32 (created by layer 'segment_ids')>]
bert_model.output
     <KerasTensor: shape=(None, 768) dtype=float32 (created by layer 'keras_layer')>
X_train_pooled_output=bert_model.predict([X_train_tokens,X_train_mask,X_train_segment])
X_test_pooled_output=bert_model.predict([X_test_tokens,X_test_mask,X_test_segment])
pickle.dump((X_train_pooled_output, X_test_pooled_output),open('/content/drive/MyDrive/Ass
X_train_pooled_output, X_test_pooled_output= pickle.load(open('/content/drive/MyDrive/Assi
def greader_output():
    assert(X_train_pooled_output.shape[1]==768)
    assert(len(y_train)==len(X_train_pooled_output))
    assert(X test pooled output.shape[1]==768)
    assert(len(y_test)==len(X_test_pooled_output))
    assert(len(y_train.shape)==1)
    assert(len(X train pooled output.shape)==2)
    assert(len(y test.shape)==1)
    assert(len(X test pooled output.shape)==2)
    return True
greader output()
     True
```

```
from keras.utils.np_utils import to_categorical
y_tr = to_categorical(y_train)
y_te = to_categorical(y_test)
from sklearn.metrics import roc_auc_score
def auc( y_true, y_pred ) :
   score = tf.py_function( lambda y_true, y_pred : roc_auc_score( y_true, y_pred, average
                   [y_true, y_pred],
                   'float32',
                   name='sklearnAUC' )
   return score
pooled_input = keras.Input(shape=(768,) , name='pooled')
model = keras.layers.Dense(512, activation='relu')(pooled_input)
model = keras.layers.Dropout(0.1)(model)
output = keras.layers.Dense(2 , activation='softmax' , name='class')(model)
model = keras.Model(inputs=[pooled_input],outputs=[output])
model.compile(optimizer=keras.optimizers.Adam(),loss='categorical_crossentropy',metrics=['
%load ext tensorboard
from keras.callbacks import TensorBoard
import datetime
logs = "logs/fit/" + datetime.datetime.now().strftime('%Y%m%d-%H%M%S')
tensor = TensorBoard(log_dir='logs',histogram_freq=1,write_graph=True,write_grads=True)
    WARNING:tensorflow:`write_grads` will be ignored in TensorFlow 2.0 for the `TensorBoat
    from keras.callbacks import EarlyStopping,ModelCheckpoint
estop = EarlyStopping( monitor='val_loss',patience=5 )
filepath='/content/drive/MyDrive/Assignment/NLP with Transfer Learning/Model.hd5f'
checkpoint = ModelCheckpoint(filepath=filepath,monitor='val_loss',save_best_only=True , mo
model.fit(X_train_pooled_output,y_tr,validation_split=0.2,batch_size=128,epochs=50,callbac
   censor from Assecs written to. /content/urive/nybrive/Assignment/NEr with fransfer
   :tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
   h 26/50
   h 27/50
   500 [============ ] - 5s 10ms/step - loss: 0.1761 - accuracy: 0.92
   ::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
   :tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
   h 29/50
   h 30/50
   h 31/50
   h 32/50
```

```
500 [========== ] - 5s 10ms/step - loss: 0.1702 - accuracy: 0.93
::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
h 33/50
h 34/50
h 35/50
::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
h 36/50
h 37/50
h 38/50
h 39/50
::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
:tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
500 [===========] - 5s 10ms/step - loss: 0.1765 - accuracy: 0.92
::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
h 41/50
h 42/50
::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
::tensorflow:Assets written to: /content/drive/MyDrive/Assignment/NLP with Transfer
h 43/50
h 44/50
500 [=========== - 4s 9ms/step - loss: 0.1743 - accuracy: 0.929
h 45/50
h 46/50
h 47/50
as.callbacks.History at 0x7ff529d77790>
```

```
print('Test AUC is: ',model.evaluate(X_test_pooled_output,y_te,batch_size=64)[2])
```

## %tensorboard--logdir /content/logs

test\_data=pd.read\_csv('\_/content/drive/MyDrive/Assignment/NLP with Transfer Learning/test.c
test data.head(2)

## Text

test\_data['Text']=remove\_html(test\_data['Text'].values)
test\_data.head(2)

## Text

- **0** Just opened Greenies Joint Care (individually ...
- 1 This product rocks :) My mom was very happy w/...

```
test_tokens,test_mask,test_segment=[],[],[]
for i in range(len(test_data)):
  tokens = tokenizer.tokenize(test data['Text'].values[i])
  if len(tokens)<(max_seq_length-2):</pre>
    tokens=tokens
  else:
    tokens = tokens[0:(max_seq_length-2)]
  tokens =['[CLS]',*tokens,'[SEP]']
  test_tokens.append(np.array(tokenizer.convert_tokens_to_ids(tokens)))
  test_mask.append(np.array([1]*len(tokens)+[0]*(max_seq_length-len(tokens))))
  test_segment.append([0]*max_seq_length)
test_token=pad_sequences(test_tokens , maxlen=max_seq_length, padding='post')
test_tokens = np.array(test_token)
test_mask=np.array(test_mask)
test_segment = np.array(test_segment)
test_pooled_output=bert_model.predict([test_tokens,test_mask,test_segment])
print(test_pooled_output.shape)
     (352, 768)
y = model.predict(test_pooled_output)
class_label=[]
for i in range(len(y)):
  if y[i][0]>y[i][1]:
    class_label.append(0)
  if y[i][1]>y[i][0]:
    class_label.append(1)
```

```
print(class label)
```

✓ 0s completed at 12:11 AM

×