

```
%%capture
!pip install transformers
```

```
import tensorflow as tf
from tensorflow.python.client import device_lib
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report

from transformers import BertTokenizer
from transformers import BertTokenizerFast
from transformers import TFBertForSequenceClassification
from transformers import DistilBertTokenizerFast, TFDistilBertForSequenceClassification

import pandas as pd
import numpy as np
import plotly.figure_factory as ff
import plotly.express as px
import plotly.graph_objects as go

import re
import nltk
from nltk.corpus import stopwords
from google.colab import drive
```

```
print(device_lib.list_local_devices())

if tf.config.list_physical_devices('GPU'):
    tf.config.experimental.set_memory_growth(tf.config.list_physical_devices('GPU')[0], True)
    tf.config.set_visible_devices(tf.config.list_physical_devices('GPU')[0], 'GPU')

[name: "/device:CPU:0"
 device_type: "CPU"
 memory_limit: 268435456
 locality {
 }
 incarnation: 16674445784876115438
 xla_global_id: -1
 , name: "/device:GPU:0"
 device_type: "GPU"
 memory_limit: 14328594432
 locality {
   bus_id: 1
   links {
   }
 }
 incarnation: 3980103699019794021
 physical_device_desc: "device: 0, name: Tesla T4, pci bus id: 0000:00:04.0, compute capability: 7.5"
 xla_global_id: 416903419
]
```

```
drive.mount('/content/drive')
nltk.download('stopwords')
```

```
Mounted at /content/drive
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
True
```

```
df = pd.read_csv('/content/drive/MyDrive/Datasets/IMDB_Movie_Review/Reviews.csv')
```

```
df = df.head(5000)
```

```
df.head(3)
```

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production. The...	positive
2	I thought this was a wonderful way to spend ti...	positive

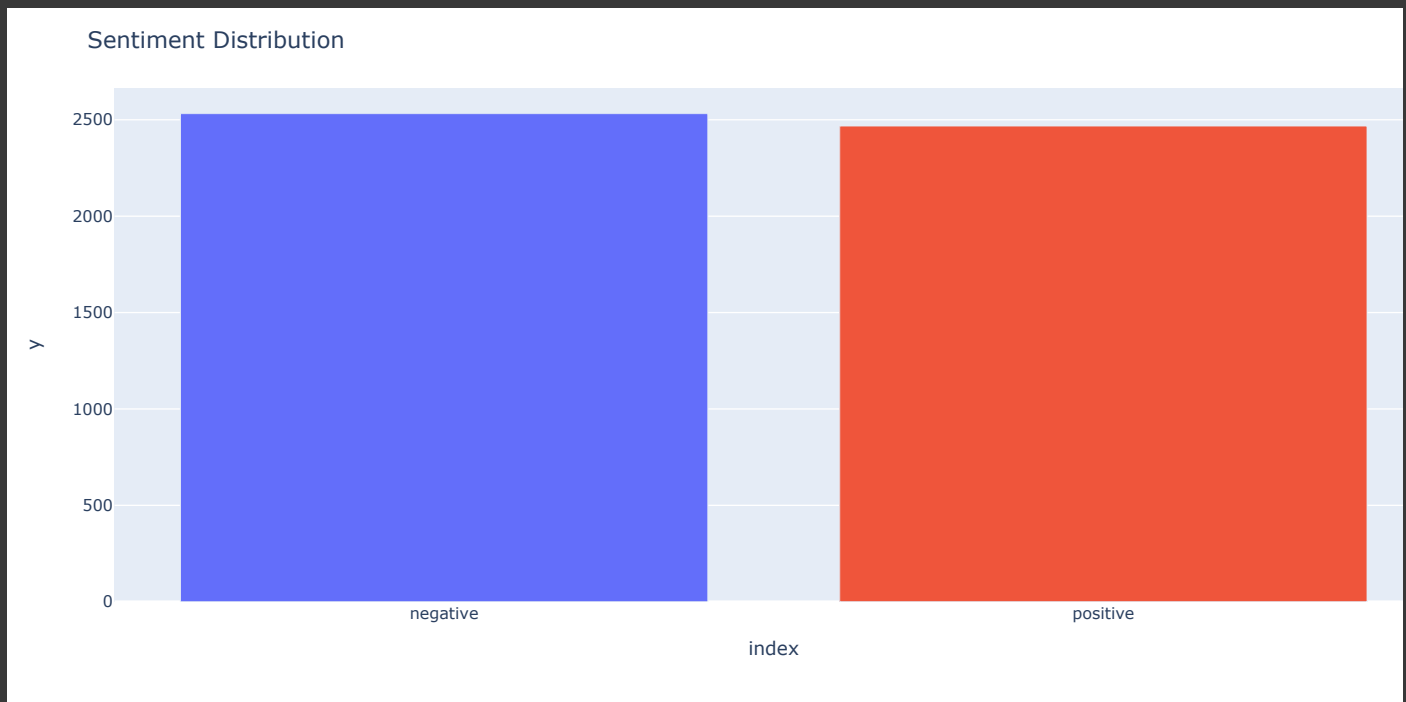
```
def preprocess_text(text):
    text = text.lower()
    text = re.sub(r"http\S+|www\S+|https\S+", "", text, flags=re.MULTILINE)
    text = re.sub(r"[^\w\s]", "", text)
    text = re.sub(r"\d+", "", text)
    stop_words = set(stopwords.words('english'))
    text = " ".join([word for word in text.split() if word not in stop_words])

    return text
```

```
num_labels = 2
```

```
X = df['review'].apply(preprocess_text)
y = df['sentiment']
```

```
sentiment_counts = y.value_counts()
fig = px.bar(sentiment_counts, x=sentiment_counts.index, y=sentiment_counts.values, color=sentiment_counts.index)
fig.update_layout(title="Sentiment Distribution")
fig.show()
```



```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
validation_num = 500
```

```
X_validation, X_train = X_train[:validation_num], X_train[validation_num:]
y_validation, y_train = y_train[:validation_num], y_train[validation_num:]
```

```
X_train.shape, X_validation.shape, X_test.shape, y_train.shape, y_validation.shape, y_test.shape
```

```
((3500,), (500,), (1000,), (3500,), (500,), (1000,))
```

```
with tf.device('GPU'):
    bert_tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
    bert_model = TFBertForSequenceClassification.from_pretrained('bert-base-uncased', num_labels=num_labels)
```

Downloading (...)solve/main/vocab.txt: 100%  232k/232k [00:00<00:00, 4.01MB/s]

Downloading (...)okenizer_config.json: 100%  28.0/28.0 [00:00<00:00, 1.64kB/s]

Downloading (...)lve/main/config.json: 100%  570/570 [00:00<00:00, 20.8kB/s]

```
def create_input_tensors(input_X, tokenizer):

    input_ids = []
    attention_masks = []

    for text in input_X:
        tokens = tokenize_text(text, tokenizer)
        input_ids.append(tokens['input_ids'][0])
        attention_masks.append(tokens['attention_mask'][0])

    input_ids = np.array(input_ids)
    attention_masks = np.array(attention_masks)

    return {'input_ids': input_ids, 'attention_mask': attention_masks}

def tokenize_text(text, tokenizer):
    tokens = tokenizer.encode_plus(text,
                                   max_length=128,
                                   truncation=True,
                                   padding='max_length',
                                   add_special_tokens=True,
                                   return_attention_mask=True,
                                   return_tensors='tf')

    return tokens

def convert_labels_to_one_hot(labels, num_classes):
    label_mapping = {'positive': 1,
                     'negative': 0}
    labels = [label_mapping[label] for label in labels]

    return tf.keras.utils.to_categorical(labels, num_classes=num_classes)
```

```
new_train_X = create_input_tensors(X_train, bert_tokenizer)
new_train_y = convert_labels_to_one_hot(y_train, num_labels)

new_validation_X = create_input_tensors(X_validation, bert_tokenizer)
new_validation_y = convert_labels_to_one_hot(y_validation, num_labels)

new_test_X = create_input_tensors(X_test, bert_tokenizer)
new_test_y = convert_labels_to_one_hot(y_test, num_labels)
```

```
with tf.device('GPU'):

    epochs = 3

    bert_model.compile(optimizer = tf.keras.optimizers.Adam(learning_rate=2e-5),
                       loss = tf.keras.losses.CategoricalCrossentropy(from_logits=True),
                       metrics = [tf.keras.metrics.CategoricalAccuracy('accuracy')])

    bert_model_history = bert_model.fit(new_train_X,
                                       new_train_y,
                                       batch_size=32,
                                       epochs=epochs,
                                       validation_data=(new_validation_X, new_validation_y))
```

```
Epoch 1/3
110/110 [=====] - 149s 860ms/step - loss: 0.4755 - accuracy: 0.7686 - val_loss: 0.3360 - val_accuracy: 0.7686
Epoch 2/3
110/110 [=====] - 92s 834ms/step - loss: 0.2545 - accuracy: 0.9040 - val_loss: 0.3506 - val_accuracy: 0.9040
Epoch 3/3
110/110 [=====] - 92s 833ms/step - loss: 0.1523 - accuracy: 0.9457 - val_loss: 0.3525 - val_accuracy: 0.9457
```

```
train_loss_history = bert_model_history.history['loss']
validation_loss_history = bert_model_history.history['val_loss']
```

```
train_acc_history = bert_model_history.history['accuracy']
validation_acc_history = bert_model_history.history['val_accuracy']
```

```
bert_model.save('/content/drive/MyDrive/Datasets/IMDB_Movie_Review/bert_model')
```

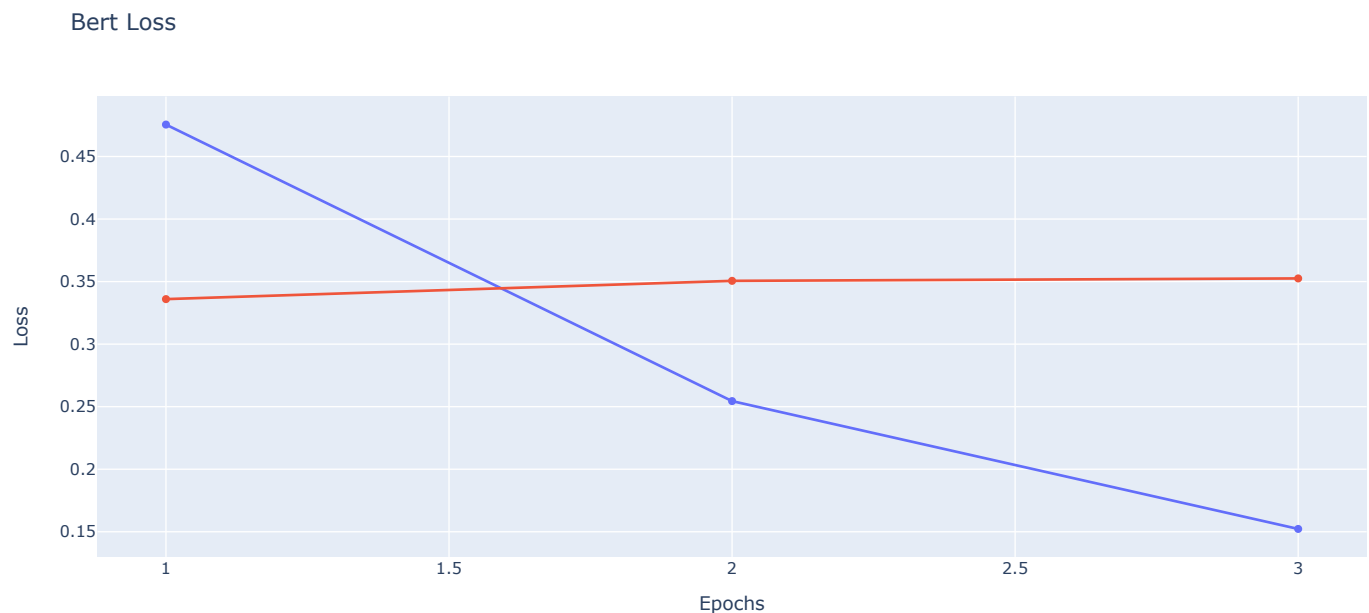
WARNING:absl:Found untraced functions such as embeddings_layer_call_fn, embeddings_layer_call_and_return_conditional_losses,

```
fig = go.Figure()
```

```
fig.add_trace(go.Scatter(x=list(range(1, epochs+1)), y=train_loss_history, mode='lines+markers', name='Train Loss'))
fig.add_trace(go.Scatter(x=list(range(1, epochs+1)), y=validation_loss_history, mode='lines+markers', name='Validation Loss'))
```

```
fig.update_layout(title="Bert Loss",
                  xaxis_title="Epochs",
                  yaxis_title="Loss")
```

```
fig.show()
```



```
fig = go.Figure()
```

```
fig.add_trace(go.Scatter(x=list(range(1, epochs+1)), y=train_acc_history, mode='lines+markers', name='Train Accuracy'))
fig.add_trace(go.Scatter(x=list(range(1, epochs+1)), y=validation_acc_history, mode='lines+markers', name='Validation Accuracy'))
```

```
fig.update_layout(title="Bert Accuracy",
                  xaxis_title="Epochs",
                  yaxis_title="Accuracy")
```

```
fig.show()
```

Bert Accuracy

0.95

```
test_predictions = bert_model.predict(new_test_X)
predicted_labels = np.argmax(test_predictions.logits, axis=1)
```

```
test_labels_ld = np.argmax(new_test_y, axis=1)
```

```
32/32 [=====] - 8s 261ms/step
```

```
print(classification_report(test_labels_ld, predicted_labels))
```

	precision	recall	f1-score	support
0	0.90	0.81	0.85	530
1	0.81	0.90	0.85	470
accuracy			0.85	1000
macro avg	0.86	0.86	0.85	1000
weighted avg	0.86	0.85	0.85	1000

```
confusion_matrix = confusion_matrix(test_labels_ld, predicted_labels)
```

```
heatmap = go.Heatmap(
    z=confusion_matrix,
    x=['Negative', 'Positive'],
    y=['Negative', 'Positive'],
    colorscale='Blues',
    reversescale=True,
)
```

```
layout = go.Layout(
    title='Confusion Matrix',
    xaxis=dict(title='Predicted'),
    yaxis=dict(title='True'),
)
```

```
fig = go.Figure(data=[heatmap], layout=layout)
fig.show()
```

