```
1 #import libraries
2 import tensorflow as tf
3 from tensorflow.keras.preprocessing.text import Tokenizer
4 from tensorflow.keras.preprocessing.sequence import pad_sequences
5 import tensorflow.keras.layers as L
6 from tensorflow.keras.losses import SparseCategoricalCrossentropy
7 from tensorflow.keras.optimizers import Adam
8
9 from sklearn.model_selection import train_test_split
10 from sklearn.metrics import confusion_matrix
11 from sklearn.metrics import classification_report
12
13 import matplotlib.pyplot as plt
14 import plotly.graph_objects as go
15 import plotly.express as px
16 import plotly.figure_factory as ff
17 import seaborn as sns
18
19
20 import numpy as np
21 import pandas as pd
22
23 import nltk
24 from nltk.stem.porter import PorterStemmer
25 from nltk.tokenize import TweetTokenizer
26 from nltk.tokenize import word_tokenize
27 from nltk.corpus import stopwords
28
29 import re
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

```
Mounted at /content/drive
```

```
1 #dataset
2 train = pd.read_csv('Corona_NLP_train.csv',encoding='latin1')
3 test = pd.read_csv('Corona_NLP_test.csv',encoding='latin1')
4
5 train.head()
```

|   | UserName | ScreenName | Location | TweetAt | OriginalTweet | Sentiment |
|---|----------|------------|----------|---------|---------------|-----------|
| 0 | 3799 | 48751 | London | 16-03-2020 | @MeNyrbie @Phil_Gahan @Chrisitv https://t.co/i... | Neutral |
| 1 | 3800 | 48752 | UK | 16-03-2020 | advice Talk to your neighbours family to excha... | Positive |
| 2 | 3801 | 48753 | Vagabonds | 16-03-2020 | Coronavirus Australia: Woolworths to give elde... | Positive |

```
1 INDEX = np.random.randint(1,train.shape[0])
2 print(INDEX)
3 print(train['OriginalTweet'][INDEX])
4 print(train['Sentiment'][INDEX])
```

```
1109
US Food Industry Scrambles To Resupply Stores Amid Apocalyptic Surge In Demand | Zero Hedge https://t.co/Pw9QbGcE1j #coronavirus #C
Extremely Negative
```

```
1 #dataset size
2 print('Examples in train data: {}'.format(len(train)))
3 print('Examples in test data: {}'.format(len(test)))
```

```
Examples in train data: 41157
Examples in test data: 3798
```

```
1 #missing value
2 train.isna().sum()
```

```
UserName            0
ScreenName          0
Location         8590
TweetAt             0
OriginalTweet       0
Sentiment           0
dtype: int64
```
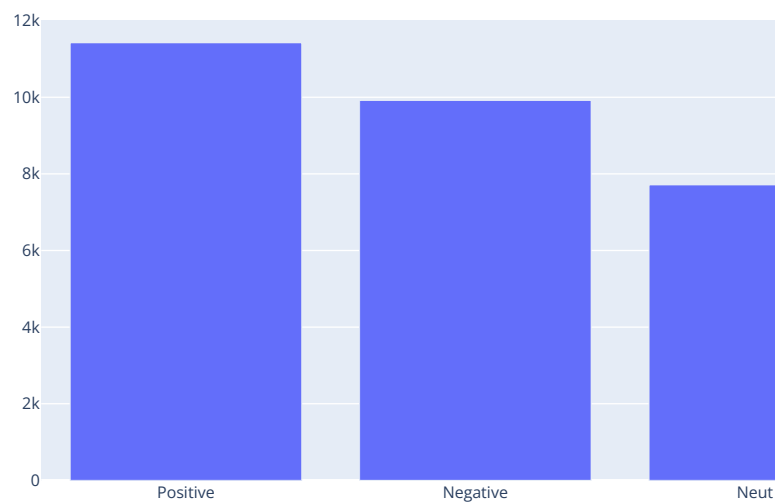
```
1 #class distribution
2 dist_train = train['Sentiment'].value_counts()
3 dist_test = test['Sentiment'].value_counts()
4
5 def ditribution_plot(x,y,name):
6     fig = go.Figure([
7         go.Bar(x=x, y=y)
8     ])
9
10     fig.update_layout(title_text=name)
11     fig.show()
```
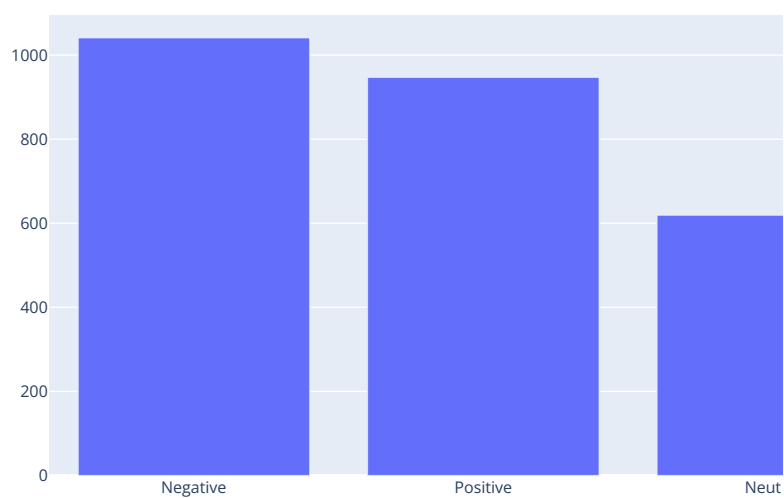
```
1 ditribution_plot(x= dist_train.index, y= dist_train.values, name= 'Class Distribution train')
```

## Class Distribution train



```
1 ditribution_plot(x= dist_test.index, y= dist_test.values, name= 'Class Distribution test')
```

## Class Distribution test



```
1 #data preprocessing
2 X = train['OriginalTweet'].copy()
3 y = train['Sentiment'].copy()
```

```
1 import nltk
2 nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
True
```

```
1 #cleaning data
2 def data_cleaner(tweet):
3
4     # remove urls
5     tweet = re.sub(r'http\S+', ' ', tweet)
6
7     # remove html tags
8     tweet = re.sub(r'<.*?>',' ', tweet)
9
10    # remove digits
11    tweet = re.sub(r'\d+',' ', tweet)
12
13    # remove hashtags
14    tweet = re.sub(r'#\w+',' ', tweet)
15
16    # remove mentions
17    tweet = re.sub(r'@\w+',' ', tweet)
18
19    #remove non words character+ underscore
20    tweet = re.sub(r'[\W_]+',' ',tweet)
21
22    #removing stop words
23    tweet = tweet.lower().split()
24    tweet = " ".join([word for word in tweet if not word in stop_words])
25
26    return tweet
27
28
29 stop_words = stopwords.words('english')
30
31 X_cleaned = X.apply(data_cleaner)
32 X_cleaned.head()
```

```
0
1     advice talk neighbours family exchange phone n...
2     coronavirus australia woolworths give elderly ...
3     food stock one empty please panic enough food ...
4     ready go supermarket outbreak paranoid food st...
Name: OriginalTweet, dtype: object
```

```
1 print(stop_words)
```

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yours
```

```
1 #tokenizing
2 tokenizer = Tokenizer()
3 tokenizer.fit_on_texts(X_cleaned)
4
5 X = tokenizer.texts_to_sequences(X_cleaned)
6
7 vocab_size = len(tokenizer.word_index)+1
8
9 print("Vocabulary size: {}".format(vocab_size))
10 print("\nExample:\n")
11 print("Sentence:\n{}".format(X_cleaned[6]))
12 print("\nAfter tokenizing :\n{}".format(X[6]))
13
14 X = pad_sequences(X, padding='post')
15 print("\nAfter padding :\n{}".format(X[6]))
```

```
Vocabulary size: 33750

Example:

Sentence:
cashier grocery store sharing insights prove credibility commented civics class know talking

After tokenizing :
[1075, 6, 3, 1167, 682, 2766, 9513, 9514, 17742, 1391, 53, 811]

After padding :
[ 1075     6     3  1167   682  2766  9513  9514 17742  1391    53   811
     0     0     0     0     0     0     0     0     0     0     0     0
```

```
          0        0        0        0        0        0        0        0        0        0        0        0
          0        0        0        0        0        0        0        0        0        0        0]
```

```python
1  #feature encoding
2  encoding = {'Extremely Negative': 0,
3              'Negative': 0,
4              'Neutral': 1,
5              'Positive':2,
6              'Extremely Positive': 2
7             }
8
9  labels = ['Negative', 'Neutral', 'Positive']
10
11
12 y.replace(encoding, inplace=True)
```

```python
1  #model building
2  tf.keras.backend.clear_session()
3
4  # hyper parameters
5  EPOCHS = 2
6  BATCH_SIZE = 32
7  embedding_dim = 16
8  units = 256
9
10 model = tf.keras.Sequential([
11     L.Embedding(vocab_size, embedding_dim, input_length=X.shape[1]),
12     L.Bidirectional(L.LSTM(units,return_sequences=True)),
13     L.GlobalMaxPool1D(),
14     L.Dense(256, activation="relu"),
15     L.Dropout(0.7),
16     L.Dense(128, activation="relu"),
17     L.Dropout(0.3),
18     L.Dense(3)
19 ])
20
21
22 model.compile(loss=SparseCategoricalCrossentropy(from_logits=True),
23               optimizer='adam',metrics=['accuracy']
24              )
25
26 model.summary()
```

```
Model: "sequential"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding (Embedding)       (None, 47, 16)            540000

 bidirectional (Bidirectiona (None, 47, 512)           559104
 l)

 global_max_pooling1d (Globa (None, 512)               0
 lMaxPooling1D)

 dense (Dense)               (None, 256)               131328

 dropout (Dropout)           (None, 256)               0

 dense_1 (Dense)             (None, 128)               32896

 dropout_1 (Dropout)         (None, 128)               0

 dense_2 (Dense)             (None, 3)                 387

=================================================================
Total params: 1,263,715
Trainable params: 1,263,715
Non-trainable params: 0
_____
```

```python
1  X[1]
```

```
array([  385,    736,   2381,    167,   2502,    715,   1084,    932,    319,
         273,    715,   1084,   2381,   1063,   3366,   3043,   4363,    368,
          12,     11,   2127,  17740,   2901,    101,    898,   2714,    139,
           0,      0,      0,      0,      0,      0,      0,      0,      0,
           0,      0,      0,      0,      0,      0,      0,      0,      0,
           0,      0], dtype=int32)
```
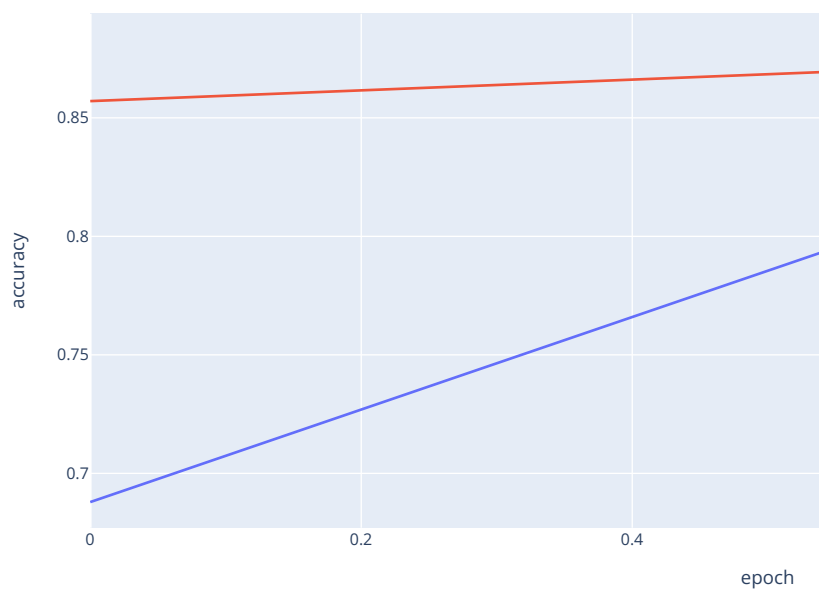
```python
1  X.shape[1]
```

```
47
```

```
1 #training
2 history = model.fit(X, y, epochs=EPOCHS, validation_split=0.12, batch_size=BATCH_SIZE)
```

```
   Epoch 1/2
   1132/1132 [==============================] - 16s 12ms/step - loss: 0.7151 - accuracy: 0.6879 - val_loss: 0.4062 - val_accuracy: 0.8
   Epoch 2/2
   1132/1132 [==============================] - 13s 11ms/step - loss: 0.3582 - accuracy: 0.8831 - val_loss: 0.3576 - val_accuracy: 0.8
```

```
1 fig = px.line(
2     history.history, y=['accuracy', 'val_accuracy'],
3     labels={'index': 'epoch', 'value': 'accuracy'}
4 )
5
6 fig.show()
```



```
1 fig = px.line(
2     history.history, y=['loss', 'val_loss'],
3     labels={'index': 'epoch', 'value': 'loss'}
4 )
5
6 fig.show()
```

```
1 #Preprocessing test data
2 X_test = test['OriginalTweet'].copy()
3 y_test = test['Sentiment'].copy()
4
5 X_test = X_test.apply(data_cleaner)
6
7 X_test = tokenizer.texts_to_sequences(X_test)
8
9 X_test = pad_sequences(X_test, padding='post',maxlen=47)
10
11 y_test.replace(encoding, inplace=True)
```

```
1 X_test.shape
```

```
(3798, 47)
```

```
1 #predict data test
2 predict_x=model.predict(X_test)
3 classes_x=np.argmax(predict_x,axis=1)
```

```
1 loss, acc = model.evaluate(X_test,y_test,verbose=0)
2 print('Test loss: {}'.format(loss))
3 print('Test Accuracy: {}'.format(acc))
```
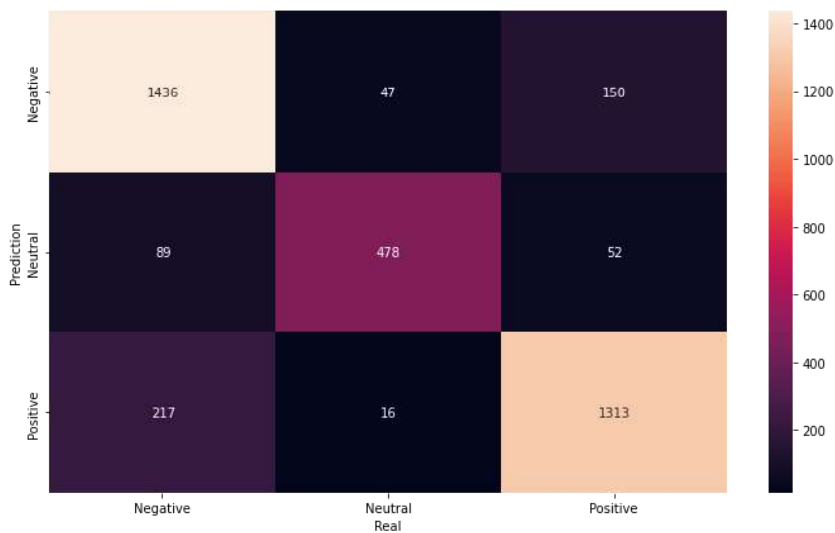
```
Test loss: 0.4198286235332489
Test Accuracy: 0.8496577143669128
```

```
1 conf = confusion_matrix(y_test, classes_x)
2
3 cm = pd.DataFrame(
4     conf, index = [i for i in labels],
5     columns = [i for i in labels]
6 )
7
8 plt.figure(figsize = (12,7))
9 sns.heatmap(cm, annot=True, fmt="d")
10 plt.xlabel('Real')
11 plt.ylabel('Prediction')
12 plt.show()
```



```
1 test['OriginalTweet']
```

```
0       TRENDING: New Yorkers encounter empty supermar...
1       When I couldn't find hand sanitizer at Fred Me...
2       Find out how you can protect yourself and love...
3       #Panic buying hits #NewYork City as anxious sh...
4       #toiletpaper #dunnypaper #coronavirus #coronav...
                              ...
3793    Meanwhile In A Supermarket in Israel -- People...
3794    Did you panic buy a lot of non-perishable item...
3795    Asst Prof of Economics @cconces was on @NBCPhi...
3796    Gov need to do somethings instead of biar je r...
3797    I and @ForestandPaper members are committed to...
Name: OriginalTweet, Length: 3798, dtype: object
```

```
1 text = [input()]
2 text = tokenizer.texts_to_sequences(text)
3 text = pad_sequences(text, padding='post',maxlen=47)
```

#Finland is boosting its #COVID-19 #vaccination campaign in preparation for the new #COVID19 #BA.5 variant first detected in #South

```
1 pred = model.predict(text)
2 output = np.argmax(pred)
3
4 if output==0:
5     print('Hasil klasifikasi : negative')
6 elif output==1:
7     print('Hasil klasifikasi : neutral')
8 else:
9     print('Hasil klasifikasi : positive')
```

Hasil klasifikasi : negative