### Reading and importing the data

```
 1 import nltk
 2
 3 # Download the stopwords resource
 4 nltk.download('stopwords')
 5
 6 # Download the tokenizer resource
 7 nltk.download('punkt')
 8
 9 import numpy as np
10 import pandas as pd
11 import matplotlib.pyplot as plt
12 import seaborn as sns
```

```
    [nltk_data] Downloading package stopwords to /root/nltk_data...
    [nltk_data]   Package stopwords is already up-to-date!
    [nltk_data] Downloading package punkt to /root/nltk_data...
    [nltk_data]   Package punkt is already up-to-date!
```

```
 1 data=pd.read_csv("/content/Twitter_Data.csv")
```

```
 1 data.info()
```

```
    <class 'pandas.core.frame.DataFrame'>
    Int64Index: 162969 entries, 0 to 162979
    Data columns (total 3 columns):
     #   Column           Non-Null Count   Dtype
    ---  ------           --------------   -----
     0   clean_text       162969 non-null  object
     1   category         162969 non-null  object
     2   sentence_length  162969 non-null  int64
    dtypes: int64(1), object(2)
    memory usage: 5.0+ MB
```

```
 1 data.describe()
```

|       | category       |
|-------|----------------|
| count | 162973.000000  |
| mean  | 0.225436       |
| std   | 0.781279       |
| min   | -1.000000      |
| 25%   | 0.000000       |
| 50%   | 0.000000       |
| 75%   | 1.000000       |
| max   | 1.000000       |

**Changing our dependent variable to categorical. (0 to "Neutral,"-1 to "Negative", 1 to "Positive")**

```
1 # Map numeric values to categorical labels
2 category_mapping = {0: 'Neutral', -1: 'Negative', 1: 'Positive'}
3 data['category'] = data['category'].map(category_mapping)
4
```

**Doing Missing value analysis and drop all null/missing values**

```
1 # Check for missing values
2 missing_values = data.isnull().sum()
3
4 # Drop rows with missing values
5 data.dropna(inplace=True)
6
```

**Dummy Variables**

```
1 # Create dummy variables for the dependent variable
2 data = pd.get_dummies(data, columns=['category'], prefix='', prefix_sep='')
3
```

**Doing text cleaning. (remove every symbol except alphanumeric, transform all words to lower case, and remove punctuationand stopwords )**

```
 1 # Remove symbols, transform to lowercase, remove punctuation, and stopwords
 2 def clean_text(text):
 3     text = re.sub(r'[^\w\s]', '', text)  # Remove symbols
 4     text = text.lower()  # Transform to lowercase
 5     text = re.sub(r'\s+', ' ', text)  # Remove extra spaces
 6     tokens = word_tokenize(text)  # Tokenize the text
 7     tokens = [word for word in tokens if word not in stopwords.words('english')]
 8     return ' '.join(tokens)
 9
10 data['cleaned_text'] = data['clean_text'].apply(clean_text)
```

**Creating a new column and finding the length of each sentence (how many words they contain)**

```
1 data['sentence_length'] = data['clean_text'].apply(lambda x: len(x.split()))
2
```

**Spliting data into dependent(X) and independent(y) dataframe**

```
1 X = data['clean_text']
2 y = data['category']
```

3

## Do operationson text data

```
 1 from tensorflow.keras.preprocessing.text import Tokenizer
 2 from tensorflow.keras.preprocessing.sequence import pad_sequences
 3 from sklearn.model_selection import train_test_split
 4
 5 # Create Tokenizer
 6 tokenizer = Tokenizer()
 7 tokenizer.fit_on_texts(X)
 8 X_encoded = tokenizer.texts_to_sequences(X)
 9 X_padded = pad_sequences(X_encoded, padding='pre')
10
```

## Build an LSTM Model:

```
 1 from tensorflow.keras.models import Sequential
 2 from tensorflow.keras.layers import Embedding, LSTM, Dense
 3
 4 # Define the model architecture
 5 model = Sequential()
 6 model.add(Embedding(input_dim=len(tokenizer.word_index)+1, output_dim=100, input_l
 7 model.add(LSTM(100))
 8 model.add(Dense(3, activation='softmax'))
 9
10 # Compile the model
11 model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accurac
12
```

## Split data into train and test

```
 1 X_train, X_test, y_train, y_test = train_test_split(X_padded, y, test_size=0.2, ra
```

## Training new model

```
 1 history = model.fit(X_train, pd.get_dummies(y_train), epochs=10, batch_size=32, va
 2
```

```
    Epoch 1/10
    3260/3260 [==============================] - 985s 301ms/step - loss: 0.3759 - ac
    Epoch 2/10
    3260/3260 [==============================] - 1013s 311ms/step - loss: 0.2170 - a
    Epoch 3/10
    3260/3260 [==============================] - 1005s 308ms/step - loss: 0.1396 - a
    Epoch 4/10
    3260/3260 [==============================] - 1027s 315ms/step - loss: 0.0854 - a
    Epoch 5/10
    3260/3260 [==============================] - 1031s 316ms/step - loss: 0.0533 - a
    Epoch 6/10
    3260/3260 [==============================] - 982s 301ms/step - loss: 0.0329 - ac
    Epoch 7/10
    3260/3260 [==============================] - 992s 304ms/step - loss: 0.0220 - ac
```

```
Epoch 8/10
3260/3260 [==============================] - 1004s 308ms/step - loss: 0.0157 - ε
Epoch 9/10
3260/3260 [==============================] - 996s 306ms/step - loss: 0.0121 - ac
Epoch 10/10
3260/3260 [==============================] - 988s 303ms/step - loss: 0.0088 - ac
```

## Normalize the Predictions

```
1 y_pred = model.predict(X_test)
2 y_pred_normalized = np.where(y_pred >= 0.5, 1, 0)
3
```

```
    1019/1019 [==============================] - 15s 15ms/step
```

```
1 y_pred_labels = np.argmax(y_pred, axis=1)
2
```

## Measure performance metricsand accuracy

```
 1 from sklearn.metrics import accuracy_score
 2 # Convert y_test to a numpy array
 3 y_test_array = np.array(y_test)
 4 # Reshape y_test if needed
 5 if len(y_test_array.shape) == 1:
 6     y_test_array = y_test_array.reshape(-1, 1)
 7
 8 y_test_labels = np.argmax(y_test_array, axis=1)
 9 accuracy = accuracy_score(y_test_labels, y_pred_labels)
10 print("Accuracy:", accuracy)
11
```

```
    Accuracy: 0.21642020003681658
```

## Print Classification report

```
1 from sklearn.metrics import classification_report
2
3 classification_rep = classification_report(y_test_labels, y_pred_labels)
4 print("Classification Report:")
5 print(classification_rep)
6
```

```
    Classification Report:
                  precision    recall  f1-score   support

               0       1.00      0.22      0.36     32594
               1       0.00      0.00      0.00         0
               2       0.00      0.00      0.00         0

        accuracy                           0.22     32594
       macro avg       0.33      0.07      0.12     32594
    weighted avg       1.00      0.22      0.36     32594
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344
  _warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344
  _warn_prf(average, modifier, msg_start, len(result))
```