

# **TWITTER SENTIMENT ANALYSIS**

REPORT

**Aarya Yogesh Pakhale**

**2nd Year Undergraduate Student  
IIT Kharagpur  
Data Analyst Intern: Nexus**

# INDEX

- **Project Goal & Context**
- **Data Sources Used**
- **Data Preprocessing Steps**
- **Exploratory Data Analysis**
- **Model Implementation**
- **Training & Evaluation**
- **Results & Insights**
- **Resources Used**

# PROJECT GOAL AND CONTEXT

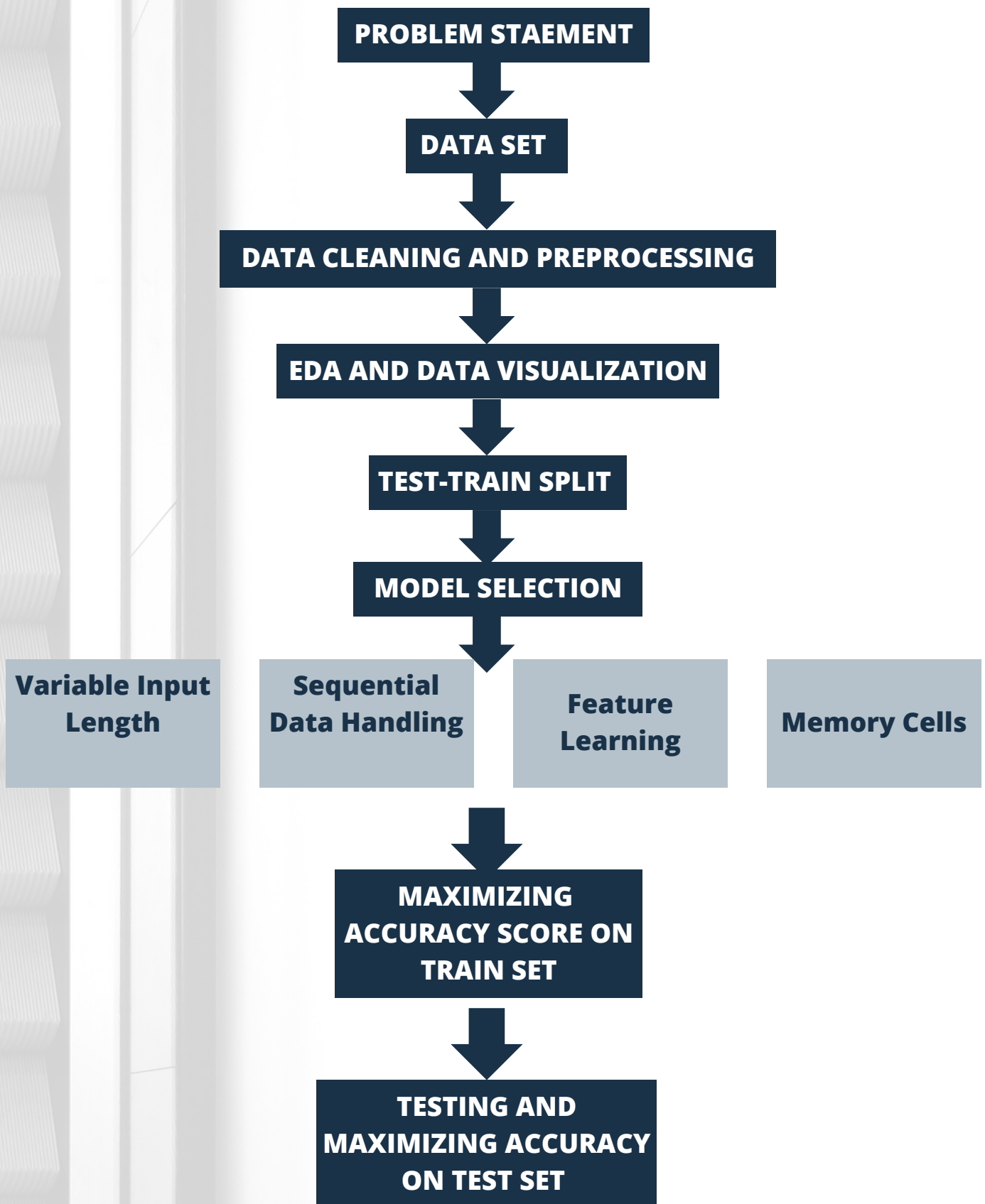
## **Project Title:**

Twitter Sentiment Analysis

## **Project Overview:**

Twitter Sentiment Analysis is a data analytics project that involves analyzing a dataset of tweets to determine the sentiment expressed in each tweet—whether it is positive, negative, or neutral. This project aims to gain insights into public opinions, trends, and sentiments shared on Twitter, utilizing data analytics techniques.

## Approach used:



# DATA SOURCES USED

(<https://www.kaggle.com/datasets/kazanova/sentiment140>)

The dataset provided is the Sentiment140 Dataset which consists of 1,600,000 tweets that have been extracted using the Twitter API. The various columns present in this Twitter data are:

- target: the polarity of the tweet (positive or negative)
- ids: Unique id of the tweet
- date: the date of the tweet
- flag: It refers to the query. If no such query exists, then it is NO QUERY.
- user: It refers to the name of the user that tweeted
- text: It refers to the text of the tweet

# DATA PREPROCESSING

Initially, I indexed the dataset and proceeded to eliminate duplicate entries. Subsequently, I employed a curated stop-word list sourced from the internet to remove irrelevant terms from the text column. I systematically cleaned the data by eliminating repeated characters, URLs, numbers, and punctuations using dedicated functions.

Additionally, I standardized the text by converting all uppercase letters to lowercase.

Following the data cleaning phase, I performed tokenization, a critical step in preparing text for sentiment analysis. Tokenization is essential for feature extraction, managing special characters, accommodating text variations, and creating a structured input for subsequent analysis.

Upon completing EDA, I encoded the text into numerical representations. This encoding step is vital for efficient training, allowing machine learning models to process the data effectively.

# DATA PREPROCESSING

## Code snippets:

```
1 #cleaning the text
2 import string
3 import re
4
5 stopwordlist = ['a', 'about', 'above', 'after', 'again', 'ain', 'all', 'am', 'an',
6                 'and', 'any', 'are', 'as', 'at', 'be', 'because', 'been', 'before',
7                 'being', 'below', 'between', 'both', 'by', 'can', 'd', 'did', 'do',
8                 'does', 'doing', 'down', 'during', 'each', 'few', 'for', 'from',
9                 'further', 'had', 'has', 'have', 'having', 'he', 'her', 'here',
10                'hers', 'herself', 'him', 'himself', 'his', 'how', 'i', 'if', 'in',
11                'into', 'is', 'it', 'its', 'itself', 'just', 'll', 'm', 'ma',
12                'me', 'more', 'most', 'my', 'myself', 'now', 'o', 'of', 'on', 'once',
13                'only', 'or', 'other', 'our', 'ours', 'ourselves', 'out', 'own', 're', 's', 'same', 'she', "shes", 'should', "shouldve", 'so', 'some', 'such',
14                't', 'than', 'that', "thatll", 'the', 'their', 'theirs', 'them',
15                'themselves', 'then', 'there', 'these', 'they', 'this', 'those',
16                'through', 'to', 'too', 'under', 'until', 'up', 've', 'very', 'was',
17                'we', 'were', 'what', 'when', 'where', 'which', 'while', 'who', 'whom',
18                'why', 'will', 'with', 'won', 'y', 'you', "you'd", "youll", "youre",
19                "youve", 'your', 'yours', 'yourself', 'yourselves']
20
21 STOPWORDS = set(stopwordlist)
22 def cleaning_stopwords(text):
23     return " ".join([word for word in str(text).split() if word not in STOPWORDS])
24 df['text'] = df['text'].apply(lambda text: cleaning_stopwords(text))
25 print(df['text'].head())
26
27 def cleaning_repeating_char(text):
28     return re.sub(r'(.1+)', r'1', text)
29 df['text'] = df['text'].apply(lambda x: cleaning_repeating_char(x))
30 print(df['text'].tail())
31
32 def cleaning_URLs(data):
33     return re.sub('((www.[^s]+)|(https?://[^s]+))', '', data)
34 df['text'] = df['text'].apply(lambda x: cleaning_URLs(x))
35 print(df['text'].tail())
```

```
36
37 def cleaning_numbers(data):
38     return re.sub('[0-9]+', '', data)
39 df['text'] = df['text'].apply(lambda x: cleaning_numbers(x))
40 print(df['text'].tail())
41
42 def cleaning_uppercase(text):
43     return text.lower()
44 df['text'] = df['text'].apply(lambda x: cleaning_uppercase(x))
45 print(df['text'].tail())
46
47 def cleaning_punctuations(text):
48     punctuations = string.punctuation + "¡" + "¿"
49     no_punc_sentence = text.translate(str.maketrans('', '', punctuations))
50     return no_punc_sentence
51 df['text'] = df['text'].apply(lambda x: cleaning_punctuations(x))
52 print(df['text'].tail())
53
```

## Code for cleaning data



# DATA PREPROCESSING

## Code snippets:

```
1 from keras.preprocessing import sequence
2 from keras.datasets.imdb import get_word_index
3
4 # tokenizing text
5 word_index=get_word_index()
6 MAXLEN= 250
7
8
9 def tokenize_text(text):
10     tokens=tf.keras.preprocessing.text.text_to_word_sequence(text)
11     tokens=[word_index[word] if word in word_index else 0 for word in tokens]
12     return sequence.pad_sequences([tokens], MAXLEN)[0]
13
```

### Tokenizing text

```
1 #test_train split
2 _, df = train_test_split(df, test_size=0.05, random_state=42)
3 train_data, test_data = train_test_split(df, test_size=0.5, random_state=42)
```

```
1 #encoding data from the dataset
2 from tqdm import tqdm
3 for x in tqdm(range(len(train_data))):
4     train_data.text.iloc[x]=tokenize_text(train_data.text.iloc[x])
5 for x in tqdm(range(len(test_data.text))):
6     test_data.text.iloc[x]=tokenize_text(test_data.text.iloc[x])
7
8 train_labels = train_data.target
9 train_text=sequence.pad_sequences(train_data.text, MAXLEN)
10
11 test_labels = test_data.target
12 test_text=sequence.pad_sequences(test_data.text, MAXLEN)
```

### Encoding text



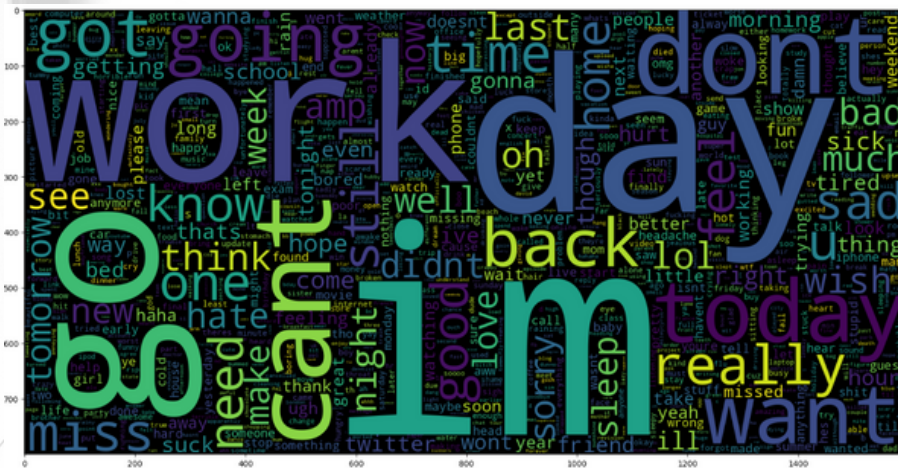
# EXPLORATORY DATA ANALYSIS

During the exploratory data analysis (EDA) phase, I generated distinct word clouds for tweets associated with negative and positive sentiments. This visual exploration provided valuable insights into the most frequently used words within each sentiment category. Words like work, day, don't , can't show a majority in negative sentiment tweets while good, day, well are majority in the positive sentiment tweets. Additionally, I created count plots for positive and negative sentiment tweets, revealing an evenly distributed dataset.

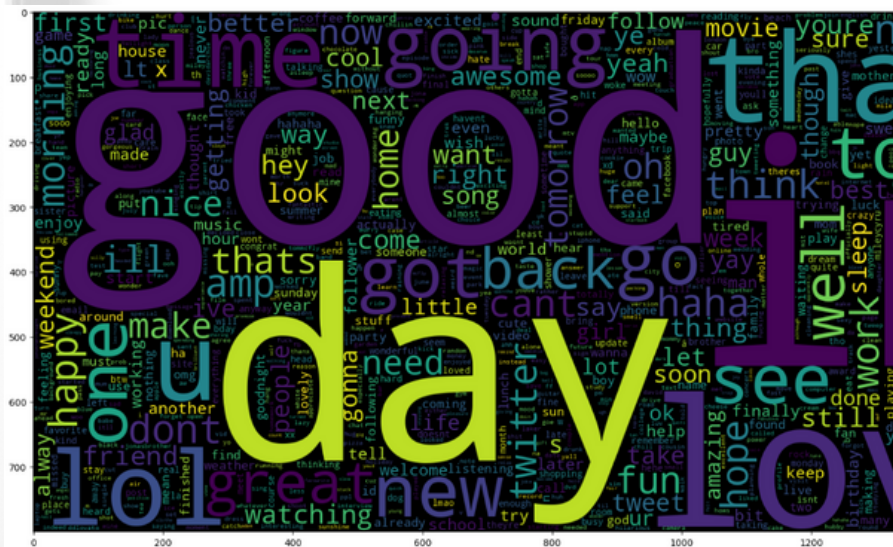
To further understand the dataset's characteristics, I plotted a histogram depicting the distribution of text lengths in the tweets. Notably, the histogram peaked at approx. 40 words and displayed a right-skewed distribution, shedding light on the varying lengths of tweets within the dataset. This information serves as a foundation for understanding the composition and structure of the text data.

# EXPLORATORY DATA ANALYSIS

**The following are the beautiful plots, that were plotted:**

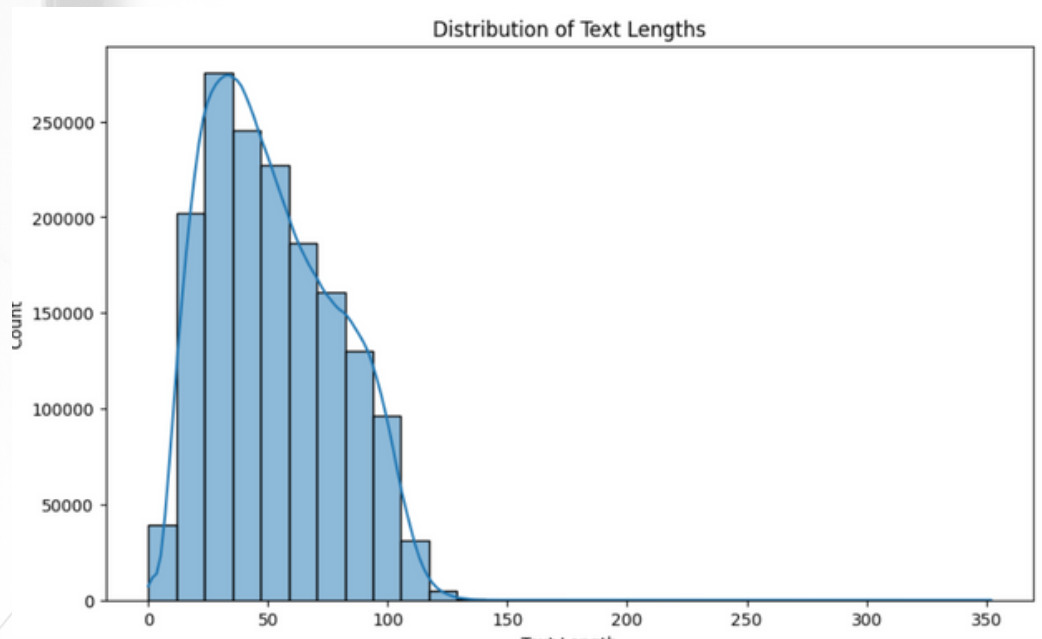


### Word cloud for negative sentiment

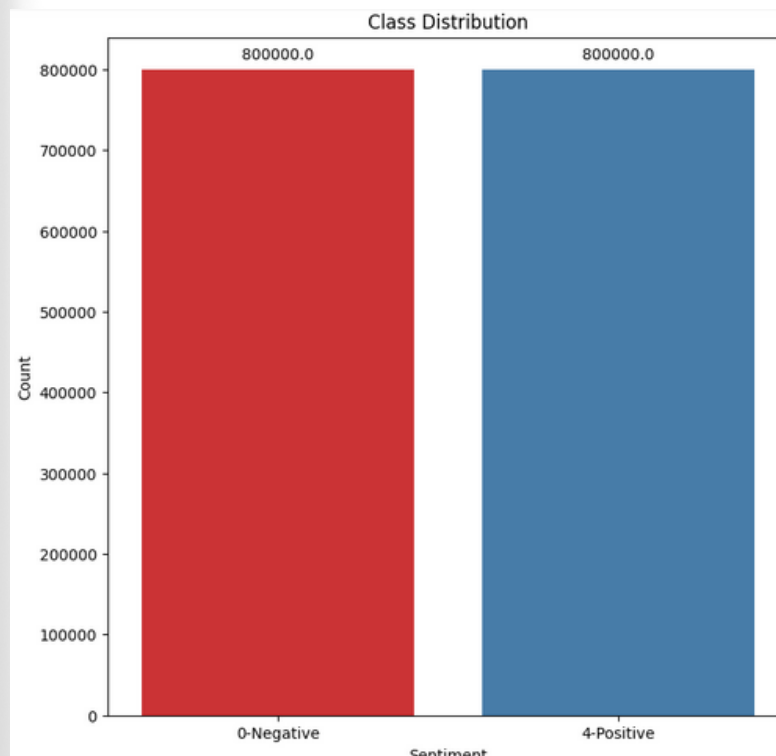


### Word cloud for positive sentiment

# EXPLORATORY DATA ANALYSIS

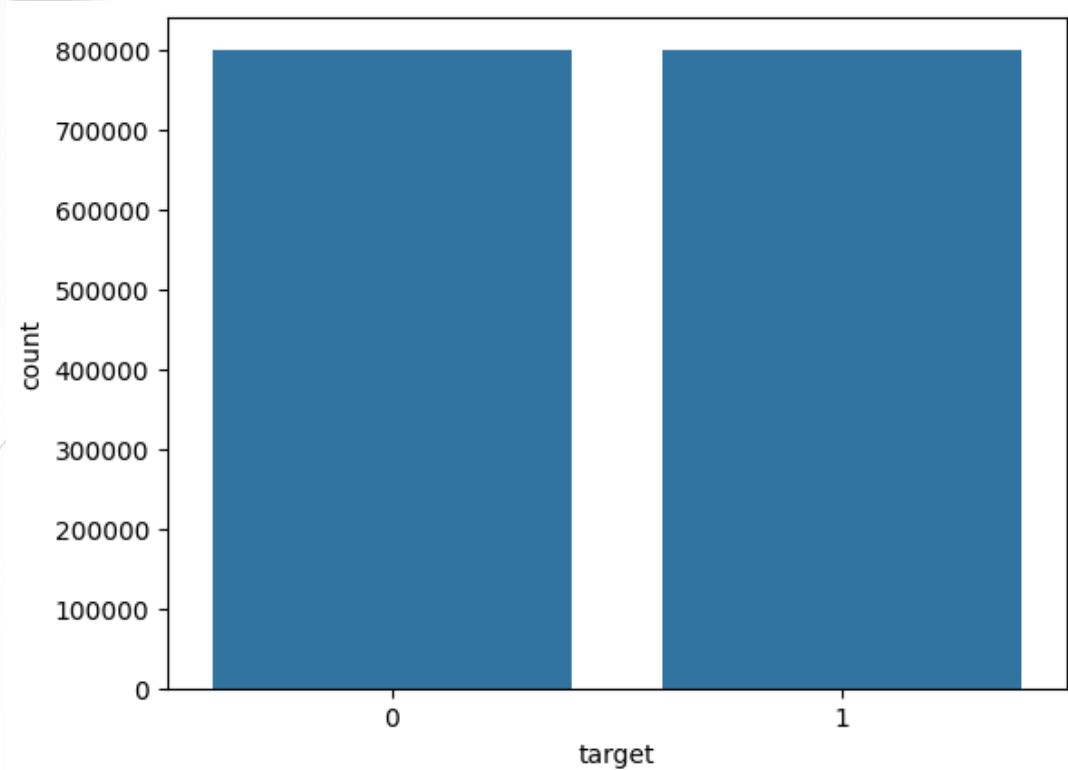


**Histogram of word lengths**



**Count plot**

# EXPLORATORY DATA ANALYSIS



**Count plot**

# MODEL IMPLEMENTATION

## Model Selection: LSTM

Long short-term memory (LSTM) is an artificial recurrent neural network (RNN) architecture. Unlike standard feedforward neural networks, LSTM has feedback connections. A standard LSTM unit is composed of a cell, an input gate, an output gate and a forget gate. The cell remembers values over arbitrary time intervals and the three gates regulate the flow of information into and out of the cell.

## Why LSTM?

- **Handling sequential data:** Twitter data consists of sentences, which are sequences of words. LSTMs are specifically designed to excel in tasks involving sequential data, unlike traditional models that struggle to capture long-term dependencies.
- **Context awareness:** LSTMs can remember information from earlier words in a sentence, allowing them to consider the context when analyzing sentiment. This is crucial for understanding sarcasm, slang, and other nuances in tweets.
- **Parameter Sharing:** LSTMs share parameters across different parts of the sequence, making them more parameter-efficient compared to some other models. This can be important when dealing with large datasets, as it helps prevent overfitting and reduces the risk of model complexity.



# MODEL IMPLEMENTATION

## ACCURACY ACHIECHED:

71.88% On test Data

98.06% On train Data

## Code Snippets:

```
1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import LSTM, Dense, Embedding, Dropout
3
4 # Define parameters
5 sequence_length = 250 # Maximum tweet length
6 embedding_dim = 100 # Dimension of word embeddings
7 vocab_size = 50000 # Number of unique words
8 lstm_units = 128 # Number of LSTM units
9
10 # Build the model
11 model = Sequential()
12 model.add(Embedding(vocab_size, embedding_dim, input_length=sequence_length))
13 model.add(LSTM(lstm_units))
14 model.add(Dropout(0.3))
15 model.add(Dense(1, activation='sigmoid'))
```



# TRAINING AND EVALUATION

## Code Snippets:

```
1 # Compile and train the model
2 from keras.callbacks import EarlyStopping
3
4 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
5 early_stopping = EarlyStopping(monitor='val_loss', patience=5, verbose=1)
6 model.fit(train_text, train_labels, epochs=20, batch_size=32, verbose=1, callbacks=[early_stopping])
```

### Compile and train the model

```
1 model.evaluate(test_text, test_labels)
2
```

```
1250/1250 [=====] - 8s 6ms/step - loss: 2.2610 - accuracy: 0.7187
[2.2610273361206055, 0.7187250256538391]
```

### Evaluation

# RESULTS AND INSIGHTS

- Cleaning and preprocessing the text data increased the accuracy significantly. This is because, the model learns better from a uniform and simple dataset rather than a complex one.
- Effective tokenization helps capture relevant features and nuances within the text that are crucial for accurately discerning sentiment. Raw text is large and complex for machines to handle directly. Tokenizing breaks it down into smaller, more manageable units, making it easier for models like LSTMs to process and analyze the information.
- Models that take variable inputs and have a good pre and post accessing memory work well with sentiment analysis.
- Using regularization techniques like Early stopping and Patience improved the model's process significantly while evaluation. During training, the model's performance on a validation set is monitored. If there is no improvement in the validation performance for a certain number of consecutive epochs (patience), the training is stopped early.
- The model has been trained using a subset of the available data, specifically 42,000 tweets out of the total 16 million provided. However, expanding the training set to encompass a larger portion of the data is currently impractical due to technical constraints such as limited memory and RAM on my PC. Executing the training on a more extensive dataset is estimated to take approximately 72 hours. It is anticipated that training on a larger dataset could potentially lead to improved accuracy, given the additional diversity and richness of information available in the expanded data.

- **FINAL ACCURACY ACHIEVED: 71.87%**

# RESOURCES USED

- ChatGPT,
- BARD,
- <https://www.analytixlabs.co.in/blog/twitter-sentiment-analysis/>
- <https://www.analyticsvidhya.com/blog/2021/06/twitter-sentiment-analysis-a-nlp-use-case-for-beginners/>