# Mini Project on DSBDA Laboratory

## Members:

### TE_A_81 Omkar Nandode (PRN No. 72034728K)
### TE_A_76 Muskan Sawant (PRN No. 72034784L)

## Title

Use the following dataset and classify tweets into positive and negative tweets.
https://www.kaggle.com/ruchi798/data-science-tweets

## Problem Definition:

Classification of tweets into positive and negative from the given csv file.

## Prerequisite:

Basic concepts of Python Programming

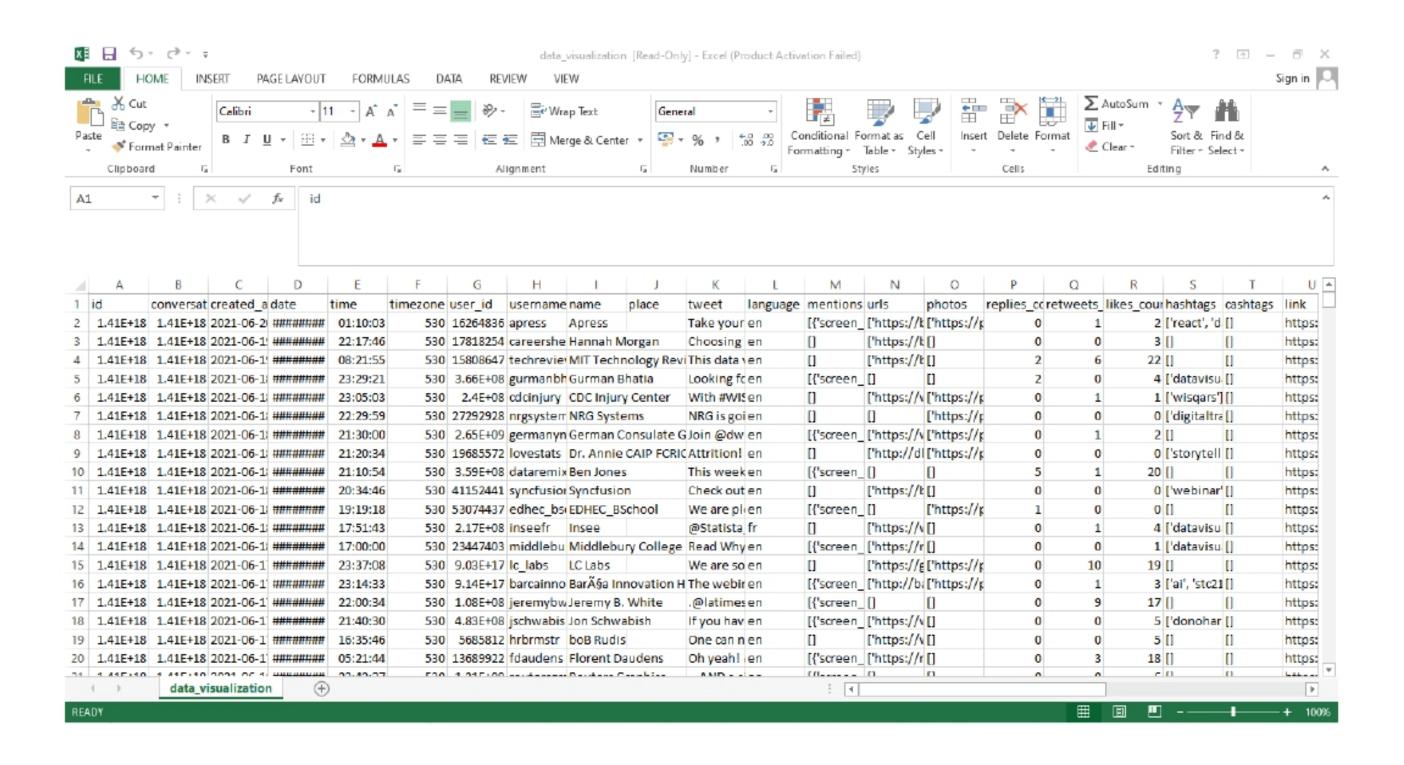## Software Requirements:

Python

## Outcomes:
Rightly segregated tweets

# Data Set Description:

- By using the given dataset we are going to classify the tweets into positive or negative.
- Perform Sentiment analysis.

# Dataset

Following is the Snapshot of Dataset :-

# Theory

## 1. Sentiment Analysis

Sentiment analysis could be performed in Python using 2 methods — i) Calculating polarity and subjectivity using the library **Textblob.** ii)
Using **sentimentIntensityAnalyzer** from the library **vader**.

The SentimentIntensityAnalyzer function relies on a dictionary that maps lexical features to emotion intensities known as sentiment scores. The sentiment score of a text can be obtained by summing up the intensity of each word in the text.

This function analyzes the text and returns the score in the form of a dictionary with the following components :negative, neutral, positive and compound. Based on the scores assigned to each component, we can define the overall sentiment of the text to be positive, negative or neutral.

## 2. Visualize the sentiment Counts

Once the sentiments are identified, we can create 3 lists for different sentiments. We can then calculate the overall percentage of each sentiment in the dataset.

# Input & Output

Input:

```
import nltk
nltk.download('vader_lexicon')
from nltk.sentiment.vader import SentimentIntensityAnalyzer
for index, row in tweets_df['Text'].iteritems():
    score = SentimentIntensityAnalyzer().polarity_scores(row)
```

```python
    if score['neg'] > score['pos']:
        tweets_df.loc[index, "Sentiment"] = "negative"
    elif score['pos'] > score['neg']:
        tweets_df.loc[index, "Sentiment"] = "positive"
    else:
        tweets_df.loc[index, "Sentiment"] = "neutral"

    tweets_df.loc[index, 'neg'] = score['neg']
    tweets_df.loc[index, 'neu'] = score['neu']
    tweets_df.loc[index, 'pos'] = score['pos']
    tweets_df.loc[index, 'compound'] = score['compound']


tweets_df.head(10)
#create new data frames for all sentiments
tweet_neg = tweets_df[tweets_df["Sentiment"] == "negative"]
tweet_neu = tweets_df[tweets_df["Sentiment"] == "neutral"]
tweet_pos = tweets_df[tweets_df["Sentiment"] == "positive"]
#function for calculating the percentage of all the sentiments
def calc_percentage(x,y):
    return x/y * 100
pos_per = calc_percentage(len(tweet_pos), len(tweets_df))
neg_per = calc_percentage(len(tweet_neg), len(tweets_df))
neu_per = calc_percentage(len(tweet_neu), len(tweets_df))
print("positive: {} {}%".format(len(tweet_pos),  format(pos_per, '.1f')))
print("negative: {} {}%".format(len(tweet_neg), format(neg_per, '.1f')))
print("neutral: {} {}%".format(len(tweet_neu), format(neu_per,
'.1f')))format(calc_percentage(len(tweet_neu), len(tweets_df)), '.1f')))
```

## Output:

```
positive: 1788 35.8%
negative: 1795 35.9%
neutral: 1417 28.3%
```

# Conclusion

Thus we have performed the twitter sentiment analysis. There are many more analysis techniques that you can apply on the twitter data like topic modelling. Other techniques are geospatial analysis, text similarity and knowledge graphs. The possibilities are endless.