Aim: Perform sentiment analysis on social media text data (e.g., Tweets). Apply NLP and text classification using ML.

## Download and Prepare the Dataset

```
# Download the Sentiment140 dataset
!wget http://cs.stanford.edu/people/alecmgo/trainingandtestdata.zip

# Unzip the dataset
!unzip -q trainingandtestdata.zip

print('Dataset downloaded and unzipped successfully.')
```

```
--2025-10-06 14:52:35--  http://cs.stanford.edu/people/alecmgo/trainingandtestdata.zip
Resolving cs.stanford.edu (cs.stanford.edu)... 171.64.64.64
Connecting to cs.stanford.edu (cs.stanford.edu)|171.64.64.64|:80... connected.
HTTP request sent, awaiting response... 301 Moved Permanently
Location: https://cs.stanford.edu/people/alecmgo/trainingandtestdata.zip [following]
--2025-10-06 14:52:36--  https://cs.stanford.edu/people/alecmgo/trainingandtestdata.zip
Connecting to cs.stanford.edu (cs.stanford.edu)|171.64.64.64|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 81363704 (78M) [application/zip]
Saving to: 'trainingandtestdata.zip'

trainingandtestdata 100%[===================>]  77.59M  29.4MB/s    in 2.6s

2025-10-06 14:52:38 (29.4 MB/s) - 'trainingandtestdata.zip' saved [81363704/81363704]

Dataset downloaded and unzipped successfully.
```

## Importing Dependencies

```
# utilities
import re
import pickle
import numpy as np
import pandas as pd

# plotting
import seaborn as sns
from wordcloud import WordCloud
import matplotlib.pyplot as plt

# nltk
import nltk
from nltk.stem import WordNetLemmatizer
# Download the 'wordnet' corpus from NLTK
nltk.download('wordnet')
nltk.download('omw-1.4')


# sklearn
from sklearn.svm import LinearSVC
from sklearn.naive_bayes import BernoulliNB
from sklearn.linear_model import LogisticRegression

from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import confusion_matrix, classification_report
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
```

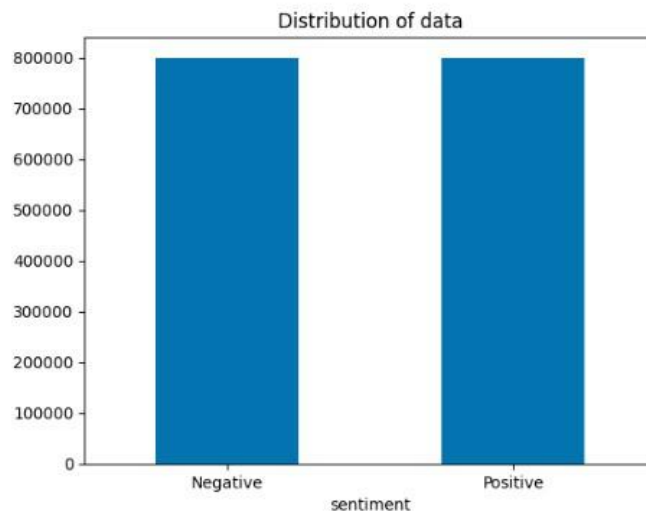## Loading and Initial Preprocessing of the Dataset

```
[3]
✓ 4s
# Importing the dataset
# Define column names and encoding
DATASET_COLUMNS = ["sentiment", "ids", "date", "flag", "user", "text"]
DATASET_ENCODING = "ISO-8859-1"

# Read the CSV file into a pandas DataFrame
# The file path is updated to the local path after unzipping
dataset = pd.read_csv('training.1600000.processed.noemoticon.csv',
                      encoding=DATASET_ENCODING , names=DATASET_COLUMNS)

# Removing the unnecessary columns.
dataset = dataset[['sentiment','text']]
# Replacing the sentiment value '4' with '1' for positive sentiment.
dataset['sentiment'] = dataset['sentiment'].replace(4,1)

# Plotting the distribution for dataset.
ax = dataset.groupby('sentiment').count().plot(kind='bar', title='Distribution of data',
                                               legend=False)
ax.set_xticklabels(['Negative','Positive'], rotation=0)

# Storing data in lists.
text, sentiment = list(dataset['text']), list(dataset['sentiment'])
```



## Define Stopwords and Emojis

```
[4]
✓ 0s
# Defining dictionary containing all emojis with their meanings.
emojis = {':)': 'smile', ':-)': 'smile', ';d': 'wink', ':-E': 'vampire', ':(': 'sad',
          ':-(': 'sad', ':-<': 'sad', ':P': 'raspberry', ':O': 'surprised',
          ':-@': 'shocked', ':@': 'shocked',':-$': 'confused', ':\\\\': 'annoyed',
          ':#': 'mute', ':X': 'mute', ':^)': 'smile', ':-&': 'confused', '$_$': 'greedy',
          '@@': 'eyeroll', ':-!': 'confused', ':-D': 'smile', ':-0': 'yell', 'O.o': 'confused',
          '<(-_-)>': 'robot', 'd[-_-]b': 'dj', ":'-)": 'sadsmile', ';)': 'wink',
          ';-)': 'wink', 'O:-)': 'angel','O*-)': 'angel','(:-D': 'gossip', '=^.^=': 'cat'}


## Defining set containing all stopwords in english.
stopwordlist = ['a', 'about', 'above', 'after', 'again', 'ain', 'all', 'am', 'an',
                'and','any','are', 'as', 'at', 'be', 'because', 'been', 'before',
                'being', 'below', 'between','both', 'by', 'can', 'd', 'did', 'do',
                'does', 'doing', 'down', 'during', 'each','few', 'for', 'from',
                'further', 'had', 'has', 'have', 'having', 'he', 'her', 'here',
                'hers', 'herself', 'him', 'himself', 'his', 'how', 'i', 'if', 'in',
                'into','is', 'it', 'its', 'itself', 'just', 'll', 'm', 'ma',
                'me', 'more', 'most','my', 'myself', 'now', 'o', 'of', 'on', 'once',
                'only', 'or', 'other', 'our', 'ours','ourselves', 'out', 'own', 're',
                's', 'same', 'she', "shes", 'should', "shouldve",'so', 'some', 'such',
                't', 'than', 'that', "thatll", 'the', 'their', 'theirs', 'them',
                'themselves', 'then', 'there', 'these', 'they', 'this', 'those',
                'through', 'to', 'too','under', 'until', 'up', 've', 'very', 'was',
                'we', 'were', 'what', 'when', 'where','which','while', 'who', 'whom',
                'why', 'will', 'with', 'won', 'y', 'you', "youd","youll", "youre",
                "youve", 'your', 'yours', 'yourself', 'yourselves']
```

## Text Preprocessing Function

```python
def preprocess(textdata):
    processedText = []

    # Create Lemmatizer.
    wordLemm = WordNetLemmatizer()

    # Defining regex patterns.
    urlPattern        = r"((http://)[^ ]*|(https://)[^ ]*|( www\.)[^ ]*)"
    userPattern       = '@[^\s]+'
    alphaPattern      = "[^a-zA-Z0-9]"
    sequencePattern   = r"(.)\1\1+"
    seqReplacePattern = r"\1\1"

    for tweet in textdata:
        tweet = tweet.lower()

        # Replace all URls with 'URL'
        tweet = re.sub(urlPattern,' URL',tweet)
        # Replace all emojis.
        for emoji in emojis.keys():
            tweet = tweet.replace(emoji, "EMOJI" + emojis[emoji])
        # Replace @USERNAME to 'USER'.
        tweet = re.sub(userPattern,' USER', tweet)
        # Replace all non alphabets.
        tweet = re.sub(alphaPattern, " ", tweet)
        # Replace 3 or more consecutive letters by 2 letters.
        tweet = re.sub(sequencePattern, seqReplacePattern, tweet)

        tweetwords = ''
        for word in tweet.split():
            # Checking if the word is a stopword is commented out in original code
            #if word not in stopwordlist:
            if len(word)>1:
                # Lemmatizing the word.
                word = wordLemm.lemmatize(word)
                tweetwords += (word+' ')

        processedText.append(tweetwords)

    return processedText
```

## Run Preprocessing

```python
import time
t = time.time()
processedtext = preprocess(text)
print(f'Text Preprocessing complete.')
print(f'Time Taken: {round(time.time()-t)} seconds')
```
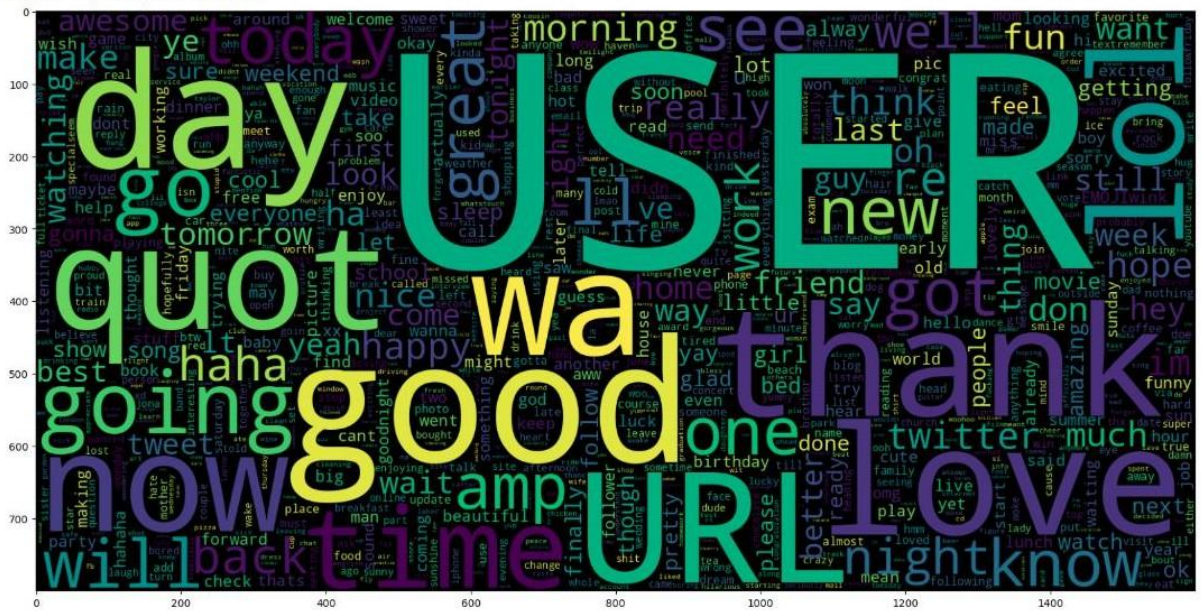
```
Text Preprocessing complete.
Time Taken: 114 seconds
```

## Word Cloud for Positive Tweets

```
[8]  data_pos = processedtext[800000:]
✓ 17s  wc = WordCloud(max_words = 1000 , width = 1600 , height = 800,
                collocations=False).generate(" ".join(data_pos))
       plt.figure(figsize = (20,20))
       plt.imshow(wc)
```
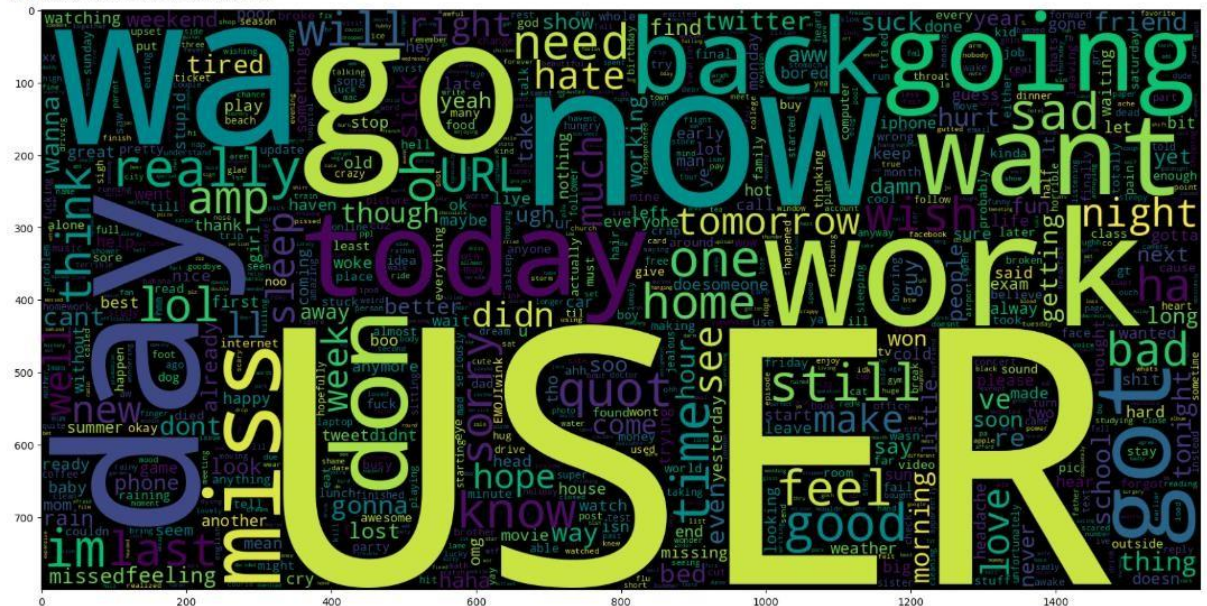
<matplotlib.image.AxesImage at 0x7dd4514e0f50>



## Word Cloud for Negative Tweets

```
[7]  data_neg = processedtext[:800000]
✓ 17s  plt.figure(figsize = (20,20))
       wc = WordCloud(max_words = 1000 , width = 1600 , height = 800,
                collocations=False).generate(" ".join(data_neg))
       plt.imshow(wc)
```

<matplotlib.image.AxesImage at 0x7dd46feeeff0>

## Splitting the Data

```
[9]
✓ 0s
    X_train, X_test, y_train, y_test = train_test_split(processedtext, sentiment,
                                                        test_size = 0.05, random_state = 0)
    print(f'Data Split done.')
```

    Data Split done.

## TF-IDF Vectorization

```
[10]
✓ 51s
    vectoriser = TfidfVectorizer(ngram_range=(1,2), max_features=500000)
    vectoriser.fit(X_train)
    print(f'Vectoriser fitted.')
    print('No. of feature_words: ', len(vectoriser.get_feature_names_out()))
```

    Vectoriser fitted.
    No. of feature_words:  500000

## Transforming the Dataset

```
[11]
✓ 38s
    X_train = vectoriser.transform(X_train)
    X_test  = vectoriser.transform(X_test)
    print(f'Data Transformed.')
```

    Data Transformed.

## Model Evaluation Function

```
[12]
✓ 0s
    def model_Evaluate(model):

        # Predict values for Test dataset
        y_pred = model.predict(X_test)

        # Print the evaluation metrics for the dataset.
        print(classification_report(y_test, y_pred))

        # Compute and plot the Confusion matrix
        cf_matrix = confusion_matrix(y_test, y_pred)

        categories  = ['Negative','Positive']
        group_names = ['True Neg','False Pos', 'False Neg','True Pos']
        group_percentages = ['{0:.2%}'.format(value) for value in cf_matrix.flatten() / np.sum(cf_matrix)]

        labels = [f'{v1}\n{v2}' for v1, v2 in zip(group_names,group_percentages)]
        labels = np.asarray(labels).reshape(2,2)

        sns.heatmap(cf_matrix, annot = labels, cmap = 'Blues',fmt = '',
                    xticklabels = categories, yticklabels = categories)

        plt.xlabel("Predicted values", fontdict = {'size':14}, labelpad = 10)
        plt.ylabel("Actual values"   , fontdict = {'size':14}, labelpad = 10)
        plt.title ("Confusion Matrix", fontdict = {'size':18}, pad = 20)
```
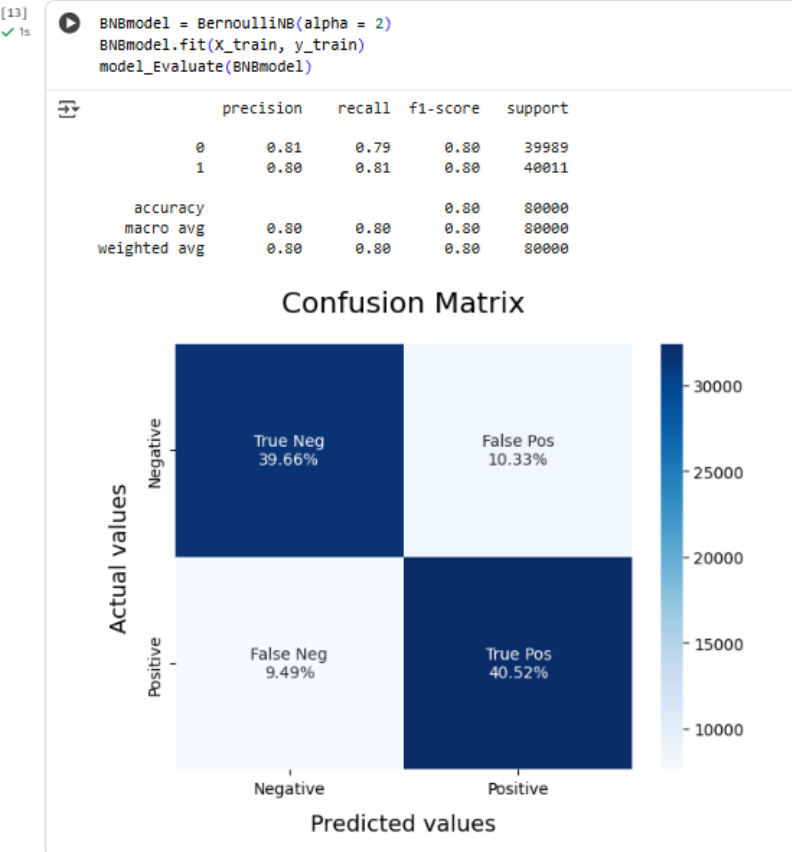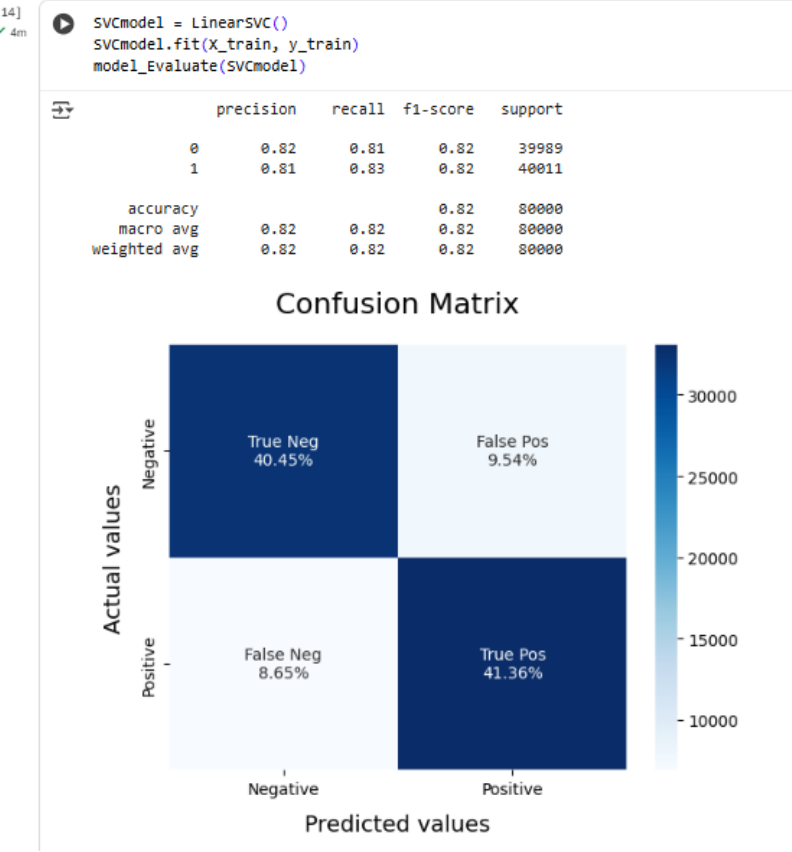
## Train and Evaluate BernoulliNB Model

```
[13]    BNBmodel = BernoulliNB(alpha = 2)
        BNBmodel.fit(X_train, y_train)
        model_Evaluate(BNBmodel)
```

```
                precision    recall  f1-score   support

            0       0.81      0.79      0.80     39989
            1       0.80      0.81      0.80     40011

     accuracy                           0.80     80000
    macro avg       0.80      0.80      0.80     80000
 weighted avg       0.80      0.80      0.80     80000
```

### Confusion Matrix



## Train and Evaluate LinearSVC Model

```
[14]    SVCmodel = LinearSVC()
        SVCmodel.fit(X_train, y_train)
        model_Evaluate(SVCmodel)
```

```
                precision    recall  f1-score   support

            0       0.82      0.81      0.82     39989
            1       0.81      0.83      0.82     40011

     accuracy                           0.82     80000
    macro avg       0.82      0.82      0.82     80000
 weighted avg       0.82      0.82      0.82     80000
```
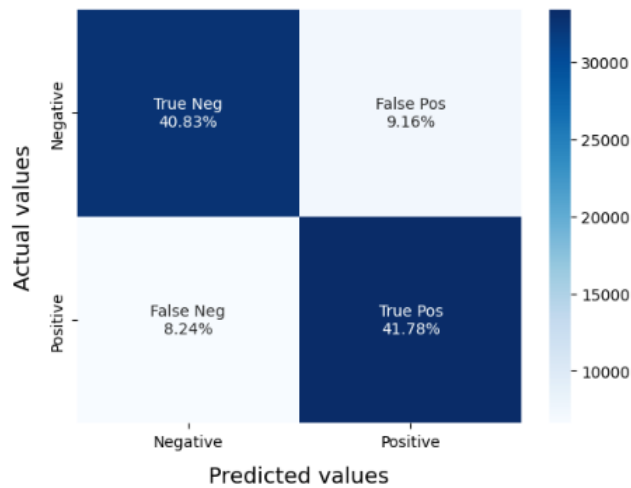
### Confusion Matrix

## Train and Evaluate Logistic Regression Model

```
[15]
✓ 24s
    LRmodel = LogisticRegression(C = 2, max_iter = 1000, n_jobs=-1)
    LRmodel.fit(X_train, y_train)
    model_Evaluate(LRmodel)
```

```
              precision    recall  f1-score   support

           0       0.83      0.82      0.82     39989
           1       0.82      0.84      0.83     40011

    accuracy                           0.83     80000
   macro avg       0.83      0.83      0.83     80000
weighted avg       0.83      0.83      0.83     80000
```

### Confusion Matrix



## Saving the Models

```
[16]
✓ 1s
    file = open('vectoriser-ngram-(1,2).pickle','wb')
    pickle.dump(vectoriser, file)
    file.close()

    file = open('Sentiment-LR.pickle','wb')
    pickle.dump(LRmodel, file)
    file.close()

    file = open('Sentiment-BNB.pickle','wb')
    pickle.dump(BNBmodel, file)
    file.close()
```

## Using the Saved Model for Prediction

```python
def load_models():
    '''
    Load the saved models from the local Colab environment.
    '''
    # Load the vectoriser.
    file = open('vectoriser-ngram-(1,2).pickle', 'rb')
    vectoriser = pickle.load(file)
    file.close()
    # Load the LR Model.
    file = open('Sentiment-LR.pickle', 'rb')
    LRmodel = pickle.load(file)
    file.close()

    return vectoriser, LRmodel

def predict(vectoriser, model, text):
    # Predict the sentiment
    textdata = vectoriser.transform(preprocess(text))
    sentiment = model.predict(textdata)

    # Make a list of text with sentiment.
    data = []
    for txt, pred in zip(text, sentiment):
        data.append((txt,pred))

    # Convert the list into a Pandas DataFrame.
    df = pd.DataFrame(data, columns = ['text','sentiment'])
    df = df.replace([0,1], ["Negative","Positive"])
    return df

# The __name__ == "__main__" block is used in scripts,
# but in a notebook we can run the code directly.

# Loading the models.
vectoriser, LRmodel = load_models()

# Text to classify should be in a list.
text = ["I hate twitter",
        "May the Force be with you.",
        "Mr. Stark, I don't feel so good"]

df = predict(vectoriser, LRmodel, text)
print(df.head())
```

```
                            text sentiment
0                 I hate twitter  Negative
1     May the Force be with you.  Positive
2  Mr. Stark, I don't feel so good  Negative
```