In [7]:

```python
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
#nltk.download('wordnet')

from nltk.stem import WordNetLemmatizer

# 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
```

In [71]:

```python
# Importing the dataset
DATASET_COLUMNS  = ["sentiment", "ids", "date", "flag", "user", "text"]
DATASET_ENCODING = "ISO-8859-1"
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 values to ease understanding.
dataset['sentiment'] = dataset['sentiment'].replace(4,1)

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

# Storing data in lists.
text, sentiment = list(dataset['text']), list(dataset['sentiment'])
dataset.groupby('sentiment').count()/dataset.shape[0]
```
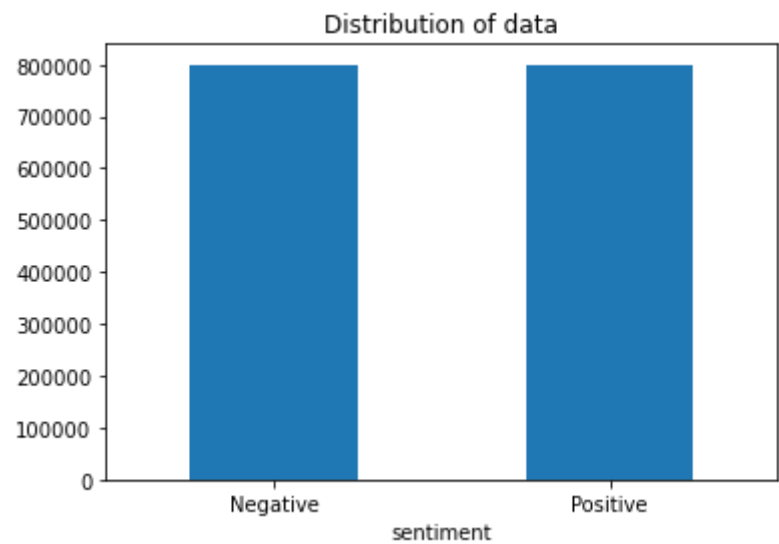
Out[71]:

| | text |
|---|---|
| **sentiment** | |
| **0** | 0.5 |
| **1** | 0.5 |

## Distribution of data



```python
In [52]:  dataset.head()
```

Out[52]:

| | sentiment | text |
|---|---|---|
| **0** | 0 | @switchfoot http://twitpic.com/2y1zl - Awww, t… |
| **1** | 0 | is upset that he can't update his Facebook by … |
| **2** | 0 | @Kenichan I dived many times for the ball. Man… |
| **3** | 0 | my whole body feels itchy and like its on fire |
| **4** | 0 | @nationwideclass no, it's not behaving at all…. |

```python
In [9]:  dataset.shape[0]
```

Out[9]:  1600000

```python
In [10]:  dataset.groupby('sentiment').count()
```

Out[10]:

| | text |
|---|---|
| **sentiment** | |
| **0** | 800000 |
| **1** | 800000 |

```python
In [11]:  dataset.head()
```

Out[11]:

| | sentiment | text |
|---|---|---|
| **0** | 0 | @switchfoot http://twitpic.com/2y1zl - Awww, t… |
| **1** | 0 | is upset that he can't update his Facebook by … |
| **2** | 0 | @Kenichan I dived many times for the ball. Man… |
| **3** | 0 | my whole body feels itchy and like its on fire |
| **4** | 0 | @nationwideclass no, it's not behaving at all…. |

```python
In [12]:  # Defining dictionary containing all emojis with their meanings.
          emojis = {':)': 'smile', ':-)': 'smile', ';d': 'wink', ':-E': 'vampire', ':(': 'sad'
                    ':-(': 'sad', ':-<': 'sad', ':P': 'raspberry', ':O': 'surprised',
```

```
        ':-@': 'shocked', ':@': 'shocked',':-$': 'confused', ':\\': 'annoyed',
        ':#': 'mute', ':X': 'mute', ':^)': 'smile', ':-&': 'confused', '$_$': 'gre
        '@@': 'eyeroll', ':-!': 'confused', ':-D': 'smile', ':-0': 'yell', 'O.o':
        '<(-_-)>': 'robot', 'd[-_-]b': 'dj', ":'-)": 'sadsmile', ';)': 'wink',
        ';-)': 'wink', 'O:-)': 'angel','O*-)': 'angel','(:-D': 'gossip', '=^.^=':

## 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']
```

# Preprocessing the data

In [66]:
```python
def preprocess(textdata):
    processedText = []

    # Create Lemmatizer and Stemmer.
    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 letter.
        tweet = re.sub(sequencePattern, seqReplacePattern, tweet)

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

```
        processedText.append(tweetwords)

    return processedText
```

In [72]:
```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: 276 seconds
```

In [75]:
```python
processedtext[1]
```

Out[75]: `'upset update facebook texting might cry result school today also blah '`

In [22]:
```python
dataset.head(10)
```

Out[22]:

| | sentiment | text |
|---|---|---|
| 0 | 0 | @switchfoot http://twitpic.com/2y1zl - Awww, t... |
| 1 | 0 | is upset that he can't update his Facebook by ... |
| 2 | 0 | @Kenichan I dived many times for the ball. Man... |
| 3 | 0 | my whole body feels itchy and like its on fire |
| 4 | 0 | @nationwideclass no, it's not behaving at all.... |
| 5 | 0 | @Kwesidei not the whole crew |
| 6 | 0 | Need a hug |
| 7 | 0 | @LOLTrish hey long time no see! Yes.. Rains a... |
| 8 | 0 | @Tatiana_K nope they didn't have it |
| 9 | 0 | @twittera que me muera ? |

In [ ]:

# Splitting the data

In [76]:
```python
X_train, X_test, y_train, y_test = train_test_split(processedtext, sentiment,test_si
print('Data Split done')
```

```
Data Split done
```

In [77]:
```python
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()))
```

```
Vectoriser fitted.
No. of feature_words:  500000
```

In [78]:
```python
vectoriser
```

Out[78]: `TfidfVectorizer(max_features=500000, ngram_range=(1, 2))`

In [79]:
```python
X_train = vectoriser.transform(X_train)
```

```
X_test  = vectoriser.transform(X_test)
print(f'Data Transformed.')
```

Data Transformed.

In [27]: `X_train`

Out[27]: 
```
<1520000x500000 sparse matrix of type '<class 'numpy.float64'>'
        with 31187017 stored elements in Compressed Sparse Row format>
```

# Evaluating Models

In [84]:
```python
from sklearn.metrics import accuracy_score
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))

    print('---------------------------')

    print(accuracy_score(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() /

    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)
```
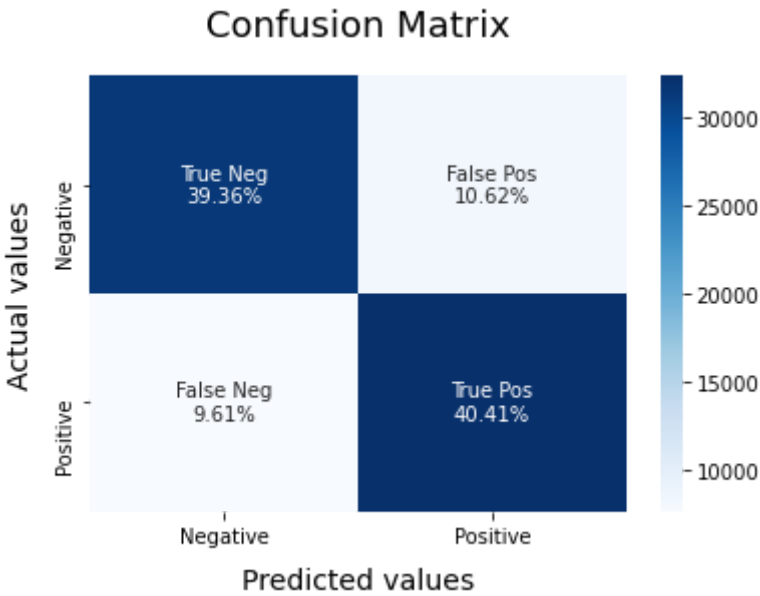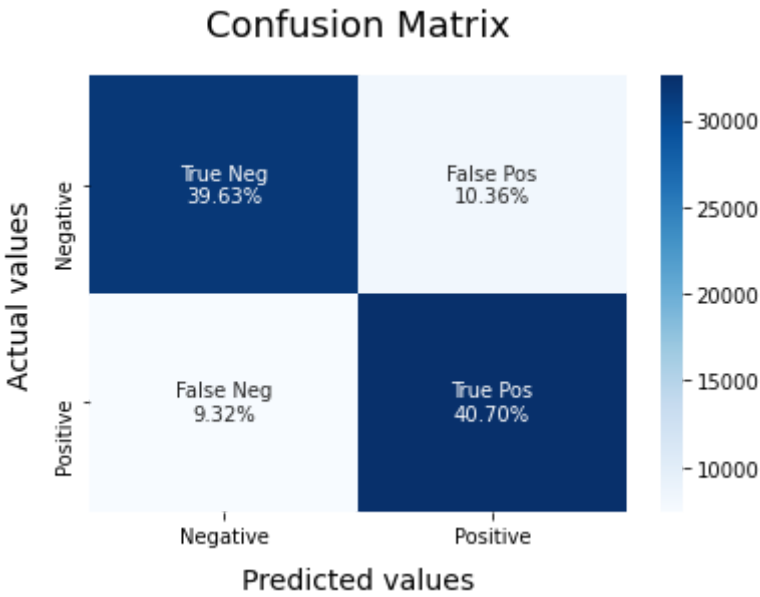
In [81]:
```python
BNBmodel = BernoulliNB(alpha = 2)
BNBmodel.fit(X_train, y_train)
model_Evaluate(BNBmodel)
```

```
              precision    recall  f1-score   support

           0       0.80      0.79      0.80     39989
           1       0.79      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
```
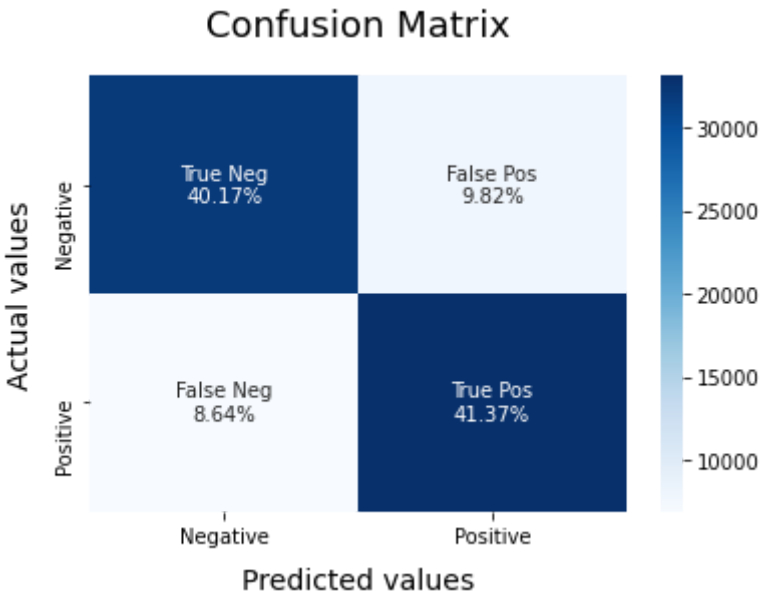
## Confusion Matrix



```
In [82]:    SVCmodel = LinearSVC()
            SVCmodel.fit(X_train, y_train)
            model_Evaluate(SVCmodel)
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.81      | 0.79   | 0.80     | 39989   |
| 1            | 0.80      | 0.81   | 0.81     | 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



```
In [83]:    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            | 0.82      | 0.80   | 0.81     | 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



```
In [85]:  model_Evaluate(LRmodel)
```

```
              precision    recall  f1-score   support

           0       0.82      0.80      0.81     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

---------------------------
0.8154125
```
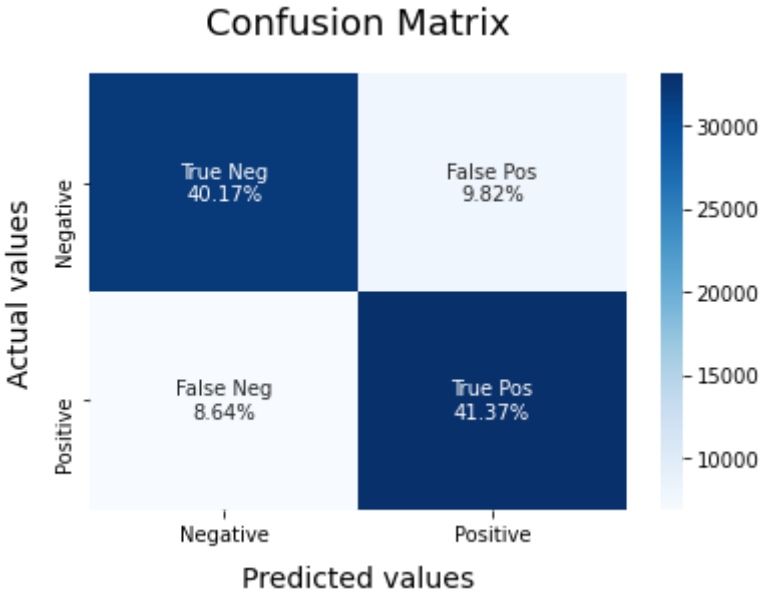
## Confusion Matrix



```
In [45]:  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()
```

In [90]:
```python
def load_models():
    '''
    Replace '..path/' by the path of the saved models.
    '''

    # Load the vectoriser.
    file = open('vectoriser-ngram-(1,2).pickle', 'rb')
    vectoriser = pickle.load(file)
    file.close()
    # Load the LR Model.
    file = open('Sentiment-LRv1.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 text, pred in zip(text, sentiment):
        data.append((text,pred))

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

if __name__=="__main__":
    # Loading the models.
    #vectoriser, LRmodel = load_models()

    # Text to classify should be in a list.
    text = ["@vijay a badddd",
            "May the Force be with you.",
            "Mr. Stark, I did not get a promotion"]

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

```
                                   text sentiment
0                       @vijay a badddd  Negative
1            May the Force be with you.  Positive
2  Mr. Stark, I did not get a promotion  Negative
```

In [ ]:
```
proce
```