	EXPERIMENT NO 6  SENTIMENT ANALYSIS  IMPORTING THE LIBRARIES
In [1]:	<pre># utilities import re import numpy as np import pandas as pd import nltk  # plotting import seaborn as sns</pre>
	<pre>from wordcloud import WordCloud import matplotlib.pyplot as plt  # nltk from nltk.stem import WordNetLemmatizer  # sklearn from sklearn.svm import LinearSVC</pre>
	<pre>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</pre> READING & PROCESSING THE DATASET
In [30]:	
	<pre># 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',</pre>
	<pre>legend=False) ax.set_xticklabels(['Negative', 'Positive'], rotation=0)  # Storing data in lists. text, sentiment = list(dataset['text']), list(dataset['sentiment'])  Distribution of data  800000</pre> <pre>Distribution of data</pre>
	700000 - 600000 - 500000 - 400000 - 300000 -
	200000 - 100000 - Negative Positive  PREPROCESSING THE DATA
	<ol> <li>Lower Casing</li> <li>Replacing URL's</li> <li>Replacing Emojis</li> <li>Replacing Usernames</li> <li>Removing Non-Alphabets</li> <li>Removing Consecutive letters"Heyyyy" to "Heyy")*</li> </ol>
In [31]:	<pre>7. Removing Short Words 8. Removing Stopwords 9. Lemmatizing  # Defining dictionary containing all emojis with their meanings. emojis = {    ':)': 'smile',</pre>
	<pre>':-)': 'smile', ';d': 'wink', ':-E': 'vampire', ':(': 'sad', ':-(': 'sad', ':-&lt;': 'sad', ':p': 'raspberry', ':0': 'surprised', ':-@': 'shocked',</pre>
	<pre>':@': 'shocked', ':-\$': 'confused', ':\\': 'annoyed', ':#': 'mute', ':X': 'mute', ':^)': 'smile', ':-&amp;': 'confused', '\$_\$': 'greedy', '@@': 'eyeroll',</pre>
	<pre>':-!': 'confused', ':-D': 'smile', ':-0': 'yell', 'O.o': 'confused', '&lt;()&gt;': 'robot', 'd[]b': 'dj', ":'-)": 'sadsmile', ';)': 'wink', ';-)': 'wink',</pre>
	<pre>'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',</pre>
	'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', 'ourselves', 'out', 'own', 're', 's', 'same', 'she', "shes", 'should', "shouldve", 'so', 'some', 'such', 't', 'than', 'that', "thatll", 'the', 'their', 'theirs',
Tn [22].	<pre>'them', 'themselves', 'then', 'there', 'these', 'they', 'this', 'those',     'through', 'to', 'under', 'until', 'up', 've', 'very', 'was', 'we',     'were', 'what', 'whene', 'which', 'while', 'who', 'whom', 'why',     'will', 'with', 'won', 'y', 'you', "youd", "youll", "youre", "youve",     'your', 'yourself', 'yourselves' ].append(nltk.corpus.stopwords.words())</pre>
In [32]:	<pre>def preprocess(textdata):     processedText = []  # Create Lemmatizer and Stemmer.     wordLemm = WordNetLemmatizer()  # Defining regex patterns.     urlPattern = r"((http://)[^ ]* (https://)[^ ]* ( www\.)[^ ]*)"     userPattern = '@[^\s]+'</pre>
	<pre>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)</pre>
	<pre># 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)</pre>
	<pre>tweetwords = '' for word in tweet.split():     # Checking if the word is a stopword.     # if word not in stopwordlist:     if len(word) &gt; 1:         # Lemmatizing the word.         word = wordLemm.lemmatize(word)         tweetwords += (word + ' ')</pre>
In [33]:	<pre>processedText.append(tweetwords)  return processedText  processedtext = preprocess(text) print(f'Text Preprocessing complete.')</pre>
In [34]:	data_neg = processedtext[.000000]
Out[34]:	<pre>plt.figure(figsize=(20, 20)) wc = WordCloud(max_words=1000, width=1600, height=800,</pre>
	month everyone trans said to set was all of the
	400 tell eat watch first bored and control of the c
	wentnextisses thing play play bear to be a sound transfer to be a so
In [35]: Out[35]:	<pre>Word-Cloud for Positive tweets.  data_pos = processedtext[800000:] wc = WordCloud(max_words=1000, width=1600, height=800,</pre>
	watching won'set getting saw bestelled need makes and the start saw bestelled need not save and the start save and the st
	take stands black stands
	beautiful everyone reading saidtend to be an interest of the said
In [36]:	SPLITTING THE DATA  X_train, X_test, y_train, y_test = train_test_split(processedtext, sentiment, test_size=0.05, prodem state=0.05, prodem state=0.05
In [37]:	<pre>print(f'Data Split done.')  Data Split done.  USING THE TF-IDF  vectoriser = TfidfVectorizer(ngram_range=(1, 2), max_features=500000) vectoriser.fit(X train)</pre>
	<pre>print(f'Vectoriser fitted.') print('No. of feature_words: ', len(vectoriser.get_feature_names()))  Vectoriser fitted.  /usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function get_feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Please use get_feature_names_out instead.     warnings.warn(msg, category=FutureWarning) No. of feature_words: 500000</pre>
In [38]:	DATA TRANSFORM (TRAIN & TEST) IN TF-IDF  X_train = vectoriser.transform(X_train) X_test = vectoriser.transform(X_test) print(f'Data Transformed.')  Data Transformed.
In [39]:	<pre>FUNCTION FOR MODEL EVALUATION  def model_Evaluate(model):     # Predict values for Test dataset     y_pred = model.predict(X_test)</pre>
	<pre># 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 = [</pre>
	<pre>'{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,</pre>
	<pre>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)</pre>
In [40]:	USING DIFFRENT MODEL AND FINDING ACCYURACY AND CONFUSION MATRIX  1. BernoulliNB Model  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 macro avg 0.80 0.80 0.80 80000 weighted avg 0.80 0.80 0.80 80000
	Confusion Matrix  - 30000  True Neg False Pos 10.35%
	Positive Pos 40.55%  False Neg 9.47%  False Neg 9.47%  False Neg 40.55%  - 15000  - 10000
In [ ]:	SVCmodel.fit(X_train, y_train)
In [42]:	<pre>model_Evaluate(SVCmodel)  3. Logistic Regression Model  LRmodel = LogisticRegression(C=2, max_iter=1000, n_jobs=-1) LRmodel.fit(X_train, y_train) model_Evaluate(LRmodel)</pre>
	precision recall f1-score support  0 0.83 0.82 0.83 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  False Pos 9.05%  - 30000 - 25000
	- 25000 - 20000 - 20000 - 15000 - 10000 - 10000
	Predicted values  We can clearly see that the Logistic Regression Model performs the best out of all the different models that we tried. It achieves nearly 82% accuracy while classifying the sentiment of a tweet.