In [1]:	<pre>import pandas as pd import numpy as np import re import matplotlib.pyplot as plt import seaborn as sns import string import nltk import warnings %matplotlib inline warnings.filterwarnings('ignore') from matplotlib import style style.use('ggplot')</pre>
In [2]: In [3]: Out[3]:	<pre>plt.style.use('fivethirtyeight')  df = pd.read_csv(r"D:\900338-SonaSelvaraj(Week 4)\Twitter Sentiments.csv")  df.head()</pre>
In [4]: Out[4]:	
In [5]:	31960 31961 1 @user #sikh #temple vandalised in in #calgary,  31961 31962 0 thank you @user for you follow  df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 31962 entries, 0 to 31961 Data columns (total 3 columns): # Column Non-Null Count Dtype</class>
In [6]:	0 id 31962 non-null int64 1 label 31962 non-null int64 2 tweet 31962 non-null object dtypes: int64(2), object(1) memory usage: 749.2+ KB  PREPROCESSING
In [7]: In [8]: Out[8]:	<pre>for word in r:     input_txt = re.sub(word, "", input_txt)     return input_txt  #remove twitter handles(@user) df['clean_tweet'] = np.vectorize(remove_pattern)(df['tweet'], "@[\w]*")  df.head()</pre>
In [9]:	thanks for #lyft credit i can't us thanks for #lyft credit i can't use cause th bihday your majesty bihday your majesty  #model i love u take with u all the time in #model i love u take with u all the time in factsguide: society now #motivation  factsguide: society now #motivation
Out[9]: In [10]:	<ul> <li>0 @user when a father is dysfunctional and is s when a father is dysfunctional and is so sel</li> <li>1 2 0 @user @user thanks for #lyft credit i can't us thanks for #lyft credit i can t use cause th</li> <li>2 3 0 bihday your majesty bihday your majesty</li> <li>3 4 0 #model i love u take with u all the time in #model i love u take with u all the time in</li> <li>4 5 0 factsguide: society now #motivation</li> <li>factsguide society now #motivation</li> </ul>
Out[10]:	<pre>df['clean_tweet'] = df['clean_tweet'].apply(lambda x: " ".join([w for w in x.split() if len(w)&gt;3])) df.head()</pre>
In [11]: Out[11]:	<pre>tokenized_tweet = df['clean_tweet'].apply(lambda x: x.split()) tokenized_tweet.head()  0    [when, father, dysfunctional, selfish, drags, 1    [thanks, #lyft, credit, cause, they, offer, wh 2</pre>
In [12]: Out[12]:	<pre>from nltk.stem.porter import PorterStemmer stemmer = PorterStemmer()  tokenized_tweet = tokenized_tweet.apply(lambda sentence: [stemmer.stem(word) for word in sentence]) tokenized_tweet.head()</pre>
<pre>In [13]: Out[13]:</pre>	<pre># combine words into single sentence for i in range(len(tokenized_tweet)):     tokenized_tweet[i] = " ".join(tokenized_tweet[i])  df['clean_tweet'] = tokenized_tweet df.head()</pre>
In [14]:	
In [15]:	Requirement already satisfied: wordcloud in c:\users\sona\anaconda3\lib\site-packages (1.6.0)Note: you may need to restart the kernel to use updated packages. Requirement already satisfied: pillow in c:\users\sona\anaconda3\lib\site-packages (from wordcloud) (9.2.0) Requirement already satisfied: numpy>=1.6.1 in c:\users\sona\anaconda3\lib\site-packages (from wordcloud) (1.19.2) Requirement already satisfied: matplotlib in c:\users\sona\anaconda3\lib\site-packages (from wordcloud) (3.3.2) Requirement already satisfied: certifi>=2020.06.20 in c:\users\sona\anaconda3\lib\site-packages (from matplotlib->wordcloud) (2.022.6.15) Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\sona\anaconda3\lib\site-packages (from matplotlib->wordcloud) (1.3.0) Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.3 in c:\users\sona\anaconda3\lib\site-packages (from matplotlib->wordcloud) (2.4.7) Requirement already satisfied: cycler>=0.10 in c:\users\sona\anaconda3\lib\site-packages (from matplotlib->wordcloud) (0.10.0) Requirement already satisfied: python-dateutil>=2.1 in c:\users\sona\anaconda3\lib\site-packages (from matplotlib->wordcloud) (2.8.1) Requirement already satisfied: six in c:\users\sona\anaconda3\lib\site-packages (from matplotlib->wordcloud) (1.15.0)  # visualize the frequent words all_words = " ".join([sentence for sentence in df['clean_tweet']])
	from wordcloud import WordCloud wordcloud = WordCloud(width=800, height=500, random_state=42, max_font_size=100).generate(all_words)  # plot the graph plt.figure(figsize=(15,8)) plt.imshow(wordcloud, interpolation='bilinear') plt.axis('off') plt.show()  think tonight read happen
	believ someth the best Ody even by Company black of the best of th
	video cute amaz food women feel will model love funni sta danc follow game true model sta danc follow game true model sta danc follow game chang of sunday done state thought noth hour good didn selfi thought now still good morn selfi thought now still good morn good didn selfi thought now still good didn selfi thought now
In [16]:	# frequent words visualization for +ve all_words = " ".join([sentence for sentence in df['clean_tweet'][df['label']==0]]) wordcloud = WordCloud(width=800, height=500, random_state=42, max_font_size=100).generate(all_words) # plot the graph plt.figure(figsize=(15,8)) plt.imshow(wordcloud, interpolation='bilinear') plt.axis('off')
	take time alway calltravel witter manipular bear calling to the state of the state
In [17]:	# frequent words visualization for -ve all_words = " ".join([sentence for sentence in df['clean_tweet'][df['label']==1]])  wordcloud = WordCloud(width=800, height=500, random_state=42, max_font_size=100).generate(all_words)
	# plot the graph plt.figure(figsize=(15,8)) plt.imshow(wordcloud, interpolation='bilinear') plt.axis('off') plt.show()   good Word true peopl look wso condemn new calgari wso find equal need world white peopl love stop love st
	In a terror retweet will problem templ vandalis hispan feel to the pleas time american girl atest fuck pleas to the people on line better latest fuck media manifamili porniti tampa miami kkk pleas tampa miami kkk pleas to the people on line better latest fuck media manifamili porniti tampa miami kkk pleas to the people on line better latest fuck media manifamili problem problem problem problem problem problem problem problem problem think problem problem think problem problem problem think problem problem think problem problem think problem problem think problem problem problem think problem problem problem think problem p
	Take thing happen talk bling suppose made that thank t
In [18]: In [19]:	<pre>def hashtag_extract(tweets):     hashtags = [] # loop words in the tweet for tweet in tweets:     ht = re.findall(r"#(\w+)", tweet)     hashtags.append(ht) return hashtags</pre>
In [21]:	<pre>[['run'], ['lyft', 'disapoint', 'getthank'], [], ['model'], ['motiv']]  ht_negative[:5]  [['cnn', 'michigan', 'tcot'],   ['australia',</pre>
In [22]:	<pre>'opkillingbay', 'seashepherd', 'helpcovedolphin', 'thecov', 'helpcovedolphin'], [], [], [meverump', 'xenophobia']]  # unnest list ht_positive = sum(ht_positive, []) ht_negative = sum(ht_negative, [])</pre>
<pre>In [23]: Out[23]: In [24]: Out[24]:</pre>	<pre>['cnn', 'michigan', 'tcot', 'australia', 'opkillingbay']  freq = nltk.FreqDist(ht_positive) d = pd.DataFrame({'Hashtag': list(freq.keys()),</pre>
In [25]:	<pre>1     lyft     2 2     disapoint     1 3     getthank     2 4     model     375  # select top 10 hashtags d = d.nlargest(columns='Count', n=10) plt.figure(figsize=(15,9)) sns.barplot(data=d, x='Hashtag', y='Count')</pre>
	1600 — 1400 — 12
	1000 - 800 - 600 -
	200 - O love posit smile healthi thank fun life affirm summer model Hashtag
In [26]: Out[26]:	<pre>d = pd.DataFrame({'Hashtag': list(freq.keys()),</pre>
In [27]:	<pre>3  australia  6 4  opkillingbay  5  # select top 10 hashtags d = d.nlargest(columns='Count', n=10) plt.figure(figsize=(15,9)) sns.barplot(data=d, x='Hashtag', y='Count') plt.show()</pre> 140
	120 -
	80 - 60 - 40 - 40 - 40 - 40 - 40 - 40 - 4
	The strump polit allahsoil liber libtard sjw retweet black miami hate liptard Hashtag
In [28]: In [29]:	<pre># feature extraction from sklearn.feature_extraction.text import CountVectorizer bow_vectorizer = CountVectorizer(max_df=0.90, min_df=2, max_features=1000, stop_words='english') bow = bow_vectorizer.fit_transform(df['clean_tweet'])</pre>
In [30]: In [31]:	<pre>from sklearn.model_selection import train_test_split  # Split dataset into training set and test set X_train, X_test, y_train, y_test = train_test_split(bow, df['label'], random_state=42, test_size=0.25)</pre>
In [36]:	<pre>#Create a svm Classifier clf = svm.SVC(kernel='linear') # Linear Kernel  #Train the model using the training sets clf.fit(X_train, y_train)  #Predict the response for test dataset y_pred = clf.predict(X_test)  #Import scikit-learn metrics module for accuracy calculation from sklearn import metrics</pre>
In [37]: Out[37]:	<pre>from sklearn.metrics import confusion_matrix, classification_report, ConfusionMatrixDisplay from sklearn.metrics import f1_score  # Model Accuracy: how often is the classifier correct? print("Accuracy:",metrics.accuracy_score(y_test, y_pred))  Accuracy: 0.9473157301964711  # training model = svm.SVC(kernel='linear', C = 1.0) model.fit(x_train, y_train)</pre>
In [38]:	<pre># testing pred = model.predict(x_test) f1_score(y_test, pred)  0.4994054696789536  print(confusion_matrix(y_test, y_pred)) print("\n") print(classification_report(y_test, y_pred))  [[7360 72]</pre>
Τ×	precision recall f1-score support  0 0.95 0.99 0.97 7432 1 0.74 0.38 0.50 559  accuracy 0.95 7991 macro avg 0.85 0.68 0.74 7991 weighted avg 0.94 0.95 0.94 7991
In [40]:	Requirement already satisfied: reportlab in c:\users\sona\anaconda3\lib\site-packages (3.6.11) Requirement already satisfied: pillow>=9.0.0 in c:\users\sona\anaconda3\lib\site-packages (from reportlab) (9.2.0) Note: you may need to restart the kernel to use updated packages.  cm = confusion_matrix(y_test, y_pred) disp = ConfusionMatrixDisplay(confusion_matrix = cm)  categories = ['Negative', 'Positive'] group_names = ['True Neg', 'False Pos', 'False Neg', 'True Pos'] group_percentages = ['{0:.2%}'.format(value) for value in cm.flatten() / np.sum(cm)]
Out[44]:	<pre>labels = [f'{v1}\n{v2}' for v1, v2 in zip(group_names, group_percentages)] labels = np.asarray(labels).reshape(2,2)  sns.heatmap(cm, annot = labels, cmap = 'Blues', fmt = '',</pre>
	True Neg 92.10% False Pos 0.90% 5000 4000 3000
In [45]:	Negative Positive Predicted values  BernoulliNB Model
<sub>[</sub> 45]:	<pre># 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']</pre>
	<pre>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 = '',</pre>
In [46]:	plt.title ("Confusion Matrix", fontdict = {'size':18}, pad = 20)
	macro avg 0.75 0.74 0.74 7991 weighted avg 0.93 0.94 0.93 7991  Confusion Matrix  True Neg False Pos 89.99% 3.02%
	False Neg True Pos 3.48% 2000 1000  Negative Positive Predicted values
In [49]: In [50]:	LRmodel.fit(X_train, y_train) model_Evaluate(LRmodel)  precision recall f1-score support  0 0.96 0.99 0.97 7432
	1 0.72 0.39 0.50 559  accuracy macro avg 0.84 0.69 0.74 7991 weighted avg 0.94 0.95 0.94 7991  Confusion Matrix  True Neg False Pos 91.95% 1.05%
	False Neg 4.29%  False Neg 2.70%  Negative  Positive
In [ ]:	Predicted values