

In [4]:

```

1 #Load the Libraries
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6
7 import nltk
8 from sklearn.feature_extraction.text import CountVectorizer
9 from sklearn.feature_extraction.text import TfidfVectorizer
10 from sklearn.preprocessing import LabelBinarizer
11 from nltk.corpus import stopwords
12 from nltk.stem.porter import PorterStemmer
13
14 from wordcloud import WordCloud,STOPWORDS
15 from nltk.stem import WordNetLemmatizer
16 from nltk.tokenize import word_tokenize,sent_tokenize
17 from bs4 import BeautifulSoup
18
19 import spacy
20 import re,string,unicodedata
21 from nltk.tokenize.toktok import ToktokTokenizer
22 from nltk.stem import LancasterStemmer,WordNetLemmatizer
23 from sklearn.linear_model import LogisticRegression,SGDClassifier
24 from sklearn.naive_bayes import MultinomialNB
25 from sklearn.svm import SVC
26 from textblob import TextBlob
27 from textblob import Word
28 from sklearn.metrics import classification_report,confusion_matrix,accuracy_score
29
30 import os
31 print(os.listdir("C:/Users/admin/Desktop/Untitled Folder"))
32 import warnings
33 warnings.filterwarnings('ignore')
34

```

['.ipynb_checkpoints', 'Assignment 2 DL.ipynb', 'desktop.ini', 'DL Binary Classification.ipynb', 'DL Binary.ipynb', 'DL Linear Regression.ipynb', 'IMDB Dataset.csv']

Sentiment Analysis of IMDB Movie Reviews

Problem Statement:

In this, we have to predict the number of positive and negative reviews based on sentiments by using different classification models.

Import necessary libraries

Import the training dataset

```
In [5]: 1 #importing the training data
2 imdb_data=pd.read_csv('IMDB Dataset.csv')
3 print(imdb_data.shape)
4 imdb_data.head(10)
```

(50000, 2)

Out[5]:

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production. The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Mattei's "Love in the Time of Money" is...	positive
5	Probably my all-time favorite movie, a story o...	positive
6	I sure would like to see a resurrection of a u...	positive
7	This show was an amazing, fresh & innovative i...	negative
8	Encouraged by the positive comments about this...	negative
9	If you like original gut wrenching laughter yo...	positive

Exploratory data analysis

```
In [6]: 1 #Summary of the dataset
2 imdb_data.describe()
```

Out[6]:

	review	sentiment
count	50000	50000
unique	49582	2
top	Loved today's show!!! It was a variety and not...	negative
freq	5	25000

Sentiment count

```
In [7]: 1 #sentiment count
2 imdb_data['sentiment'].value_counts()
```

```
Out[7]: negative    25000
positive     25000
Name: sentiment, dtype: int64
```

We can see that the dataset is balanced.

Splitting the training dataset

In [8]:

```

1 #split the dataset
2 #train dataset
3 train_reviews=imdb_data.review[:40000]
4 train_sentiments=imdb_data.sentiment[:40000]
5 #test dataset
6 test_reviews=imdb_data.review[40000:]
7 test_sentiments=imdb_data.sentiment[40000:]
8 print(train_reviews.shape,train_sentiments.shape)
9 print(test_reviews.shape,test_sentiments.shape)

```

(40000,) (40000,
(10000,) (10000,)

Text normalization

In [9]:

```

1 #Tokenization of text
2
3
4 tokenizer=ToktokTokenizer()
5 #Setting English stopwords
6 stopword_list=nltk.corpus.stopwords.words('english')

```

Removing html strips and noise text

In [10]:

```

1 #Removing the html strips
2 def strip_html(text):
3     soup = BeautifulSoup(text, "html.parser")
4     return soup.get_text()
5
6 #Removing the square brackets
7 def remove_between_square_brackets(text):
8     return re.sub('\[\[^]*\]', '', text)
9
10 #Removing the noisy text
11 def denoise_text(text):
12     text = strip_html(text)
13     text = remove_between_square_brackets(text)
14     return text
15 #Apply function on review column
16 imdb_data['review']=imdb_data['review'].apply(denoise_text)

```

Removing special characters

In [11]:

```
1 #Define function for removing special characters
2 def remove_special_characters(text, remove_digits=True):
3     pattern=r'[^a-zA-Z0-9\s]'
4     text=re.sub(pattern,'',text)
5     return text
6 #Apply function on review column
7 imdb_data['review']=imdb_data['review'].apply(remove_special_characters)
```

*Text stemming *

In [14]:

```
1 #Stemming the text
2 def simple_stemmer(text):
3     ps=nltk.porter.PorterStemmer()
4     text= ' '.join([ps.stem(word) for word in text.split()])
5     return text
6 #Apply function on review column
7 imdb_data['review']=imdb_data['review'].apply(simple_stemmer)
```

Removing stopwords

In [15]:

```

1 #set stopwords to english
2 stop=set(stopwords.words('english'))
3 print(stop)
4
5 #removing the stopwords
6 def remove_stopwords(text, is_lower_case=False):
7     tokens = tokenizer.tokenize(text)
8     tokens = [token.strip() for token in tokens]
9     if is_lower_case:
10         filtered_tokens = [token for token in tokens if token not in stop]
11     else:
12         filtered_tokens = [token for token in tokens if token.lower() not
13     filtered_text = ' '.join(filtered_tokens)
14     return filtered_text
15 #Apply function on review column
16 imdb_data['review']=imdb_data['review'].apply(remove_stopwords)

```

{'i', 'she', 'theirs', 'was', 'yours', "she's", "doesn't", 'they', 'ourselves', 'from', 'too', "won't", "you're", 't', 'her', 'at', 'having', 'against', 'few', 'couldn', 'mustn', 'am', 'have', 'above', 'now', 'until', 'an', 'some', 'is', "should've", 'how', 'we', 'but', 'these', 'when', 'you', 'needn', 'didn', 'weren', 'can', 'doing', "you'd", 'same', 'while', 'nor', 'that', "it's", 'all', 'shan', 'more', 'any', 'each', 'this', 'before', "wasn't", 'to', "didn't", "you'll", 'than', 'what', 'does', 'aren', 'off', 'not', 'don', 'only', "you've", 'those', 'both', 'did', 'and', 'the', 'with', 're', 'being', "shan't", 'should', 'where', 'hers', 'their', "mustn't", 'do', 'a', 'which', 'own', 'of', 'because', 'yourselves', "hasn't", 'will', 's', 'once', 'haven', 'why', 'or', 'its', 'ma', 'out', 'such', "aren't", 'under', "needn't", 'won', 'whom', 'y', 'in', 'm', 'itself', 'below', 'if', 'had', 'after', 'he', 'd', 'be', 'for', 'them', 'into', "couldn't", "hadn't", "mightn't", 'further', 'ain', 'ours', 'him', 'me', 'just', 'then', 'here', 'about', 'my', 'myself', 'hadn', 'so', "weren't", 'has', 'himself', 'there', 'are', 've', 'been', 'up', 'o', "don't", "haven't", "isn't", 'll', 'very', 'wasn', 'mightn', 'wouldn', 'on', 'during', 'themselves', 'were', 'who', 'shouldn', 'by', 'through', 'over', 'your', 'between', 'doesn', "wouldn't", 'as', "shouldn't", 'herself', 'it', 'again', 'other', 'isn', 'his', 'down', 'yourself', 'most', "that'll", 'hasn', 'our', 'no'}

Normalized train reviews

In [16]:

```
1 #normalized train reviews
2 norm_train_reviews=imdb_data.review[:40000]
3 norm_train_reviews[0]
4 #convert dataframe to string
5 #norm_train_string=norm_train_reviews.to_string()
6 #Spelling correction using Textblob
7 #norm_train_spelling=TextBlob(norm_train_string)
8 #norm_train_spelling.correct()
9 #Tokenization using Textblob
10 #norm_train_words=norm_train_spelling.words
11 #norm_train_words
```

Out[16]:

'one review ha mention watch 1 oz episod youll hook right thi exactli happen
meth first thing struck oz wa brutal unflinch scene violenc set right word go
trust thi show faint heart timid thi show pull punch regard drug sex violenc
hardcor classic use wordit call oz nicknam given oswald maximum secur state p
enitentari focus mainli emerald citi experiment section prison cell glass fro
nt face inward privaci high agenda em citi home manyaryan muslim gangsta lati
no christian italian irish moreso scuffl death stare dodgi deal shadi agreeme
nt never far awayi would say main appeal show due fact goe show woudnt dare
forget pretti pictur paint mainstream audienc forget charm forget romanceoz d
oesnt mess around first episod ever saw struck nasti wa surreal couldnt say w
a readi watch develop tast oz got accustom high level graphic violenc violenc
injustic crook guard wholl sold nickel inmat wholl kill order get away well m
anner middl class inmat turn prison bitch due lack street skill prison experi
watch oz may becom comfort uncomfornt viewingthat get touch darker side'

Normalized test reviews

In [17]:

```

1 #Normalized test reviews
2 norm_test_reviews=imdb_data.review[40000:]
3 norm_test_reviews[45005]
4 ##convert dataframe to string
5 #norm_test_string=norm_test_reviews.to_string()
6 #spelling correction using Textblob
7 #norm_test_spelling=TextBlob(norm_test_string)
8 #print(norm_test_spelling.correct())
9 #Tokenization using Textblob
10 #norm_test_words=norm_test_spelling.words
11 #norm_test_words

```

Out[17]:

'read review watch thi piec cinemat garbag took least 2 page find somebody el s didnt think thi appallingli unfunni montag wasnt acm humour 70 inde ani era thi isnt least funni set sketch comed i ive ever seen itll till come along hal f skit alreadi done infinit better act monti python woodi allen wa say nice p iec anim last 90 second highlight thi film would still get close sum mindless drivelridden thi wast 75 minut semin comed i onli world semin realli doe mean semen scatalog humour onli world scat actual fece precursor joke onli mean th i handbook comed i tit bum odd beaver niceif pubesc boy least one hand free ha vent found playboy exist give break becaus wa earli 70 way sketch comed i go b ack least ten year prior onli way could even forgiv thi film even made wa gun point retro hardli sketch clown subtli pervert children may cut edg circl cou ld actual funni come realli quit sad kept go throughout entir 75 minut sheer belief may save genuin funni skit end gave film 1 becaus wa lower scoreand on li recommend insomniac coma patientsor perhap peopl suffer lockjawtheir jaw w ould final drop open disbelief'

**Bags of words model **

It is used to convert text documents to numerical vectors or bag of words.

In [18]:

```

1 #Count vectorizer for bag of words
2 cv=CountVectorizer(min_df=0,max_df=1,binary=False,ngram_range=(1,3))
3 #transformed train reviews
4 cv_train_reviews=cv.fit_transform(norm_train_reviews)
5 #transformed test reviews
6 cv_test_reviews=cv.transform(norm_test_reviews)
7
8 print('BOW_cv_train:',cv_train_reviews.shape)
9 print('BOW_cv_test:',cv_test_reviews.shape)
10 #vocab=cv.get_feature_names()-toget feature names

```

BOW_cv_train: (40000, 6209089)

BOW_cv_test: (10000, 6209089)

Term Frequency-Inverse Document Frequency model (TFIDF)

It is used to convert text documents to matrix of tfidf features.

```
In [19]: 1 #Tfidf vectorizer
2 tv=TfidfVectorizer(min_df=0,max_df=1,use_idf=True,ngram_range=(1,3))
3 #transformed train reviews
4 tv_train_reviews=tv.fit_transform(norm_train_reviews)
5 #transformed test reviews
6 tv_test_reviews=tv.transform(norm_test_reviews)
7 print('Tfidf_train:',tv_train_reviews.shape)
8 print('Tfidf_test:',tv_test_reviews.shape)
```

Tfidf_train: (40000, 6209089)
Tfidf_test: (10000, 6209089)

Labeling the sentiment text

```
In [20]: 1 #Labeling the sentient data
2 lb=LabelBinarizer()
3 #transformed sentiment data
4 sentiment_data=lb.fit_transform(imdb_data['sentiment'])
5 print(sentiment_data.shape)
```

(50000, 1)

Split the sentiment tdata

```
In [21]: 1 #Splitting the sentiment data
2 train_sentiments=sentiment_data[:40000]
3 test_sentiments=sentiment_data[40000:]
4 print(train_sentiments)
5 print(test_sentiments)
```

[[1]
[1]
[1]
...
[1]
[0]
[0]]
[[0]
[0]
[0]
...
[0]
[0]
[0]]]

Modelling the dataset

Let us build logistic regression model for both bag of words and tfidf features

In [39]:

```

1 #training the model
2 lr=LogisticRegression(penalty='l2',max_iter=500,C=1,random_state=42)
3 #Fitting the model for Bag of words
4 lr_bow=lr.fit(cv_train_reviews,train_sentiments)
5 print(lr_bow)
6 #Fitting the model for tfidf features
7 lr_tfidf=lr.fit(tv_train_reviews,train_sentiments)
8 print(lr_tfidf)

```

```
LogisticRegression(C=1, max_iter=500, random_state=42)
LogisticRegression(C=1, max_iter=500, random_state=42)
```

Logistic regression model performance on test dataset

In [23]:

```

1 #Predicting the model for bag of words
2 lr_bow_predict=lr.predict(cv_test_reviews)
3 print(lr_bow_predict)
4 ##Predicting the model for tfidf features
5 lr_tfidf_predict=lr.predict(tv_test_reviews)
6 print(lr_tfidf_predict)

```

```
[0 0 0 ... 0 1 1]
[0 0 0 ... 0 1 1]
```

Accuracy of the model

In [24]:

```

1 #Accuracy score for bag of words
2 lr_bow_score=accuracy_score(test_sentiments,lr_bow_predict)
3 print("lr_bow_score :",lr_bow_score)
4 #Accuracy score for tfidf features
5 lr_tfidf_score=accuracy_score(test_sentiments,lr_tfidf_predict)
6 print("lr_tfidf_score :",lr_tfidf_score)

```

```
lr_bow_score : 0.7512
lr_tfidf_score : 0.75
```

Print the classification report

In [25]:

```

1 #Classification report for bag of words
2 lr_bow_report=classification_report(test_sentiments,lr_bow_predict,target_
3 print(lr_bow_report)
4
5 #Classification report for tfidf features
6 lr_tfidf_report=classification_report(test_sentiments,lr_tfidf_predict,tar_
7 print(lr_tfidf_report)

```

	precision	recall	f1-score	support
Positive	0.75	0.75	0.75	4993
Negative	0.75	0.75	0.75	5007
accuracy			0.75	10000
macro avg	0.75	0.75	0.75	10000
weighted avg	0.75	0.75	0.75	10000
	precision	recall	f1-score	support
Positive	0.74	0.77	0.75	4993
Negative	0.76	0.73	0.75	5007
accuracy			0.75	10000
macro avg	0.75	0.75	0.75	10000
weighted avg	0.75	0.75	0.75	10000

Confusion matrix

In [26]:

```

1 #confusion matrix for bag of words
2 cm_bow=confusion_matrix(test_sentiments,lr_bow_predict,labels=[1,0])
3 print(cm_bow)
4 #confusion matrix for tfidf features
5 cm_tfidf=confusion_matrix(test_sentiments,lr_tfidf_predict,labels=[1,0])
6 print(cm_tfidf)

```

```

[[3768 1239]
 [1249 3744]]
[[3663 1344]
 [1156 3837]]

```

Stochastic gradient descent or Linear support vector machines for bag of words and tfidf features

In [27]:

```

1 #training the Linear svm
2 svm=SGDClassifier(loss='hinge',max_iter=500,random_state=42)
3 #fitting the svm for bag of words
4 svm_bow=svm.fit(cv_train_reviews,train_sentiments)
5 print(svm_bow)
6 #fitting the svm for tfidf features
7 svm_tfidf=svm.fit(tv_train_reviews,train_sentiments)
8 print(svm_tfidf)

```

```

SGDClassifier(max_iter=500, random_state=42)
SGDClassifier(max_iter=500, random_state=42)

```

Model performance on test data

In [28]:

```

1 #Predicting the model for bag of words
2 svm_bow_predict=svm.predict(cv_test_reviews)
3 print(svm_bow_predict)
4 #Predicting the model for tfidf features
5 svm_tfidf_predict=svm.predict(tv_test_reviews)
6 print(svm_tfidf_predict)

```

```

[1 1 0 ... 1 1 1]
[1 1 1 ... 1 1 1]

```

Accuracy of the model

In [29]:

```

1 #Accuracy score for bag of words
2 svm_bow_score=accuracy_score(test_sentiments,svm_bow_predict)
3 print("svm_bow_score :",svm_bow_score)
4 #Accuracy score for tfidf features
5 svm_tfidf_score=accuracy_score(test_sentiments,svm_tfidf_predict)
6 print("svm_tfidf_score :",svm_tfidf_score)

```

```

svm_bow_score : 0.5829
svm_tfidf_score : 0.5112

```

Print the classification report

In [30]:

```

1 #Classification report for bag of words
2 svm_bow_report=classification_report(test_sentiments,svm_bow_predict,target_names=['0','1'])
3 print(svm_bow_report)
4 #Classification report for tfidf features
5 svm_tfidf_report=classification_report(test_sentiments,svm_tfidf_predict,target_names=['0','1'])
6 print(svm_tfidf_report)

```

	precision	recall	f1-score	support
Positive	0.94	0.18	0.30	4993
Negative	0.55	0.99	0.70	5007
accuracy			0.58	10000
macro avg	0.74	0.58	0.50	10000
weighted avg	0.74	0.58	0.50	10000
	precision	recall	f1-score	support
Positive	1.00	0.02	0.04	4993
Negative	0.51	1.00	0.67	5007
accuracy			0.51	10000
macro avg	0.75	0.51	0.36	10000
weighted avg	0.75	0.51	0.36	10000

Plot the confusion matrix

In [31]:

```

1 #confusion matrix for bag of words
2 cm_bow=confusion_matrix(test_sentiments,svm_bow_predict,labels=[1,0])
3 print(cm_bow)
4 #confusion matrix for tfidf features
5 cm_tfidf=confusion_matrix(test_sentiments,svm_tfidf_predict,labels=[1,0])
6 print(cm_tfidf)

```

```

[[4948  59]
 [4112 881]]
[[5007   0]
 [4888 105]]

```

Multinomial Naive Bayes for bag of words and tfidf features

In [32]:

```

1 #training the model
2 mnb=MultinomialNB()
3 #fitting the svm for bag of words
4 mnb_bow=mnb.fit(cv_train_reviews,train_sentiments)
5 print(mnb_bow)
6 #fitting the svm for tfidf features
7 mnb_tfidf=mnb.fit(tv_train_reviews,train_sentiments)
8 print(mnb_tfidf)

```

```

MultinomialNB()
MultinomialNB()

```

Model performance on test data

```
In [33]: 1 #Predicting the model for bag of words
2 mnb_bow_predict=mnb.predict(cv_test_reviews)
3 print(mnb_bow_predict)
4 #Predicting the model for tfidf features
5 mnb_tfidf_predict=mnb.predict(tv_test_reviews)
6 print(mnb_tfidf_predict)
```

[0 0 0 ... 0 1 1]
[0 0 0 ... 0 1 1]

Accuracy of the model

```
In [34]: 1 #Accuracy score for bag of words
2 mnb_bow_score=accuracy_score(test_sentiments,mnb_bow_predict)
3 print("mnb_bow_score :",mnb_bow_score)
4 #Accuracy score for tfidf features
5 mnb_tfidf_score=accuracy_score(test_sentiments,mnb_tfidf_predict)
6 print("mnb_tfidf_score :",mnb_tfidf_score)
```

mnb_bow_score : 0.751
mnb_tfidf_score : 0.7509

Print the classification report

```
In [35]: 1 #Classification report for bag of words
2 mnb_bow_report=classification_report(test_sentiments,mnb_bow_predict,target_names=['Positive','Negative'])
3 print(mnb_bow_report)
4 #Classification report for tfidf features
5 mnb_tfidf_report=classification_report(test_sentiments,mnb_tfidf_predict,target_names=['Positive','Negative'])
6 print(mnb_tfidf_report)
```

	precision	recall	f1-score	support
Positive	0.75	0.76	0.75	4993
Negative	0.75	0.75	0.75	5007
accuracy			0.75	10000
macro avg	0.75	0.75	0.75	10000
weighted avg	0.75	0.75	0.75	10000
	precision	recall	f1-score	support
Positive	0.75	0.76	0.75	4993
Negative	0.75	0.74	0.75	5007
accuracy			0.75	10000
macro avg	0.75	0.75	0.75	10000
weighted avg	0.75	0.75	0.75	10000

Plot the confusion matrix

```
In [36]: 1 #confusion matrix for bag of words
          2 cm_bow=confusion_matrix(test_sentiments,mnb_bow_predict,labels=[1,0])
          3 print(cm_bow)
          4 #confusion matrix for tfidf features
          5 cm_tfidf=confusion_matrix(test_sentiments,mnb_tfidf_predict,labels=[1,0])
          6 print(cm_tfidf)
```

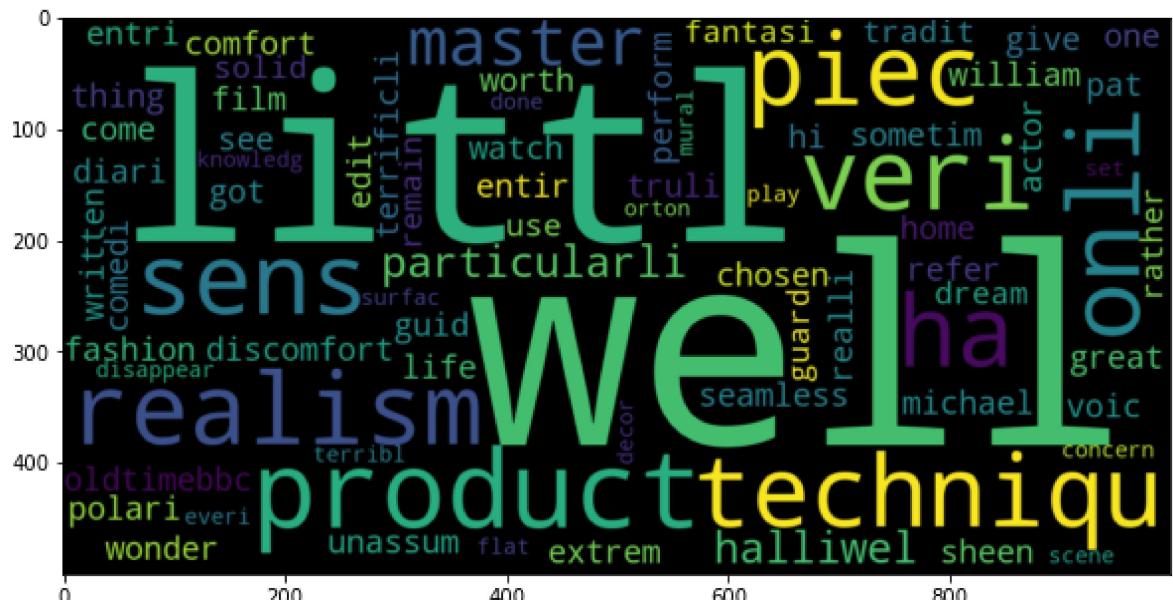
```
[ [ 3736 1271 ]
  [ 1219 3774 ] ]
[[ 3729 1278 ]
 [ 1213 3780 ]]
```

Let us see positive and negative words by using WordCloud.

Word cloud for positive review words

```
In [37]: 1 #word cloud for positive review words
          2 plt.figure(figsize=(10,10))
          3 positive_text=norm_train_reviews[1]
          4 WC=WordCloud(width=1000,height=500,max_words=500,min_font_size=5)
          5 positive_words=WC.generate(positive_text)
          6 plt.imshow(positive_words,interpolation='bilinear')
          7 plt.show
```

Out[37]: <function matplotlib.pyplot.show(close=None, block=None)>



Word cloud for negative review words

```
In [38]: 1 #Word cloud for negative review words
2 plt.figure(figsize=(10,10))
3 negative_text=norm_train_reviews[8]
4 WC=WordCloud(width=1000,height=500,max_words=500,min_font_size=5)
5 negative_words=WC.generate(negative_text)
6 plt.imshow(negative_words,interpolation='bilinear')
7 plt.show
```

Out[38]: <function matplotlib.pyplot.show(close=None, block=None)>

