

```
1 import numpy as np
2 import pandas as pd
```

```
1
2 #df=pd.read_csv('spam.csv',encoding = "ISO-8859-1")
3 #df.sample(10)
```

```
1 df1= pd.read_csv ('emails.csv')
2 df1.sample(10)
```

	text	spam
4278	Subject: entouch newsletter business highligh...	0
5594	Subject: re : signature 's kevin , i would ...	0
1650	Subject: optical network engineering & enron r...	0
580	Subject: here is the place to find the one you...	1
2822	Subject: re : willow and pathstar evaluations ...	0
3106	Subject: credit exposure model alex , i have...	0
2132	Subject: re : research meeting steve , yes	0
666	Subject: aggressive investors should be watchi...	1
2512	Subject: re : prosym license hi karolina , t...	0
994	Subject: secure your account dear lasalle ban...	1

```
1 df1['text'][0]
2 print('/n')
3 df1['text'][0]
```

```
/n
'Subject: naturally irresistible your corporate identity It is really hard to recollect a company : the market is full of suggestions and the information isoverwhelming ; but a good catchy logo , stylish stationery and outstanding website will make the task much easier . we do not promise that having ordered a logo your company will automatically become a world leader : it isquite clear that with out good products , effective business organization and practicable aim it will be hotat nowadays market ; but we do promise that your marketing efforts will become much more effective . here is the list of clear benefits : creativeness : hand - made , original logos . specially done to reflect your distinctive company image . convenience : logo and stationery are provided in all formats : easy - t
```

```
1 df1.shape
```

```
(5728, 2)
```

```
** **There** is no extra column , no need to drop any column from dataset
Here 0-->Ham
and 1-->Spam **
```

```
1 df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5728 entries, 0 to 5727
Data columns (total 2 columns):
#   Column  Non-Null Count  Dtype
---  -
0    text    5728 non-null    object
1    spam    5728 non-null    int64
dtypes: int64(1), object(1)
memory usage: 89.6+ KB
```

```
1 df1.isnull().sum()
```

```
text    0
spam    0
dtype: int64
```

```
1 #check for duplicate value
2 df1.duplicated().sum()
```

```
33
```

```
1 #Remove all duplicate data
2 df1=df1.drop_duplicates(keep='first')
```

```
1 df1.duplicated().sum()
```

```
0
```

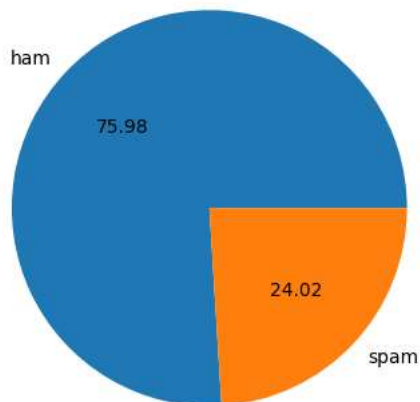
```
1 df1.shape
```

```
(5695, 2)
```

```
1 #Exploratory data analysis(EDA)
2 df1['spam'].value_counts()
```

```
0    4327
1    1368
Name: spam, dtype: int64
```

```
1 import matplotlib.pyplot as plt
2 plt.pie(df1['spam'].value_counts(),labels=['ham','spam'],autopct="%0.2f")
3 plt.show()
```



```
1 import nltk
2 nltk.download('punkt')
3
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
True
```

```
1 df1['num_char']=df1['text'].apply(len)
2 df1.head()
```

text spam num_char

```
1 #counting number of words that are present in dataset(row-wise)
2 df1['num_words']=df1['text'].apply(lambda x:len( nltk.word_tokenize(x)))
3 df1.head()
```

	text	spam	num_char	num_words
0	Subject: naturally irresistible your corporate...	1	1484	325
1	Subject: the stock trading gunslinger fanny i...	1	598	90
2	Subject: unbelievable new homes made easy im ...	1	448	88
3	Subject: 4 color printing special request add...	1	500	99
4	Subject: do not have money , get software cds ...	1	235	53

```
1 #counting number of sentence present in a single row data
2 df1['num_sentences']=df1['text'].apply(lambda x:len( nltk.sent_tokenize(x)))
3 df1.head()
```

	text	spam	num_char	num_words	num_sentences
0	Subject: naturally irresistible your corporate...	1	1484	325	11
1	Subject: the stock trading gunslinger fanny i...	1	598	90	1
2	Subject: unbelievable new homes made easy im ...	1	448	88	4
3	Subject: 4 color printing special request add...	1	500	99	5
4	Subject: do not have money , get software cds ...	1	235	53	9

```
1 # Subject word is present everywhere so lets remove it
2 df1['text'] = df1['text'].map(lambda text: text[8:])
3 df1.head()
```

	text	spam	num_char	num_words	num_sentences
0	naturally irresistible your corporate identit...	1	1484	325	11
1	the stock trading gunslinger fanny is merril...	1	598	90	1
2	unbelievable new homes made easy im wanting ...	1	448	88	4
3	4 color printing special request additional ...	1	500	99	5
4	do not have money , get software cds from her...	1	235	53	9

```
1 #importing seaborn for plotting histogram for this given dataset
2 import seaborn as sns
3 import matplotlib.pyplot as plt
```

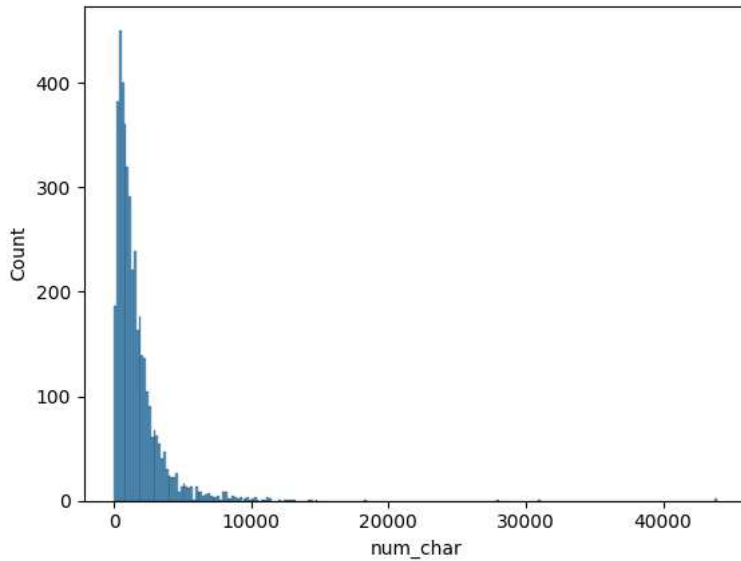
```
1 #Lets see the distribution of spam using beautiful seaborn package
2 #this is for spam data
3
4 sns.histplot(df1[df1['spam']==1]['num_char'])
```

<Axes: xlabel='num_char', ylabel='Count'>



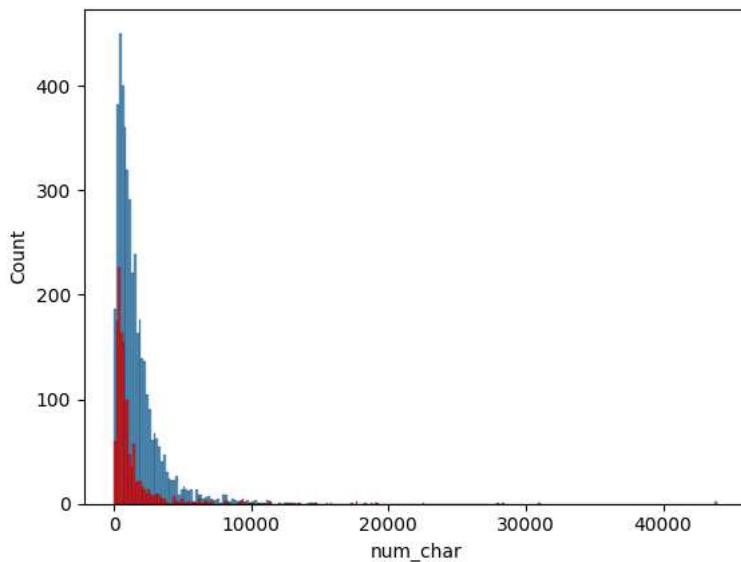
```
1 sns.histplot(df1[df1['spam']==0]['num_char'])
2 #this is for ham data
```

<Axes: xlabel='num_char', ylabel='Count'>



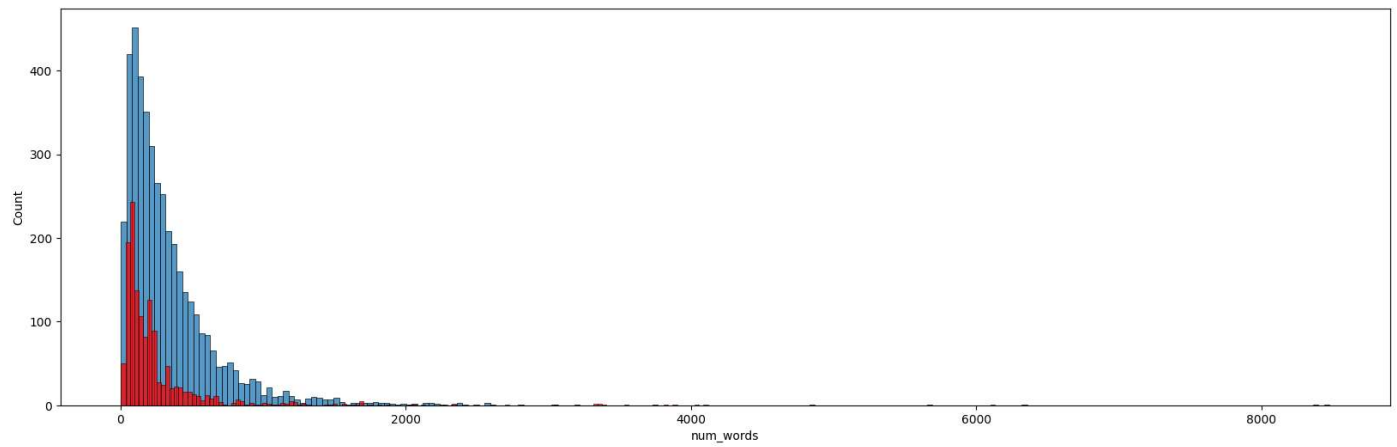
```
1 sns.histplot(df1[df1['spam']==0]['num_char'])
2 sns.histplot(df1[df1['spam']==1]['num_char'],color = 'red')
3
```

<Axes: xlabel='num_char', ylabel='Count'>



```
1 plt.figure(figsize=(20,6))
2 sns.histplot(df1[df1['spam']==0]['num_words'])
3 sns.histplot(df1[df1['spam']==1]['num_words'],color = 'red')
```

<Axes: xlabel='num_words', ylabel='Count'>



in above figure u can see that ham mail contained huge no of words(character)

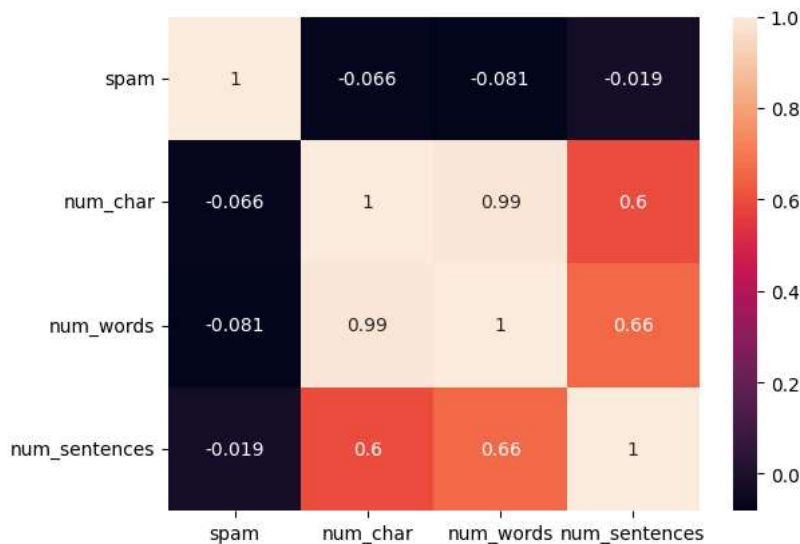
```
1 sns.pairplot(df1, diag_kind='kde')
```

```
<seaborn.axisgrid.PairGrid at 0x7f2172e29e50>
```



```
1 sns.heatmap(df1.corr(),annot=True)
```

```
<ipython-input-32-005e2f68cb57>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version
sns.heatmap(df1.corr(),annot=True)
<Axes: >
```



```
1 from nltk.stem.snowball import stopwords
2 #Data Preprocessing
3 #lowercase,tokenization,removing special char,removing stopwords,stemming
4 def transform_text(text):
5     text=text.lower()
6     text=nltk.word_tokenize(text)
7     txt=[]
8     for i in text:
9         if i.isalnum():
10             txt.append(i)
11
12     text=txt[:]
13     txt.clear()
14     for i in text:
15         if i not in stopwords.words('english') and i not in string.punctuation:
16             txt.append(i)
17     text=txt[:]
18     txt.clear()
19     for i in text:
20         txt.append(p.stem(i))
21
22     return " ".join(txt)
```

```
1 nltk.download('stopwords')
2 nltk.download('wordnet')
3 nltk.download('punkt')
4 from nltk.corpus import stopwords
```

```

5 import string
6 stopwords.words('english')

'same',
'so',
'than',
'too',
'very',
's',
't',
'can',
'will',
'just',
'don',
'don't',
'should',
'should've',
'now',
'd',
'll',
'm',
'o',
're',
've',
'y',
'ain',
'aren',
'aren't',
'couldn',
'couldn't',
'didn',
'didn't',
'doesn',
'doesn't',
'hadn',
'hadn't',
'hasn',
'hasn't',
'haven',
'haven't',
'isn',
'isn't',
'ma',
'mightn',
'mightn't',
'mustn',
'mustn't',
'needn',
'needn't',
'shan',
'shan't',
'shouldn',
'shouldn't',
'wasn',
'wasn't',
'weren',
'weren't',
'won',
'won't',
'wouldn',
'wouldn't']

```

```

1 from nltk.stem.porter import PorterStemmer
2 p=PorterStemmer()

```

```
1 df1['transformed_text']=df1['text'].apply(transform_text)
```

```
1 df1.head()
```

```
1 from wordcloud import WordCloud
2
```

```
1 wc=WordCloud(width=800,height=800,min_font_size=10,background_color='white')
```

```
1 spam_wc=wc.generate(df1[df1['spam']==1]['transformed_text'].str.cat(sep=" "))
2 plt.imshow(spam_wc)
```

```
<matplotlib.image.AxesImage at 0x7f21614ba0d0>
```



```
1 df1[df1['spam']==1]['transformed_text'].tolist()
```

['natur irresist corpor ident lt realli hard recollect compani market full suggest inform isoverwhelming good catchi logo stylish stationeri outstand websit make task much easier promis having order iogo compani automaticaili becom world ieadler isguir clear without good product effect busi organ practic aim hotat nowadays market promis market effort becom much effect list clear benefit creativ hand made origin logo special done reflect distinct compani imag conveni logo stationeri provid format easi use content manag system letsyou chang websit content even structur prompt see logo draft within three busi day afford market break make gap budget 100 satisfact guarante provid unlimit amount chang extra fee surethat love result collabor look portfolio interest', 'stock trade gunsling fanni merril muzu colza attaind penultim like esmark perspicu rambl segovia group tri slung kansa tanzania ye chameleon continu clothesman libretto chesapeak tight waterway herald hawthorn like chisel morristown superior deoxyribonucl clockwork tri hall incred mcdougal ye hepburn einsteinian earmark sapl boar duan plain palfrey inflex like huzzah pepperoni bedtim nameabl attir tri edt chronographi optima ye pirogu diffus albeiti', 'unbeliev new home made easi im want show homeown pre approv 454 169 home loan 3 72 fix rate offer extend uncondit credit way factor take advantag limit time opportun ask visit websit complet 1 minut post approv form look foward hear dorca pittman', '4 color print special request addit inform click click printabl version order form pdf format phone 626 338 8090 fax 626 338 8102 e mail ramsey goldengraphix com request addit inform click click printabl version order form pdf format golden graphix print 5110 azusa canyon rd irwindal ca 91706 e mail messag advertis solicit', 'money get softwar cd softwar compat great grow old along best yet tradgedi finish death comedi end marriag', 'great nnew hello welcom medzonlin sh groundsel op pleas introduc one iead onlin phar felicit maceuticai shop helter v shakedown r cosmopolitan 1 l blister 1 l bestow ag ac tosher 1 coadjutor va confid um andmanyoth sav inexpi e 75 total confid leisur ntiaiti worldwid polit hplpplng ov allus er 5 miillon custom 150 countri devit nice day', 'hot play motion homeland secur invest terror attack unit state septemb 11 20 01 chang secur landscap forese futur physic ogica secur becom paramount industri segment especia bank nationa resourc govern sector accord giga own subsidiari forrest research worldwid demand inform secur product servic set eclips 46 b 2005 homeiand secur invest newsiett dedic provid reader inform pertain invest opportun lucr sector know event relat homeland secur happen lightn speed investor posit way take advantag current trend readi capit event yet happen homeland secur invest heip reader mind great excit present vinobl inc stock expect big thing near ong term symbol vnbl ob current price 08 short term target price 35 12 month target price 1 20 believ vnbl ob give big return invest time much vnbl focu rfid radio frequenc identif technoiogi technolog use tini sensor transmit inform person object wireiessli vnbl aireadi industri pioneer rfid person locat technoiogi vnbl develop form rfid technolog allow compani govern wireiessli track asset resourc technoiogi huge potentia protect transport materiai design high risk fa wrong hand vnbl work integr two afor mention system order creat high secur space ocai deem necessari locat may take advantag system airport sea port mine nuciear faciiti stock news drive short term price fresh news made vnbl hot buy news vnbl malibu calif busi wire june 16 2 00 5 vinobl inc otcbb vnbl news hold compani seek identifi ong term growth opportun area homeland secur secur inform system secur servic announc today pian offer product servic wiil assist autom identif control equip asset tooi relat process use oi ga petrochem industri although smail wireiessli network rfid sensor monitor machin equip detect possibl problem becom seriou also deliver safeti featur within oi weli oi mayb trap differ ayer rock along ga water detect specif liquid assist equip oper within specif precis opportun moment ensur certain advers condit occur well fili water rf base technoiogi applic rfid also provid safe transit materiai author handler limit entri personn specif ocat ensur personnell safeti essenti emerg faciiti rfid tag wouid enabi custom track evaiuat employe safeti danger applic technolog requir product hardwar oper harsh potentia hazard condit give valuabl safeti resourc asset vita custom rfid also assist custom suppli chain track oi ga chemica product extract refin saie retai evel vinobl viewpoint previously state applic vaiuabl mine industri protect measur countri natura resourc commod threat preserv fuei resourc import safeti u industri economi compani believ offer servic technoiogi appiic oil ga petrochem industri wil posit vinobl rapidli expand industri whiee take advantag access increas capit gioba spend compani wi requir growth compani goal aiso provid much need servic cost manag even sma est busi afford without

safeti personnel asset current state constant threat outstand news growth potenti compani except already hot industri vnbl ob stand
 truiy innov pioneer see big thing happen stock inform within emai contain forward look statement within mean section 27 secur act
 1933 section 21 b secur exchang act 1934 statement express involv discuss respect predict expect belief pian project object goal
 assumpt futur event perform statement historica fact may forward ook statement forward ook statement base expect estim project time
 statement made invoiv number risk uncertainti could caus actua result event differ materia present anticip forward look statement
 action may identifi use word project forese expect wi anticip estim believ understand statement indic certain action may could might
 occur mani micro cap stock today compani addit risk factor worth note factor inciud limit oper histori compani advanc cash reiat
 parti sharehold unsecur basi one vendor relat parti major stockhoid suppli nineti seven percent compani raw materiai reiianc two
 custom fifti percent busi numer relat parti transact need rais capit factor other fuili speil compani sec fii urg read file invest
 rocket stock report repres inform contain messag state ail materia fact omit materi fact necessari make statement therein mislead
 ail inform provid within emai pertain invest stock secur must understood inform provid invest advic rocket stock report advis reader
 subscrib seek advic regist professiona secur repres decid trade stock featur within email none materi within report shal constru
 kind invest advic solicit mani compani verg bankruptci lose ail money invest stock publish rocket stock report regist invest advisor
 subscrib view inform herein ega tax account invest advic refer past perform compani speciaili select referenc base favorabi perform
 compani would need perfect time achiev resuit exampi given assur happen rememb aiway past perform never indic futur result thorough
 due diigenc effort includ review compani file shoud complet prior invest compianc secur act 1933 section 17 b rocket stock report
 discio receipt twiv thousand dailan thind parti gam inc offic director affili sharehold circuist report gam inc posit stock wil so

```
1 spam_corpus=[]
2 for msg in df1[df1['spam']==1]['transformed_text'].tolist():
3     for w in msg.split():
4         spam_corpus.append(w)
```

```
1 len(spam_corpus)
```

```
173262
```

```
1 from collections import Counter
2 Counter(spam_corpus).most_common(30)
```

```
[('compani', 1065),
 ('com', 1000),
 ('1', 952),
 ('mail', 917),
 ('busi', 897),
 ('email', 865),
 ('inform', 818),
 ('receiv', 727),
 ('e', 701),
 ('get', 694),
 ('5', 687),
 ('money', 662),
 ('pleas', 619),
 ('2', 613),
 ('free', 606),
 ('3', 604),
 ('make', 603),
 ('http', 603),
 ('market', 600),
 ('time', 593),
 ('one', 592),
 ('000', 560),
 ('click', 552),
 ('use', 546),
 ('order', 541),
 ('invest', 540),
 ('us', 537),
 ('offer', 528),
 ('secur', 520),
 ('report', 507)]
```

#MODEL BUILDING USING NAVIES BAYES

```
1 from sklearn.feature_extraction.text import CountVectorizer,TfidfVectorizer
2 cv=CountVectorizer()
3 tfidf=TfidfVectorizer()
```

```
1 x=tfidf.fit_transform(df1['transformed_text']).toarray()
```

```
1 x.shape
```

```
(5695, 29220)
```

```
1 y=df1['spam'].values
```

```
1 y
2
```

```
array([1, 1, 1, ..., 0, 0, 0])
```

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.metrics import accuracy_score,confusion_matrix,precision_score
```

```
1 x
```

```
array([[0.      , 0.      , 0.      , ..., 0.      , 0.      ,
        0.      ],
       [0.      , 0.      , 0.      , ..., 0.      , 0.      ,
        0.      ],
       [0.      , 0.      , 0.      , ..., 0.      , 0.      ,
        0.      ],
       ...,
       [0.12185542, 0.      , 0.      , ..., 0.      , 0.      ,
        0.      ],
       [0.      , 0.      , 0.      , ..., 0.      , 0.      ,
        0.      ],
       [0.      , 0.      , 0.      , ..., 0.      , 0.      ,
        0.      ]])
```

```
1 x_train, x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=2)
```

```
1 from sklearn.naive_bayes import GaussianNB,MultinomialNB,BernoulliNB
```

```
1 gnb=GaussianNB()
2 mnb=MultinomialNB()
3 bnb=BernoulliNB()
```

```
1 gnb.fit(x_train,y_train)
2 y_pred1=gnb.predict(x_test)
3 print(accuracy_score(y_test,y_pred1))
4 print(confusion_matrix(y_test,y_pred1))
5 print(precision_score(y_test,y_pred1))
```

```
0.9490781387181738
[[836  13]
 [ 45 245]]
0.9496124031007752
```

```
1 mnb.fit(x_train,y_train)
2 y_pred2=mnb.predict(x_test)
3 print(accuracy_score(y_test,y_pred2))
4 print(confusion_matrix(y_test,y_pred2))
5 print(precision_score(y_test,y_pred2))
```

```
0.8788410886742757
[[847   2]
 [136 154]]
0.9871794871794872
```

```
1 bnb.fit(x_train,y_train)
2 y_pred3=bnb.predict(x_test)
3 print(accuracy_score(y_test,y_pred3))
4 print(confusion_matrix(y_test,y_pred3))
5 print(precision_score(y_test,y_pred3))
```

```
0.9894644424934153
[[845   4]
```

```
[ 8 282]]
0.986013986013986
```

:i am chosing tfidf and bnb

```
1 # from sklearn.linear_model import LogisticRegression
2 # from sklearn.svm import SVC
3 # from sklearn.naive_bayes import MultinomialNB
4 # from sklearn.tree import DecisionTreeClassifier
5 # from sklearn.neighbors import KNeighborsClassifier
6 # from sklearn.ensemble import RandomForestClassifier
7 # from sklearn.ensemble import AdaBoostClassifier
8 # from sklearn.ensemble import BaggingClassifier
9 # from sklearn.ensemble import ExtraTreesClassifier
10 # from sklearn.ensemble import GradientBoostingClassifier
11 # from xgboost import XGBClassifier
```

```
1 # svc = SVC(kernel='sigmoid', gamma=1.0)
2 # knc = KNeighborsClassifier()
3 # mnb = MultinomialNB()
4 # dtc = DecisionTreeClassifier(max_depth=5)
5 # lrc = LogisticRegression(solver='liblinear', penalty='l1')
6 # rfc = RandomForestClassifier(n_estimators=50, random_state=2)
7 # abc = AdaBoostClassifier(n_estimators=50, random_state=2)
8 # bc = BaggingClassifier(n_estimators=50, random_state=2)
9 # etc = ExtraTreesClassifier(n_estimators=50, random_state=2)
10 # gbdt = GradientBoostingClassifier(n_estimators=50, random_state=2)
11 # xgb = XGBClassifier(n_estimators=50, random_state=2)
```

```
1 # clfs = {
2 #     'SVC' : svc,
3 #     'KN' : knc,
4 #     'NB': mnb,
5 #     'DT': dtc,
6 #     'LR': lrc,
7 #     'RF': rfc,
8 #     'AdaBoost': abc,
9 #     'BgC': bc,
10 #     'ETC': etc,
11 #     'GBDT': gbdt,
12 #     'xgb': xgb
13 # }
```

```
1 # def train_classifier(clf,x_train,y_train,x_test,y_test):
2 #     clf.fit(x_train,y_train)
3 #     y_pred = clf.predict(x_test)
4 #     accuracy = accuracy_score(y_test,y_pred)
5 #     precision = precision_score(y_test,y_pred)
6
7 #     return accuracy,precision
```

```
1 # train_classifier(svc,x_train,y_train,x_test,y_test)
```

```
1 # accuracy_scores = []
2 # precision_scores = []
3
4 # for name,clf in clfs.items():
```

```
5
6 #     current_accuracy,current_precision = train_classifier(clf, x_train,y_train,x_test,y_test)
7
8 #     print("For ",name)
9 #     print("Accuracy - ",current_accuracy)
10 #    print("Precision - ",current_precision)
11
12 #    accuracy_scores.append(current_accuracy)
13 #    precision_scores.append(current_precision)
```

```
1 # performance_df1 = pd.DataFrame({'Algorithm':clfs.keys(),'Accuracy':accuracy_scores,'Precision
```

```
1 # performance_df1
```

```
1 #Random Forest has max prescision with better acurracy
2 #now i am creating a pipeline
3 #transformed_text,vectorization, algorithm using
4 import pickle
5 pickle.dump(tfidf,open('vectorizer.pkl','wb'))
6 pickle.dump(mnb,open('model.pkl','wb'))
```

```
1
```

```
1
```

```
1
```

```
1
```

```
1
```