

```
In [2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

df=pd.read_csv("C:\Users\umari\OneDrive\Documents\Desktop\email_spam_classification\spam2.csv",encoding="utf-8")
```

Data Cleaning

```
In [4]: df

Out[4]:
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy.. Available only in	NaN	NaN	NaN
1	ham	Ok lar... Joking wif u oni...	NaN	NaN	NaN
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	NaN	NaN	NaN
3	ham	U dun say so early hor... U c already then say...	NaN	NaN	NaN
4	ham	Nah I don't think he goes to usf, he lives aro...	NaN	NaN	NaN
...
5567	spam	This is the 2nd time we have tried 2 contact u...	NaN	NaN	NaN
5568	ham	Will ♠. b going to explanade fr home?	NaN	NaN	NaN
5569	ham	Pty. * was in mood for that. So..any other s...	NaN	NaN	NaN
5570	ham	The guy did some bitching but I acted like i'd...	NaN	NaN	NaN
5571	ham	Rofl. Its true to its name	NaN	NaN	NaN

5572 rows × 5 columns

```
In [5]: df.rename(columns={'v1':'spam','v2':'email'},inplace=True)

In [6]: df

Out[6]:
```

	spam	email	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy.. Available only in	NaN	NaN	NaN
1	ham	Ok lar... Joking wif u oni...	NaN	NaN	NaN
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	NaN	NaN	NaN
3	ham	U dun say so early hor... U c already then say...	NaN	NaN	NaN
4	ham	Nah I don't think he goes to usf, he lives aro...	NaN	NaN	NaN
...
5567	spam	This is the 2nd time we have tried 2 contact u...	NaN	NaN	NaN
5568	ham	Will ♠. b going to explanade fr home?	NaN	NaN	NaN
5569	ham	Pty. * was in mood for that. So..any other s...	NaN	NaN	NaN
5570	ham	The guy did some bitching but I acted like i'd...	NaN	NaN	NaN
5571	ham	Rofl. Its true to its name	NaN	NaN	NaN

5572 rows × 5 columns

```
In [7]: df=df.drop(['Unnamed: 2','Unnamed: 3','Unnamed: 4'],axis=1)

In [8]: df

Out[8]:
```

	spam	email
0	ham	Go until jurong point, crazy.. Available only in
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...
...
5567	spam	This is the 2nd time we have tried 2 contact u...
5568	ham	Will ♠. b going to explanade fr home?
5569	ham	Pty. * was in mood for that. So..any other s...
5570	ham	The guy did some bitching but I acted like i'd...
5571	ham	Rofl. Its true to its name

5572 rows × 2 columns

```
In [9]: from sklearn.preprocessing import LabelEncoder
encoder=LabelEncoder()

In [10]: df['spam']=encoder.fit_transform(df['spam'])

In [11]: df

Out[11]:
```

	spam	email
0	0	Go until jurong point, crazy.. Available only in
1	0	Ok lar... Joking wif u oni...
2	1	Free entry in 2 a wkly comp to win FA Cup fina...
3	0	U dun say so early hor... U c already then say...
4	0	Nah I don't think he goes to usf, he lives aro...
...
5567	1	This is the 2nd time we have tried 2 contact u...
5568	0	Will ♠. b going to explanade fr home?
5569	0	Pty. * was in mood for that. So..any other s...
5570	0	The guy did some bitching but I acted like i'd...
5571	0	Rofl. Its true to its name

5572 rows × 2 columns

```
In [12]: df.isnull().sum()
#There are no null values

Out[12]:
spam      0
email     0
dtype: int64

In [13]: df.duplicated().sum()

Out[13]:
403

In [14]: df.drop_duplicates(inplace=True)

In [15]: df.duplicated().sum()

Out[15]:
0

In [16]: df.shape

Out[16]:
(5169, 2)
```

Exploratory Data Analysis

```
In [17]: df.head()

Out[17]:
```

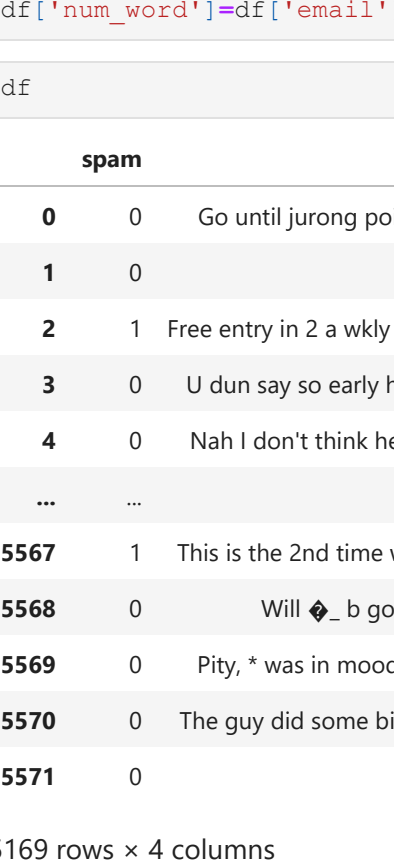
	spam	email
0	0	Go until jurong point, crazy.. Available only in
1	0	Ok lar... Joking wif u oni...
2	1	Free entry in 2 a wkly comp to win FA Cup fina...
3	0	U dun say so early hor... U c already then say...
4	0	Nah I don't think he goes to usf, he lives aro...

```
In [18]: df['spam'].value_counts()

Out[18]:
0    4516
1     653
Name: spam, dtype: int64

In [19]: # Now here we see that the proportion of ham is much more than that of the spam--- data is imbalanced
plt.pie(df['spam'].value_counts(),labels=['ham','spam'])

Out[19]:
```



```
In [20]: import nltk

In [21]: nltk.download('punkt')

[nltk_data] Error loading punkt: <urlopen error [WinError 10060] A
[nltk_data] connection attempt failed because the connected party
[nltk_data] did not properly respond after a period of time, or
[nltk_data] established connection failed because connected host
[nltk_data] has failed to respond>

Out[21]:

In [22]: #no. of characters present is determined here.
df['num_characters']=df['email'].apply(len)

In [23]: df['num_word']=df['email'].apply(lambda x:len(nltk.word_tokenize(x)))

In [24]: df

Out[24]:
```

	spam	email	num_characters	num_word
0	0	Go until jurong point, crazy.. Available only in	111	24
1	0	Ok lar... Joking wif u oni...	29	8
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37
3	0	U dun say so early hor... U c already then say...	49	13
4	0	Nah I don't think he goes to usf, he lives aro...	61	15
...
5567	1	This is the 2nd time we have tried 2 contact u...	160	35
5568	0	Will ♠. b going to explanade fr home?	37	9
5569	0	Pty. * was in mood for that. So..any other s...	57	15
5570	0	The guy did some bitching but I acted like i'd...	125	27
5571	0	Rofl. Its true to its name	26	7

5169 rows × 4 columns

```
In [25]: df[['num_characters','num_word','num_sentences']].describe()

Out[27]:
```

	num_characters	num_word	num_sentences
count	5169.000000	5169.000000	5169.000000
mean	78.923776	18.454246	1.948152
std	58.174846	13.325668	1.363792
min	2.000000	1.000000	1.000000
25%	36.000000	9.000000	1.000000
50%	60.000000	15.000000	1.000000
75%	117.000000	26.000000	2.000000
max	910.000000	220.000000	28.000000

```
In [28]: df(df['spam']==1)[['num_characters','num_word','num_sentences']].describe()

Out[28]:
```

	num_characters	num_word	num_sentences
count	653.000000	653.000000	653.000000
mean	137.479326	27.675345	2.975498
std	30.014336	7.011513	1.487993
min	13.000000	2.000000	1.000000
25%	131.000000	25.000000	2.000000
50%	148.000000	29.000000	3.000000
75%	157.000000	32.000000	4.000000
max	223.000000	46.000000	8.000000


```
In [29]: df(df['spam']==0)[['num_characters','num_word','num_sentences']].describe()

Out[29]:
```

	num_characters	num_word	num_sentences
count	4516.000000	4516.000000	4516.000000
mean	70.456802	17.120903	1.799601
std	56.356802	13.493725	1.278465
min	2.000000	1.000000	1.000000
25%	34.000000	8.000000	1.000000
50%	52.000000	13.000000	1.000000
75%	90.000000	22.000000	2.000000
max	910.000000	220.000000	28.000000


```
In [30]: # Now let us visualize the cases of number of words and characters for the two cases of spam and ham
fig, ax = plt.subplots(figsize=(20, 5))
sns.histplot(df[df['spam']==0]['num_characters'],color='red')
sns.histplot(df[df['spam']==1]['num_characters'],color='black')
<AxesSubplot: xlabel='num_characters', ylabel='Count'>

Out[30]:
```




```
In [31]: fig, ax = plt.subplots(figsize=(20, 5))
sns.histplot(df[df['spam']==0]['num_word'],color='red')
sns.histplot(df[df['spam']==1]['num_word'],color='black')
<AxesSubplot: xlabel='num_word', ylabel='Count'>

Out[31]:
```



```
In [32]: sns.pairplot(df,hue='spam') ## Tells us the relation between each variables.

Out[32]:
```



```
In [33]: sns.heatmap(df.corr(),annot=True)

Out[33]:
```

	spam	num_characters	num_word	num_sentences
spam	1	0.38	0.26	0.29
num_characters	0.38	1	0.97	0.64
num_word	0.26	0.97	1	0.68
num_sentences	0.29	0.64	0.68	1

#Here we can see that the correlation between the coefficient therefore multilinearity is present , Therefore

Text Preprocessing

```
In [35]: nltk.download('stopwords')
from nltk.corpus import stopwords
import string
from string import punctuation
from nltk.stem import PorterStemmer
ps=PorterStemmer()

[nltk_data] Error loading stopwords: <urlopen error [WinError 10060] A
[nltk_data] connection attempt failed because the connected party
[nltk_data] did not properly respond after a period of time, or
[nltk_data] established connection failed because connected host
[nltk_data] has failed to respond>

In [36]: def transform_text(text):
    text=text.lower()
    text=nltk.word_tokenize(text)
    y=[]
    for i in text:
        if i.isalnum():
            y.append(i)
        texty=[]
        y.clear()
        for i in text:
            if i not in stopwords.words('english') and i not in string.punctuation:
                y.append(i)
        texty=y[:]
        y.clear()
        for i in text:
            y.append(ps.stem(i))
    return " ".join(y)

In [37]: transform_text("Hi How are You Miss i love dancing")

Out[37]:
'hi love danc'

In [38]: df['transformed_text']=df['email'].apply(transform_text)

In [39]: df

Out[39]:
```

	spam	email	num_characters	num_word	num_sentences	transformed_text
0	0	Go until jurong point, crazy.. Available only in	111	24	2	go jurong point craci avail bugi n great world...
1	0	Ok lar... Joking wif u oni...	29	8	2	ok lar joke wif u oni
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37	2	free entrn 2 wkli comp win fa cup final txt 21...
3	0	U dun say so early hor... U c already then say...	49	13	1	u dun say earli hor u c already say
4	0	Nah I don't think he goes to usf, he lives aro...	61	15	1	nah think goe usf live around thogh
...
5567	1	This is the 2nd time we have tried 2 contact u...	160	35	4	2nd time tri 2 contact u pound prize 2 claim e...
5568	0	Will ♠. b going to explanade fr home?	37	9	1	b go explanad fr home
5569	0	Pty. * was in mood for that. So..any other s...	57	15	2	ptly mood suggest
5570	0	The guy did some bitching but I acted like i'd...	125	27	1	guy bitch act like interest buy someth els nex...
5571	0	Rofl. Its true to its name	26	7	2	rofl true name

5169 rows × 6 columns

```
In [40]: #Shows the top 30 most frequent values in case of ham messages
for i in df[df['spam']==0]['transformed_text'].to_list():
    for word in i.split():
        spam_corpus.append(word)
spam_corpus
```


Word	Frequency
from	850
collect	450
to	420
the	380
in	350
a	320
for	300
word	280
an	250
if	220
and	200
open_corpus	180


```

'phone',
'cell',
'mobilesdirect',
'08000938767',
'updat',
'for2stopxt',
'free',
'rington',
'text',
'first',
'87131',
'87131',
'text',
'get',
'87131',
'true',
'tone',
'0845',
'2814032',
'16',
'list',
'free',
'tone',
'txt',
'stop',
'100',
'date',
'servic',
'text',
'1',
'0904012103',
'pobox349438ch',
'free',
'entri',
'weekl',
'compet',
'text',
'ward',
'win',
'80006',
'18',
'c',
'send',
'10p',
'2',
'ur',
'love',
'2',
'name',
'join',
'heart',
'txt',
'love',
'name',
'name2',
'mob',
'eg',
'love',
'adam',
'eve',
'07123456789',
'87077',
'yaho',
'pobox36504w45wq',
'text',
'4',
'ad',
'15p',
'someos',
'contact',
'date',
'servic',
'enter',
'phone',
'fanci',
'find',
'text',
'landin',
'0911032124',
'pobox1246tf150p',
'urgent',
'mobil',
'mob',
'award',
'prize',
'guarante',
'call',
'09058094455',
'land',
'line',
'claim',
'radio',
'12hr',
'congrat',
'mob',
'3650',
'video',
...

In [43]:
from collections import Counter
import numpy as np
from sklearn.feature_extraction.text import CountVecorizer,TfidfVectorizer
tfidf=TfidfVectorizer(max_features=3000)
cvc=CountVecorizer()

In [44]:
#bag of words
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2)

In [45]:
X=ttfidf.fit_transform(df['transformed_text']).toarray()

In [46]:
X.shape

Out[46]:
(5169, 3000)

In [47]:
Y=ndf['spam'].values

In [48]:
from sklearn.model_selection import train_test_split
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2)

In [49]:
from sklearn.naive_bayes import GaussianNB
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import BernoulliNB

In [50]:
gmb=GaussianNB()
mnb=MultinomialNB()
bnb=BernoulliNB()

In [51]:
from sklearn.metrics import accuracy_score,confusion_matrix,precision_score

In [52]:
gmb.fit(X_train,Y_train)

Out[52]:
GaussianNB()

In [53]:
Y_pred=gmb.predict(X_test)

In [54]:
for i in range(X_test.shape[0]):
    print('Accuracy: ',accuracy_score(Y_test,Y_pred),confusion_matrix(Y_test,Y_pred),precision_score(Y_test,Y_pred))

0.8655705996131529
[[783 122]
 [ 17 112]]
0.47863247863247865

In [55]:
mnb.fit(X_train,Y_train)
Y_pred=mnb.predict(X_test)
for i in range(X_test.shape[0]):
    print('Accuracy: ',accuracy_score(Y_test,Y_pred),confusion_matrix(Y_test,Y_pred),precision_score(Y_test,Y_pred))

0.9758220502901354
[[905  0]
 [ 25 194]]
1.0

In [56]:
bnb.fit(X_train,Y_train)
Y_pred=bnb.predict(X_test)
for i in range(X_test.shape[0]):
    print('Accuracy: ',accuracy_score(Y_test,Y_pred),confusion_matrix(Y_test,Y_pred),precision_score(Y_test,Y_pred))

0.976905222437138
[[963  2]
 [ 19 110]]
0.9621428571428571

We can see that in this case the multinomial naive bayes gives the precision of 1 and accuracy is also higher than the other two naive bayes
```

Model building followed by Bag of words

```

In [57]:
Xlcv=fit_transform(df['transformed_text']).toarray()

In [58]:
Yl=ndf['spam']
Xl_train,Xl_test,Yl_train,Yl_test=train_test_split(Xl,Yl,test_size=0.2)

In [59]:
gmb.fit(Xl_train,Y_train)
Yl_pred=gmb.predict(Xl_test)
print(accuracy_score(Yl_test,Yl_pred))
print(confusion_matrix(Yl_test,Yl_pred))
print(precision_score(Yl_test,Yl_pred))

0.4990328820116054
[[487 411]
 [107 291]]
0.065909090909090909

In [60]:
mnb.fit(Xl_train,Y_train)
Yl_pred=mnb.predict(Xl_test)
print(accuracy_score(Yl_test,Yl_pred))
print(confusion_matrix(Yl_test,Yl_pred))
print(precision_score(Yl_test,Yl_pred))

0.8210831721470019
[[836 62]
 [123 131]]
0.17333333333333334

In [61]:
bnb.fit(Xl_train,Y_train)
Yl_pred=bnb.predict(Xl_test)
print(accuracy_score(Yl_test,Yl_pred))
print(confusion_matrix(Yl_test,Yl_pred))
print(precision_score(Yl_test,Yl_pred))

0.8646034816247582
[[879 191]
 [121 151]]
0.4411146705882353

Therefore we see that the vectorization by TFIDF gives a better result than the bag of words, Hence we use a model of TF_IDF followed by MultinomialNaiveBayes

In [62]:
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.naive_bayes import MultinomialNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import BaggingClassifier
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.ensemble import GradientBoostingClassifier

In [63]:
svc = SVC(kernel='sigmoid', gamma=1.0)
knc = KNeighborsClassifier()
mnb = MultinomialNB()
dtc = DecisionTreeClassifier(max_depth=5)
lrc = LogisticRegression(solver='liblinear', penalty='l1')
rfc = RandomForestClassifier(n_estimators=50, random_state=2)
abc = AdaBoostClassifier(n_estimators=50, random_state=2)
etc = BaggingClassifier(n_estimators=50, random_state=2)
etc2 = ExtraTreesClassifier(n_estimators=50, random_state=2)
gbd = GradientBoostingClassifier(n_estimators=50, random_state=2)

In [64]:
#lets make a dictionary of the algorithms
clfs={}
clfs['SVC']=svc,
clfs['KNC']=knc,
clfs['MNB']=mnb,
clfs['DTC']=dtc,
clfs['ABC']=abc,
clfs['RFC']=rfc,
clfs['ETC']=etc,
clfs['ETC2']=etc2,
clfs['GBD']=gbd

In [65]:
def training_data(clf,X_train,X_test,Y_train,Y_test):
    clf.fit(X_train,Y_train)
    Y_pred=clf.predict(X_test)
    acc_score=accuracy_score(Y_test,Y_pred)
    prec_score=precision_score(Y_test,Y_pred)
    return acc_score,prec_score

In [66]:
accuracy=[]
precision=[]
for name,clf in clfs.items():
    current_accuracy,current_precision = training_data(clf, X_train,X_test,Y_train,Y_test)

    print("For ",name)
    print("Accuracy = ",current_accuracy)
    print("Precision = ",current_precision)
    accuracy.append(current_accuracy)
    precision.append(current_precision)

For SVC
Accuracy = 0.97678916827653
Precision = 0.9646017699115044
For KNC
Accuracy = 0.9177949709864603
Precision = 1.0
For MNB
Accuracy = 0.9758220502901354
Precision = 1.0
For DTC
Accuracy = 0.9352030947775629
Precision = 0.82978723042425532
For LRC
Accuracy = 0.9574468085106383
Precision = 0.9292929292929293
For RFC
Accuracy = 0.9758220502901354
Precision = 0.9905660377358491
For ABC
Accuracy = 0.960348162475822
Precision = 0.9313725490196079
For BC
Accuracy = 0.9574468085106383
Precision = 0.84
For ETC
Accuracy = 0.9738978143133463
Precision = 0.9636363636363636
For GBDT
Accuracy = 0.9448742746615088
Precision = 0.9285714285714286

In [67]:
performance_df=pd.DataFrame({'Algorithm':clfs.keys(), 'Accuracy':accuracy, 'Precision':precision}).sort_values('P

In [68]:
performance_df

Out[68]:
  Algorithm  Accuracy  Precision
0  KNC  Accuracy  0.917795
1  MNB  Accuracy  0.975822
2  RFC  Accuracy  0.973888
0  SVC  Accuracy  0.964602
8  ETC  Accuracy  0.963636
6  ABC  Accuracy  0.931373
4  LRC  Accuracy  0.929293
9  GBDT  Accuracy  0.928571
7  BC  Accuracy  0.840000
3  DTC  Accuracy  0.829787

In [69]:
performance_df=pd.melt(performance_df,id_vars='Algorithm')

In [70]:
performance_df

Out[70]:
  Algorithm  Accuracy  Precision
0  KNC  Accuracy  0.917795
1  MNB  Accuracy  0.975822
2  RFC  Accuracy  0.973888
3  SVC  Accuracy  0.964602
4  ETC  Accuracy  0.963636
5  ABC  Accuracy  0.960348
6  LRC  Accuracy  0.957447
7  GBDT  Accuracy  0.944874
8  BC  Accuracy  0.957447
9  DTC  Accuracy  0.928571
10 KNC  Precision  1.000000
11 MNB  Precision  1.000000
12 RFC  Precision  0.990566
13 SVC  Precision  0.964602
14 ETC  Precision  0.963636
15 ABC  Precision  0.931373
16 LRC  Precision  0.929293
17 GBDT  Precision  0.928571
18 BC  Precision  0.840000
19 DTC  Precision  0.829787

In [71]:
sns.barplot(x='Algorithm',y='value',hue='variable',data=performance_df_1,kind='bar',height=0)
plt.ylim(0.5,1)

Out[71]:
(0.5, 1.0)

0.5 0.6 0.7 0.8 0.9 1.0
KNC MNB RFC SVC ETC ABC LRC GBDT BC DTC
variable
Accuracy
Precision

THEREFORE WE SEE THAT THE MULTINOMIAL NAIVE BAYES GIVES THE BEST RESULT BASED ON ACCURACY AND ALSO ACCURACY WHENEVER WE ARE USING IT AFTER VECTORIZING WITH TF-IDF
```

```

We see that scaling does not make much of a difference on the accuracy and precision , Therefore We go ahead with about 3000 max_features

In [74]:
import pickle
pickle.dump(tfidf,open('new_vectorizer.pkl','wb'))
pickle.dump(mnb,open('model_email.pkl','wb'))

In [ ]:
```