

```

import numpy as np
import pandas as pd

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

/kaggle/input/sms-spam-collection-dataset/spam.csv

import chardet
with open('/kaggle/input/sms-spam-collection-dataset/spam.csv', 'rb') as rawdata:
    result = chardet.detect(rawdata.read(100000))
result

{'encoding': 'Windows-1252', 'confidence': 0.7272080023536335, 'language': ''}

df = pd.read_csv('/kaggle/input/sms-spam-collection-dataset/spam.csv', encoding = 'Windows-1252')

df.sample(5)

```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
<b>1271</b>	ham	If you still havent collected the dough pls le...	NaN	NaN	NaN
<b>2416</b>	ham	Could you not read me, my Love ? I answered you	NaN	NaN	NaN
<b>2745</b>	ham	R Ì_ going 4 today's meeting?	NaN	NaN	NaN
<b>3332</b>	spam	You are being contacted by our dating service ...	NaN	NaN	NaN
<b>4926</b>	ham	Wanna do some art?! :D	NaN	NaN	NaN

```
df.shape
```

```
(5572, 5)
```

## ▼ 1. Data Cleaning

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0   v1               5572 non-null   object
1   v2               5572 non-null   object
2   Unnamed: 2       50 non-null     object
3   Unnamed: 3       12 non-null     object
4   Unnamed: 4       6 non-null      object
dtypes: object(5)
memory usage: 217.8+ KB
```

```
# drop last 3 cols
```

```
df.drop(columns=['Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], inplace=True)
```

```
df.sample(5)
```

**v1****v2**

```
# renaming the cols
df.rename(columns={'v1':'target','v2':'text'},inplace=True)
df.sample(5)
```

	target	text
<b>1818</b>	ham	Am i that much dirty fellow?
<b>530</b>	spam	PRIVATE! Your 2003 Account Statement for 07815...
<b>377</b>	ham	Well there's not a lot of things happening in ...
<b>4905</b>	ham	Will you come online today night
<b>3508</b>	ham	Hi Petey!noiãÕm ok just wanted 2 chat coz aven...

```
from sklearn.preprocessing import LabelEncoder
encoder = LabelEncoder()
```

```
df['target'] = encoder.fit_transform(df['target'])
```

```
df.head()
```

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

```
# missing values
df.isnull().sum()
```

```
target    0
text      0
dtype: int64
```

```
# check for duplicate values
df.duplicated().sum()
```

```
403
```

```
# remove duplicates
df = df.drop_duplicates(keep='first')
```

```
df.duplicated().sum()
```

```
0
```

```
df.shape
```

```
(5169, 2)
```

## ▼ 2.EDA

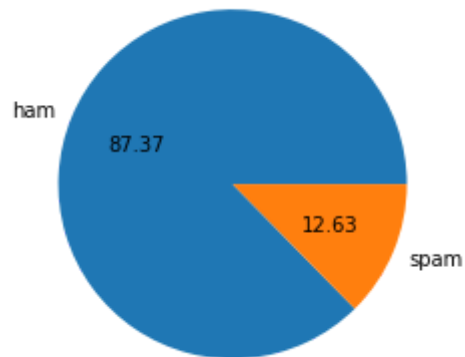
```
df.head()
```

	target	text
0	0	Go until jurong point, crazy.. Available only ...
1	0	Ok lar... Joking wif u oni...

```
df['target'].value_counts()
```

```
0    4516
1     653
Name: target, dtype: int64
```

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



```
# Data is imbalanced
```

```
import nltk
```

```
df['num_characters'] = df['text'].apply(len)
```

```
df.head()
```

	target	text	num_characters
0	0	Go until jurong point, crazy.. Available only ...	111
1	0	Ok lar... Joking wif u oni...	29
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155
3	0	U dun say so early hor... U c already then say...	49
4	0	Nah I don't think he goes to usf, he lives aro...	61

```
# num of words
```

```
df['num_words'] = df['text'].apply(lambda x:len(nltk.word_tokenize(x)))
```

```
df.head()
```

	target	text	num_characters	num_words
0	0	Go until jurong point, crazy.. Available only ...	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

```
df['num_sentences'] = df['text'].apply(lambda x:len(nltk.sent_tokenize(x)))
```

```
df.head()
```

	target	text	num_characters	num_words	num_sentences
0	0	Go until jurong point, crazy.. Available only ...	111	24	2
1	0	Ok lar... Joking wif u oni...	29	8	2
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37	2
3	0	U dun say so early hor... U c already then say...	49	13	1
4	0	Nah I don't think he goes to usf, he lives aro...	61	15	1

```
df[['num_characters', 'num_words', 'num_sentences']].describe()
```

	num_characters	num_words	num_sentences
<b>count</b>	5169.000000	5169.000000	5169.000000
<b>mean</b>	78.977945	18.453279	1.947185
<b>std</b>	58.236293	13.324793	1.362406
<b>min</b>	2.000000	1.000000	1.000000
<b>25%</b>	36.000000	9.000000	1.000000
<b>50%</b>	60.000000	15.000000	1.000000
<b>75%</b>	117.000000	26.000000	2.000000
<b>max</b>	910.000000	220.000000	28.000000

```
# ham
```

```
df[df['target'] == 0][['num_characters', 'num_words', 'num_sentences']].describe()
```

	num_characters	num_words	num_sentences
<b>count</b>	4516.000000	4516.000000	4516.000000
<b>mean</b>	70.459256	17.120903	1.799601
<b>std</b>	56.358207	13.493725	1.278465
<b>min</b>	2.000000	1.000000	1.000000
<b>25%</b>	34.000000	8.000000	1.000000

```
#spam
```

```
df[df['target'] == 1][['num_characters', 'num_words', 'num_sentences']].describe()
```

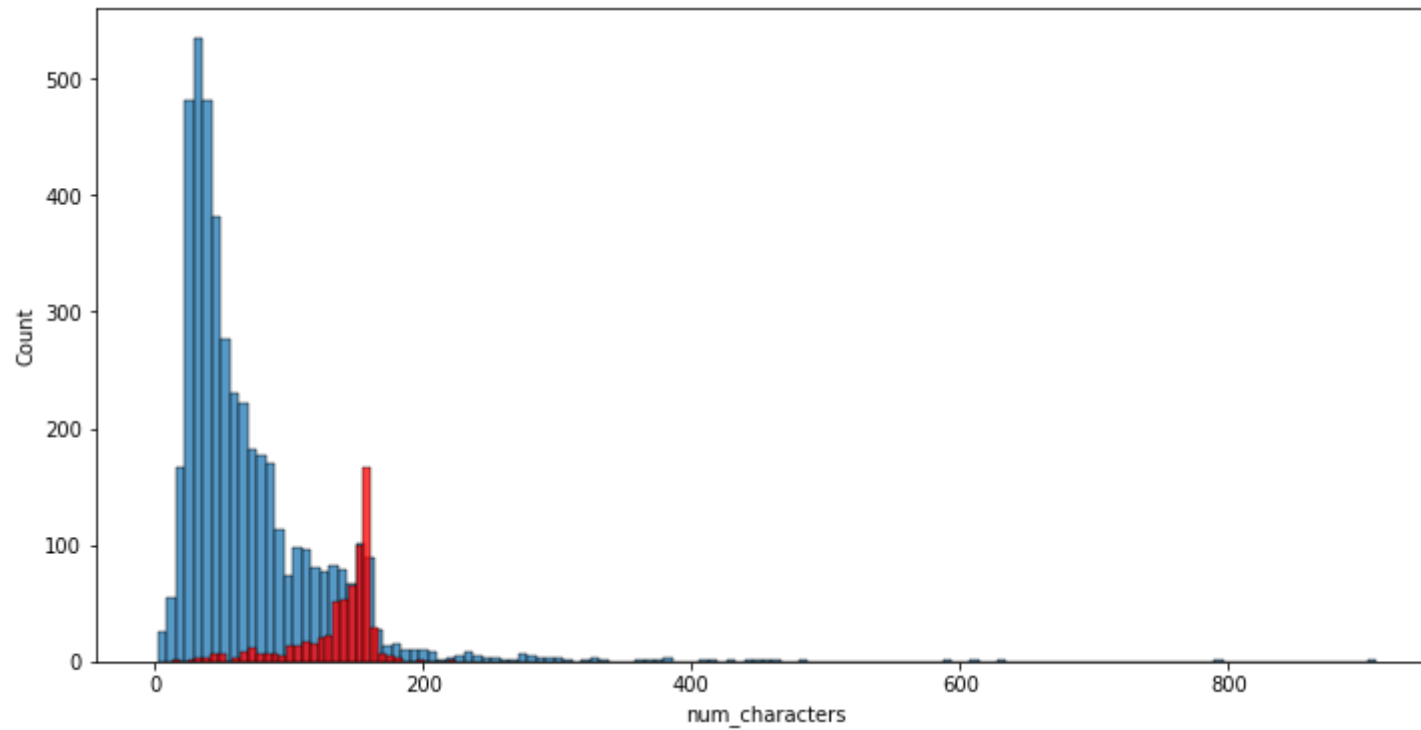
	num_characters	num_words	num_sentences
<b>count</b>	653.000000	653.000000	653.000000
<b>mean</b>	137.891271	27.667688	2.967841
<b>std</b>	30.137753	7.008418	1.483201
<b>min</b>	13.000000	2.000000	1.000000
<b>25%</b>	132.000000	25.000000	2.000000
<b>50%</b>	149.000000	29.000000	3.000000
<b>75%</b>	157.000000	32.000000	4.000000
<b>max</b>	224.000000	46.000000	8.000000

```
import seaborn as sns
```

```
plt.figure(figsize=(12,6))
sns.histplot(df[df['target'] == 0]['num_characters'])
sns.histplot(df[df['target'] == 1]['num_characters'],color='red')
```



<AxesSubplot:xlabel='num\_characters', ylabel='Count'>



```
plt.figure(figsize=(12,6))
sns.histplot(df[df['target'] == 0]['num_words'])
sns.histplot(df[df['target'] == 1]['num_words'],color='red')
```

<AxesSubplot:xlabel='num\_words', ylabel='Count'>



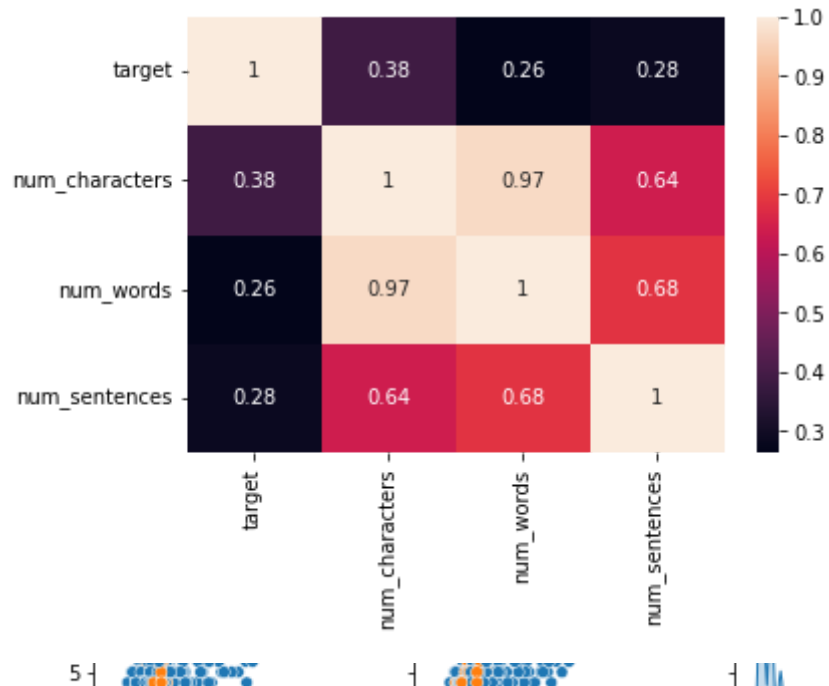
```
sns.pairplot(df,hue='target')
```

```
<seaborn.axisgrid.PairGrid at 0x7fbd221a9850>
```



```
sns.heatmap(df.corr(),annot=True)
```

```
<AxesSubplot:>
```



### ▼ 3. Data Preprocessing

- Lower case
- Tokenization
- Removing special characters
- Removing stop words and punctuation
- Stemming

```
nltk.download('punkt')
from nltk.corpus import stopwords
import string
```

```
[nltk_data] Downloading package punkt to /usr/share/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

```
from nltk.stem.porter import PorterStemmer
ps = PorterStemmer()
ps.stem('loving')
```

```
'love'
```

```
def transform_text(text):
```

```
    text = text.lower()
    text = nltk.word_tokenize(text)
```

```
    y = []
    for i in text:
        if i.isalnum():
            y.append(i)
```

```
    text = y[:]
    y.clear()
```

```
    for i in text:
        if i not in stopwords.words('english') and i not in string.punctuation:
            y.append(i)
```

```
    text = y[:]
    y.clear()
```

```
    for i in text:
        y.append(ps.stem(i))
```

```
return " ".join(y)
```

```
df['text'][10]
```

```
"I'm gonna be home soon and i don't want to talk about this stuff anymore tonight, k? I've cried enough today."
```

```
transform_text("I'm gonna be home soon and i don't want to talk about this stuff anymore tonight, k? I've cried enough today.")
```

```
'gon na home soon want talk stuff anymor tonight k cri enough today'
```

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

```
df.head()
```

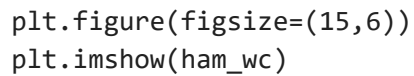
	target	text	num_characters	num_words	num_sentences	transformed_text
0	0	Go until jurong point, crazy.. Available only ...	111	24	2	go jurong point crazi avail bugi n great world...
1	0	Ok lar... Joking wif u oni...	29	8	2	ok lar joke wif u oni
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37	2	free entri 2 wkli comp win fa cup final tkt 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

```
from wordcloud import WordCloud
```

```
wc = WordCloud(width=500,height=500,min_font_size=10,background_color='white')
```

```
spam_wc = wc.generate(df[df['target'] == 1]['transformed_text'].str.cat(sep=" "))
```

```
<matplotlib.image.AxesImage at 0x7fbd13ed6210>
```



```
<matplotlib.image.AxesImage at 0x7fbd13ecad90>
```



```
df.head()
```

	target	text	num_characters	num_words	num_sentences	transformed_text
0	0	Go until jurong point, crazy.. Available only ...	111	24	2	go jurong point crazi 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 entri 2 wkli comp win fa cup final tkt 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

```
spam_corpus = []
for msg in df[df['target'] == 1]['transformed_text'].tolist():
    for word in msg.split():
        spam_corpus.append(word)
```

```
len(spam_corpus)
```

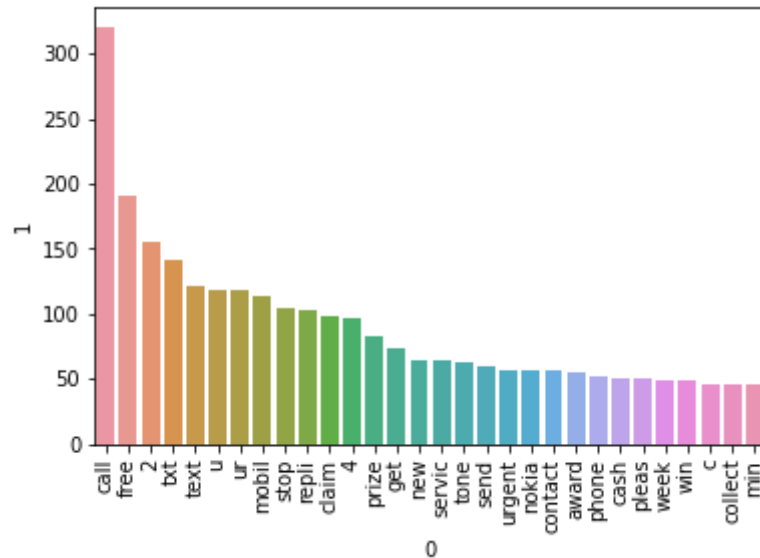
9939

```

from collections import Counter
sns.barplot(pd.DataFrame(Counter(spam_corpus).most_common(30))[0],pd.DataFrame(Counter(spam_corpus).most_common(30))[1])
plt.xticks(rotation='vertical')
plt.show()

```

/opt/conda/lib/python3.7/site-packages/seaborn/\_decorators.py:43: FutureWarning: Pass the following variables as keyword args: FutureWarning



```

ham_corpus = []
for msg in df[df['target'] == 0]['transformed_text'].tolist():
    for word in msg.split():
        ham_corpus.append(word)

```

```
len(ham_corpus)
```

35394

```

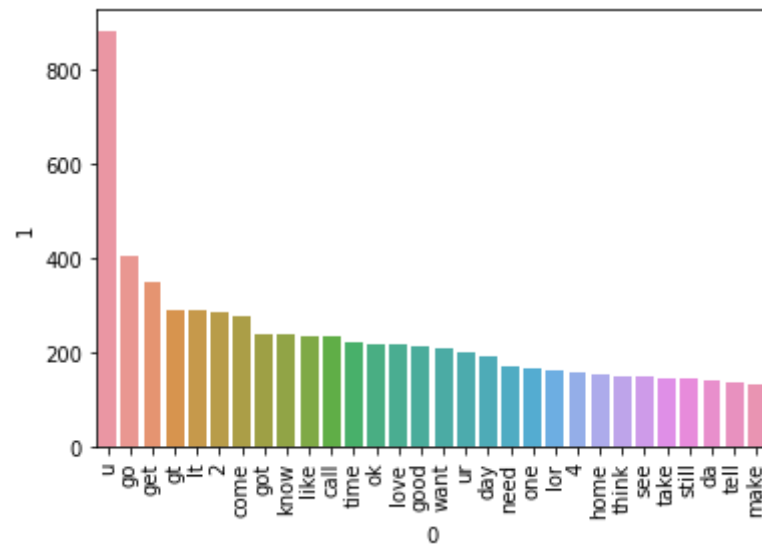
from collections import Counter
sns.barplot(pd.DataFrame(Counter(ham_corpus).most_common(30))[0],pd.DataFrame(Counter(ham_corpus).most_common(30))[1])

```



```
plt.xticks(rotation='vertical')  
plt.show()
```

/opt/conda/lib/python3.7/site-packages/seaborn/\_decorators.py:43: FutureWarning: Pass the following variables as keyword args: FutureWarning



```
# Text Vectorization  
# using Bag of Words  
df.head()
```

## 4. Model Building

```
from sklearn.feature_extraction.text import CountVectorizer,TfidfVectorizer
cv = CountVectorizer()
tfidf = TfidfVectorizer(max_features=3000)

# U dun say so early hor... U c already then
X = tfidf.fit_transform(df['transformed_text']).toarray()

X.shape

(5169, 3000)

y = df['target'].values

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,random_state=2)

from sklearn.naive_bayes import GaussianNB,MultinomialNB,BernoulliNB
from sklearn.metrics import accuracy_score,confusion_matrix,precision_score

gnb = GaussianNB()
mnb = MultinomialNB()
bnb = BernoulliNB()

gnb.fit(X_train,y_train)
y_pred1 = gnb.predict(X_test)
print(accuracy_score(y_test,y_pred1))
```

```
print(confusion_matrix(y_test,y_pred1))
print(precision_score(y_test,y_pred1))
```

```
0.8694390715667312
[[788 108]
 [ 27 111]]
0.5068493150684932
```

```
mnb.fit(X_train,y_train)
y_pred2 = mnb.predict(X_test)
print(accuracy_score(y_test,y_pred2))
print(confusion_matrix(y_test,y_pred2))
print(precision_score(y_test,y_pred2))
```

```
0.9709864603481625
[[896  0]
 [ 30 108]]
1.0
```

```
bnb.fit(X_train,y_train)
y_pred3 = bnb.predict(X_test)
print(accuracy_score(y_test,y_pred3))
print(confusion_matrix(y_test,y_pred3))
print(precision_score(y_test,y_pred3))
```

```
0.9835589941972921
[[895  1]
 [ 16 122]]
0.991869918699187
```

```
# tfidf --> MNB
```

```
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
from xgboost import XGBClassifier

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)
bc = BaggingClassifier(n_estimators=50, random_state=2)
etc = ExtraTreesClassifier(n_estimators=50, random_state=2)
gbdt = GradientBoostingClassifier(n_estimators=50, random_state=2)
xgb = XGBClassifier(n_estimators=50, random_state=2)

clfs = {
    'SVC' : svc,
    'KN' : knc,
    'NB': mnb,
    'DT': dtc,
    'LR': lrc,
    'RF': rfc,
    'AdaBoost': abc,
    'BgC': bc,
    'ETC': etc,
    'GBDT': gbdt,
    'xgb': xgb
}
```

```
def train_classifier(clf,X_train,y_train,X_test,y_test):  
    clf.fit(X_train,y_train)  
    y_pred = clf.predict(X_test)  
    accuracy = accuracy_score(y_test,y_pred)  
    precision = precision_score(y_test,y_pred)  
  
    return accuracy,precision
```

```
train_classifier(svc,X_train,y_train,X_test,y_test)
```

```
(0.9758220502901354, 0.9747899159663865)
```

```
accuracy_scores = []  
precision_scores = []
```

```
for name,clf in clfs.items():
```

```
    current_accuracy,current_precision = train_classifier(clf, X_train,y_train,X_test,y_test)
```

```
    print("For ",name)  
    print("Accuracy - ",current_accuracy)  
    print("Precision - ",current_precision)
```

```
    accuracy_scores.append(current_accuracy)  
    precision_scores.append(current_precision)
```

```
For SVC  
Accuracy - 0.9758220502901354  
Precision - 0.9747899159663865  
For KN  
Accuracy - 0.9052224371373307  
Precision - 1.0  
For NB  
Accuracy - 0.9709864603481625  
Precision - 1.0  
For DT  
Accuracy - 0.9294003868471954
```

```
Precision - 0.8282828282828283
For LR
Accuracy - 0.9584139264990329
Precision - 0.9702970297029703
For RF
Accuracy - 0.9748549323017408
Precision - 0.9827586206896551
For AdaBoost
Accuracy - 0.960348162475822
Precision - 0.9292035398230089
For BgC
Accuracy - 0.9574468085106383
Precision - 0.8671875
For ETC
Accuracy - 0.9748549323017408
Precision - 0.9745762711864406
For GBDT
Accuracy - 0.9477756286266924
Precision - 0.92
For xgb
Accuracy - 0.971953578336557
Precision - 0.943089430894309
```

```
performance_df = pd.DataFrame({'Algorithm':clfs.keys(),'Accuracy':accuracy_scores,'Precision':precision_scores}).sort_values('Precisi
```

```
performance_df
```

	Algorithm	Accuracy	Precision
1	KN	0.905222	1.000000
2	NB	0.970986	1.000000
5	RF	0.974855	0.982759
0	SVC	0.975822	0.974790
8	ETC	0.974855	0.974576

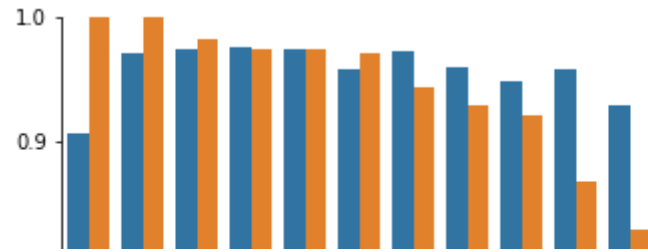
```
performance_df1 = pd.melt(performance_df, id_vars = "Algorithm")
```

```
-- . - - - - -  
performance_df1
```

	Algorithm	variable	value
0	KN	Accuracy	0.905222
1	NB	Accuracy	0.970986
2	RF	Accuracy	0.974855
3	SVC	Accuracy	0.975822
4	ETC	Accuracy	0.974855
5	LR	Accuracy	0.958414
6	xgb	Accuracy	0.971954
7	AdaBoost	Accuracy	0.960348
8	GBDT	Accuracy	0.947776

```
sns.catplot(x = 'Algorithm', y='value',  
            hue = 'variable',data=performance_df1, kind='bar',height=5)  
plt.ylim(0.5,1.0)  
plt.xticks(rotation='vertical')  
plt.show()
```





```
# Voting Classifier
```

```
svc = SVC(kernel='sigmoid', gamma=1.0, probability=True)
```

```
mnb = MultinomialNB()
```

```
etc = ExtraTreesClassifier(n_estimators=50, random_state=2)
```

```
from sklearn.ensemble import VotingClassifier
```



```
voting = VotingClassifier(estimators=[('svm', svc), ('nb', mnb), ('et', etc)], voting='soft')
```



```
voting.fit(X_train, y_train)
```

```
VotingClassifier(estimators=[('svm',
                             SVC(gamma=1.0, kernel='sigmoid',
                                probability=True)),
                             ('nb', MultinomialNB()),
                             ('et',
                              ExtraTreesClassifier(n_estimators=50,
                                                    random_state=2))],
                 voting='soft')
```

```
y_pred = voting.predict(X_test)
```

```
print("Accuracy", accuracy_score(y_test, y_pred))
```

```
print("Precision", precision_score(y_test, y_pred))
```

```
Accuracy 0.9816247582205029
```

```
Precision 0.9917355371900827
```

```
# Applying stacking
```

```
estimators=[('svm', svc), ('nb', mnb), ('et', etc)]
final_estimator=RandomForestClassifier()

from sklearn.ensemble import StackingClassifier

clf = StackingClassifier(estimators=estimators, final_estimator=final_estimator)

clf.fit(X_train,y_train)
y_pred = clf.predict(X_test)
print("Accuracy",accuracy_score(y_test,y_pred))
print("Precision",precision_score(y_test,y_pred))

    Accuracy 0.9796905222437138
    Precision 0.9398496240601504

import pickle
pickle.dump(tfidf,open('vectorizer.pkl','wb'))
pickle.dump(mnb,open('model.pkl','wb'))
```

Selected TFIDF Vectorizer= max\_features parameter of -->3000

With Multinomial Naive Bayes algorithm which gives highest precision value

[Colab paid products](#) - [Cancel contracts here](#)

