

```

1 import numpy as np
2 import pandas as pd

1 df = pd.read_csv('/content/drive/MyDrive/Bharat Intern Internship/SMS spam detection/encoded-sp

1 df.sample(5)

```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
5318	ham	Good morning, my Love ... I go to sleep now an...	NaN	NaN	NaN
711	ham	It just seems like weird timing that the night...	NaN	NaN	NaN
		Ok I din get ur			

```

1 df.shape

(5572, 5)

```

```

1 # 1. Data cleaning
2 # 2. EDA
3 # 3. Text Preprocessing
4 # 4. Model building
5 # 5. Evaluation
6 # 6. Improvement
7 # 7. Website
8 # 8. Deploy

```

✓ 1.Data Cleaning

```

1 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

1 # drop last 3 cols
2 df.drop(columns=['Unnamed: 2','Unnamed: 3','Unnamed: 4'],inplace=True)

1 df.sample(5)

```

	v1	v2	
1464	ham	Wat makes u thk i'll fall down. But actually i...	
288	ham	My life Means a lot to me, Not because I love ...	
4840	ham	Hmmm.... Mayb can try e shoppin area one, but ...	
3973	ham	Roger that. We%Űre probably going to rem in a...	
3375	ham	Good afternon, my love. How are today? I hope ...	

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

	target	text	
1236	ham	How much are we getting?	
4493	ham	Man this bus is so so so slow. I think you're ...	
2087	ham	\alright babe	
2501	ham	No da..today also i forgot..	
575	spam	You have won ?1,000 cash or a ?2,000 prize! To...	

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

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

```
1 df.head()
```

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

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

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

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

```
403
```

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

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

```
0
```

```
1 df.shape  
  
(5169, 2)
```

✓ 2.EDA

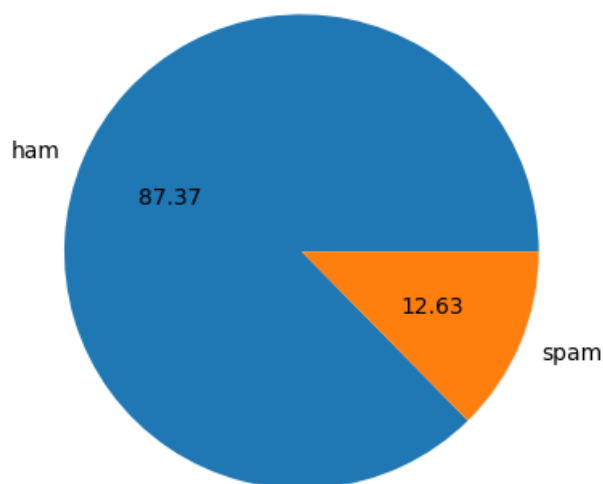
```
1 df.head()
```

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

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

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

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



```
1 # Data is imbalanced
```

```
1 !pip install nltk
```

```
Requirement already satisfied: nltk in /usr/local/lib/python3.10/dist-packages (3.8.1)  
Requirement already satisfied: click in /usr/local/lib/python3.10/dist-packages (from nltk) (8.1.7)  
Requirement already satisfied: joblib in /usr/local/lib/python3.10/dist-packages (from nltk) (1.3.2)  
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.10/dist-packages (from nltk) (2023  
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from nltk) (4.66.1)
```



```
1 import nltk
```

```
1 nltk.download('punkt')
```

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

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



```
1 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	

```
1 # num of words
```



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

```
1 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	

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



```
1 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	



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

	num_characters	num_words	num_sentences	
count	5169.000000	5169.000000	5169.000000	
mean	78.977945	18.455794	1.965564	
std	58.236293	13.324758	1.448541	
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	38.000000	

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

	num_characters	num_words	num_sentences	
count	4516.000000	4516.000000	4516.000000	
mean	70.459256	17.123782	1.820195	
std	56.358207	13.493970	1.383657	
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	38.000000	

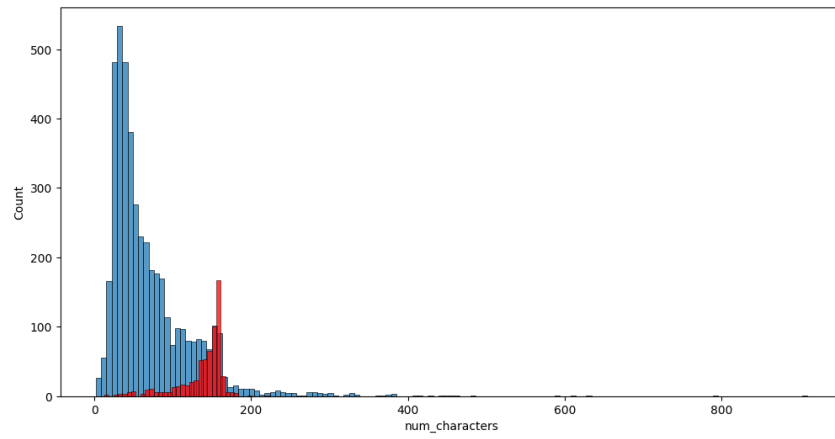
```
1 #spam
2 df[df['target'] == 1][['num_characters', 'num_words', 'num_sentences']].describe()
```

	num_characters	num_words	num_sentences	
count	653.000000	653.000000	653.000000	
mean	137.891271	27.667688	2.970904	
std	30.137753	7.008418	1.488425	
min	13.000000	2.000000	1.000000	
25%	132.000000	25.000000	2.000000	
50%	149.000000	29.000000	3.000000	
75%	157.000000	32.000000	4.000000	
max	224.000000	46.000000	9.000000	

```
1 import seaborn as sns
```

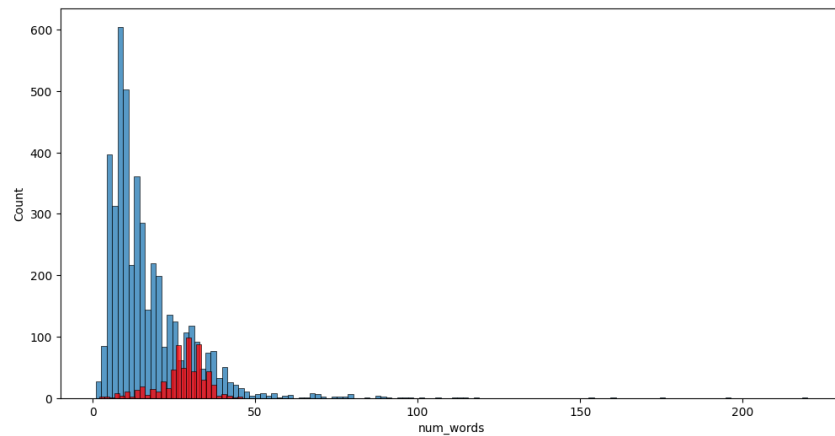
```
1 plt.figure(figsize=(12,6))
2 sns.histplot(df[df['target'] == 0]['num_characters'])
3 sns.histplot(df[df['target'] == 1]['num_characters'],color='red')
```

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



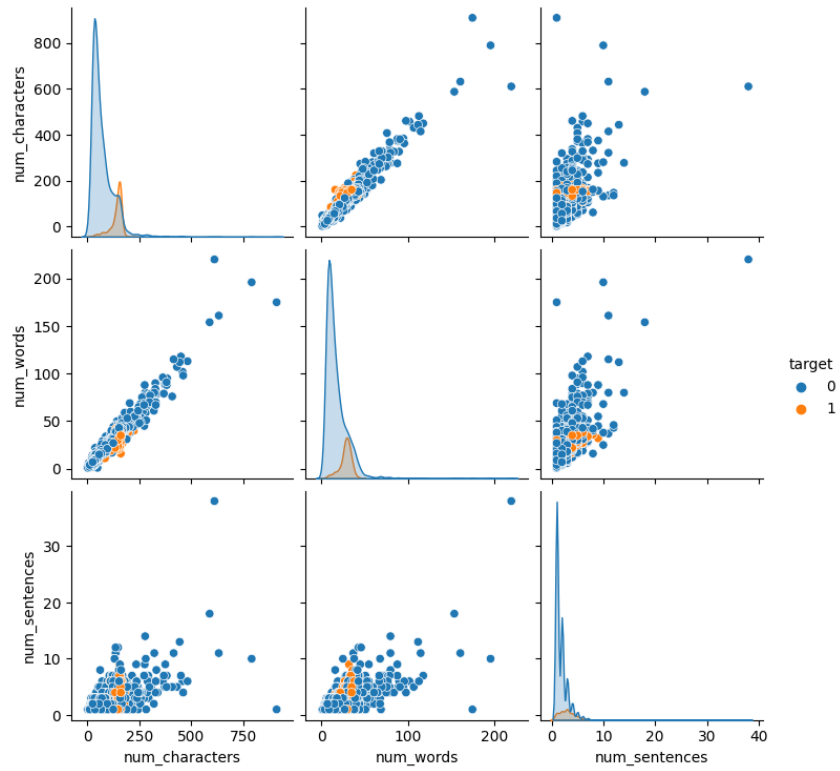
```
1 plt.figure(figsize=(12,6))
2 sns.histplot(df[df['target'] == 0]['num_words'])
3 sns.histplot(df[df['target'] == 1]['num_words'],color='red')
```

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



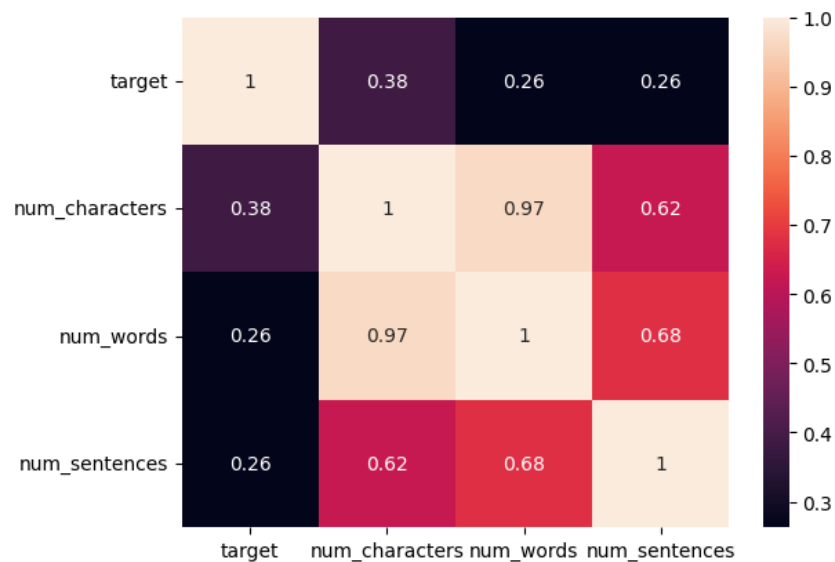
```
1 sns.pairplot(df,hue='target')
```

```
<seaborn.axisgrid.PairGrid at 0x780ed062bbb0>
```



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

```
<ipython-input-81-8df7bcac526d>:1: FutureWarning: The default value c
sns.heatmap(df.corr(),annot=True)
<Axes: >
```



✓ 3.Data Preprocessing

1. Lower case
 2. Tokenization
 3. Removing special characters
 4. Removing stop words and punctuation
 5. Stemming
-

```
1 from nltk.stem.porter import PorterStemmer
2 ps = PorterStemmer()
3 import nltk
4 import string
5 nltk.download('stopwords')
6 ps.stem('loving')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
'love'
```



```

1 def transform_text(text):
2     text = text.lower()
3     text = nltk.word_tokenize(text)
4
5     y = []
6     for i in text:
7         if i.isalnum():
8             y.append(i)
9
10    text = y[:]
11    y.clear()
12
13    for i in text:
14        if i not in stopwords.words('english') and i not in string.punctuation:
15            y.append(i)
16
17    text = y[:]
18    y.clear()
19
20    for i in text:
21        y.append(ps.stem(i))
22
23
24    return " ".join(y)

```

```

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

```

```

1 df.head()

```

	target	text	num_characters	num_words	num_sentences	transformed_text
0	0	Go until jurong point, crazy.. Available only ... Ok lar... ...	111	24	2	go crazi g

```

1 from wordcloud import WordCloud
2 wc = WordCloud(width=500,height=500,min_font_size=10,background_color='white')

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

1 plt.figure(figsize=(15,6))
2 plt.imshow(spam_wc)

```

[illegible]

```
1 plt.figure(figsize=(15,6))
2 plt.imshow(ham_wc)
```

[illegible]

```
1 df.head()
```

	target	text	num_characters	num_words	num_sentences	transf
0	0	Go until jurong point, crazy.. Available only ... Ok lar...	111	24	2	go crazi g

```

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

```

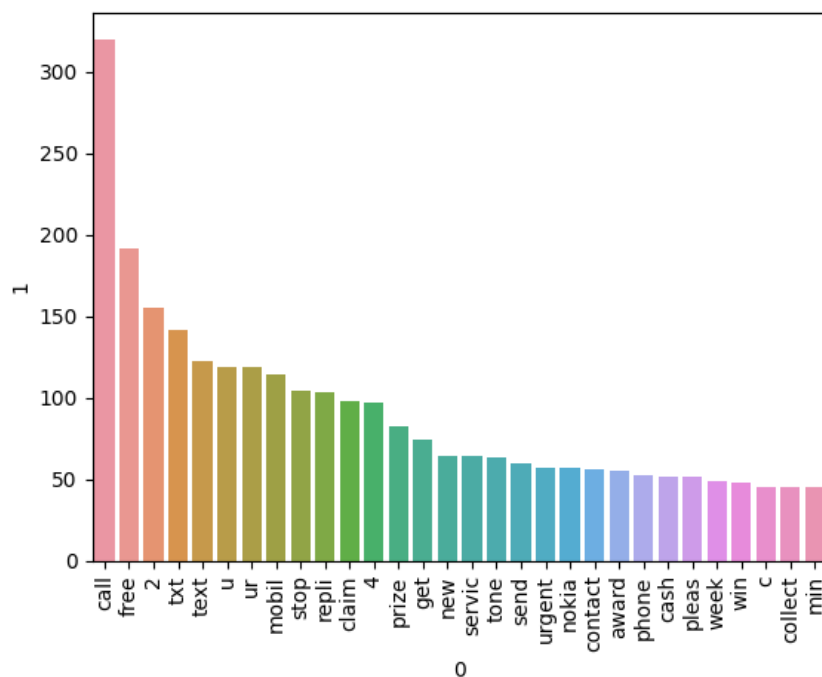
```
1 len(spam_corpus)
```

9939

```

1 from collections import Counter
2 import matplotlib.pyplot as plt
3 sns.barplot(x=pd.DataFrame(Counter(spam_corpus).most_common(30))[0],y=pd.DataFrame(Counter(span
4 plt.xticks(rotation='vertical')
5 plt.show()

```



```

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

```

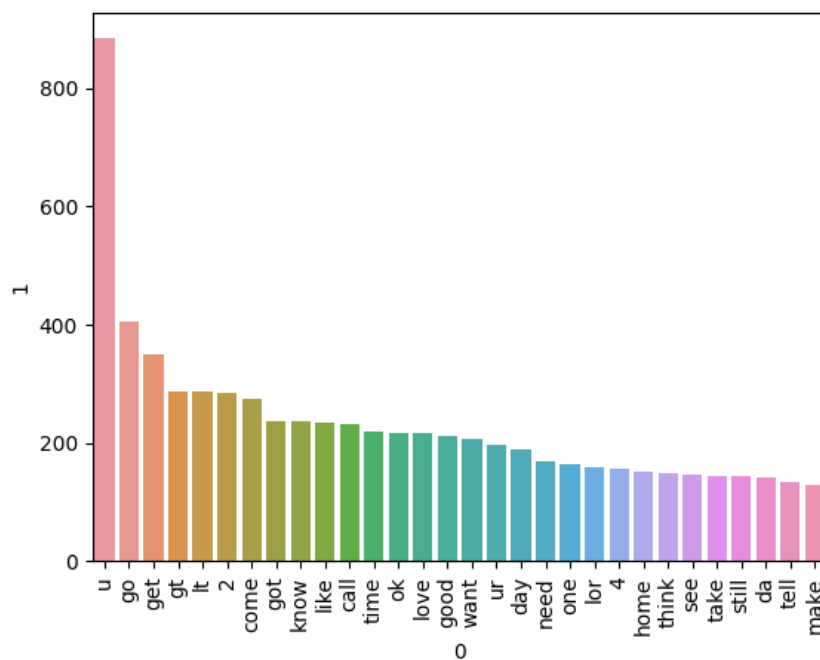
```
1 len(ham_corpus)
```

35404

```

1 from collections import Counter
2 sns.barplot(x=pd.DataFrame(Counter(ham_corpus).most_common(30))[0],y=pd.DataFrame(Counter(ham_c
3 plt.xticks(rotation='vertical')
4 plt.show()

```



```

1 # Text Vectorization
2 # using Bag of Words
3 df.head()

```

	target	text	num_characters	num_words	num_sentences	transfo
0	0	Go until jurong point, crazy.. Available only in lar...	111	24	2	go crazy g
1	1	Ok lar... I prefer the sea view!	41	10	1	Ok lar... I prefer the sea view!

4.Model Building

```

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

```

```

1 X = tfidf.fit_transform(df['transformed_text']).toarray()

```

```

1 #from sklearn.preprocessing import MinMaxScaler
2 #scaler = MinMaxScaler()
3 #X = scaler.fit_transform(X)

```

```

1 # appending the num_character col to X
2 #X = np.hstack((X,df['num_characters'].values.reshape(-1,1)))

```

```

1 X.shape

(5169, 3000)

1 y = df['target'].values

1 from sklearn.model_selection import train_test_split

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
2 from sklearn.metrics import accuracy_score,confusion_matrix,precision_score

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))

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

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))

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

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))

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

1 # tfidf --> MNB

```

```

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)

(0.9758220502901354, 0.9747899159663865)

```

```

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)

```

```

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 - 0.9303675048355899
Precision - 0.8173076923076923
For LR
Accuracy - 0.9584139264990329
Precision - 0.9702970297029703
For RF
Accuracy - 0.9758220502901354
Precision - 0.9829059829059829
For AdaBoost
Accuracy - 0.960348162475822
Precision - 0.9292035398230089
For BgC
Accuracy - 0.9584139264990329
Precision - 0.8682170542635659
For ETC
Accuracy - 0.9748549323017408
Precision - 0.9745762711864406
For GBDT
Accuracy - 0.9468085106382979
Precision - 0.9191919191919192
For xgb
Accuracy - 0.9671179883945842
Precision - 0.9262295081967213

```

```

1 performance_df = pd.DataFrame({'Algorithm':clfs.keys(),'Accuracy':accuracy_scores,'Precision':p
2 performance_df

```

	Algorithm	Accuracy	Precision	
1	KN	0.905222	1.000000	
2	NB	0.970986	1.000000	
5	RF	0.975822	0.982906	
0	SVC	0.975822	0.974790	
8	ETC	0.974855	0.974576	
4	LR	0.958414	0.970297	
6	AdaBoost	0.960348	0.929204	
10	xgb	0.967118	0.926230	
9	GBDT	0.946809	0.919192	
7	BgC	0.958414	0.868217	
3	DT	0.930368	0.817308	

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

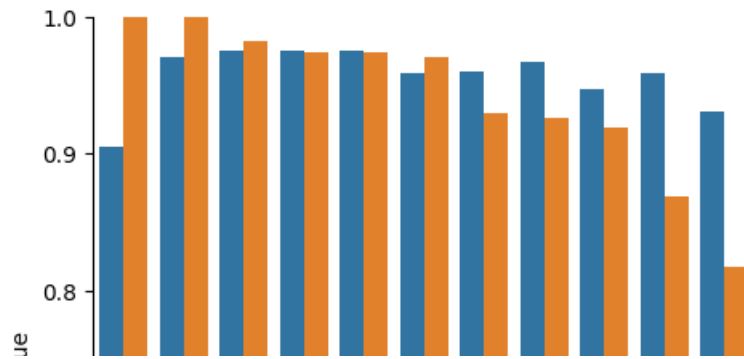
```
1 performance_df1
```

	Algorithm	variable	value	
0	KN	Accuracy	0.905222	
1	NB	Accuracy	0.970986	
2	RF	Accuracy	0.975822	
3	SVC	Accuracy	0.975822	
4	ETC	Accuracy	0.974855	
5	LR	Accuracy	0.958414	
6	AdaBoost	Accuracy	0.960348	
7	xgb	Accuracy	0.967118	
8	GBDT	Accuracy	0.946809	
9	BgC	Accuracy	0.958414	
10	DT	Accuracy	0.930368	
11	KN	Precision	1.000000	
12	NB	Precision	1.000000	
13	RF	Precision	0.982906	
14	SVC	Precision	0.974790	
15	ETC	Precision	0.974576	
16	LR	Precision	0.970297	
17	AdaBoost	Precision	0.929204	
18	xgb	Precision	0.926230	
19	GBDT	Precision	0.919192	
20	BgC	Precision	0.868217	
21	DT	Precision	0.817308	

```
1 sns.catplot(x = 'Algorithm', y='value',
2             hue = 'variable',data=performance_df1, kind='bar',height=5)
```



```
3 plt.ylim(0.5,1.0)
4 plt.xticks(rotation='vertical')
5 plt.show()
```



Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.