```
import numpy as np
import pandas as pd
```

```
!iconv -f ISO-8859-1 -t UTF-8 spam.csv > spam_utf8.csv
```

```
df = pd.read_csv('spam_utf8.csv')
```

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

|  | v1 | v2 | Unnamed: 2 | Unnamed: 3 | Unnamed: 4 |
|---|---|---|---|---|---|
| 2876 | ham | twenty past five he said will this train have ... | NaN | NaN | NaN |
| 3807 | ham | Mm you ask him to come its enough :-) | NaN | NaN | NaN |
| 1769 | ham | Ha... Both of us doing e same thing. But i got... | NaN | NaN | NaN |
| 4012 | ham | Ok. | NaN | NaN | NaN |
| 373 | ham | I cant keep talking to people if am not sure i... | NaN | NaN | NaN |

```
df.shape
```

```
(5572, 5)
```

```
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
```

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

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

1 to 5 of 5 entries    Filter

| index | v1 | v2 |
|---|---|---|
| 4465 | ham | Hey u still at the gym? |
| 5252 | ham | You do your studies alone without anyones help. If you cant no need to study. |
| 3676 | ham | Great! So what attracts you to the brothas? |
| 3194 | ham | Great. P diddy is my neighbor and comes for toothpaste every morning |
| 890 | ham | Why do you ask princess? |

Show 25 ▾ per page

Like what you see? Visit the data table notebook to learn more about interactive tables.

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

|  | target | text |
|---|---|---|
| 170 | ham | Sir, I need AXIS BANK account no and bank addr... |
| 989 | ham | Ugh. Gotta drive back to sd from la. My butt i... |
| 2534 | ham | Ok enjoy . R u there in home. |
| 225 | ham | Would really appreciate if you call me. Just n... |
| 756 | ham | Cant think of anyone with * spare room off * t... |

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

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

```
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... | |

```
df.isnull().sum()
```

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

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

```
403
```

```
df = df.drop_duplicates(keep='first')
```

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

```
0
```
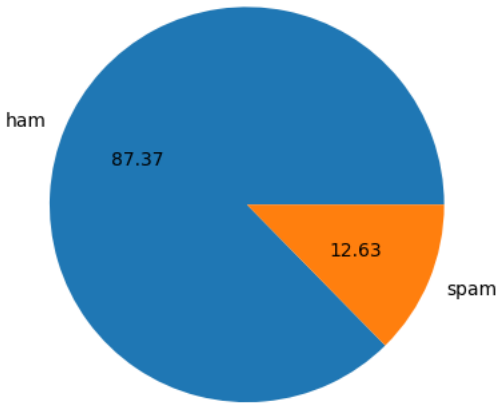
```
df.shape
```

```
(5169, 2)
```

```
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... | |

```
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()
```

```
import nltk
```

```
!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) (2023.6.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from nltk) (4.66.1)
```

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

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

```
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 |

```
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 |

```
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 |

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

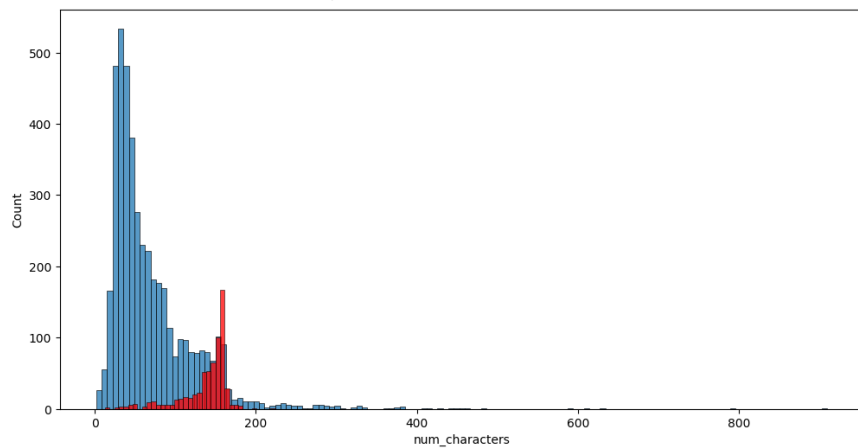| | num_characters | num_words | num_sentences |
|---|---|---|---|
| **count** | 4516.000000 | 4516.000000 | 4516.000000 |
| **mean** | 70.459256 | 17.123782 | 1.820195 |
| **std** | 56.358207 | 13.493970 | 1.383657 |
| **min** | 2.000000 | 1.000000 | 1.000000 |
| **25%** | 34.000000 | 8.000000 | 1.000000 |
| **50%** | 52.000000 | 13.000000 | 1.000000 |
| **75%** | 90.000000 | 22.000000 | 2.000000 |
| **max** | 910.000000 | 220.000000 | 38.000000 |

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

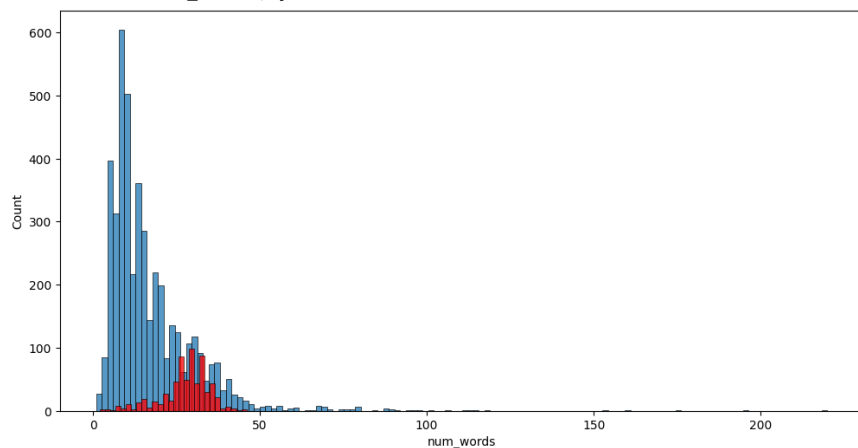| | num_characters | num_words | num_sentences |
|---|---|---|---|
| **count** | 653.000000 | 653.000000 | 653.000000 |
| **mean** | 137.891271 | 27.667688 | 2.970904 |
| **std** | 30.137753 | 7.008418 | 1.488425 |
| **min** | 13.000000 | 2.000000 | 1.000000 |
| **25%** | 132.000000 | 25.000000 | 2.000000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

```
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')
```

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



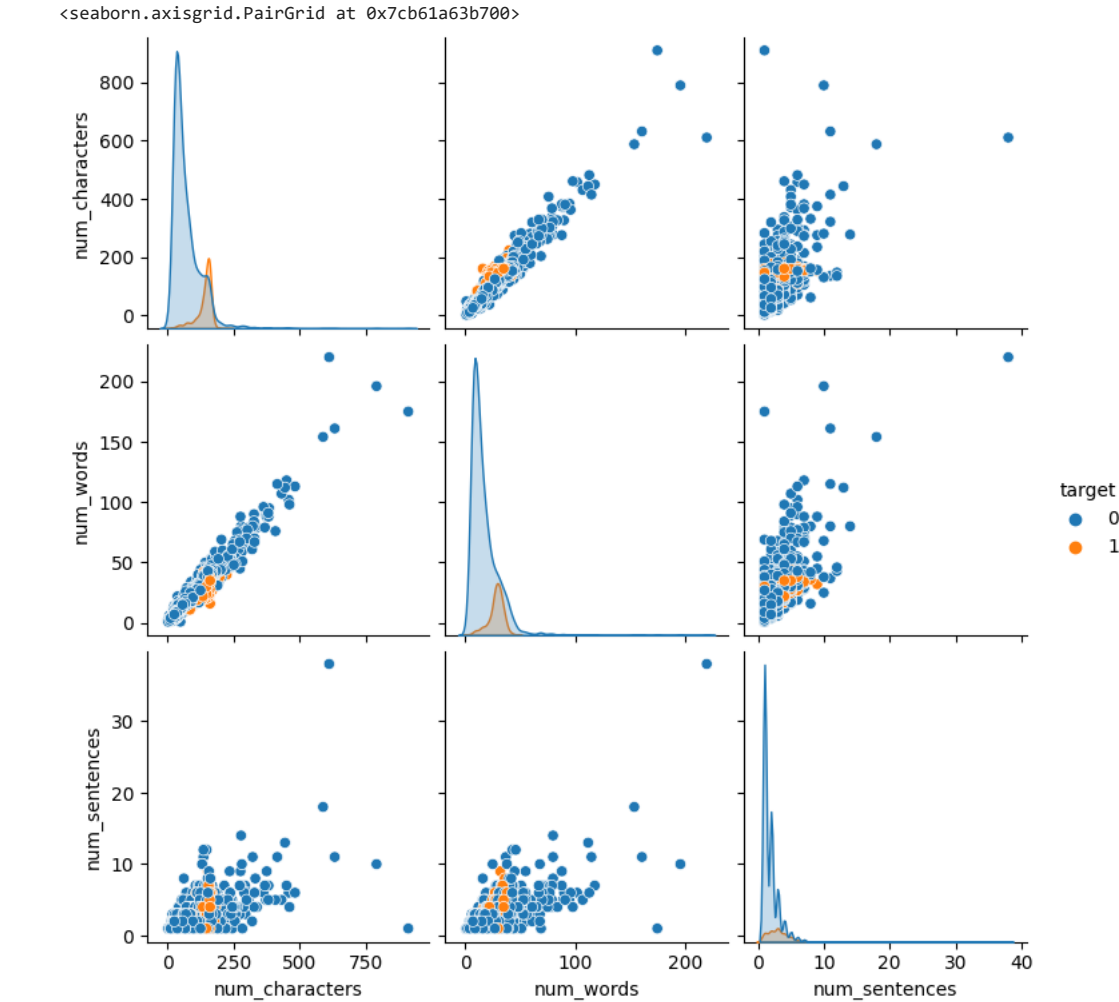```
plt.figure(figsize=(12,6))
sns.histplot(df[df['target'] == 0]['num_words'])
sns.histplot(df[df['target'] == 1]['num_words'],color='red')
```
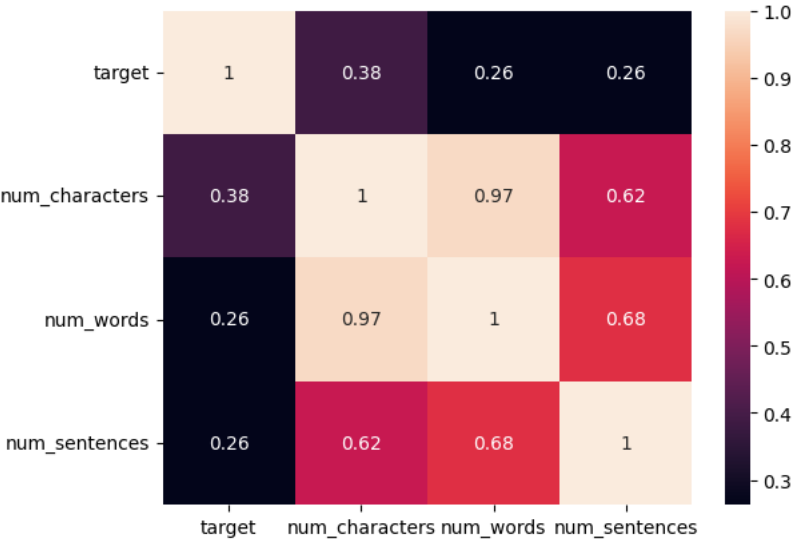
<Axes: xlabel='num_words', ylabel='Count'>



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

```
<seaborn.axisgrid.PairGrid at 0x7cb61a63b700>
```



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

```
<ipython-input-38-8df7bcac526d>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future ver
  sns.heatmap(df.corr(),annot=True)
<Axes: >
```

```python
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)
```

```python
!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) (2023.6.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from nltk) (4.66.1)
```

```python
!nltk.download('stopwords')
```

```
/bin/bash: -c: line 1: syntax error near unexpected token `'stopwords''
/bin/bash: -c: line 1: `nltk.download('stopwords')'
```

```python
from nltk.corpus import stopwords
stopwords.words('english')
```

```
 'needn',
 "needn't",
 'shan',
 "shan't",
 'shouldn',
 "shouldn't",
 'wasn',
 "wasn't",
 'weren',
 "weren't",
 'won',
 "won't",
 'wouldn',
 "wouldn't"]
```

```
import string
string.punctuation
```

```
'!"#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~'
```

```
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.'
```

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

```
'love'
```

```
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 alreadi 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=" "))
```

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

```
<matplotlib.image.AxesImage at 0x7cb60fd75420>
```



```
ham_wc = wc.generate(df[df['target'] == 0]['transformed_text'].str.cat(sep=" "))
```

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

```
<matplotlib.image.AxesImage at 0x7cb60fdd90f0>
```



```
df.head()
```

1 to 5 of 5 entries  Filter

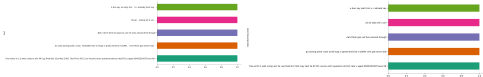| index | target | text | num_characters | num_words | num_sentences | transformed_text |
|-------|--------|------|----------------|-----------|---------------|------------------|
| 0 | 0 | Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat... | 111 | 24 | 2 | go jurong point crazi avail bugi n great world la e buffet cine got amor wat |
| 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 final tkts 21st May 2005. Text FA to 87121 to receive entry question(std txt rate)T&C's apply 08452810075over18's | 155 | 37 | 2 | free entri 2 wkli comp win fa cup final tkt 21st may text fa 87121 receiv entri question std txt rate c appli 08452810075over18 |
| 3 | 0 | U dun say so early hor... U c already then say... | 49 | 13 | 1 | u dun say earli hor u c alreadi say |
| 4 | 0 | Nah I don't think he goes to usf, he lives around here though | 61 | 15 | 1 | nah think goe usf live around though |

Show 25 ▾ per page

Like what you see? Visit the [data table notebook](#) to learn more about interactive tables.
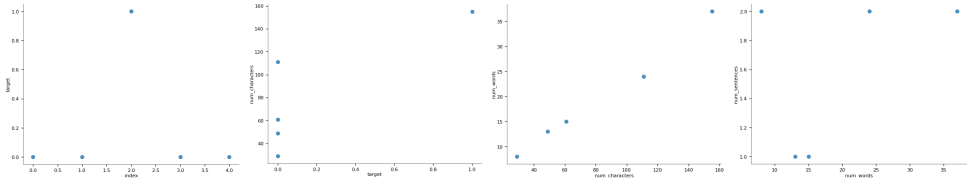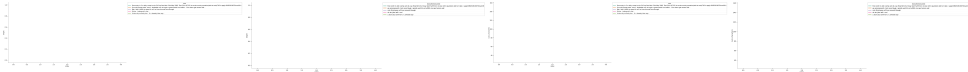
**Distributions**



**Categorical distributions**
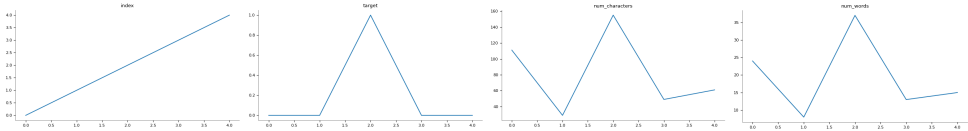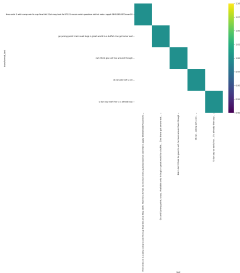


**2-d distributions**



**Time series**



**Values**



**2-d categorical distributions**



**Faceted distributions**

```
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
```

```
# Text Vectorization
# using Bag of Words
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 sav... | 49 | 13 | 1 | u dun say earli hor u c alreadi say |

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

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

```python
X.shape
```

```
(5169, 3000)
```

```python
y = df['target'].values
```

```python
from sklearn.model_selection import train_test_split
```

```python
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,random_state=2)
```

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

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

```python
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
```

```python
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
```

```python
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
```

```python
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
```

```python
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)
```

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

```python
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
```

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

```
    (0.9758220502901354, 0.9747899159663865)
```

```python
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 -  0.9303675048355899
    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
```

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

```
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  |

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

```
performance_df1
```

| | Algorithm | variable | value |
|---|---|---|---|
| **0** | KN | Accuracy | 0.905222 |
| **1** | NB | Accuracy | 0.970986 |

```
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()
```