```
!pip install https://github.com/scikit-learn/scikit-learn/archive/master.zip
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting https://github.com/scikit-learn/scikit-learn/archive/master.zip
  Using cached https://github.com/scikit-learn/scikit-learn/archive/master.zip
  Installing build dependencies ... done
  Getting requirements to build wheel ... done
  Preparing metadata (pyproject.toml) ... done
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.8/dist-packages (from scikit-learn==1.3.dev0) (1.2.0)
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.8/dist-packages (from scikit-learn==1.3.dev0) (1.22.4)
Requirement already satisfied: scipy>=1.3.2 in /usr/local/lib/python3.8/dist-packages (from scikit-learn==1.3.dev0) (1.10.1)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.8/dist-packages (from scikit-learn==1.3.dev0) (3.1.0)
```

```
pip install -U scikit-learn
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.8/dist-packages (1.3.dev0)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.8/dist-packages (from scikit-learn) (1.2.0)
Requirement already satisfied: scipy>=1.3.2 in /usr/local/lib/python3.8/dist-packages (from scikit-learn) (1.10.1)
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.8/dist-packages (from scikit-learn) (1.22.4)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.8/dist-packages (from scikit-learn) (3.1.0)
```

```python
import io
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import nltk
from nltk.corpus import stopwords
import string
import math
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix, accuracy_score, roc_auc_score, roc_curve
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV, train_test_split
%matplotlib inline
```

```python
from google.colab import files
uploaded = files.upload()
```

```
Choose Files  yelp.csv
  • yelp.csv(text/csv) - 8091185 bytes, last modified: 3/3/2023 - 100% done
  Saving yelp.csv to yelp (2).csv
```

```python
# LOAD THE DATASET AND SEEING THE DETAILS
data = pd.read_csv(io.BytesIO(uploaded['yelp.csv']))
# SHAPE OF THE DATASET
print("Shape of the dataset:")
print(data.shape)
# COLUMN NAMES
print("Column names:")
print(data.columns)
# DATATYPE OF EACH COLUMN
print("Datatype of each column:")
print(data.dtypes)
# SEEING FEW OF THE ENTRIES
print("Few dataset entries:")
print(data.head())
# DATASET SUMMARY
data.describe(include='all')
```

```
Shape of the dataset:
(10000, 10)
Column names:
Index(['business_id', 'date', 'review_id', 'stars', 'text', 'type', 'user_id',
       'cool', 'useful', 'funny'],
      dtype='object')
Datatype of each column:
business_id    object
date           object
review_id      object
stars           int64
text           object
type           object
user_id        object
cool            int64
useful          int64
funny           int64
dtype: object
Few dataset entries:
       business_id        date              review_id  stars  \
0  9yKzy9PApeiPPOUJEtnvkg  2011-01-26  fWKvX83p0-ka4JS3dc6E5A      5
1  ZRJwVLyzEJq1VAihDhYiow  2011-07-27  IjZ33sJrzXqU-0X6U8NwyA      5
2  6oRAC4uyJCsJl1X0WZpVSA  2012-06-14  IESLBzqUCLdSzSqm0eCSxQ      4
3  _1QQZuf4zZOyFCvXc0o6Vg  2010-05-27  G-WvGaISbqqaMHlNnByodA      5
4  6ozycU1RpktNG2-1BroVtw  2012-01-05  1uJFq2r5QfJG_6ExMRCaGw      5


                                              text    type  \
0  My wife took me here on my birthday for breakf...  review
1  I have no idea why some people give bad review...  review
2  love the gyro plate. Rice is so good and I als...  review
3  Rosie, Dakota, and I LOVE Chaparral Dog Park!!...  review
4  General Manager Scott Petello is a good egg!!!...  review


              user_id  cool  useful  funny
0  rLtl8ZkDX5vH5nAx9C3q5Q     2       5      0
1  0a2KyEL0d3Yb1V6aivbIuQ     0       0      0
2  0hT2KtfLiobPvh6cDC8JQg     0       1      0
3  uZetl9T0NcROGOyFfughhg     1       2      0
4  vYmM4KTsC8ZfQBg-j5MWkw     0       0      0
```

|       | business_id | date | review_id | stars | text | type | user_id | |
|-------|-------------|------|-----------|-------|------|------|---------|---|
| count | 10000 | 10000 | 10000 | 10000.000000 | 10000 | 10000 | 10000 | 100 |
| unique | 4174 | 1995 | 10000 | NaN | 9998 | 1 | 6403 | |
| top | JokKtdXU7zXHcr20Lrk29A | 2011-03-28 | fWKvX83p0-ka4JS3dc6E5A | NaN | Great service | review | fczQCSmaWF78toLEmb0Zsw | |
| freq | 37 | 21 | 1 | NaN | 2 | 10000 | 38 | |

```
#CREATING A NEW COLUMN IN THE DATASET FOR THE NUMBER OF WORDS IN THE REVIEW
data['length'] = data['text'].apply(len)
data.head()
```

|   | business_id | date | review_id | stars | text | type | user_id | |
|---|-------------|------|-----------|-------|------|------|---------|---|
| 0 | 9yKzy9PApeiPPOUJEtnvkg | 2011-01-26 | fWKvX83p0-ka4JS3dc6E5A | 5 | My wife took me here on my birthday for breakf... | review | rLtl8ZkDX5vH5nAx9C3q5Q | |
| 1 | ZRJwVLyzEJq1VAihDhYiow | 2011-07-27 | IjZ33sJrzXqU-0X6U8NwyA | 5 | I have no idea why some people give bad review... | review | 0a2KyEL0d3Yb1V6aivbIuQ | |
| | | | | | love the gyro | | | |

```
# COMPARING TEXT LENGTH TO STARS
graph = sns.FacetGrid(data=data,col='stars')
graph.map(plt.hist,'length',bins=50,color='blue')
```

```
<seaborn.axisgrid.FacetGrid at 0x7f6c7c0b70a0>
```



```
# GETTING THE MEAN VALUES OF THE VOTE COLUMNS WRT THE STARS ON THE REVIEW
stval = data.groupby('stars').mean()
stval
```

|  stars | cool | useful | funny | length |
|---|---|---|---|---|
| 1 | 0.576769 | 1.604806 | 1.056075 | 826.515354 |
| 2 | 0.719525 | 1.563107 | 0.875944 | 842.256742 |
| 3 | 0.788501 | 1.306639 | 0.694730 | 758.498289 |
| 4 | 0.954623 | 1.395916 | 0.670448 | 712.923142 |
| 5 | 0.944261 | 1.381780 | 0.608631 | 624.999101 |

```
stval.corr()
```

|  | cool | useful | funny | length |
|---|---|---|---|---|
| cool | 1.000000 | -0.743329 | -0.944939 | -0.857664 |
| useful | -0.743329 | 1.000000 | 0.894506 | 0.699881 |
| funny | -0.944939 | 0.894506 | 1.000000 | 0.843461 |
| length | -0.857664 | 0.699881 | 0.843461 | 1.000000 |

```
# CLASSIFICATION
data_classes = data[(data['stars']==1) | (data['stars']==3) | (data['stars']==5)]
data_classes.head()
print(data_classes.shape)

# Seperate the dataset into X and Y for prediction
x = data_classes['text']
y = data_classes['stars']
print(x.head())
print(y.head())
```

```
    (5547, 11)
    0    My wife took me here on my birthday for breakf...
    1    I have no idea why some people give bad review...
    3    Rosie, Dakota, and I LOVE Chaparral Dog Park!!...
    4    General Manager Scott Petello is a good egg!!!...
    6    Drop what you're doing and drive here. After I...
    Name: text, dtype: object
    0    5
    1    5
    3    5
    4    5
    6    5
    Name: stars, dtype: int64
```

```
# CLEANING THE REVIEWS - REMOVAL OF STOPWORDS AND PUNCTUATION
def text_process(text):
    nopunc = [char for char in text if char not in string.punctuation]
    nopunc = ''.join(nopunc)
    return [word for word in nopunc.split() if word.lower() not in stopwords.words('english')]
```

```
X = ["Mr. Green killed Colonel Mustard in the study with the candlestick. \
Mr. Green is not a very nice fellow.",
    "Professor Plum has a green plant in his study.",
    "Miss Scarlett watered Professor Plum's green plant while he was away \
from his office last week."]
```

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer()
vectorizer.fit(X)
```

```
▾ CountVectorizer
CountVectorizer()
```

```
import nltk
nltk.download("popular")
```

```
[nltk_data] Downloading collection 'popular'
[nltk_data]    |
[nltk_data]    | Downloading package cmudict to /root/nltk_data...
[nltk_data]    |   Package cmudict is already up-to-date!
[nltk_data]    | Downloading package gazetteers to /root/nltk_data...
[nltk_data]    |   Package gazetteers is already up-to-date!
[nltk_data]    | Downloading package genesis to /root/nltk_data...
[nltk_data]    |   Package genesis is already up-to-date!
[nltk_data]    | Downloading package gutenberg to /root/nltk_data...
[nltk_data]    |   Package gutenberg is already up-to-date!
[nltk_data]    | Downloading package inaugural to /root/nltk_data...
[nltk_data]    |   Package inaugural is already up-to-date!
[nltk_data]    | Downloading package movie_reviews to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package movie_reviews is already up-to-date!
[nltk_data]    | Downloading package names to /root/nltk_data...
[nltk_data]    |   Package names is already up-to-date!
[nltk_data]    | Downloading package shakespeare to /root/nltk_data...
[nltk_data]    |   Package shakespeare is already up-to-date!
[nltk_data]    | Downloading package stopwords to /root/nltk_data...
[nltk_data]    |   Package stopwords is already up-to-date!
[nltk_data]    | Downloading package treebank to /root/nltk_data...
[nltk_data]    |   Package treebank is already up-to-date!
[nltk_data]    | Downloading package twitter_samples to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package twitter_samples is already up-to-date!
[nltk_data]    | Downloading package omw to /root/nltk_data...
[nltk_data]    |   Package omw is already up-to-date!
[nltk_data]    | Downloading package omw-1.4 to /root/nltk_data...
[nltk_data]    |   Package omw-1.4 is already up-to-date!
[nltk_data]    | Downloading package wordnet to /root/nltk_data...
[nltk_data]    |   Package wordnet is already up-to-date!
[nltk_data]    | Downloading package wordnet2021 to /root/nltk_data...
[nltk_data]    |   Package wordnet2021 is already up-to-date!
[nltk_data]    | Downloading package wordnet31 to /root/nltk_data...
[nltk_data]    |   Package wordnet31 is already up-to-date!
[nltk_data]    | Downloading package wordnet_ic to /root/nltk_data...
[nltk_data]    |   Package wordnet_ic is already up-to-date!
[nltk_data]    | Downloading package words to /root/nltk_data...
[nltk_data]    |   Package words is already up-to-date!
[nltk_data]    | Downloading package maxent_ne_chunker to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package maxent_ne_chunker is already up-to-date!
[nltk_data]    | Downloading package punkt to /root/nltk_data...
[nltk_data]    |   Package punkt is already up-to-date!
[nltk_data]    | Downloading package snowball_data to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package snowball_data is already up-to-date!
[nltk_data]    | Downloading package averaged_perceptron_tagger to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package averaged_perceptron_tagger is already up-
[nltk_data]    |       to-date!
[nltk_data]    |
[nltk_data]  Done downloading collection popular
True
```

```python
vocab = CountVectorizer(analyzer=text_process).fit(x)
print(len(vocab.vocabulary_))
r0 = x[0]
print(r0)
vocab0 = vocab.transform([r0])
print(vocab0)
"""
    Now the words in the review number 78 have been converted into a vector.
    The data that we can see is the transformed words.
    If we now get the feature's name - we can get the word back!
"""
print("Getting the words back:")
print(vocab.get_feature_names()[19648])
print(vocab.get_feature_names()[10643])
```

```
While EVERYTHING on the menu looks excellent, I had the white truffle scrambled e

Anyway, I can't wait to go back!
  (0, 292)      1
  (0, 1213)     1
  (0, 1811)     1
  (0, 3537)     1
  (0, 5139)     1
  (0, 5256)     2
  (0, 6275)     1
  (0, 8521)     1
  (0, 10646)    1
  (0, 10647)    1
  (0, 11128)    2
  (0, 11479)    1
  (0, 11779)    1
  (0, 12206)    2
  (0, 12221)    1
  (0, 12297)    1
  (0, 12386)    1
  (0, 12675)    1
  (0, 12689)    1
  (0, 13135)    1
  (0, 13186)    1
  (0, 14247)    1
  (0, 15385)    1
  (0, 16292)    1
  (0, 16412)    1
  :        :
  (0, 23318)    1
  (0, 23801)    1
  (0, 23902)    1
  (0, 23976)    1
  (0, 24080)    1
  (0, 24177)    1
  (0, 24544)    2
  (0, 24972)    2
  (0, 26383)    1
  (0, 26543)    1
  (0, 26978)    1
  (0, 27029)    1
  (0, 27068)    1
  (0, 28403)    1
  (0, 28735)    1
  (0, 29230)    1
  (0, 29313)    1
  (0, 29620)    1
  (0, 30135)    1
  (0, 30240)    1
  (0, 30471)    1
  (0, 30488)    1
  (0, 30672)    1
  (0, 30854)    1
  (0, 30900)    1
Getting the words back:
```

```
x = vocab.transform(x)
#Shape of the matrix:
print("Shape of the sparse matrix: ", x.shape)
#Non-zero occurences:
print("Non-Zero occurences: ",x.nnz)

# DENSITY OF THE MATRIX
density = (x.nnz/(x.shape[0]*x.shape[1]))*100
print("Density of the matrix = ",density)
```

```
Shape of the sparse matrix:  (5547, 31336)
Non-Zero occurences:  312457
Density of the matrix =  0.17975812697942373
```

```
# SPLITTING THE DATASET INTO TRAINING SET AND TESTING SET
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.2,random_state=101)


#random forest
from sklearn.ensemble import RandomForestClassifier
rmfr = RandomForestClassifier()
rmfr.fit(x_train,y_train)
predrmfr = rmfr.predict(x_test)
print("Confusion Matrix for Random Forest Classifier:")
print(confusion_matrix(y_test,predrmfr))
```

```
print("Score:",round(accuracy_score(y_test,predrmfr)*100,2))
print(Confusiofication FeportRandomForestclassif_report(y_test,predrmfr))
    [[ 30  29 103]
     [  3 100 189]
     [  1  22 633]]
    Score: 68.74
    Classification Report:                 precision    recall  f1-score   support

               1       0.88      0.19      0.31       162
               3       0.66      0.34      0.45       292
               5       0.68      0.96      0.80       656

        accuracy                           0.69      1110
       macro avg       0.74      0.50      0.52      1110
    weighted avg       0.71      0.69      0.64      1110
```

```
# Support Vector Machine
from sklearn.svm import SVC
svm = SVC(random_state=101)
svm.fit(x_train,y_train)
predsvm = svm.predict(x_test)
print("Confusion Matrix for Support Vector Machines:")
print(confusion_matrix(y_test,predsvm))
print("Score:",round(accuracy_score(y_test,predsvm)*100,2))
print("Classification Report:",classification_report(y_test,predsvm))
```

```
    Confusion Matrix for Support Vector Machines:
    [[ 31  23 108]
     [  5 122 165]
     [  1  19 636]]
    Score: 71.08
    Classification Report:                 precision    recall  f1-score   support

               1       0.84      0.19      0.31       162
               3       0.74      0.42      0.54       292
               5       0.70      0.97      0.81       656

        accuracy                           0.71      1110
       macro avg       0.76      0.53      0.55      1110
    weighted avg       0.73      0.71      0.67      1110
```

Double-click (or enter) to edit

```
# Decision Tree
from sklearn.tree import DecisionTreeClassifier
dt = DecisionTreeClassifier()
dt.fit(x_train,y_train)
preddt = dt.predict(x_test)
print("Confusion Matrix for Decision Tree:")
print(confusion_matrix(y_test,preddt))
print("Score:",round(accuracy_score(y_test,preddt)*100,2))
print("Classification Report:",classification_report(y_test,preddt))
```

```
    Confusion Matrix for Decision Tree:
    [[ 60  48  54]
     [ 37 139 116]
     [ 44 110 502]]
    Score: 63.15
    Classification Report:                 precision    recall  f1-score   support

               1       0.43      0.37      0.40       162
               3       0.47      0.48      0.47       292
               5       0.75      0.77      0.76       656

        accuracy                           0.63      1110
       macro avg       0.55      0.54      0.54      1110
    weighted avg       0.63      0.63      0.63      1110
```

Double-click (or enter) to edit

```
#Gradient Boosting Classifier
from sklearn.ensemble import GradientBoostingClassifier
"""# parameter evaluation
gbe = GradientBoostingClassifier(random_state=0)
parameters = {
    'learning_rate': [0.05, 0.1, 0.5],
    'max_features': [0.5, 1],
    'max_depth': [3, 4, 5]}
gridsearch=GridSearchCV(gbe,parameters,cv=100,scoring='roc_auc')
```

```
gridsearch.fit(x,y)
print(gridsearch.best_params_)
print(gridsearch.best_score_)"""
#Boosting
gbi = GradientBoostingClassifier(learning_rate=0.1,max_depth=5,max_features=0.5,random_state=999999)
gbi.fit(x_train,y_train)
predgbi = gbi.predict(x_test)
print("Confusion Matrix for Gradient Boosting Classifier:")
print(confusion_matrix(y_test,predgbi))
print("Score:",round(accuracy_score(y_test,predgbi)*100,2))
print("Classification Report:",classification_report(y_test,predgbi))
```

```
    Confusion Matrix for Gradient Boosting Classifier:
    [[ 61  33  68]
     [  9 139 144]
     [  4  33 619]]
    Score: 73.78
    Classification Report:               precision    recall  f1-score   support

               1       0.82      0.38      0.52       162
               3       0.68      0.48      0.56       292
               5       0.74      0.94      0.83       656

        accuracy                           0.74      1110
       macro avg       0.75      0.60      0.64      1110
    weighted avg       0.74      0.74      0.71      1110
```

```
# K Nearest Neighbour Algorithm
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=10)
knn.fit(x_train,y_train)
predknn = knn.predict(x_test)
print("Confusion Matrix for K Neighbors Classifier:")
print(confusion_matrix(y_test,predknn))
print("Score: ",round(accuracy_score(y_test,predknn)*100,2))
print("Classification Report:")
print(classification_report(y_test,predknn))
```

```
    Confusion Matrix for K Neighbors Classifier:
    [[ 12  10 140]
     [  3  33 256]
     [  8  12 636]]
    Score:  61.35
    Classification Report:
                  precision    recall  f1-score   support

               1       0.52      0.07      0.13       162
               3       0.60      0.11      0.19       292
               5       0.62      0.97      0.75       656

        accuracy                           0.61      1110
       macro avg       0.58      0.39      0.36      1110
    weighted avg       0.60      0.61      0.51      1110
```

```
# MULTILAYER PERCEPTRON CLASSIFIER
from sklearn.neural_network import MLPClassifier
mlp = MLPClassifier()
mlp.fit(x_train,y_train)
predmlp = mlp.predict(x_test)
print("Confusion Matrix for Multilayer Perceptron Classifier:")
print(confusion_matrix(y_test,predmlp))
print("Score:",round(accuracy_score(y_test,predmlp)*100,2))
print("Classification Report:")
print(classification_report(y_test,predmlp))
```

```
    Confusion Matrix for Multilayer Perceptron Classifier:
    [[ 95  36  31]
     [ 24 184  84]
     [ 12  63 581]]
    Score: 77.48
    Classification Report:
                  precision    recall  f1-score   support

               1       0.73      0.59      0.65       162
               3       0.65      0.63      0.64       292
               5       0.83      0.89      0.86       656

        accuracy                           0.77      1110
       macro avg       0.74      0.70      0.72      1110
    weighted avg       0.77      0.77      0.77      1110
```