

```
# utilities
import re
import numpy as np
import pandas as pd
# plotting
import seaborn as sns
from wordcloud import WordCloud
import matplotlib.pyplot as plt
# nltk
from nltk.stem import WordNetLemmatizer
# sklearn
from sklearn.svm import LinearSVC
from sklearn.naive_bayes import BernoulliNB
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import confusion_matrix, classification_report

# Importing the dataset
DATASET_COLUMNS=['target','ids','date','flag','user','text']
DATASET_ENCODING = "ISO-8859-1"
df = pd.read_csv('/content/training.1600000.processed.noemoticon.csv', encoding=DATASET_ENCODING)
df.sample(8)
```

	target	ids	date	flag	user	text
207247	0	1973433334	Sat May 30 11:08:17 PDT 2009	NO_QUERY	Reynolds_x	Want some warburtons
195764	0	1970565860	Sat May 30 04:22:37 PDT 2009	NO_QUERY	jesssicababesss	@mariedancerr next thursday :/ idk if i will s...
1111750	4	1972277915	Sat May 30 08:55:28 PDT 2009	NO_QUERY	JadeLittish	@saaaaaarah15 its cool
1527292	4	2177165765	Mon Jun 15 05:35:07 PDT 2009	NO_QUERY	katienaas	@SmashMe_EraseMe yes ma'am... It was a good mo...

```
df.head()
```

	target	ids	date	flag	user	text
			Mon Apr 06			@switchfoot
0	0	1467810369	22:19:45 PDT	NO QUERY	TheSpecialOne	http://twitpic.com/2y1zl

```
df.columns
```

```
Index(['target', 'ids', 'date', 'flag', 'user', 'text'], dtype='object')
```

```
2009
```

```
by ...
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1600000 entries, 0 to 1599999
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   target      1600000 non-null  int64
1   ids         1600000 non-null  int64
2   date        1600000 non-null  object
3   flag        1600000 non-null  object
4   user        1600000 non-null  object
5   text        1600000 non-null  object
dtypes: int64(2), object(4)
memory usage: 73.2+ MB
```

```
df.dtypes
```

```
target      int64
ids          int64
date         object
flag         object
user         object
text         object
dtype: object
```

```
np.sum(df.isnull().any(axis=1))
```

```
0
```

```
print('Count of columns in the data is: ', len(df.columns))
```

```
print('Count of rows in the data is: ', len(df))
```

```
Count of columns in the data is: 6
Count of rows in the data is: 1600000
```

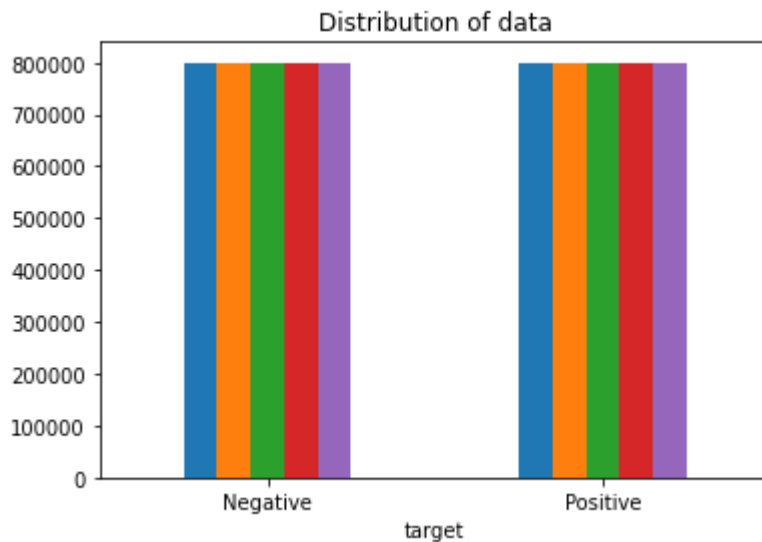
```
df['target'].unique()
```

```
array([0, 4])
```

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

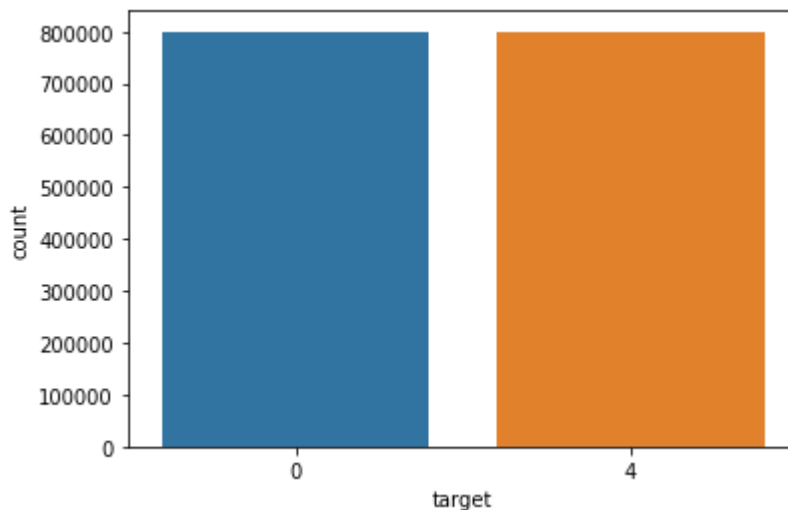
2

```
# Plotting the distribution for dataset.
ax = df.groupby('target').count().plot(kind='bar', title='Distribution of data', legend=False)
ax.set_xticklabels(['Negative', 'Positive'], rotation=0)
# Storing data in lists.
text, sentiment = list(df['text']), list(df['target'])
```



```
import seaborn as sns
sns.countplot(x='target', data=df)
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f3308e0fb50>



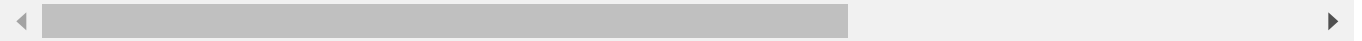
```
data=df[['text', 'target']]
```

```
data['target'] = data['target'].replace(4,1)
```

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.
Try using `.loc[row_indexer,col_indexer] = value` instead

See the caveats in the documentation: <https://pandas.pydata.org/pandas-docs/stable/user>
""Entry point for launching an IPython kernel.



```
data['target'].unique()
```

```
array([0, 1])
```

```
data_pos = data[data['target'] == 1]
data_neg = data[data['target'] == 0]
```

```
data_pos = data_pos.iloc[:int(20000)]
data_neg = data_neg.iloc[:int(20000)]
```

```
dataset = pd.concat([data_pos, data_neg])
```

```
dataset['text']=dataset['text'].str.lower()
dataset['text'].tail()
```

```
19995    not much time off this weekend, work trip to m...
19996                                one more day of holidays
19997    feeling so down right now .. i hate you damn h...
19998    geez,i hv to read the whole book of personalit...
19999    i threw my sign at donnie and he bent over to ...
Name: text, dtype: object
```

```
stopwordlist = ['a', 'about', 'above', 'after', 'again', 'ain', 'all', 'am', 'an',
                'and', 'any', 'are', 'as', 'at', 'be', 'because', 'been', 'before',
                'being', 'below', 'between', 'both', 'by', 'can', 'd', 'did', 'do',
                'does', 'doing', 'down', 'during', 'each', 'few', 'for', 'from',
                'further', 'had', 'has', 'have', 'having', 'he', 'her', 'here',
                'hers', 'herself', 'him', 'himself', 'his', 'how', 'i', 'if', 'in',
                'into', 'is', 'it', 'its', 'itself', 'just', 'll', 'm', 'ma',
                'me', 'more', 'most', 'my', 'myself', 'now', 'o', 'of', 'on', 'once',
                'only', 'or', 'other', 'our', 'ours', 'ourselves', 'out', 'own', 're', 's', 'same',
                't', 'than', 'that', 'thatll', 'the', 'their', 'theirs', 'them',
                'themselves', 'then', 'there', 'these', 'they', 'this', 'those',
                'through', 'to', 'too', 'under', 'until', 'up', 've', 'very', 'was',
                'we', 'were', 'what', 'when', 'where', 'which', 'while', 'who', 'whom',
                'why', 'will', 'with', 'won', 'y', 'you', 'youd', 'youll', 'youre',
                'youve', 'your', 'yours', 'yourself', 'yourselves']
```

```
STOPWORDS = set(stopwordlist)
def cleaning_stopwords(text):
```

```

return " ".join([word for word in str(text).split() if word not in STOPWORDS])
dataset['text'] = dataset['text'].apply(lambda text: cleaning_stopwords(text))
dataset['text'].head()

```

```

800000          love @health4uandpets u guys r best!!
800001    im meeting one besties tonight! cant wait!! - ...
800002    @darealsunisakim thanks twitter add, sunisa! g...
800003    sick really cheap hurts much eat real food plu...
800004          @lovesbrooklyn2 effect everyone
Name: text, dtype: object

```

```

import string
english_punctuations = string.punctuation
punctuations_list = english_punctuations
def cleaning_punctuations(text):
    translator = str.maketrans('', '', punctuations_list)
    return text.translate(translator)
dataset['text'] = dataset['text'].apply(lambda x: cleaning_punctuations(x))
dataset['text'].tail()

```

```

19995    not much time off weekend work trip malmi½ fr...
19996          one day holidays
19997          feeling right hate damn humprey
19998    geezi hv read whole book personality types emb...
19999    threw sign donnie bent over get but thingee ma...
Name: text, dtype: object

```

```

def cleaning_repeating_char(text):
    return re.sub(r'(. )1+', r'1', text)
dataset['text'] = dataset['text'].apply(lambda x: cleaning_repeating_char(x))
dataset['text'].tail()

```

```

19995    not much time off weekend work trip malmi½ fr...
19996          one day holidays
19997          feeling right hate damn humprey
19998    geezi hv read whole book personality types emb...
19999    threw sign donnie bent over get but thingee ma...
Name: text, dtype: object

```

```

def cleaning_URLs(data):
    return re.sub('((www.[^s]+)|(https?://[^\s]+))', ' ', data)
dataset['text'] = dataset['text'].apply(lambda x: cleaning_URLs(x))
dataset['text'].tail()

```

```

19995    not much time off weekend work trip malmi½ fr...
19996          one day holidays
19997          feeling right hate damn humprey
19998    geezi hv read whole book personality types emb...
19999    threw sign donnie bent over get but thingee ma...
Name: text, dtype: object

```

```
def cleaning_numbers(data):
    return re.sub('[0-9]+', '', data)
dataset['text'] = dataset['text'].apply(lambda x: cleaning_numbers(x))
dataset['text'].tail()
```

```
19995    not much time off weekend work trip malmi;½ fr...
19996                                one day holidays
19997                    feeling right  hate damn humprey
19998    geezi hv read whole book personality types emb...
19999    threw sign donnie bent over get but thingee ma...
Name: text, dtype: object
```

```
from nltk.tokenize import RegexpTokenizer
tokenizer = RegexpTokenizer(r'w+')
dataset['text'] = dataset['text'].apply(tokenizer.tokenize)
dataset['text'].head()
```

```
800000    []
800001    [w]
800002    [w, w, w]
800003    []
800004    []
Name: text, dtype: object
```

```
import nltk
st = nltk.PorterStemmer()
def stemming_on_text(data):
    text = [st.stem(word) for word in data]
    return data
dataset['text'] = dataset['text'].apply(lambda x: stemming_on_text(x))
dataset['text'].head()
```

```
800000    []
800001    [w]
800002    [w, w, w]
800003    []
800004    []
Name: text, dtype: object
```

```
nltk.download('wordnet')
nltk.download('omw-1.4')
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
True
```

```
lm = nltk.WordNetLemmatizer()
def lemmatizer_on_text(data):
    text = [lm.lemmatize(word) for word in data]
```

```
return data
dataset['text'] = dataset['text'].apply(lambda x: lemmatizer_on_text(x))
dataset['text'].head()

800000      []
800001      [w]
800002    [w, w, w]
800003      []
800004      []
Name: text, dtype: object
```

```
X=data.text
y=data.target
```

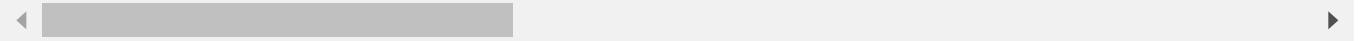
```
data_neg = data['text'][:800000]
plt.figure(figsize = (20,20))
wc = WordCloud(max_words = 1000 , width = 1600 , height = 800,
               collocations=False).generate(" ".join(data_neg))
plt.imshow(wc)
```



```

/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: F
warnings.warn(msg, category=FutureWarning)
No. of feature_words: 500000

```



```
X_train = vectoriser.transform(X_train)
```

```
X_test = vectoriser.transform(X_test)
```

```
def model_Evaluate(model):
```

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

```
    print(classification_report(y_test, y_pred))
```

```
    cf_matrix = confusion_matrix(y_test, y_pred)
```

```
    categories = ['Negative', 'Positive']
```

```
    group_names = ['True Neg', 'False Pos', 'False Neg', 'True Pos']
```

```
    group_percentages = ['{0:.2%}'.format(value) for value in cf_matrix.flatten() / np.sum(cf_m
```

```
    labels = [f'{v1}n{v2}' for v1, v2 in zip(group_names, group_percentages)]
```

```
    labels = np.asarray(labels).reshape(2,2)
```

```
    sns.heatmap(cf_matrix, annot = labels, cmap = 'Blues', fmt = '', xticklabels = categories, yti
```

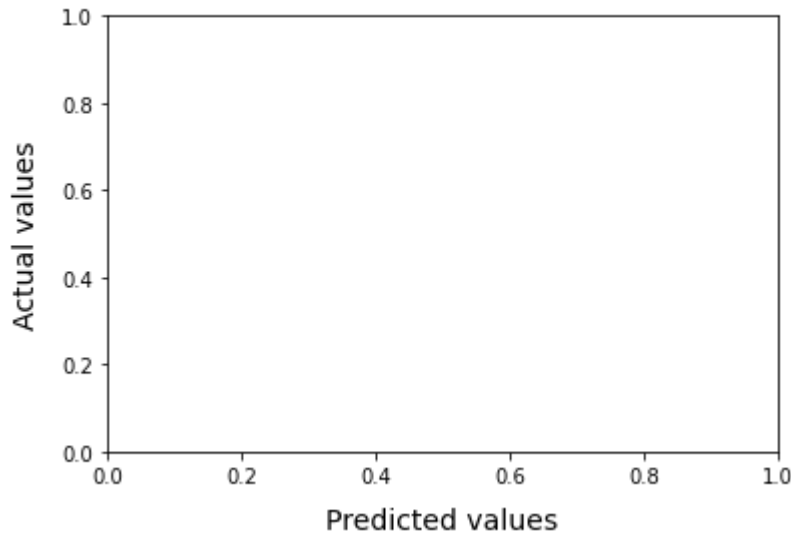
```
    plt.xlabel("Predicted values", fontdict = {'size':14}, labelpad = 10)
```

```
    plt.ylabel("Actual values", fontdict = {'size':14}, labelpad = 10)
```

```
    plt.title ("Confusion Matrix", fontdict = {'size':18}, pad = 20)
```

```
Text(0.5, 1.0, 'Confusion Matrix')
```

Confusion Matrix



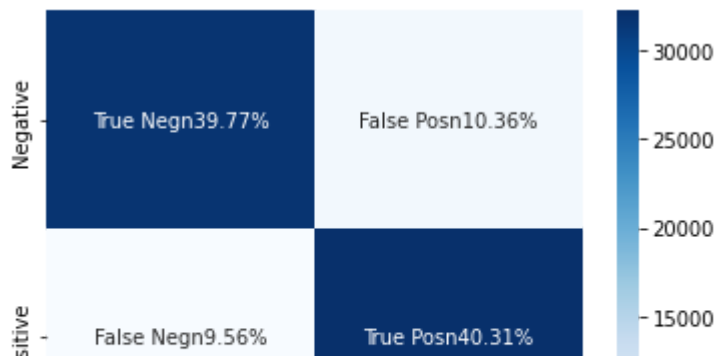
```
BNBmodel = BernoulliNB()
```

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

```
model_Evaluate(BNBmodel)
```

```
y_pred1 = BNBmodel.predict(X_test)
```

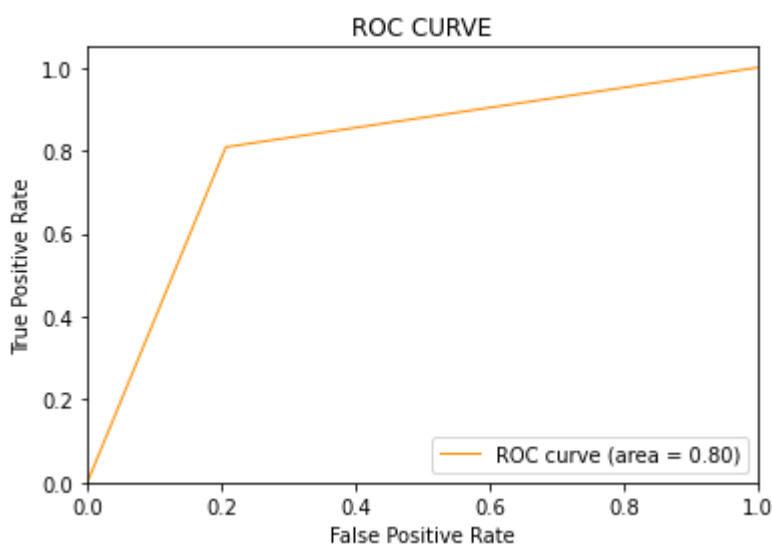
	precision	recall	f1-score	support
0	0.81	0.79	0.80	40100
1	0.80	0.81	0.80	39900
accuracy			0.80	80000
macro avg	0.80	0.80	0.80	80000
weighted avg	0.80	0.80	0.80	80000



```

from sklearn.metrics import roc_curve, auc
fpr, tpr, thresholds = roc_curve(y_test, y_pred1)
roc_auc = auc(fpr, tpr)
plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=1, label='ROC curve (area = %0.2f)' % roc_auc)
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC CURVE')
plt.legend(loc="lower right")
plt.show()

```



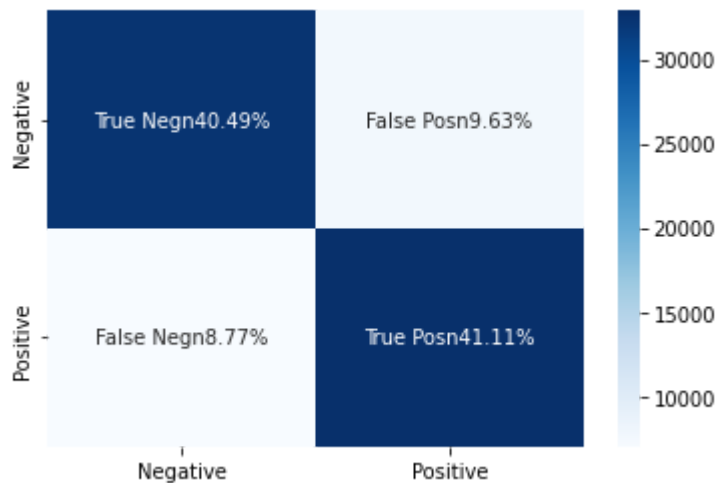
```

SVCmodel = LinearSVC()
SVCmodel.fit(X_train, y_train)

```

```
model_Evaluate(SVCmodel)
y_pred2 = SVCmodel.predict(X_test)
```

	precision	recall	f1-score	support
0	0.82	0.81	0.81	40100
1	0.81	0.82	0.82	39900
accuracy			0.82	80000
macro avg	0.82	0.82	0.82	80000
weighted avg	0.82	0.82	0.82	80000

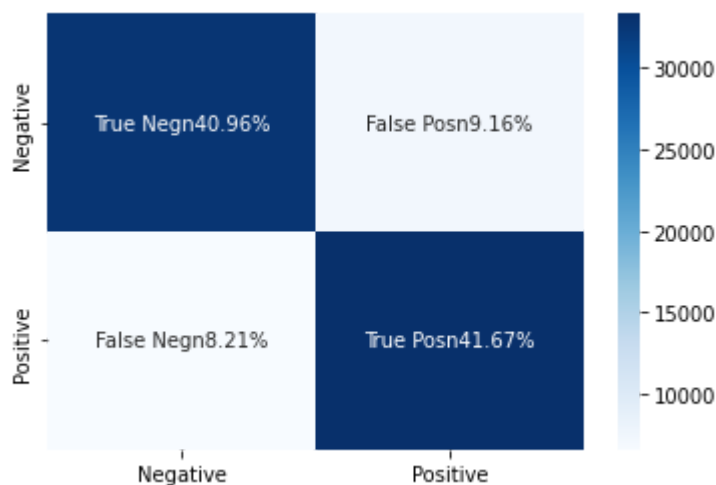


```
from sklearn.metrics import roc_curve, auc
fpr, tpr, thresholds = roc_curve(y_test, y_pred2)
roc_auc = auc(fpr, tpr)
plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=1, label='ROC curve (area = %0.2f)' % roc_auc)
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC CURVE')
plt.legend(loc="lower right")
plt.show()
```

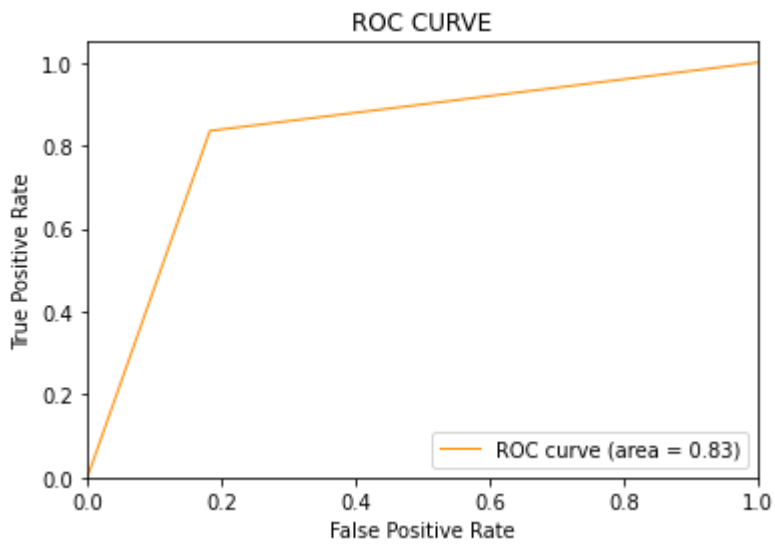
ROC CURVE

```
LRmodel = LogisticRegression(C = 2, max_iter = 1000, n_jobs=-1)
LRmodel.fit(X_train, y_train)
model_Evaluate(LRmodel)
y_pred3 = LRmodel.predict(X_test)
```

	precision	recall	f1-score	support
0	0.83	0.82	0.83	40100
1	0.82	0.84	0.83	39900
accuracy			0.83	80000
macro avg	0.83	0.83	0.83	80000
weighted avg	0.83	0.83	0.83	80000



```
from sklearn.metrics import roc_curve, auc
fpr, tpr, thresholds = roc_curve(y_test, y_pred3)
roc_auc = auc(fpr, tpr)
plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=1, label='ROC curve (area = %0.2f)' % roc_auc)
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC CURVE')
plt.legend(loc="lower right")
plt.show()
```



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

