

## ✓ Phishing URL Detection

The Internet has become an indispensable part of our life, However, It also has provided opportunities to anonymously perform malicious activities like Phishing. Phishers try to deceive their victims by social engineering or creating mockup websites to steal information such as account ID, username, password from individuals and organizations. Although many methods have been proposed to detect phishing websites, Phishers have evolved their methods to escape from these detection methods. One of the most successful methods for detecting these malicious activities is Machine Learning. This is because most Phishing attacks have some common characteristics which can be identified by machine learning methods.

The steps demonstrated in this notebook are:

1. Loading the data
2. Familiarizing with data & EDA
3. Visualizing the data
4. Splitting the data
5. Training the data
6. Comparison of Model
7. Conclusion

```
#importing required libraries

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
from sklearn import metrics
import warnings
warnings.filterwarnings('ignore')
```

### ✓ 1. Loading Data:

The dataset is borrowed from Kaggle, <https://www.kaggle.com/code/hasibur013/phishing-url-detection>.

A collection of website URLs for 11000+ websites. Each sample has 30 website parameters and a class label identifying it as a phishing website or not (1 or -1).

The overview of this dataset is, it has 11054 samples with 32 features. Download the dataset from the link provided.

```
#Loading data into dataframe
```

```
data = pd.read_csv("/content/phishing.csv")
data.head()
```

```
↗
```

	Index	UsingIP	LongURL	ShortURL	Symbol@	Redirecting//	PrefixSuffix-	SubDomains	HTTPS	DomainRegLen	...	UsingPopupWindow
0	0	1	1	1	1	1	-1	0	1	-1	...	1
1	1	1	0	1	1	1	-1	-1	-1	-1	...	1
2	2	1	0	1	1	1	-1	-1	-1	1	...	1
3	3	1	0	-1	1	1	-1	1	1	-1	...	-1
4	4	-1	0	-1	1	-1	-1	1	1	-1	...	1

5 rows × 32 columns

### ✓ 2. Familiarizing with Data & EDA:

In this step, few dataframe methods are used to look into the data and its features.

```
data.shape
```

```
↗ (11054, 32)
```

data.columns

```
Index(['Index', 'UsingIP', 'LongURL', 'ShortURL', 'Symbol@', 'Redirecting//',
      'PrefixSuffix-', 'SubDomains', 'HTTPS', 'DomainRegLen', 'Favicon',
      'NonStdPort', 'HTTPSDomainURL', 'RequestURL', 'AnchorURL',
      'LinksInScriptTags', 'ServerFormHandler', 'InfoEmail', 'AbnormalURL',
      'WebsiteForwarding', 'StatusBarCust', 'DisableRightClick',
      'UsingPopupWindow', 'IframeRedirection', 'AgeofDomain', 'DNSRecording',
      'WebsiteTraffic', 'PageRank', 'GoogleIndex', 'LinksPointingToPage',
      'StatsReport', 'class'],
      dtype='object')
```

#Information about the dataset

data.info()

Show hidden output

# nunique value in columns

data.nunique()

	0
Index	11054
UsingIP	2
LongURL	3
ShortURL	2
Symbol@	2
Redirecting//	2
PrefixSuffix-	2
SubDomains	3
HTTPS	3
DomainRegLen	2
Favicon	2
NonStdPort	2
HTTPSDomainURL	2
RequestURL	2
AnchorURL	3
LinksInScriptTags	3
ServerFormHandler	3
InfoEmail	2
AbnormalURL	2
WebsiteForwarding	2
StatusBarCust	2
DisableRightClick	2
UsingPopupWindow	2
IframeRedirection	2
AgeofDomain	2
DNSRecording	2
WebsiteTraffic	3
PageRank	2
GoogleIndex	2
LinksPointingToPage	3
StatsReport	2
class	2

```
#dropping index column

data = data.drop(['Index'],axis = 1)

#description of dataset

data.describe().T
```

 [Show hidden output](#)

data\_set.append(9 OBSERVATIONS:

1. There are 11054 instances and 31 features in dataset.
2. Out of which 30 are independent features where as 1 is dependent feature.
3. Each feature is in int datatype, so there is no need to use LabelEncoder.
4. There is no outlier present in dataset.
5. There is no missing value in dataset.

### ✓ 3. Visualizing the data:


Few plots and graphs are displayed to find how the data is distributed and the how features are related to each other.

```
plt.figure(figsize=(15,15))
correlation_matrix = data.corr() # Assign the result of data.corr() to correlation_matrix
sns.heatmap(correlation_matrix, annot=True)
plt.show()
```

 [Show hidden output](#)

```
# Filter correlations above a certain threshold
```

```
high_corr_features = correlation_matrix[(correlation_matrix > 0.8) | (correlation_matrix < -0.8)]
print(high_corr_features)
```

 [Show hidden output](#)

```
#pairplot for particular features
```

```
df = data[['PrefixSuffix-', 'SubDomains', 'HTTPS', 'AnchorURL', 'WebsiteTraffic', 'ShortURL', 'DomainRegLen', 'class']]
sns.pairplot(data = df,hue="class",corner=True);
```

 [Show hidden output](#)

```
# Phishing Count in pie chart
```

```
data['class'].value_counts().plot(kind='pie',autopct='%1.2f%%')
plt.title("Phishing Count")
plt.show()
```

 [Show hidden output](#)

### ✓ 4. Splitting the Data:

The data is split into train & test sets, 80-20 split.


```
# Splitting the dataset into dependant and independant fetature
```

```
X = data.drop(["class"],axis =1)
y = data["class"]
```

```
# Splitting the dataset into train and test sets: 80-20 split
```

```
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 = 42)
X_train.shape, y_train.shape, X_test.shape, y_test.shape
```

 ((8843, 30), (8843,), (2211, 30), (2211,))

## ✓ 5. Model Building & Training:

Supervised machine learning is one of the most commonly used and successful types of machine learning. Supervised learning is used whenever we want to predict a certain outcome/label from a given set of features, and we have examples of features-label pairs. We build a machine learning model from these features-label pairs, which comprise our training set. Our goal is to make accurate predictions for new, never-before-seen data.

There are two major types of supervised machine learning problems, called classification and regression. Our data set comes under regression problem, as the prediction of suicide rate is a continuous number, or a floating-point number in programming terms. The supervised machine learning models (regression) considered to train the dataset in this notebook are:

1. Logistic Regression
2. k-Nearest Neighbors
3. Support Vector Classifier
4. Naive Bayes
5. Decision Tree
6. Random Forest
7. Gradient Boosting

The metrics considered to evaluate the model performance are Accuracy & F1 score.

```
# Creating holders to store the model performance results
ML_Model = []
accuracy = []
f1_score = []
recall = []
precision = []

#function to call for storing the results
def storeResults(model, a,b,c,d):
    ML_Model.append(model)
    accuracy.append(round(a, 3))
    f1_score.append(round(b, 3))
    recall.append(round(c, 3))
    precision.append(round(d, 3))

from sklearn.model_selection import cross_val_score, StratifiedKFold
from sklearn.model_selection import GridSearchCV
import numpy as np
import pandas as pd

def perform_cross_validation(model, X, y, cv=5):

    skf = StratifiedKFold(n_splits=cv, shuffle=True, random_state=42)
    cv_scores = cross_val_score(model, X, y, cv=skf, scoring='accuracy')
    print(f"Cross-validation scores: {cv_scores}")
    print(f"Mean CV score: {cv_scores.mean():.4f} (+/- {cv_scores.std() * 2:.4f})")

    return cv_scores
```

### ✓ 5.1. Logistic Regression

Logistic regression predicts the output of a categorical dependent variable. Therefore the outcome must be a categorical or discrete value. Logistic Regression is much similar to the Linear Regression except that how they are used. Linear Regression is used for solving Regression problems, whereas Logistic regression is used for solving the classification problems.

```
from sklearn.linear_model import LogisticRegression
log = LogisticRegression(max_iter=1000, random_state=42)
print("Logistic Regression Cross-Validation:")
log_cv_scores = perform_cross_validation(log, X, y)
log.fit(X_train, y_train)
```

🔍 Logistic Regression Cross-Validation:  
Cross-validation scores: [0.92265943 0.92492085 0.9348711 0.92672999 0.92714932]  
Mean CV score: 0.9273 (+/- 0.0082)

▼ LogisticRegression ⓘ ?

LogisticRegression(max\_iter=1000, random\_state=42)

```
y_train_log = log.predict(X_train)
```

```

y_test_log = log.predict(X_test)


acc_train_log = metrics.accuracy_score(y_train,y_train_log)
acc_test_log = metrics.accuracy_score(y_test,y_test_log)
print("Logistic Regression : Accuracy on training Data: {:.3f}".format(acc_train_log))
print("Logistic Regression : Accuracy on test Data: {:.3f}".format(acc_test_log))
print()

f1_score_train_log = metrics.f1_score(y_train,y_train_log)
f1_score_test_log = metrics.f1_score(y_test,y_test_log)
print("Logistic Regression : f1_score on training Data: {:.3f}".format(f1_score_train_log))
print("Logistic Regression : f1_score on test Data: {:.3f}".format(f1_score_test_log))
print()

recall_score_train_log = metrics.recall_score(y_train,y_train_log)
recall_score_test_log = metrics.recall_score(y_test,y_test_log)
print("Logistic Regression : Recall on training Data: {:.3f}".format(recall_score_train_log))
print("Logistic Regression : Recall on test Data: {:.3f}".format(recall_score_test_log))
print()

precision_score_train_log = metrics.precision_score(y_train,y_train_log)
precision_score_test_log = metrics.precision_score(y_test,y_test_log)
print("Logistic Regression : precision on training Data: {:.3f}".format(precision_score_train_log))
print("Logistic Regression : precision on test Data: {:.3f}".format(precision_score_test_log))

```


 Logistic Regression : Accuracy on training Data: 0.927  
 Logistic Regression : Accuracy on test Data: 0.934

Logistic Regression : f1\_score on training Data: 0.935  
 Logistic Regression : f1\_score on test Data: 0.941

Logistic Regression : Recall on training Data: 0.943  
 Logistic Regression : Recall on test Data: 0.953

Logistic Regression : precision on training Data: 0.927  
 Logistic Regression : precision on test Data: 0.930

```
print(metrics.classification_report(y_test, y_test_log))
```



	precision	recall	f1-score	support
-1	0.94	0.91	0.92	976
1	0.93	0.95	0.94	1235
accuracy			0.93	2211
macro avg	0.93	0.93	0.93	2211
weighted avg	0.93	0.93	0.93	2211

```

storeResults('Logistic Regression',acc_test_log,f1_score_test_log,
            recall_score_train_log,precision_score_train_log)

```


## ✓ 5.2. K-Nearest Neighbors : Classifier



K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique. K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.

```

from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X_train,y_train)

```



KNeighborsClassifier
 


```

KNeighborsClassifier(n_neighbors=1)

```

```

y_train_knn = knn.predict(X_train)
y_test_knn = knn.predict(X_test)

```

```
#computing the accuracy,f1_score,Recall,precision of the model performance
```

```

acc_train_knn = metrics.accuracy_score(y_train,y_train_knn)
acc_test_knn = metrics.accuracy_score(y_test,y_test_knn)
print("K-Nearest Neighbors : Accuracy on training Data: {:.3f}".format(acc_train_knn))
print("K-Nearest Neighbors : Accuracy on test Data: {:.3f}".format(acc_test_knn))

```

```

print()

f1_score_train_knn = metrics.f1_score(y_train,y_train_knn)
f1_score_test_knn = metrics.f1_score(y_test,y_test_knn)
print("K-Nearest Neighbors : f1_score on training Data: {:.3f}".format(f1_score_train_knn))
print("K-Nearest Neighbors : f1_score on test Data: {:.3f}".format(f1_score_test_knn))
print()

recall_score_train_knn = metrics.recall_score(y_train,y_train_knn)
recall_score_test_knn = metrics.recall_score(y_test,y_test_knn)
print("K-Nearest Neighbors : Recall on training Data: {:.3f}".format(recall_score_train_knn))
print("Logistic Regression : Recall on test Data: {:.3f}".format(recall_score_test_knn))
print()

precision_score_train_knn = metrics.precision_score(y_train,y_train_knn)
precision_score_test_knn = metrics.precision_score(y_test,y_test_knn)
print("K-Nearest Neighbors : precision on training Data: {:.3f}".format(precision_score_train_knn))
print("K-Nearest Neighbors : precision on test Data: {:.3f}".format(precision_score_test_knn))

```

```

↗ K-Nearest Neighbors : Accuracy on training Data: 0.989
K-Nearest Neighbors : Accuracy on test Data: 0.956

```

```

K-Nearest Neighbors : f1_score on training Data: 0.990
K-Nearest Neighbors : f1_score on test Data: 0.961

```

```

K-Nearest Neighbors : Recall on training Data: 0.991
Logistic Regression : Recall on test Data: 0.962

```

```

K-Nearest Neighbors : precision on training Data: 0.989
K-Nearest Neighbors : precision on test Data: 0.960

```

```

print(metrics.classification_report(y_test, y_test_knn))

```

```

↗

```

	precision	recall	f1-score	support
-1	0.95	0.95	0.95	976
1	0.96	0.96	0.96	1235
accuracy			0.96	2211
macro avg	0.96	0.96	0.96	2211
weighted avg	0.96	0.96	0.96	2211

```

training_accuracy = []
test_accuracy = []

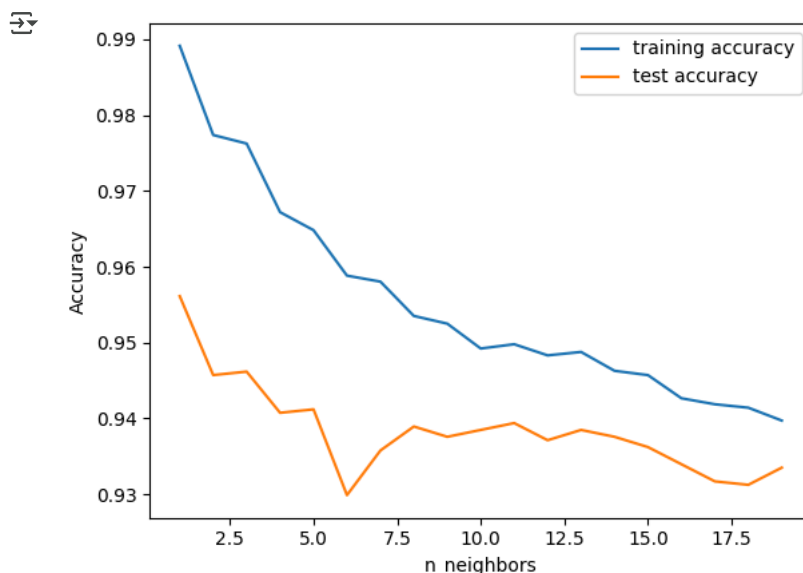
```

```

depth = range(1,20)
for n in depth:
    knn = KNeighborsClassifier(n_neighbors=n)

    knn.fit(X_train, y_train)
    training_accuracy.append(knn.score(X_train, y_train))
    test_accuracy.append(knn.score(X_test, y_test))
plt.plot(depth, training_accuracy, label="training accuracy")
plt.plot(depth, test_accuracy, label="test accuracy")
plt.ylabel("Accuracy")
plt.xlabel("n_neighbors")
plt.legend();

```



```
storeResults('K-Nearest Neighbors',acc_test_knn,f1_score_test_knn,
            recall_score_train_knn,precision_score_train_knn)
```

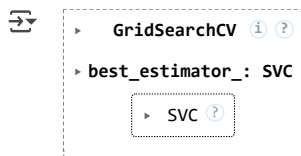
### 5.3. Support Vector Machine : Classifier

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future.

```
from sklearn.svm import SVC
from sklearn.model_selection import GridSearchCV

param_grid = {'gamma': [0.1], 'kernel': ['rbf', 'linear']}

svc = GridSearchCV(SVC(), param_grid)
svc.fit(X_train, y_train)
```



```
y_train_svc = svc.predict(X_train)
y_test_svc = svc.predict(X_test)
```

```
acc_train_svc = metrics.accuracy_score(y_train,y_train_svc)
acc_test_svc = metrics.accuracy_score(y_test,y_test_svc)
print("Support Vector Machine : Accuracy on training Data: {:.3f}".format(acc_train_svc))
print("Support Vector Machine : Accuracy on test Data: {:.3f}".format(acc_test_svc))
print()
```

```
f1_score_train_svc = metrics.f1_score(y_train,y_train_svc)
f1_score_test_svc = metrics.f1_score(y_test,y_test_svc)
print("Support Vector Machine : f1_score on training Data: {:.3f}".format(f1_score_train_svc))
print("Support Vector Machine : f1_score on test Data: {:.3f}".format(f1_score_test_svc))
print()
```

```
recall_score_train_svc = metrics.recall_score(y_train,y_train_svc)
recall_score_test_svc = metrics.recall_score(y_test,y_test_svc)
print("Support Vector Machine : Recall on training Data: {:.3f}".format(recall_score_train_svc))
print("Support Vector Machine : Recall on test Data: {:.3f}".format(recall_score_test_svc))
print()
```

```
precision_score_train_svc = metrics.precision_score(y_train,y_train_svc)
precision_score_test_svc = metrics.precision_score(y_test,y_test_svc)
print("Support Vector Machine : precision on training Data: {:.3f}".format(precision_score_train_svc))
print("Support Vector Machine : precision on test Data: {:.3f}".format(precision_score_test_svc))
```

```
Support Vector Machine : Accuracy on training Data: 0.969
Support Vector Machine : Accuracy on test Data: 0.964
```

```
Support Vector Machine : f1_score on training Data: 0.973
Support Vector Machine : f1_score on test Data: 0.968
```

```
Support Vector Machine : Recall on training Data: 0.980
Support Vector Machine : Recall on test Data: 0.980
```

```
Support Vector Machine : precision on training Data: 0.965
Support Vector Machine : precision on test Data: 0.957
```

```
print(metrics.classification_report(y_test, y_test_svc))
```

	precision	recall	f1-score	support
-1	0.97	0.94	0.96	976
1	0.96	0.98	0.97	1235
accuracy			0.96	2211
macro avg	0.97	0.96	0.96	2211
weighted avg	0.96	0.96	0.96	2211

```
storeResults('Support Vector Machine',acc_test_svc,f1_score_test_svc,
            recall_score_train_svc,precision_score_train_svc)
```

## ✓ 5.4. Naive Bayes : Classifier

Naïve Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems. It is mainly used in text, image classification that includes a high-dimensional training dataset. Naïve Bayes Classifier is one of the simple and most effective Classification algorithms which helps in building the fast machine learning models that can make quick predictions.

```
from sklearn.naive_bayes import GaussianNB
from sklearn.pipeline import Pipeline
nb= GaussianNB()
nb.fit(X_train,y_train)
```



```
#predicting the target value from the model for the samples
y_train_nb = nb.predict(X_train)
y_test_nb = nb.predict(X_test)
```

```
#computing the accuracy, f1_score, Recall, precision of the model performance
```

```
acc_train_nb = metrics.accuracy_score(y_train,y_train_nb)
acc_test_nb = metrics.accuracy_score(y_test,y_test_nb)
print("Naive Bayes Classifier : Accuracy on training Data: {:.3f}".format(acc_train_nb))
print("Naive Bayes Classifier : Accuracy on test Data: {:.3f}".format(acc_test_nb))
print()
```

```
f1_score_train_nb = metrics.f1_score(y_train,y_train_nb)
f1_score_test_nb = metrics.f1_score(y_test,y_test_nb)
print("Naive Bayes Classifier : f1_score on training Data: {:.3f}".format(f1_score_train_nb))
print("Naive Bayes Classifier : f1_score on test Data: {:.3f}".format(f1_score_test_nb))
print()
```

```
recall_score_train_nb = metrics.recall_score(y_train,y_train_nb)
recall_score_test_nb = metrics.recall_score(y_test,y_test_nb)
print("Naive Bayes Classifier : Recall on training Data: {:.3f}".format(recall_score_train_nb))
print("Naive Bayes Classifier : Recall on test Data: {:.3f}".format(recall_score_test_nb))
print()
```

```
precision_score_train_nb = metrics.precision_score(y_train,y_train_nb)
precision_score_test_nb = metrics.precision_score(y_test,y_test_nb)
print("Naive Bayes Classifier : precision on training Data: {:.3f}".format(precision_score_train_nb))
print("Naive Bayes Classifier : precision on test Data: {:.3f}".format(precision_score_test_nb))
```

```
Naive Bayes Classifier : Accuracy on training Data: 0.605
Naive Bayes Classifier : Accuracy on test Data: 0.605
```

```
Naive Bayes Classifier : f1_score on training Data: 0.451
Naive Bayes Classifier : f1_score on test Data: 0.454
```

```
Naive Bayes Classifier : Recall on training Data: 0.292
Naive Bayes Classifier : Recall on test Data: 0.294
```

```
Naive Bayes Classifier : precision on training Data: 0.997
Naive Bayes Classifier : precision on test Data: 0.995
```

```
#computing the classification report of the model
```

```
print(metrics.classification_report(y_test, y_test_svc))
```

```
precision    recall  f1-score   support

-1          0.97      0.94      0.96         976
1           0.96      0.98      0.97        1235

accuracy          0.96         2211
macro avg         0.97         0.96         0.96         2211
weighted avg      0.96         0.96         0.96         2211
```

```
#storing the results. The below mentioned order of parameter passing is important.
```



```
storeResults('Naive Bayes Classifier',acc_test_nb,f1_score_test_nb,
            recall_score_train_nb,precision_score_train_nb)
```

## 5.5. Decision Trees : Classifier

Decision Tree is a Supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome.

```
from sklearn.tree import DecisionTreeClassifier
```

```
tree = DecisionTreeClassifier(max_depth=30)
```

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



```
DecisionTreeClassifier(max_depth=30)
```

```
y_train_tree = tree.predict(X_train)
y_test_tree = tree.predict(X_test)
```

```
acc_train_tree = metrics.accuracy_score(y_train,y_train_tree)
acc_test_tree = metrics.accuracy_score(y_test,y_test_tree)
print("Decision Tree : Accuracy on training Data: {:.3f}".format(acc_train_tree))
print("Decision Tree : Accuracy on test Data: {:.3f}".format(acc_test_tree))
print()
```

```
f1_score_train_tree = metrics.f1_score(y_train,y_train_tree)
f1_score_test_tree = metrics.f1_score(y_test,y_test_tree)
print("Decision Tree : f1_score on training Data: {:.3f}".format(f1_score_train_tree))
print("Decision Tree : f1_score on test Data: {:.3f}".format(f1_score_test_tree))
print()
```

```
recall_score_train_tree = metrics.recall_score(y_train,y_train_tree)
recall_score_test_tree = metrics.recall_score(y_test,y_test_tree)
print("Decision Tree : Recall on training Data: {:.3f}".format(recall_score_train_tree))
print("Decision Tree : Recall on test Data: {:.3f}".format(recall_score_test_tree))
print()
```

```
precision_score_train_tree = metrics.precision_score(y_train,y_train_tree)
precision_score_test_tree = metrics.precision_score(y_test,y_test_tree)
print("Decision Tree : precision on training Data: {:.3f}".format(precision_score_train_tree))
print("Decision Tree : precision on test Data: {:.3f}".format(precision_score_test_tree))
```

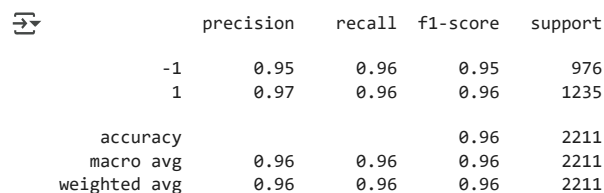
```
Decision Tree : Accuracy on training Data: 0.991
Decision Tree : Accuracy on test Data: 0.959
```

```
Decision Tree : f1_score on training Data: 0.992
Decision Tree : f1_score on test Data: 0.963
```

```
Decision Tree : Recall on training Data: 0.991
Decision Tree : Recall on test Data: 0.961
```

```
Decision Tree : precision on training Data: 0.993
Decision Tree : precision on test Data: 0.966
```

```
print(metrics.classification_report(y_test, y_test_tree))
```



```
precision    recall  f1-score   support

-1          0.95      0.96      0.95        976
 1          0.97      0.96      0.96       1235

 accuracy          0.96      0.96      0.96       2211
 macro avg          0.96      0.96      0.96       2211
 weighted avg        0.96      0.96      0.96       2211
```

```
training_accuracy = []
test_accuracy = []
```

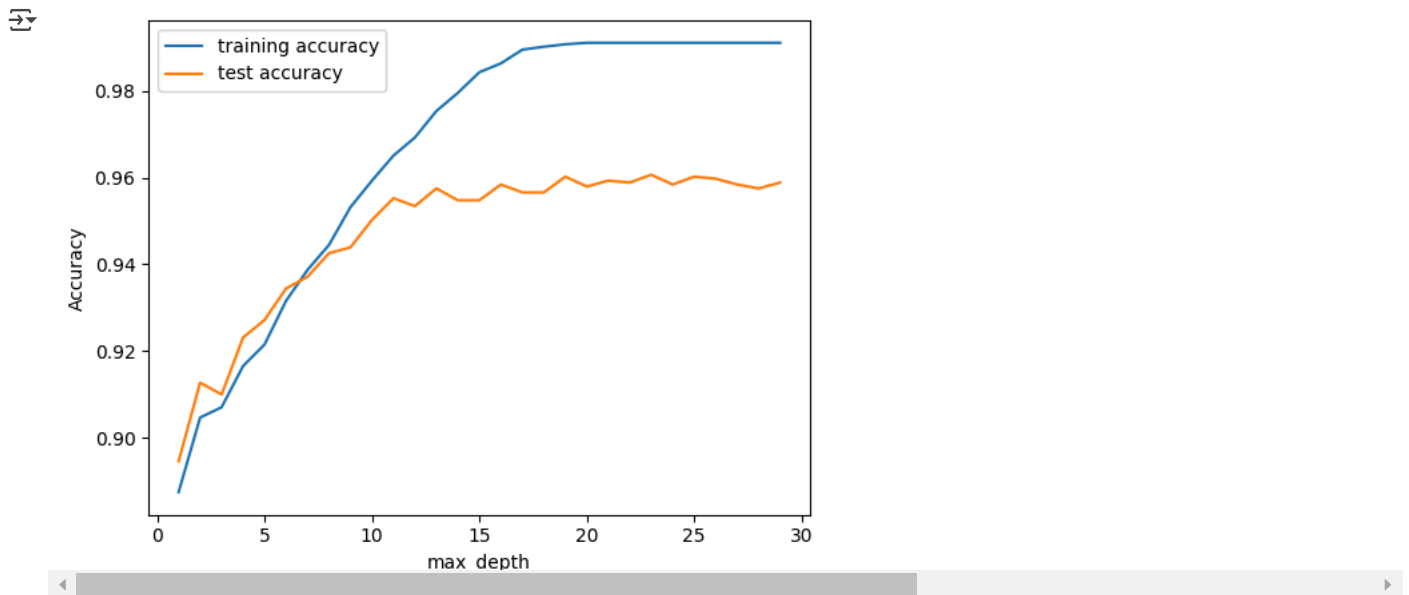
```
depth = range(1,30)
```

```

for n in depth:
    tree_test = DecisionTreeClassifier(max_depth=n)
    tree_test.fit(X_train, y_train)
    training_accuracy.append(tree_test.score(X_train, y_train))
    test_accuracy.append(tree_test.score(X_test, y_test))

plt.plot(depth, training_accuracy, label="training accuracy")
plt.plot(depth, test_accuracy, label="test accuracy")
plt.ylabel("Accuracy")
plt.xlabel("max_depth")
plt.legend();

```



```

storeResults('Decision Tree',acc_test_tree,f1_score_test_tree,
            recall_score_train_tree,precision_score_train_tree)

```

## 5.6. Random Forest : Classifier

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

```

from sklearn.ensemble import RandomForestClassifier
forest = RandomForestClassifier(n_estimators=20)
forest.fit(X_train,y_train)

```

```

RandomForestClassifier
RandomForestClassifier(n_estimators=20)

```

```

y_train_forest = forest.predict(X_train)
y_test_forest = forest.predict(X_test)

```

```

acc_train_forest = metrics.accuracy_score(y_train,y_train_forest)
acc_test_forest = metrics.accuracy_score(y_test,y_test_forest)
print("Random Forest : Accuracy on training Data: {:.3f}".format(acc_train_forest))
print("Random Forest : Accuracy on test Data: {:.3f}".format(acc_test_forest))
print()

```

```

f1_score_train_forest = metrics.f1_score(y_train,y_train_forest)
f1_score_test_forest = metrics.f1_score(y_test,y_test_forest)
print("Random Forest : f1_score on training Data: {:.3f}".format(f1_score_train_forest))
print("Random Forest : f1_score on test Data: {:.3f}".format(f1_score_test_forest))
print()

```

```

recall_score_train_forest = metrics.recall_score(y_train,y_train_forest)
recall_score_test_forest = metrics.recall_score(y_test,y_test_forest)
print("Random Forest : Recall on training Data: {:.3f}".format(recall_score_train_forest))
print("Random Forest : Recall on test Data: {:.3f}".format(recall_score_test_forest))
print()

```

```
precision_score_train_forest = metrics.precision_score(y_train,y_train_forest)
precision_score_test_forest = metrics.precision_score(y_test,y_test_forest)
print("Random Forest : precision on training Data: {:.3f}".format(precision_score_train_forest))
print("Random Forest : precision on test Data: {:.3f}".format(precision_score_test_forest))
```

```
Random Forest : Accuracy on training Data: 0.991
Random Forest : Accuracy on test Data: 0.968
```

```
Random Forest : f1_score on training Data: 0.992
Random Forest : f1_score on test Data: 0.972
```

```
Random Forest : Recall on training Data: 0.993
Random Forest : Recall on test Data: 0.977
```

```
Random Forest : precision on training Data: 0.990
Random Forest : precision on test Data: 0.966
```

```
print(metrics.classification_report(y_test, y_test_forest))
```

```

              precision    recall  f1-score   support

     -1         0.97         0.96         0.96         976
         1         0.97         0.98         0.97        1235

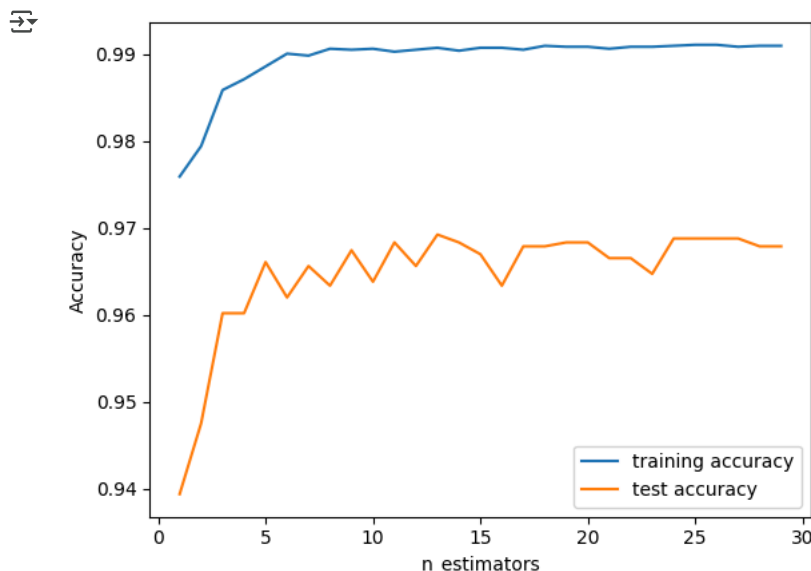
 accuracy                   0.97         2211
 macro avg          0.97         0.97         0.97         2211
 weighted avg       0.97         0.97         0.97         2211
```

```
training_accuracy = []
test_accuracy = []
```

```
depth = range(1,30)
for n in depth:
    forest_test = RandomForestClassifier(n_estimators=n)
```

```
    forest_test.fit(X_train, y_train)
    training_accuracy.append(forest_test.score(X_train, y_train))
    test_accuracy.append(forest_test.score(X_test, y_test))
```

```
plt.figure(figsize=None)
plt.plot(depth, training_accuracy, label="training accuracy")
plt.plot(depth, test_accuracy, label="test accuracy")
plt.ylabel("Accuracy")
plt.xlabel("n_estimators")
plt.legend();
```



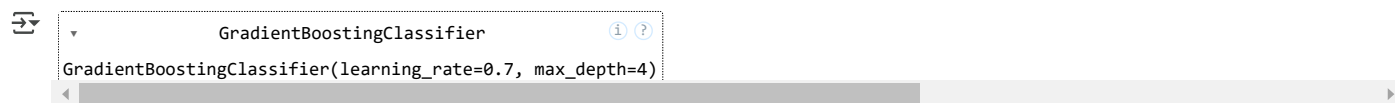
```
storeResults('Random Forest',acc_test_forest,f1_score_test_forest,
            recall_score_train_forest,precision_score_train_forest)
```

## 5.7.Gradient Boosting Classifier

Gradient boosting classifiers are a group of machine learning algorithms that combine many weak learning models together to create a strong predictive model. Decision trees are usually used when doing gradient boosting. Boosting algorithms play a crucial role in dealing with

bias variance trade-off. Unlike bagging algorithms, which only controls for high variance in a model, boosting controls both the aspects (bias & variance), and is considered to be more effective.

```
from sklearn.ensemble import GradientBoostingClassifier
gbc = GradientBoostingClassifier(max_depth=4, learning_rate=0.7)
gbc.fit(X_train, y_train)
```



```
y_train_gbc = gbc.predict(X_train)
y_test_gbc = gbc.predict(X_test)
```

```
acc_train_gbc = metrics.accuracy_score(y_train, y_train_gbc)
acc_test_gbc = metrics.accuracy_score(y_test, y_test_gbc)
print("Gradient Boosting Classifier : Accuracy on training Data: {:.3f}".format(acc_train_gbc))
print("Gradient Boosting Classifier : Accuracy on test Data: {:.3f}".format(acc_test_gbc))
print()

f1_score_train_gbc = metrics.f1_score(y_train, y_train_gbc)
f1_score_test_gbc = metrics.f1_score(y_test, y_test_gbc)
print("Gradient Boosting Classifier : f1_score on training Data: {:.3f}".format(f1_score_train_gbc))
print("Gradient Boosting Classifier : f1_score on test Data: {:.3f}".format(f1_score_test_gbc))
print()

recall_score_train_gbc = metrics.recall_score(y_train, y_train_gbc)
recall_score_test_gbc = metrics.recall_score(y_test, y_test_gbc)
print("Gradient Boosting Classifier : Recall on training Data: {:.3f}".format(recall_score_train_gbc))
print("Gradient Boosting Classifier : Recall on test Data: {:.3f}".format(recall_score_test_gbc))
print()

precision_score_train_gbc = metrics.precision_score(y_train, y_train_gbc)
precision_score_test_gbc = metrics.precision_score(y_test, y_test_gbc)
print("Gradient Boosting Classifier : precision on training Data: {:.3f}".format(precision_score_train_gbc))
print("Gradient Boosting Classifier : precision on test Data: {:.3f}".format(precision_score_test_gbc))
```

```
↗ Gradient Boosting Classifier : Accuracy on training Data: 0.989
  Gradient Boosting Classifier : Accuracy on test Data: 0.974

  Gradient Boosting Classifier : f1_score on training Data: 0.990
  Gradient Boosting Classifier : f1_score on test Data: 0.977

  Gradient Boosting Classifier : Recall on training Data: 0.994
  Gradient Boosting Classifier : Recall on test Data: 0.989

  Gradient Boosting Classifier : precision on training Data: 0.986
  Gradient Boosting Classifier : precision on test Data: 0.966
```

```
print(metrics.classification_report(y_test, y_test_gbc))
```

```
↗
```

	precision	recall	f1-score	support
-1	0.99	0.96	0.97	976
1	0.97	0.99	0.98	1235
accuracy			0.97	2211
macro avg	0.98	0.97	0.97	2211
weighted avg	0.97	0.97	0.97	2211

```
training_accuracy = []
test_accuracy = []

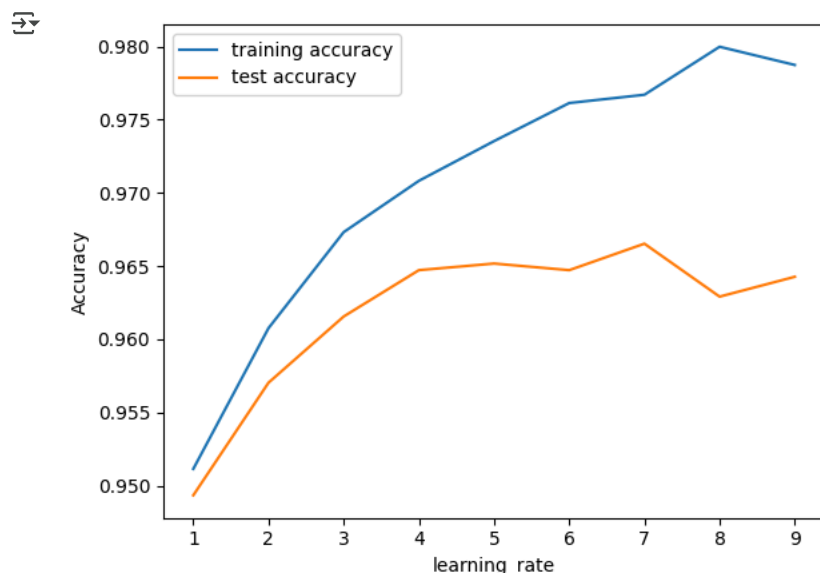
depth = range(1,10)
for n in depth:
    forest_test = GradientBoostingClassifier(learning_rate = n*0.1)

    forest_test.fit(X_train, y_train)

    training_accuracy.append(forest_test.score(X_train, y_train))

    test_accuracy.append(forest_test.score(X_test, y_test))
plt.figure(figsize=None)
plt.plot(depth, training_accuracy, label="training accuracy")
plt.plot(depth, test_accuracy, label="test accuracy")
```

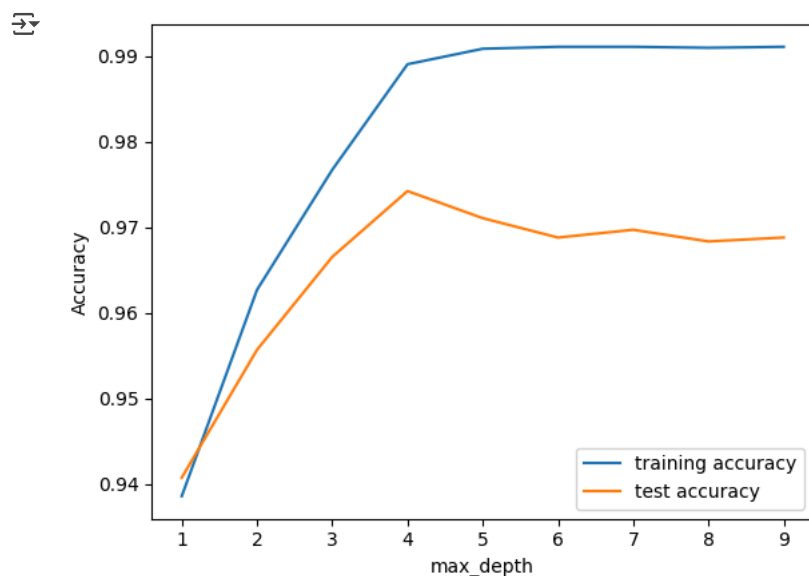
```
plt.ylabel("Accuracy")
plt.xlabel("learning_rate")
plt.legend();
```



```
training_accuracy = []
test_accuracy = []

depth = range(1,10,1)
for n in depth:
    forest_test = GradientBoostingClassifier(max_depth=n,learning_rate = 0.7)

    forest_test.fit(X_train, y_train)
    training_accuracy.append(forest_test.score(X_train, y_train))
    test_accuracy.append(forest_test.score(X_test, y_test))
plt.figure(figsize=None)
plt.plot(depth, training_accuracy, label="training accuracy")
plt.plot(depth, test_accuracy, label="test accuracy")
plt.ylabel("Accuracy")
plt.xlabel("max_depth")
plt.legend();
```



```
#storing the results. The below mentioned order of parameter passing is important.
storeResults('Gradient Boosting Classifier',acc_test_gbc,f1_score_test_gbc,
            recall_score_train_gbc,precision_score_train_gbc)
```

## ✓ 5.8. XGBoost Classifier

XGBoost is an implementation of gradient boosted decision trees designed for speed and performance that is dominative competitive machine learning. In this post you will discover how you can install and create your first XGBoost model in Python





	ML Model	Accuracy	f1_score	Recall	Precision
0	Logistic Regression	0.934	0.941	0.943	0.927
1	K-Nearest Neighbors	0.956	0.961	0.991	0.989
2	Support Vector Machine	0.964	0.968	0.980	0.965
3	Naive Bayes Classifier	0.605	0.454	0.292	0.997
4	Decision Tree	0.959	0.963	0.991	0.993
5	Random Forest	0.968	0.972	0.993	0.990
6	Gradient Boosting Classifier	0.974	0.977	0.994	0.986
7	XGBoost Classifier	0.549	0.544	0.553	0.548

```
#Sorting the datafram on accuracy
sorted_result=result.sort_values(by=['Accuracy', 'f1_score'],ascending=False).reset_index(drop=True)
```

```
# dispalying total result
sorted_result
```



	ML Model	Accuracy	f1_score	Recall	Precision
0	Gradient Boosting Classifier	0.974	0.977	0.994	0.986
1	Random Forest	0.968	0.972	0.993	0.990
2	Support Vector Machine	0.964	0.968	0.980	0.965
3	Decision Tree	0.959	0.963	0.991	0.993
4	K-Nearest Neighbors	0.956	0.961	0.991	0.989
5	Logistic Regression	0.934	0.941	0.943	0.927
6	Naive Bayes Classifier	0.605	0.454	0.292	0.997
7	XGBoost Classifier	0.549	0.544	0.553	0.548

## ✓ Storing Best Model

```
# XGBoost Classifier Model
from xgboost import XGBClassifier

# instantiate the model
gbc = GradientBoostingClassifier(max_depth=4,learning_rate=0.7)

# fit the model
gbc.fit(X_train,y_train)
```



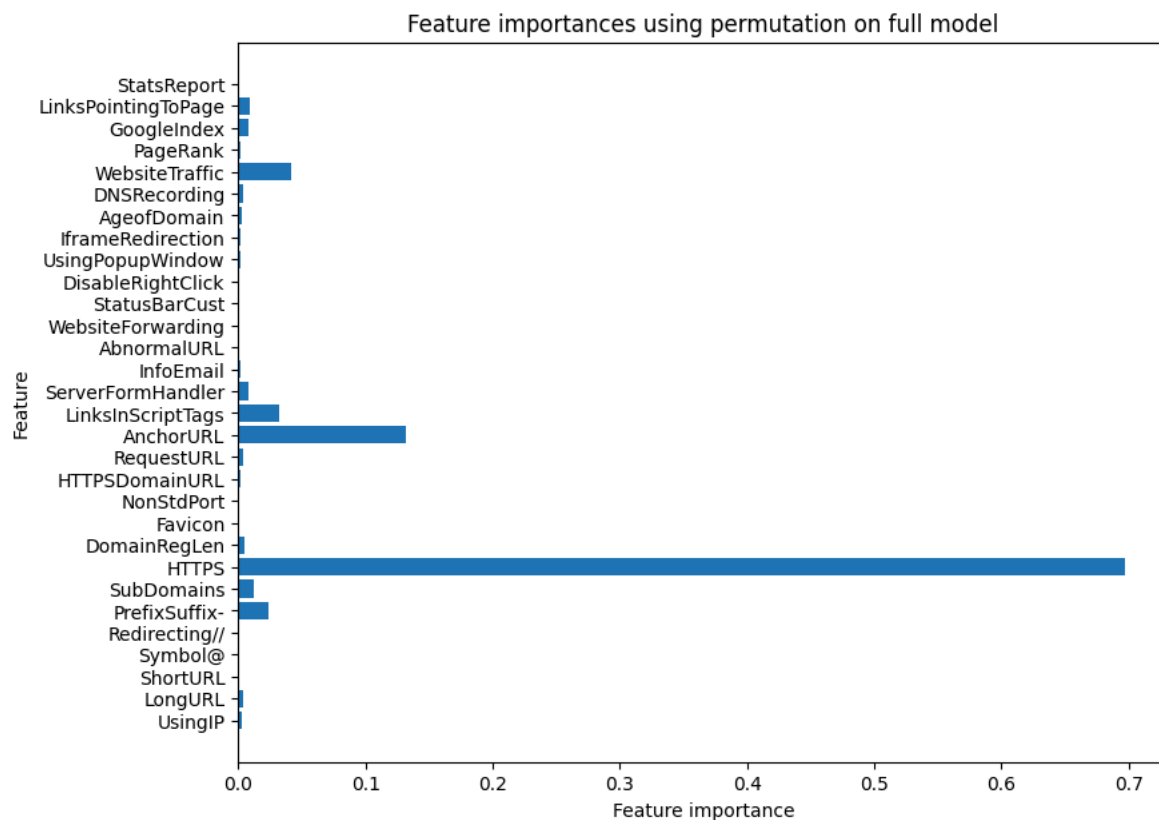
```
GradientBoostingClassifier
GradientBoostingClassifier(learning_rate=0.7, max_depth=4)
```

```
import os
import pickle

# Create the 'pickle' directory if it doesn't exist
os.makedirs('pickle', exist_ok=True)

# dump information to that file
pickle.dump(gbc, open('pickle/model.pkl', 'wb'))

plt.figure(figsize=(9,7))
n_features = X_train.shape[1]
plt.barh(range(n_features), gbc.feature_importances_, align='center')
plt.yticks(np.arange(n_features), X_train.columns)
plt.title("Feature importances using permutation on full model")
plt.xlabel("Feature importance")
plt.ylabel("Feature")
plt.show()
```



## 7. Conclusion

1. The final take away from this project is to explore various machine learning models, perform Exploratory Data Analysis on phishing dataset and understanding their features.
2. Creating this notebook helped me to learn a lot about the features affecting the models to detect whether URL is safe or not, also I came to know how to tune model and how they affect the model performance.
3. The final conclusion on the Phishing dataset is that some features like "HTTPS", "AnchorURL", "WebsiteTraffic" have more importance to classify if URL is phishing URL or not.