```python
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/eswarchandt/phishing-website-detector .

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.

```python
#Loading data into dataframe
import os, types
import pandas as pd
from botocore.client import Config
import ibm_boto3

def __iter__(self): return 0

# @hidden_cell
# The following code accesses a file in your IBM Cloud Object Storage. It includes your credentials.
# You might want to remove those credentials before you share the notebook.
cos_client = ibm_boto3.client(service_name='s3',
    ibm_api_key_id='OhdKjFpiXr7DU5QWMb68AMArQw7XueIuEhpVQqmvTHzY',
    ibm_auth_endpoint="https://iam.cloud.ibm.com/oidc/token",
    config=Config(signature_version='oauth'),
    endpoint_url='https://s3.private.eu.cloud-object-storage.appdomain.cloud')

bucket = 'webphishingdeployment-donotdelete-pr-sxj3w6vd0aq0yn'
object_key = 'Phishing.csv'

body = cos_client.get_object(Bucket=bucket,Key=object_key)['Body']
# add missing __iter__ method, so pandas accepts body as file-like object
if not hasattr(body, "__iter__"): body.__iter__ = types.MethodType( __iter__, body )

data = pd.read_csv(body)
```

# 2. Familiarizing with Data & EDA:

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

In [3]:
```python
#Shape of dataframe

data.shape
```

Out[3]: (11054, 32)

In [4]:
```python
#Listing the features of the dataset

data.columns
```

Out[4]: 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')

In [5]:

```python
#Information about the dataset

data.info()
```

```
RangeIndex: 11054 entries, 0 to 11053
Data columns (total 32 columns):
 #   Column             Non-Null Count  Dtype
---  ------             --------------  -----
 0   Index              11054 non-null  int64
 1   UsingIP            11054 non-null  int64
 2   LongURL            11054 non-null  int64
 3   ShortURL           11054 non-null  int64
 4   Symbol@            11054 non-null  int64
 5   Redirecting//      11054 non-null  int64
 6   PrefixSuffix-      11054 non-null  int64
 7   SubDomains         11054 non-null  int64
 8   HTTPS              11054 non-null  int64
 9   DomainRegLen       11054 non-null  int64
 10  Favicon            11054 non-null  int64
 11  NonStdPort         11054 non-null  int64
 12  HTTPSDomainURL     11054 non-null  int64
 13  RequestURL         11054 non-null  int64
 14  AnchorURL          11054 non-null  int64
 15  LinksInScriptTags  11054 non-null  int64
 16  ServerFormHandler  11054 non-null  int64
 17  InfoEmail          11054 non-null  int64
 18  AbnormalURL        11054 non-null  int64
 19  WebsiteForwarding  11054 non-null  int64
 20  StatusBarCust      11054 non-null  int64
 21  DisableRightClick  11054 non-null  int64
 22  UsingPopupWindow   11054 non-null  int64
 23  IframeRedirection  11054 non-null  int64
 24  AgeofDomain        11054 non-null  int64
 25  DNSRecording       11054 non-null  int64
 26  WebsiteTraffic     11054 non-null  int64
 27  PageRank           11054 non-null  int64
 28  GoogleIndex        11054 non-null  int64
 29  LinksPointingToPage 11054 non-null int64
 30  StatsReport        11054 non-null  int64
 31  class              11054 non-null  int64
dtypes: int64(32)
memory usage: 2.7 MB
```
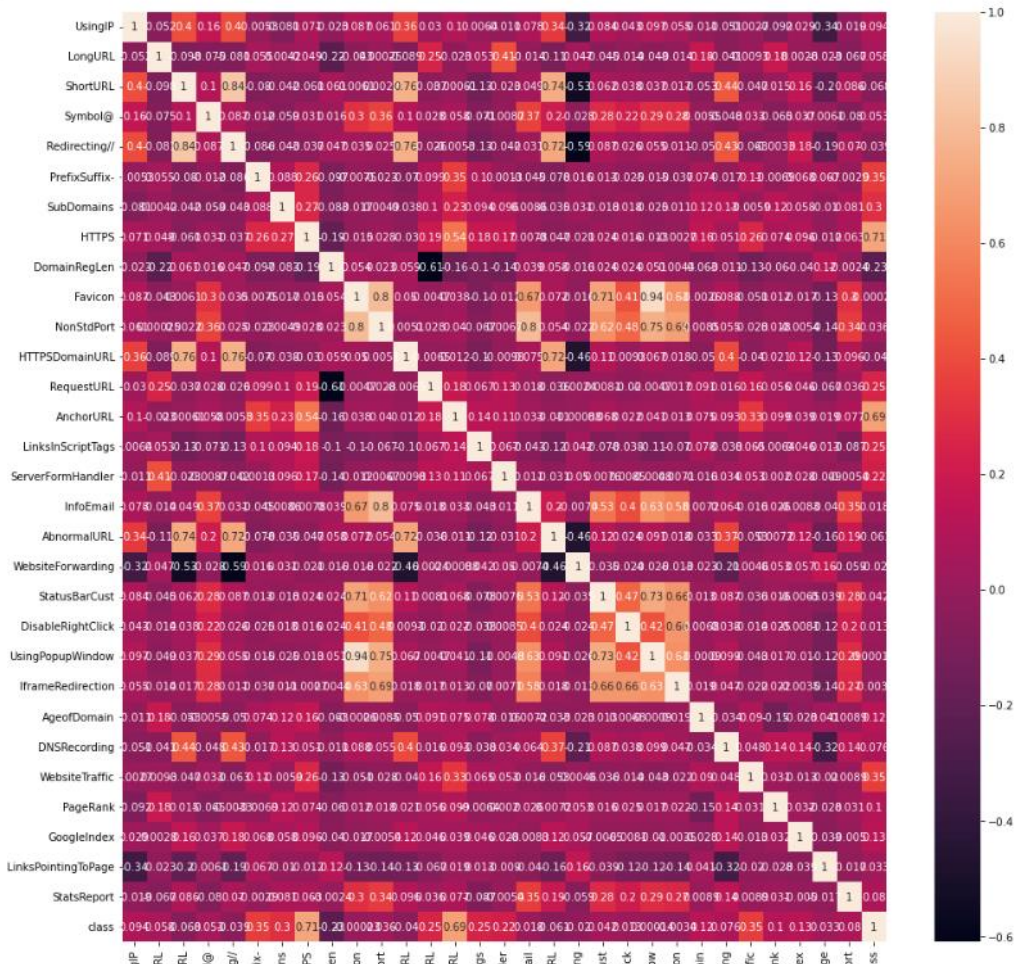
## 3. Visualizing the data:

Few plots and graphs are displayed to find how the data is distributed and the how features are related to each other.
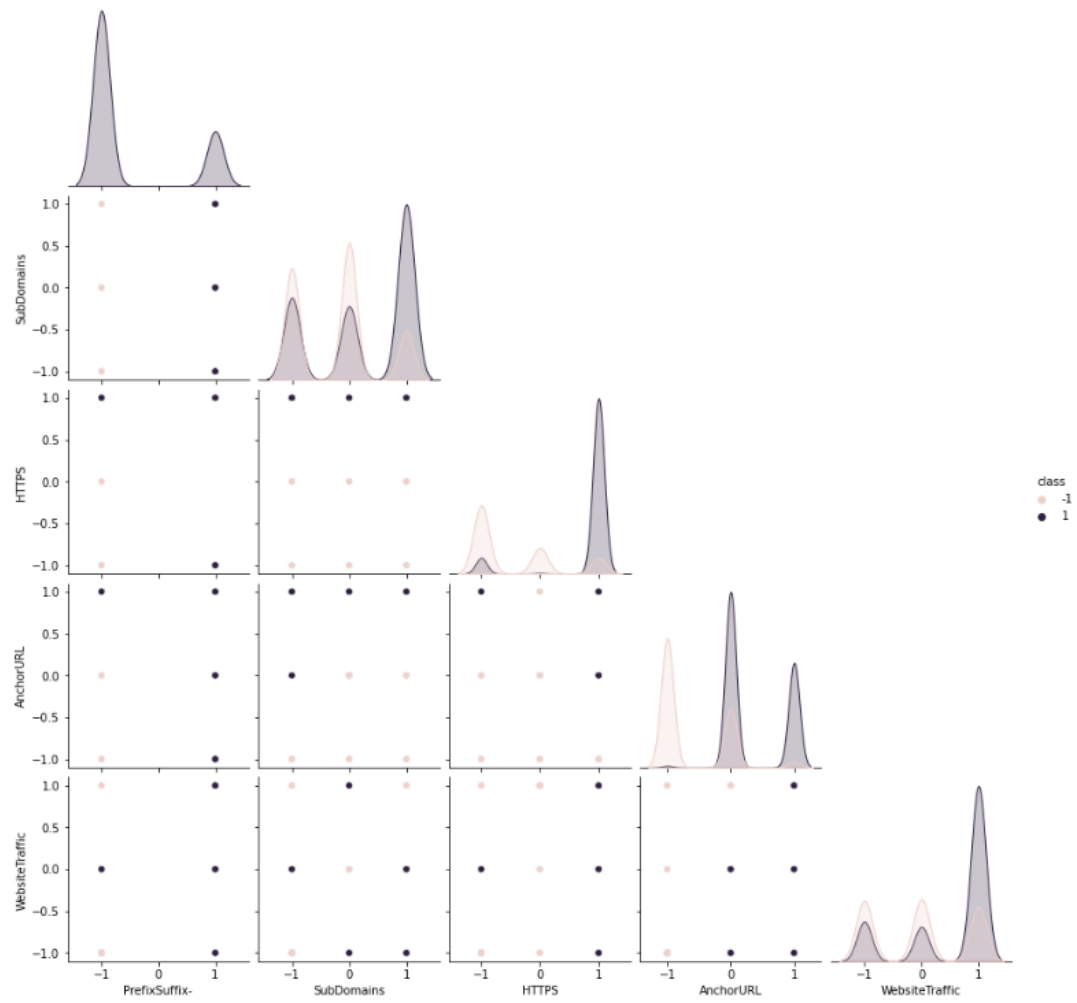
```
In [10]:   #Correlation heatmap

           plt.figure(figsize=(16,16))
           sns.heatmap(data.corr(), annot=True)
           plt.show()
```
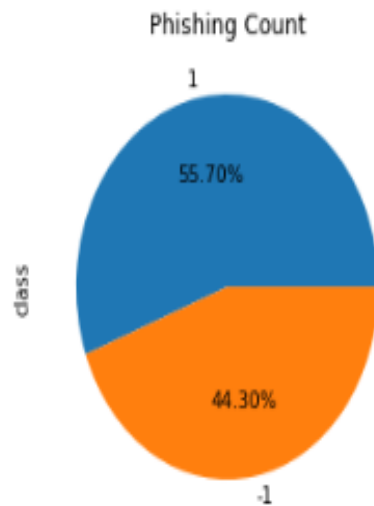
```
#pairplot for particular features

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

```python
# Phishing Count in pie chart

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

**Phishing Count**



## 4. Splitting the Data:

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

[13]:
```python
# Splitting the dataset into dependant and independant fetature

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

[14]:
```python
# 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
```

[14]: ((8843, 30), (8843,), (2211, 30), (2211,))

## 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.

```python
# Linear regression model
from sklearn.linear_model import LogisticRegression
#from sklearn.pipeline import Pipeline

# instantiate the model
log = LogisticRegression()

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

16]: LogisticRegression()

```python
#predicting the target value from the model for the samples

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

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

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))
```

## 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.

In [21]:
```python
# K-Nearest Neighbors Classifier model
from sklearn.neighbors import KNeighborsClassifier

# instantiate the model
knn = KNeighborsClassifier(n_neighbors=1)

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

Out[21]: KNeighborsClassifier(n_neighbors=1)

In [22]:
```python
#predicting the target value from the model for the samples
y_train_knn = knn.predict(X_train)
y_test_knn = knn.predict(X_test)
```

In [23]:
```python
#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 Neighborsn : 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 Neighborsn : 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
```
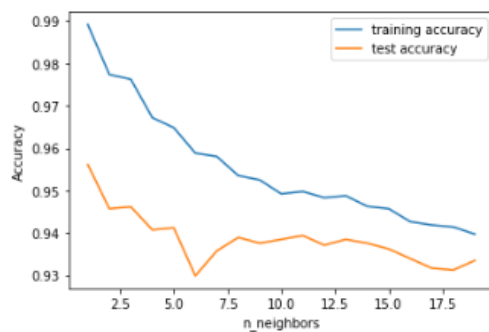
```
[24]:   #computing the classification report of the model

         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
```

```python
[25]:   training_accuracy = []
        test_accuracy = []
        # try max_depth from 1 to 20
        depth = range(1,20)
        for n in depth:
            knn = KNeighborsClassifier(n_neighbors=n)

            knn.fit(X_train, y_train)
            # record training set accuracy
            training_accuracy.append(knn.score(X_train, y_train))
            # record generalization accuracy
            test_accuracy.append(knn.score(X_test, y_test))


        #plotting the training & testing accuracy for n_estimators from 1 to 20
        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();
```



```python
[26]:   #storing the results. The below mentioned order of parameter passing is important.

        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.

[27]:
```python
# Support Vector Classifier model
from sklearn.svm import SVC
from sklearn.model_selection import GridSearchCV

# defining parameter range
param_grid = {'gamma': [0.1],'kernel': ['rbf','linear']}

svc = GridSearchCV(SVC(), param_grid)

# fitting the model for grid search
svc.fit(X_train, y_train)
```

[27]:
```
GridSearchCV(estimator=SVC(),
             param_grid={'gamma': [0.1], 'kernel': ['rbf', 'linear']})
```

[28]:
```python
#predicting the target value from the model for the samples
y_train_svc = svc.predict(X_train)
y_test_svc = svc.predict(X_test)
```

[29]:
```python
#computing the accuracy, f1_score, Recall, precision of the model performance

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 : 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
```

[30]:
```python
#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
```

[31]:
```python
#storing the results. The below mentioned order of parameter passing is important.

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.

```python
# Naive Bayes Classifier Model
from sklearn.naive_bayes import GaussianNB
from sklearn.pipeline import Pipeline

# instantiate the model
nb=  GaussianNB()

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

```
GaussianNB()
```

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

```python
#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
```

```python
#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    |

```python
#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)
```
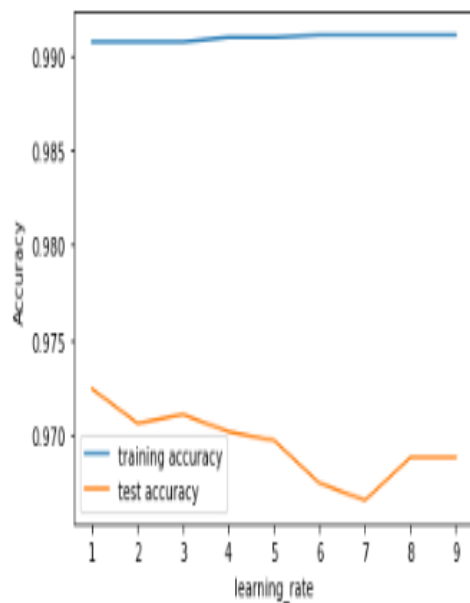
```
In [63]:  training_accuracy = []
          test_accuracy = []
          # try learning_rate from 0.1 to 0.9
          depth = range(1,10)
          for n in depth:
              forest_test =  CatBoostClassifier(learning_rate = n*0.1)

              forest_test.fit(X_train, y_train)
              # record training set accuracy
              training_accuracy.append(forest_test.score(X_train, y_train))
              # record generalization accuracy
              test_accuracy.append(forest_test.score(X_test, y_test))
```

```
0:      learn: 0.5487232        total: 7.17ms   remaining: 7.17s
1:      learn: 0.4349357        total: 16.7ms   remaining: 8.33s
2:      learn: 0.3609236        total: 27ms     remaining: 8.98s
3:      learn: 0.3050829        total: 36.4ms   remaining: 9.07s
4:      learn: 0.2766620        total: 45.9ms   remaining: 9.14s
5:      learn: 0.2475476        total: 56ms     remaining: 9.28s
6:      learn: 0.2286637        total: 66.2ms   remaining: 9.38s
7:      learn: 0.2138754        total: 76.2ms   remaining: 9.44s
8:      learn: 0.2013643        total: 86.3ms   remaining: 9.5s
9:      learn: 0.1896378        total: 95.8ms   remaining: 9.49s
10:     learn: 0.1819539        total: 106ms    remaining: 9.49s
11:     learn: 0.1767867        total: 115ms    remaining: 9.49s
12:     learn: 0.1727735        total: 125ms    remaining: 9.5s
13:     learn: 0.1682578        total: 135ms    remaining: 9.48s
14:     learn: 0.1641759        total: 144ms    remaining: 9.48s
15:     learn: 0.1614218        total: 154ms    remaining: 9.49s
16:     learn: 0.1558968        total: 164ms    remaining: 9.49s
17:     learn: 0.1535881        total: 174ms    remaining: 9.49s
18:     learn: 0.1514228        total: 184ms    remaining: 9.48s
19:     learn: 0.1482580        total: 195ms    remaining: 9.56s
20:     learn: 0.1452536        total: 206ms    remaining: 9.61s
21:     learn: 0.1426992        total: 218ms    remaining: 9.71s
22:     learn: 0.1405068        total: 229ms    remaining: 9.73s
23:     learn: 0.1381617        total: 239ms    remaining: 9.7s
24:     learn: 0.1363558        total: 249ms    remaining: 9.71s
25:     learn: 0.1341378        total: 259ms    remaining: 9.7s
26:     learn: 0.1323241        total: 269ms    remaining: 9.68s
27:     learn: 0.1305175        total: 278ms    remaining: 9.66s
28:     learn: 0.1289123        total: 288ms    remaining: 9.64s
29:     learn: 0.1278445        total: 298ms    remaining: 9.62s
30:     learn: 0.1256489        total: 308ms    remaining: 9.61s
31:     learn: 0.1239303        total: 317ms    remaining: 9.6s
32:     learn: 0.1216332        total: 327ms    remaining: 9.59s
33:     learn: 0.1207898        total: 337ms    remaining: 9.58s
34:     learn: 0.1197948        total: 347ms    remaining: 9.56s
35:     learn: 0.1187666        total: 358ms    remaining: 9.58s
36:     learn: 0.1175524        total: 369ms    remaining: 9.6s
37:     learn: 0.1164566        total: 381ms    remaining: 9.64s
38:     learn: 0.1156115        total: 392ms    remaining: 9.67s
39:     learn: 0.1143547        total: 404ms    remaining: 9.7s
40:     learn: 0.1131558        total: 414ms    remaining: 9.69s
41:     learn: 0.1122965        total: 424ms    remaining: 9.68s
42:     learn: 0.1113744        total: 435ms    remaining: 9.67s
43:     learn: 0.1103635        total: 444ms    remaining: 9.64s
44:     learn: 0.1093154        total: 453ms    remaining: 9.62s
45:     learn: 0.1085949        total: 463ms    remaining: 9.6s
46:     learn: 0.1078806        total: 472ms    remaining: 9.58s
47:     learn: 0.1071933        total: 483ms    remaining: 9.57s
48:     learn: 0.1061508        total: 492ms    remaining: 9.56s
49:     learn: 0.1051555        total: 502ms    remaining: 9.54s
50:     learn: 0.1044416        total: 511ms    remaining: 9.51s
51:     learn: 0.1035176        total: 521ms    remaining: 9.49s
```

In [64]:

```python
#plotting the training & testing accuracy for n_estimators from 1 to 50
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();
```



In [65]:

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

storeResults('CatBoost Classifier',acc_test_cat,f1_score_test_cat,
             recall_score_train_cat,precision_score_train_cat)
```

## 5.9. 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

```
In [67]:   from sklearn.preprocessing import LabelEncoder
           le = LabelEncoder()
           y_train = le.fit_transform(y_train)
```

```
In [68]:   # XGBoost Classifier Model
           from xgboost import XGBClassifier

           # instantiate the model
           xgb = XGBClassifier()

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

```
Out[68]:  XGBClassifier(base_score=0.5, booster='gbtree', callbacks=None,
                        colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1,
                        early_stopping_rounds=None, enable_categorical=False,
                        eval_metric=None, gamma=0, gpu_id=-1, grow_policy='depthwise',
                        importance_type=None, interaction_constraints='',
                        learning_rate=0.300000012, max_bin=256, max_cat_to_onehot=4,
                        max_delta_step=0, max_depth=6, max_leaves=0, min_child_weight=1,
                        missing=nan, monotone_constraints='()', n_estimators=100,
                        n_jobs=0, num_parallel_tree=1, predictor='auto', random_state=0,
                        reg_alpha=0, reg_lambda=1, ...)
```

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**
**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

```
In [69]:   #predicting the target value from the model for the samples
           y_train_xgb = xgb.predict(X_train)
           y_test_xgb = xgb.predict(X_test)
```

```
In [74]:   #computing the accuracy, f1_score, Recall, precision of the model performance

           acc_train_xgb = metrics.accuracy_score(y_train,y_train_xgb)
           acc_test_xgb = metrics.accuracy_score(y_test,y_test_xgb)
           print("XGBoost Classifier : Accuracy on training Data: {:.3f}".format(acc_train_xgb))
           print("XGBoost Classifier : Accuracy on test Data: {:.3f}".format(acc_test_xgb))
           print()

           f1_score_train_xgb = metrics.f1_score(y_train,y_train_xgb)
           f1_score_test_xgb = metrics.f1_score(y_test,y_test_xgb,average='micro')
           print("XGBoost Classifier : f1_score on training Data: {:.3f}".format(f1_score_train_xgb))
           print("XGBoost Classifier : f1_score on test Data: {:.3f}".format(f1_score_test_xgb))
           print()

           recall_score_train_xgb = metrics.recall_score(y_train,y_train_xgb)
           recall_score_test_xgb = metrics.recall_score(y_test,y_test_xgb,average='micro')
           print("XGBoost Classifier : Recall on training Data: {:.3f}".format(recall_score_train_xgb))
           print("XGBoost Classifier : Recall on test Data: {:.3f}".format(recall_score_train_xgb))
           print()

           precision_score_train_xgb = metrics.precision_score(y_train,y_train_xgb)
           precision_score_test_xgb = metrics.precision_score(y_test,y_test_xgb,average='micro')
           print("XGBoost Classifier : precision on training Data: {:.3f}".format(precision_score_train_xgb))
           print("XGBoost Classifier : precision on test Data: {:.3f}".format(precision_score_train_xgb))
```

```
XGBoost Classifier : Accuracy on training Data: 0.997
XGBoost Classifier : Accuracy on test Data: 0.698

XGBoost Classifier : f1_score on training Data: 0.998
XGBoost Classifier : f1_score on test Data: 0.698

XGBoost Classifier : Recall on training Data: 0.998
XGBoost Classifier : Recall on test Data: 0.998

XGBoost Classifier : precision on training Data: 0.998
XGBoost Classifier : precision on test Data: 0.998
```

```
In [75]:   #storing the results. The below mentioned order of parameter passing is important.

           storeResults('XGBoost Classifier',acc_test_xgb,f1_score_test_xgb,
                        recall_score_train_xgb,precision_score_train_xgb)
```

## 5.10. Multi-layer Perceptron classifier

MLPClassifier stands for Multi-layer Perceptron classifier which in the name itself connects to a Neural Network. Unlike other classification algorithms such as Support Vectors or Naive Bayes Classifier, MLPClassifier relies on an underlying Neural Network to perform the task of classification.

```python
# Multi-layer Perceptron Classifier Model
from sklearn.neural_network import MLPClassifier

# instantiate the model
mlp = MLPClassifier()
#mlp = GridSearchCV(mlpc, parameter_space)

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

```
MLPClassifier()
```
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```python
#predicting the target value from the model for the samples
y_train_mlp = mlp.predict(X_train)
y_test_mlp = mlp.predict(X_test)
```

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

acc_train_mlp  = metrics.accuracy_score(y_train,y_train_mlp)
acc_test_mlp = metrics.accuracy_score(y_test,y_test_mlp)
print("Multi-layer Perceptron : Accuracy on training Data: {:.3f}".format(acc_train_mlp))
print("Multi-layer Perceptron : Accuracy on test Data: {:.3f}".format(acc_test_mlp))
print()

f1_score_train_mlp = metrics.f1_score(y_train,y_train_mlp)
f1_score_test_mlp = metrics.f1_score(y_test,y_test_mlp,average="micro")
print("Multi-layer Perceptron : f1_score on training Data: {:.3f}".format(f1_score_train_mlp))
print("Multi-layer Perceptron : f1_score on test Data: {:.3f}".format(f1_score_train_mlp))
print()

recall_score_train_mlp = metrics.recall_score(y_train,y_train_mlp)
recall_score_test_mlp = metrics.recall_score(y_test,y_test_mlp,average="micro")
print("Multi-layer Perceptron : Recall on training Data: {:.3f}".format(recall_score_train_mlp))
print("Multi-layer Perceptron : Recall on test Data: {:.3f}".format(recall_score_test_mlp))
print()

precision_score_train_mlp = metrics.precision_score(y_train,y_train_mlp)
precision_score_test_mlp = metrics.precision_score(y_test,y_test_mlp,average="micro")
print("Multi-layer Perceptron : precision on training Data: {:.3f}".format(precision_score_train_mlp))
print("Multi-layer Perceptron : precision on test Data: {:.3f}".format(precision_score_test_mlp))
```

```
Multi-layer Perceptron : Accuracy on training Data: 0.985
Multi-layer Perceptron : Accuracy on test Data: 0.543

Multi-layer Perceptron : f1_score on training Data: 0.986
Multi-layer Perceptron : f1_score on test Data: 0.986

Multi-layer Perceptron : Recall on training Data: 0.989
Multi-layer Perceptron : Recall on test Data: 0.543

Multi-layer Perceptron : precision on training Data: 0.983
Multi-layer Perceptron : precision on test Data: 0.543
```

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

storeResults('Multi-layer Perceptron',acc_test_mlp,f1_score_test_mlp,
             recall_score_train_mlp,precision_score_train_mlp)
```

# 6. Comparision of Models

To compare the models performance, a dataframe is created. The columns of this dataframe are the lists created to store the results of the model.

```
In [80]:   #creating dataframe
           result = pd.DataFrame({ 'ML Model' : ML_Model,
                                   'Accuracy' : accuracy,
                                   'f1_score' : f1_score,
                                   'Recall'   : recall,
                                   'Precision': precision,
                                  })
```

```
In [81]:   # dispalying total result
           result
```

Out[81]:

|   | 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 | 0.958 | 0.962 | 0.991 | 0.993 |
| 5 | Random Forest | 0.969 | 0.972 | 0.992 | 0.991 |
| 6 | Gradient Boosting Classifier | 0.974 | 0.977 | 0.994 | 0.986 |
| 7 | CatBoost Classifier | 0.972 | 0.975 | 0.994 | 0.989 |
| 8 | XGBoost Classifier | 0.548 | 0.548 | 0.993 | 0.984 |
| 9 | Multi-layer Perceptron | 0.543 | 0.543 | 0.989 | 0.983 |

```
In [82]:   #Sorting the datafram on accuracy
           sorted_result=result.sort_values(by=['Accuracy', 'f1_score'],ascending=False).reset_index(drop=True)
```

```
In [83]:   # dispalying total result
           sorted_result
```

Out[83]:

|   | ML Model | Accuracy | f1 score | Recall | Precision |
|---|----------|----------|----------|--------|-----------|
| 0 | Gradient Boosting Classifier | 0.974 | 0.977 | 0.994 | 0.986 |
| 1 | CatBoost Classifier | 0.972 | 0.975 | 0.994 | 0.989 |
| 2 | Random Forest | 0.969 | 0.972 | 0.992 | 0.991 |
| 3 | Support Vector Machine | 0.964 | 0.968 | 0.980 | 0.965 |
| 4 | Decision Tree | 0.958 | 0.962 | 0.991 | 0.993 |
| 5 | K-Nearest Neighbors | 0.956 | 0.961 | 0.991 | 0.989 |
| 6 | Logistic Regression | 0.934 | 0.941 | 0.943 | 0.927 |
| 7 | Naive Bayes Classifier | 0.605 | 0.454 | 0.292 | 0.997 |
| 8 | XGBoost Classifier | 0.548 | 0.548 | 0.993 | 0.984 |
| 9 | Multi-layer Perceptron | 0.543 | 0.543 | 0.989 | 0.983 |

## Storing Best Model

```python
#  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(learning_rate=0.7, max_depth=4)
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```python
import pickle

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

```python
#checking the feature improtance in the model
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()
```



Feature importances using permutation on full model