

Cyber Security

DESCRIPTION

Problem Statement:

BookMyShow will enable the ads on their website, but they are also very cautious about their user privacy and information who visit their website. Some ads URL could contain a malicious link that can trick any recipient and lead to a malware installation, freezing the system as part of a ransomware attack or revealing sensitive information. BookMyShow now wants to analyze that whether the particular URL is prone to phishing (malicious) or not.

Dataset Details:

The input dataset contains an 11k sample corresponding to the 11k URL. Each sample contains 32 features that give a different and unique description of URL ranging from -1,0,1.

1: Phishing
0: Suspicious

1: Legitimate

The sample could be either legitimate or phishing.

Project Task: Week 1

Exploratory Data Analysis:

Each sample has 32 features ranging from -1,0,1. Explore the data using histogram, heatmaps.

Determine the number of samples present in the data, unique elements in all the features.

Check if there is any null value in any features.

Correlation of features and feature selection:

Next, we have to find if there are any correlated features present in the data. Remove the feature which might be correlated with some threshold.

Project Task: Week 2

Building Classification Model

Finally, build a robust classification system that classifies whether the URL sample is a phishing site or not.

Build classification models using a binary classifier to detect malicious or phishing URLs.

Illustrate the diagnostic ability of the binary classifier by plotting the ROC curve.

Validate the accuracy of data by the K-Fold cross-validation technique.

The final output consists of the model, which will give maximum accuracy on the validation dataset with selected attributes.

Import necessary modules

```
In [13]: import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import os
import seaborn as sns

import warnings
warnings.filterwarnings('ignore')

pd.options.display.max_columns = 999

# pre-processing
from sklearn.decomposition import PCA

# import libraries for model validation
from sklearn.model_selection import StratifiedKFold
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import train_test_split
from sklearn.model_selection import leave_one_out

# import the ML algorithm
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.metrics import precision_score
from sklearn.metrics import recall_score
from sklearn.metrics import f1_score
from sklearn.metrics import classification_report
from sklearn import metrics

# import libraries for metrics and reporting
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.metrics import accuracy_score
from sklearn.metrics import precision_score
from sklearn.metrics import recall_score
from sklearn.metrics import f1_score
from sklearn import metrics

# import libraries for metrics and reporting
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.metrics import accuracy_score
from sklearn.metrics import precision_score
from sklearn.metrics import recall_score
from sklearn.metrics import f1_score
from sklearn import metrics

# Grid searching key hyperparameters for logistic regression
from sklearn.model_selection import GridSearchCV
from sklearn.feature_selection import SelectKBest
from sklearn.pipeline import make_pipeline
from sklearn.pipeline import Pipeline

# Grid searching key hyperparameters for logistic regression
from sklearn.model_selection import GridSearchCV
from sklearn.feature_selection import SelectKBest
from sklearn.pipeline import make_pipeline
from sklearn.pipeline import Pipeline
```

Load Data

```
In [36]: import os
import os
for filename in os.listdir('input'):
    print(os.path.join('input', filename))

In [37]: file = 'input/dataset.csv'

In [38]: df = pd.read_csv(file, header=0)
print("Training Data has ", df.shape[0], "Rows and ", df.shape[1], "Columns", )

Training Data has 11055 Rows and 32 Columns
```

Exploratory Data Analysis

```
In [9]: df.sample(5)

Out[9]:
```

	index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Prefix_Suffix	having
7750	7751	-1	-1	1	1	1	1	-1
4485	4486	1	-1	1	1	1	1	-1
747	748	-1	-1	-1	1	-1	-1	-1
9628	9629	1	-1	1	1	1	1	-1
6462	6463	-1	-1	1	1	1	1	-1

```
In [10]: df.dtypes

Out[10]:
```

	index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Prefix_Suffix	having
count	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000
mean	5528.000000	0.313795	-0.633198	0.738761	0.700588	0.741474	-0.734962	
std	3191.447947	0.949534	0.766095	0.673998	0.713598	0.671011	0.678139	
min	1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	
25%	2764.500000	-1.000000	-1.000000	1.000000	1.000000	1.000000	-1.000000	
50%	5528.000000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
75%	8291.500000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
max	11055.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	

```
In [41]: df.describe(include='all')

Out[11]:
```

	index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Prefix_Suffix	having
count	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000
mean	5528.000000	0.313795	-0.633198	0.738761	0.700588	0.741474	-0.734962	
std	3191.447947	0.949534	0.766095	0.673998	0.713598	0.671011	0.678139	
min	1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	
25%	2764.500000	-1.000000	-1.000000	1.000000	1.000000	1.000000	-1.000000	
50%	5528.000000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
75%	8291.500000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
max	11055.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	

```
In [41]: df.dtypes

Out[41]:
```

	index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Prefix_Suffix	having
count	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000
mean	5528.000000	0.313795	-0.633198	0.738761	0.700588	0.741474	-0.734962	
std	3191.447947	0.949534	0.766095	0.673998	0.713598	0.671011	0.678139	
min	1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	
25%	2764.500000	-1.000000	-1.000000	1.000000	1.000000	1.000000	-1.000000	
50%	5528.000000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
75%	8291.500000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
max	11055.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	

```
In [41]: df.dtypes

Out[41]:
```

	index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Prefix_Suffix	having
count	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000
mean	5528.000000	0.313795	-0.633198	0.738761	0.700588	0.741474	-0.734962	
std	3191.447947	0.949534	0.766095	0.673998	0.713598	0.671011	0.678139	
min	1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	
25%	2764.500000	-1.000000	-1.000000	1.000000	1.000000	1.000000	-1.000000	
50%	5528.000000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
75%	8291.500000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
max	11055.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	

Removing Index column as this column will not have any significance in training and testing/prediction

```
In [42]: df.drop(['index', 'having_IPhaving_IP_Address', 'URLURL_Length',
'Shortining_Service', 'having_At_Symbol', 'double_slash_redirecting',
'Prefix_Suffix', 'having_Sub_Domain', 'SSLfinal_State',
'Domain_registration_length', 'Favicon', 'port', 'HTTPS_token',
'Request_URL', 'URL_of_Anchor', 'Links_in_tags', 'SFB',
'Submitting_to_email', 'Abnormal_URL', 'Redirect', 'on_mouseover',
'RightClick', 'popupWindow', 'Iframe', 'age_of_domain', 'DNSRecord',
'web_traffic', 'Page_Rank', 'Google_Index', 'Links_pointing_to_page',
'Statistical_report', 'Result'])

In [43]: df.dtypes

Out[43]:
```

	index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Prefix_Suffix	having
count	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000
mean	5528.000000	0.313795	-0.633198	0.738761	0.700588	0.741474	-0.734962	
std	3191.447947	0.949534	0.766095	0.673998	0.713598	0.671011	0.678139	
min	1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	
25%	2764.500000	-1.000000	-1.000000	1.000000	1.000000	1.000000	-1.000000	
50%	5528.000000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
75%	8291.500000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
max	11055.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	

Each sample has 30 features ranging from -1,0,1. Explore the data using histogram, heatmaps.

Histogram

```
In [15]: cols = ['having_IPhaving_IP_Address', 'URLURL_Length',
'Shortining_Service', 'having_At_Symbol', 'double_slash_redirecting',
'Prefix_Suffix', 'having_Sub_Domain', 'SSLfinal_State',
'Domain_registration_length', 'Favicon', 'port', 'HTTPS_token',
'Request_URL', 'URL_of_Anchor', 'Links_in_tags', 'SFB',
'Submitting_to_email', 'Abnormal_URL', 'Redirect', 'on_mouseover',
'RightClick', 'popupWindow', 'Iframe', 'age_of_domain', 'DNSRecord',
'web_traffic', 'Page_Rank', 'Google_Index', 'Links_pointing_to_page',
'Statistical_report', 'Result']

for col in cols:
    plt.figure(figsize=(18,5))
    plt.hist(df[col], bins=3, color = 'green', density=False)
    plt.xlabel(col, fontsize=10)
    plt.ylabel('Count', fontsize=10)
    plt.xticks(fontsize=8)
    plt.yticks(fontsize=8)
```

Heatmap

```
In [16]: df.columns

Out[16]:
```

	index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Prefix_Suffix	having
count	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000
mean	5528.000000	0.313795	-0.633198	0.738761	0.700588	0.741474	-0.734962	
std	3191.447947	0.949534	0.766095	0.673998	0.713598	0.671011	0.678139	
min	1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	
25%	2764.500000	-1.000000	-1.000000	1.000000	1.000000	1.000000	-1.000000	
50%	5528.000000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
75%	8291.500000	1.000000	-1.000000	1.000000	1.000000	1.000000	1.000000	
max	11055.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	

```
In [17]: corr = df.corr()
corr
```



```
[1487 128]
[1481 134]
Accuracy for pcca no of component: 1 is 0.205128205128206
[1465 150]
[126 129]
Accuracy for pcca no of component: 2 is 0.6102564102564103
[1469 126]
[158 57]
Accuracy for pcca no of component: 4 is 0.7931623931623931
[1558 57]
Accuracy for pcca no of component: 5 is 0.9213675213675213
[1555 60]
Accuracy for pcca no of component: 6 is 0.9213675213675213
[1555 60]
Accuracy for pcca no of component: 7 is 0.9213675213675213
[1555 60]
Accuracy for pcca no of component: 8 is 0.9213675213675213
[1556 59]
Accuracy for pcca no of component: 9 is 0.9213675213675213
[1556 59]
Accuracy for pcca no of component: 10 is 0.9213675213675213
[1556 59]
Accuracy for pcca no of component: 11 is 0.9213675213675213
[1557 58]
[127 129]
Accuracy for pcca no of component: 12 is 0.923931623931624
[1557 58]
[127 129]
Accuracy for pcca no of component: 13 is 0.923931623931624
[1556 59]
[126 129]
Accuracy for pcca no of component: 14 is 0.9205128205128205
[1556 59]
[126 129]
Accuracy for pcca no of component: 15 is 0.9230769230769231
[1561 54]
Accuracy for pcca no of component: 16 is 0.9282051282051282
[1562 53]
Accuracy for pcca no of component: 17 is 0.9333333333333333
[1569 46]
[127 129]
Accuracy for pcca no of component: 18 is 0.9358974358974359
[1570 45]
[127 129]
Accuracy for pcca no of component: 19 is 0.9376068376068376
[1569 46]
[126 129]
Accuracy for pcca no of component: 20 is 0.9384615384615385
[1568 47]
[127 129]
Accuracy for pcca no of component: 21 is 0.9376068376068376
[1567 48]
[126 129]
Accuracy for pcca no of component: 22 is 0.9367521367521368
[1566 49]
[127 129]
Accuracy for pcca no of component: 23 is 0.9341880341880342
[1567 48]
[127 129]
Accuracy for pcca no of component: 24 is 0.9333333333333333
[1568 47]
[126 129]
Accuracy for pcca no of component: 25 is 0.9341880341880342
[1570 45]
[127 129]
Accuracy for pcca no of component: 26 is 0.9384615384615385
[1570 45]
[127 129]
Accuracy for pcca no of component: 27 is 0.9376068376068376
[1569 46]
[126 129]
Accuracy for pcca no of component: 28 is 0.9384615384615385
[1568 47]
[127 129]
Accuracy for pcca no of component: 29 is 0.9367521367521368
[1568 47]
[126 129]
Accuracy for pcca no of component: 30 is 0.9367521367521368
```

It seems if we keep no of components as 20 & 26 then the accuracy is 93.84%.

Still we would go for 26 components as prescribed by gridsearch as True positive is better (570) if we use 26 components.

Perform Logistic Regression

```
In [130]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
pca = PCA(n_components=26)
X_train = pca.fit_transform(X_train)
X_test = pca.transform(X_test)
# Training and Making Predictions
classifier = LogisticRegression()
classifier.fit(X_train, y_train)

# Predicting the Test set results
y_pred = classifier.predict(X_test)

cm = confusion_matrix(y_test, y_pred)
print(cm)
print('Accuracy : ', metrics.accuracy_score(y_test, y_pred))

[570 45]
[127 128]
Accuracy : 0.9384615384615385
```

Illustrate the diagnostic ability of this binary classifier by plotting the ROC curve.

What Are ROC Curves?

A useful tool when predicting the probability of a binary outcome is the Receiver Operating Characteristic curve or ROC curve.

It is a plot of the false positive rate (x-axis) versus the true positive rate (y-axis) for a number of different candidate threshold values between 0.0 and 1.0. Put another way, it plots the false alarm rate versus the hit rate.

The true positive rate is calculated as the number of true positives divided by the sum of the number of true positives and the number of false negatives. It describes how good the model is at predicting the positive class when the actual outcome is positive.

True Positive Rate = True Positives / (True Positives + False Negatives)

The true positive rate is also referred to as sensitivity.

Sensitivity = True Positives / (True Positives + False Negatives)

The false positive rate is calculated as the number of false positives divided by the sum of the number of false positives and the number of true negatives. It is also called the false alarm rate as it summarizes how often a positive class is predicted when the actual outcome is negative.

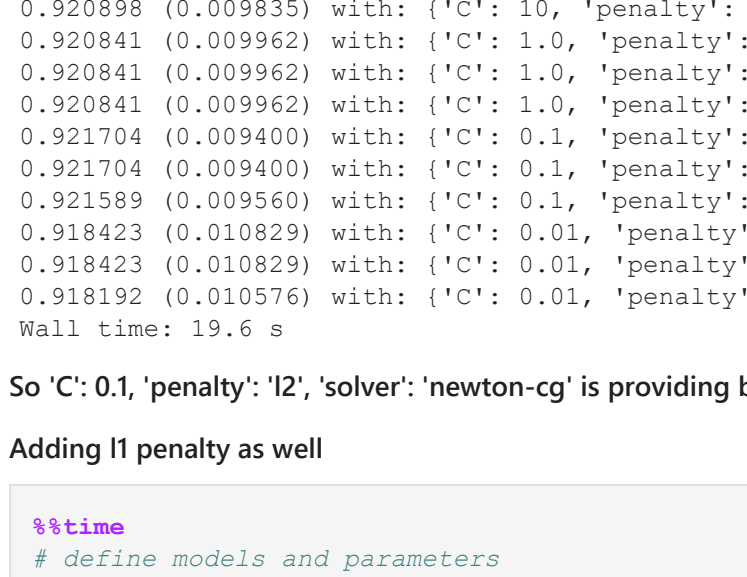
False Positive Rate = False Positives / (False Positives + True Negatives) The false positive rate is also referred to as the inverted specificity where specificity is the total number of true negatives divided by the sum of the number of true negatives and false positives. Specificity = True Negatives / (True Negatives + False Positives) False Positive Rate = 1 - Specificity

The ROC curve is a useful tool for a few reasons:

The curves of different models can be compared directly in general or for different thresholds. The area under the curve (AUC) can be used as a summary of the model skill. The shape of the curve contains a lot of information, including what we might care about most for a problem, the expected false positive rate, and the false negative rate.

To make this clear:

Smaller values on the x-axis of the plot indicate lower false positives and higher true negatives. Larger values on the y-axis of the plot indicate higher true positives and lower false negatives.



As we can observe from the ROC curve the x-axis values are Smaller which indicate lower false positives and higher true negatives.

Also y-axis values are > .9 which shows higher true positives and lower false negatives.

Below are the list of parameters to do hyper parameter tuning for logistic regression

Solver

solver in ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']

Regularization (penalty)

penalty in ['none', 'l1', 'l2', 'elasticnet'] Note: not all solvers support all regularization terms.

The C parameter controls the penalty strength, which can also be effective.

C in [100, 10, 1.0, 0.1, 0.01]

```
In [132]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.01, random_state=17)
pca = PCA(n_components=26)
X_train = pca.fit_transform(X_train)
X_test = pca.transform(X_test)
```

```
In [136]: # Define models and parameters
model = LogisticRegression()
solvers = ['newton-cg', 'lbfgs', 'liblinear']
penalty = ['l2']
c_values = [100, 10, 1.0, 0.1, 0.01]
# Define grid search
grid = dict(solver=solvers, penalty=penalty, C=c_values)
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
grid_search = GridSearchCV(estimator=model, param_grid=grid, n_jobs=-1, cv=cv, scoring='accuracy', error_score='raise')
grid_result = grid_search.fit(X_train, y_train)

# Summarize results
print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))
means = grid_result.cv_results_['mean_test_score']
stds = grid_result.cv_results_['std_test_score']
params = grid_result.cv_results_['params']
for mean, stdev, param in zip(means, stds, params):
    print("%f (%f) with: %s" % (mean, stdev, param))

Best: 0.921704 using {'C': 0.1, 'penalty': 'l2', 'solver': 'newton-cg'}
0.920455 (0.010267) with: {'C': 100, 'penalty': 'l1', 'solver': 'newton-cg'}
0.920265 (0.010267) with: {'C': 100, 'penalty': 'l2', 'solver': 'lbfgs'}
0.920495 (0.010267) with: {'C': 100, 'penalty': 'l1', 'solver': 'liblinear'}
0.920265 (0.010267) with: {'C': 10, 'penalty': 'l2', 'solver': 'newton-cg'}
0.920956 (0.009818) with: {'C': 10, 'penalty': 'l2', 'solver': 'lbfgs'}
0.920898 (0.009835) with: {'C': 10, 'penalty': 'l2', 'solver': 'liblinear'}
0.920841 (0.009862) with: {'C': 1.0, 'penalty': 'l2', 'solver': 'newton-cg'}
0.920841 (0.009862) with: {'C': 1.0, 'penalty': 'l2', 'solver': 'liblinear'}
0.920841 (0.009862) with: {'C': 1.0, 'penalty': 'l2', 'solver': 'lbfgs'}
0.921704 (0.009400) with: {'C': 0.1, 'penalty': 'l2', 'solver': 'liblinear'}
0.921589 (0.009560) with: {'C': 0.1, 'penalty': 'l2', 'solver': 'newton-cg'}
0.921589 (0.009560) with: {'C': 0.1, 'penalty': 'l2', 'solver': 'lbfgs'}
0.918423 (0.010829) with: {'C': 0.01, 'penalty': 'l2', 'solver': 'newton-cg'}
0.918423 (0.010829) with: {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}
0.918192 (0.010576) with: {'C': 0.01, 'penalty': 'l2', 'solver': 'liblinear'}
Wall time: 19.6 s
```

So 'C': 0.1, 'penalty': 'l2', 'solver': 'newton-cg' is providing best accuracy score of 92.17%

Adding l1 penalty as well

```
In [137]: # Define models and parameters
model = LogisticRegression()
solvers = ['newton-cg', 'lbfgs', 'liblinear']
penalty = ['l1', 'l2']
c_values = [100, 10, 1.0, 0.1, 0.01]
# Define grid search
grid = dict(solver=solvers, penalty=penalty, C=c_values)
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
grid_search = GridSearchCV(estimator=model, param_grid=grid, n_jobs=-1, cv=cv, scoring='accuracy', error_score='raise')
grid_result = grid_search.fit(X_train, y_train)

# Summarize results
print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))
means = grid_result.cv_results_['mean_test_score']
stds = grid_result.cv_results_['std_test_score']
params = grid_result.cv_results_['params']
for mean, stdev, param in zip(means, stds, params):
    print("%f (%f) with: %s" % (mean, stdev, param))

Best: 0.921704 using {'C': 0.1, 'penalty': 'l2', 'solver': 'newton-cg'}
0.000000 (0.000000) with: {'C': 100, 'penalty': 'l1', 'solver': 'newton-cg'}
0.000000 (0.000000) with: {'C': 100, 'penalty': 'l1', 'solver': 'lbfgs'}
0.920495 (0.010267) with: {'C': 100, 'penalty': 'l1', 'solver': 'liblinear'}
0.920265 (0.010267) with: {'C': 100, 'penalty': 'l2', 'solver': 'newton-cg'}
0.920265 (0.010267) with: {'C': 100, 'penalty': 'l1', 'solver': 'lbfgs'}
0.920265 (0.010267) with: {'C': 10, 'penalty': 'l2', 'solver': 'liblinear'}
0.000000 (0.000000) with: {'C': 10, 'penalty': 'l1', 'solver': 'newton-cg'}
0.000000 (0.000000) with: {'C': 10, 'penalty': 'l1', 'solver': 'lbfgs'}
0.920783 (0.009848) with: {'C': 10, 'penalty': 'l1', 'solver': 'liblinear'}
0.920956 (0.009818) with: {'C': 10, 'penalty': 'l2', 'solver': 'newton-cg'}
0.920898 (0.009835) with: {'C': 10, 'penalty': 'l2', 'solver': 'liblinear'}
0.000000 (0.000000) with: {'C': 1.0, 'penalty': 'l1', 'solver': 'newton-cg'}
0.000000 (0.000000) with: {'C': 1.0, 'penalty': 'l1', 'solver': 'lbfgs'}
0.920725 (0.010005) with: {'C': 1.0, 'penalty': 'l1', 'solver': 'liblinear'}
0.920841 (0.009862) with: {'C': 1.0, 'penalty': 'l2', 'solver': 'newton-cg'}
0.920841 (0.009862) with: {'C': 1.0, 'penalty': 'l2', 'solver': 'liblinear'}
0.000000 (0.000000) with: {'C': 0.1, 'penalty': 'l1', 'solver': 'newton-cg'}
0.000000 (0.000000) with: {'C': 0.1, 'penalty': 'l1', 'solver': 'lbfgs'}
0.920380 (0.010509) with: {'C': 0.1, 'penalty': 'l1', 'solver': 'liblinear'}
0.921704 (0.009400) with: {'C': 0.1, 'penalty': 'l2', 'solver': 'newton-cg'}
0.921589 (0.009560) with: {'C': 0.1, 'penalty': 'l2', 'solver': 'liblinear'}
0.000000 (0.000000) with: {'C': 0.01, 'penalty': 'l1', 'solver': 'newton-cg'}
0.000000 (0.000000) with: {'C': 0.01, 'penalty': 'l1', 'solver': 'lbfgs'}
0.901900 (0.013323) with: {'C': 0.01, 'penalty': 'l1', 'solver': 'liblinear'}
0.918423 (0.010829) with: {'C': 0.01, 'penalty': 'l2', 'solver': 'newton-cg'}
0.918423 (0.010829) with: {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}
0.918192 (0.010576) with: {'C': 0.01, 'penalty': 'l2', 'solver': 'liblinear'}
Wall time: 59.9 s
```

Still 'C': 0.1, 'penalty': 'l2', 'solver': 'newton-cg' is providing best accuracy score of 92.17%

```
In [138]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
pca = PCA(n_components=26)
X_train = pca.fit_transform(X_train)
X_test = pca.transform(X_test)
# Training and Making Predictions
classifier = LogisticRegression(C=0.1, penalty='l2', solver='newton-cg')
classifier.fit(X_train, y_train)
```

```
# Predicting the Test set results
y_pred = classifier.predict(X_test)

cm = confusion_matrix(y_test, y_pred)
print(cm)
print('Accuracy : ', metrics.accuracy_score(y_test, y_pred))

[570 45]
[126 129]
Accuracy : 0.9393162393162393
```

Validate the accuracy of data by the K-Fold cross-validation technique.

```
In [142]: # Import libraries for model validation
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_validate
```

```
In [145]: kfold = KFold(n_splits=11, shuffle=True, random_state=20)
```

```
In [146]: %timeit (accuracy, 'accuracy')
# Instantiate the lin reg model
classifier = LogisticRegression(C=0.1, penalty='l2', solver='newton-cg')
scores = cross_validate(estimator=classifier, X=X_train, y=y_train, cv=kfold, scoring='accuracy', return_train_score=False)

print('Score keys : \n', scores.keys())
print('Training accuracy : {}'.format(scores['train_accuracy'].mean()))
print('Testing accuracy : {}'.format(scores['test_accuracy'].mean()))

Score keys :
Training accuracy : 0.91927606901751
Testing accuracy : 0.917290525846702
Wall time: 677 ms
```

The KFold cross validation result shows the training and testing scores are almost similar, no overfitting.

Final Logistic Regression Model

```
In [155]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=7)
pca = PCA(n_components=26)
X_train = pca.fit_transform(X_train)
X_test = pca.transform(X_test)
# Training and Making Predictions
classifier = LogisticRegression(C=0.1, penalty='l2', solver='newton-cg')
lr.fit(X_train, y_train)

# Predicting the Test set results
y_pred = lr.predict(X_test)

cm = confusion_matrix(y_test, y_pred)
from sklearn.metrics import plot_confusion_matrix

fig, ax = plot_confusion_matrix(conf_mat=cm, figsize=(8, 6), cmap=plt.cm.Greens)
plt.xlabel('Predictions', fontsize=18)
plt.ylabel('Actuals', fontsize=18)
plt.title('Confusion Matrix', fontsize=18)
plt.show()
print('Accuracy : ', metrics.accuracy_score(y_test, y_pred))
```



```
Accuracy : 0.9253561253561253

from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))

precision    recall  f1-score   support

-1           0.94         0.91         0.93         910
 1           0.91         0.94         0.92         845

 accuracy          0.93         0.93         0.93        1755
 macro avg         0.93         0.93         0.93        1755
 weighted avg         0.93         0.93         0.93        1755
```

In [] :