



COMP534: Applied Artificial Intelligence

First Assignment

A Report

on

*“Performance analysis of multiple supervised learning methods for solving
binary classification problem”*

Submitted to: **Dr. Pamela Bezerra**

Submitted by: **Ananya Swami**

Rajesh Adepu

Introduction

The performance analysis has been executed using python programming language. A code.py file has been submitted along with this report containing the code for performance analysis.

The code begins with importing libraries. Various libraries are used within the code namely; numpy, matplotlib, pandas, scikit-learn (sklearn).

Loading the dataset:

To read the “diabetes.csv” file, we use pandas library and divide the data into input X and output Y.

Splitting the dataset:

Here, we have split the data into 3:1 ratio using train_test_split from sklearn library i.e. train(75%) and test(25%) data.

Cleaning the dataset:

To clean and preprocess the data we have used StandardScaler from sklearn library.

Classification:

The classification has been done using three different supervised learning methods namely; **k-nearest neighbors**, **random forest** and **logistic regression**.

To get this classification done we have directly used KNearestClassifier(), RandomForestClassifier() and LogisticRegression() functions from sklearn library.

Analysis:

The analysis is done by creating tables for each classification method and graph comparing accuracies of all three methods, out of which “*logistic regression was found to be the best*”.

Evaluation

k-Nearest Neighbors Classifier: (with k=10)

```
The KNN Classification details are as follows:
```

	precision	recall	f1-score	support
0.0	0.74	0.87	0.80	121
1.0	0.69	0.49	0.57	71
accuracy			0.73	192
macro avg	0.72	0.68	0.69	192
weighted avg	0.72	0.73	0.72	192

The Confusion Matrix for KNN:
[[105 16]
[36 35]]
The Accuracy for KNN = 0.7291666666666666

Random Forest Classifier: (with 50 estimators)

```
The Random Forest Classification details are as follows:
```

	precision	recall	f1-score	support
0.0	0.78	0.88	0.83	121
1.0	0.74	0.56	0.64	71
accuracy			0.77	192
macro avg	0.76	0.72	0.73	192
weighted avg	0.76	0.77	0.76	192

The Confusion Matrix for Random Forest:
[[107 14]
[31 40]]
The Accuracy for Random Forest = 0.765625

Logistic Regression Classifier: (best among all three)

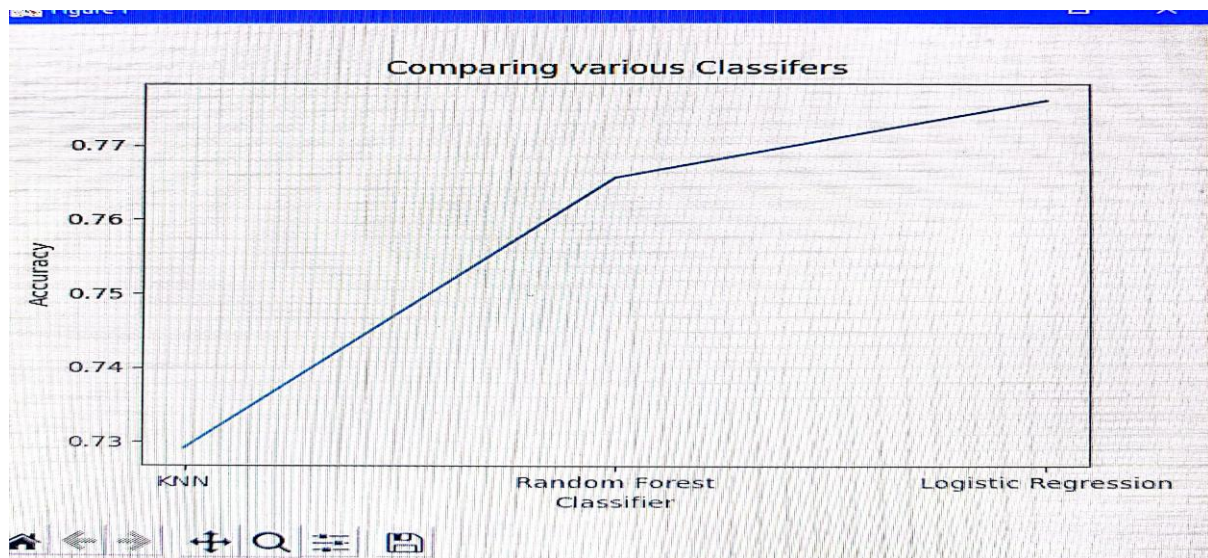
	precision	recall	f1-score	support
0.0	0.77	0.92	0.84	121
1.0	0.79	0.54	0.64	71
accuracy			0.78	192
macro avg	0.78	0.73	0.74	192
weighted avg	0.78	0.78	0.76	192

The Confusion Matrix for Logistic Regression:
[[111 10]
[33 38]]
The Accuracy for Logistic Regression = 0.7760416666666666

0 A 4

Results and Conclusion

As mentioned, a graph was also plotted using matplotlib library, comparing the accuracies of all three methods and it was found that logistic regression performed best among three as the number of instances were very large as compared to the number of features which were just 8. It generally does not makes any assumptions in distributing classes in a feature spce. The graph can be seen as follows:



Challenges faced in creating this project are as follows:

- There was a conflict of ideas wherein choosing three supervised learning methods out of the list of methods.
- Also, finding the best version for each method which includes finding appropriate k value and estimators was tedious.
- Time management was also one of the problems as we had to put in equal efforts in taking the assignment to its accomplishment stage.

Task Allocation:

Choosing 3 methods was done on a mutual agreement basis.

Programming and Code Creation: Ananya

Report Preparation: Rajesh