# AHSANULLAH UNIVERSITY OF SCIENCE AND TECHNOLOGY DHAKA-1208, BANGLADESH.



Department of Computer Science and Engineering
Spring 2019

Program: Bachelor of Science in Computer Science and Engineering
Course No: CSE 4108
Course Title: Artificial Intelligence Lab

Term Project: 03
Topic: K-Nearest Neighbor Classification and Decision Tree Classification

Date of Submission: October 6, 2019

Submitted to:
Dr. S.M. Abdullah Al-Mamun Professor, Department of
CSE, AUST.

Md. Siam Ansary
Adjunct Faculty, Department of CSE, AUST.

Submitted By:
Name: Syed Sanzam
ID: 16.01.04.042
Group: A2

**Preface:**
In machine learning and statistics, classification is a **supervised learning** approach in which the computer program learns from the data input given to it and then uses this learning to classify new observation. It is a tool widely used in the classification in many domains such as in credit approval, medical diagnosis, target marketing etc.

**Classification Algorithms:**
There are a number of classification algorithms used in various sectors, depending on the nature of a specific task. In this case, the **K-Nearest Neighbor** and **Decision Tree** classification algorithms are applied in order to find a predict a certain result.

1. K-Nearest Neighbor (KNN): A supervised classification algorithm, which takes a bunch of labelled points and uses them to learn how to label other points. To label a new point, it looks at the labelled points closest to that new point (those are its nearest neighbors), and has those neighbors vote, so whichever label the most of the neighbors have is the label for the new point (the "k" is the number of neighbors it checks).

2. Decision Tree: Builds classification models in the form of a tree structure. It breaks down a data set into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes. A decision node has two or more branches and a leaf node represents a classification or decision. The topmost decision node in a tree which corresponds to the best predictor called root node. Decision trees can handle both categorical and numerical data.

**Tools Used:**
- Python 3
- Sci-kit learn
- Pandas
- Numpy
- Matplotlib

**Dataset:**

*LiverPatients.csv*
Associated Task: Classification
Number of Instances: 583
Number of attributes: 11
Attribute Information:

Contents:
1. age: Age of the patient in years
2. gender: Patient Gender: Male or Female
3. TB: Total Bilirubin
4. DB: Direct Bilirubin
5. alkphos: Alkaline Phosphotase
6. sgpt: Alamine Aminotransferase
7. sgot: Aspartate Aminotransferase
8. TP: Total Proteins
9. ALB: Albumin
10. A_G: Ratio of Albumin and Globulin
11. class: Predictor Class: 1 if patient has Liver Disease and 2 if they do not


## Major Steps of Processing:

1. The dataset is divided into X and Y where X is represented as following columns from the LiverPatients.csv dataset,
   *'age','gender','TB','DB','alkphos','sgpt','sgot','TP','ALB','A_G'*
   And Y as, *'class'*

2. The dataset is then split into 10 folds (groups) for *Cross Validation*. For each group, take the group as **train dataset** and the remaining groups as the **test dataset.**
3. Fit a model on the training set and evaluate it on the test set
4. Performance Matrices (Accuracy, Mean Absolute Error, Mean Squared Absolute Error and Root Mean Squared Error) have been calculated using the testing data and predicted data for both models.
5. Finally, results have been visualized with a *Grouped Bar Chart* to illustrate the performance of both models.


## Python Implementation:
*kNNVsDecisionTree.py*

```
1.  import pandas as pd
2.  from pandas import DataFrame as df
3.  import numpy as np
4.  import matplotlib.pyplot as plot
5.  from sklearn.neighbors import KNeighborsClassifier
6.  from sklearn.model_selection import KFold
7.  from sklearn import metrics
8.  from sklearn.metrics import accuracy_score
9.  from sklearn import tree
10. from sklearn.preprocessing import LabelEncoder
11.
```

```python
12. def readDataset():
13.     global df
14.     df = pd.read_csv("LiverPatient.csv")
15.     #print(df.head())
16.
17. def allocateValues():
18.     enc = LabelEncoder()
19.     enc.fit(df['gender'])
20.     df['gender'] = enc.transform(df['gender'])
21.     # X = df[['age','gender','TB','DB','alkphos','sgpt','sgot','TP','ALB','A_G']]
22.     # Y = df['class']
23.     global X
24.     global Y
25.     X = df.iloc[:, [0,9]].values
26.     Y = df.iloc[:, 10].values
27.
28. def crossValidation(classifier):
29.     global acc
30.     global mae
31.     global mse
32.     global rmse
33.
34.     acc = np.array([ ])
35.     mae = 0
36.     mse = 0
37.     rmse = 0
38.
39.     ns = 10
40.     crossValidation = KFold(n_splits=ns)
41.
42.
43.     for train_index, test_index in crossValidation.split(X):
44.         X_train, X_test = X[train_index], X[test_index]
45.         Y_train, Y_test = Y[train_index], Y[test_index]
46.         classifier.fit(X_train,Y_train)
47.         prediction = classifier.predict(X_test)
48.         acc = np.append(acc, [accuracy_score(Y_test, prediction)], axis = 0)
49.
50.         mae = mae + metrics.mean_absolute_error(Y_test, prediction)
51.         mse  = mse + metrics.mean_squared_error(Y_test, prediction)
52.         rmse = rmse + np.sqrt(metrics.mean_squared_error(Y_test, prediction))
53.
54.
55.     acc = np.mean(acc)
56.     mae = mae/ns
57.     mse = mse/ns
58.     rmse = rmse/ns
59.
60.     print('Accuracy',acc)
61.     print('Mean Absolute Error: ',mae)
62.     print('Mean Squared Error',mse)
63.     print('Root Mean Squared Error',rmse)
64.     print('\n')
65.
66. def testClassifiers():
67.     knn = KNeighborsClassifier(n_neighbors=20)
68.     dct = tree.DecisionTreeClassifier()
69.     print('For KNN : ')
70.     crossValidation(knn)
71.
72.     global k_acc
```

```python
73.        global k_mae
74.        global k_mse
75.        global k_rmse
76.
77.        global d_acc
78.        global d_mae
79.        global d_mse
80.        global d_rmse
81.
82.        k_acc = acc
83.        k_mae = mae
84.        k_mse = mse
85.        k_rmse = rmse
86.
87.        print('For Decision Tree : ')
88.        crossValidation(dct)
89.
90.        d_acc = acc
91.        d_mae = mae
92.        d_mse = mse
93.        d_rmse = rmse
94.
95. def illustrateResult():
96.        n_groups = 4
97.        knn_bar = (k_acc*100, k_mae*100, k_mse*100, k_rmse*100)
98.        dct_bar = (d_acc*100, d_mae*100, d_mse*100, d_rmse*100)
99.
100.          # create plot
101.             fig, ax = plot.subplots()
102.             index = np.arange(n_groups)
103.             bar_width = 0.35
104.             opacity = 0.8
105.
106.             rects1 = plot.bar(index, knn_bar, bar_width,
107.             alpha=opacity,
108.             color='b',
109.             label='KNN')
110.
111.             rects2 = plot.bar(index + bar_width, dct_bar, bar_width,
112.             alpha=opacity,
113.             color='g',
114.             label='Decision Tree')
115.
116.             plot.xlabel('Classifiers')
117.             plot.ylabel('')
118.             plot.title('KNN vs Decision Tree')
119.             plot.xticks(index + bar_width, ('Accuracy', 'MAE', 'MSE', 'RMSE'))
120.             plot.legend()
121.
122.             plot.tight_layout()
123.             plot.show()
124.
125.         #Main
126.         readDataset()
127.         allocateValues()
128.         testClassifiers()
129.         illustrateResult()
```
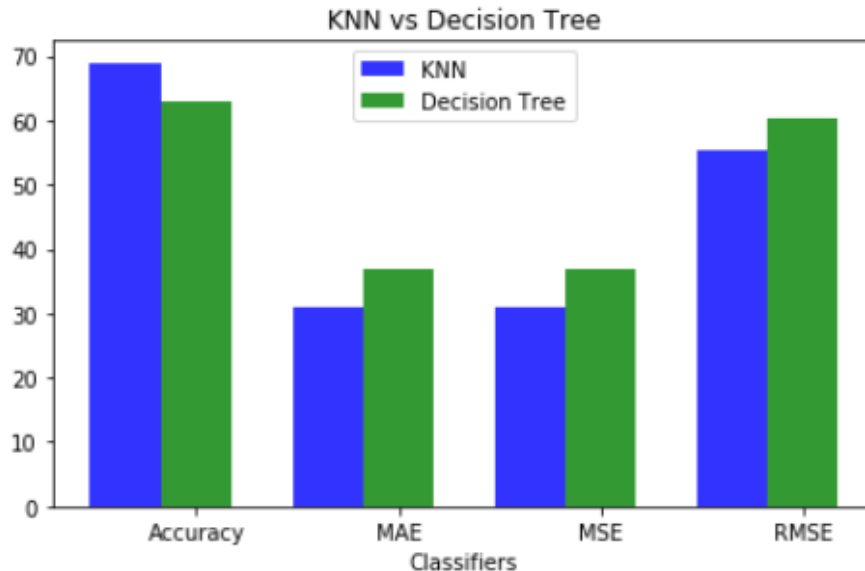
**Output:**

```
For KNN :
Accuracy 0.6893921683226184
Mean Absolute Error:  0.31060783167738165
Mean Squared Error 0.31060783167738165
Root Mean Squared Error 0.5542988792254294


For Decision Tree :
Accuracy 0.6310052600818234
Mean Absolute Error:  0.3689947399181766
Mean Squared Error 0.3689947399181766
Root Mean Squared Error 0.6044338437309891
```



**References:**

1.  https://machinelearningmastery.com/k-fold-cross-validation/
2.  https://towardsdatascience.com/machine-learning-classifiers-a5cc4e1b0623
3.  https://medium.com/@Mandysidana/machine-learning-types-of-classification-9497bd4f2e14

**My GitHub:** https://github.com/sanzamsyed/Artificial-Intelligence