

ASSIGNMENT 2

FIRST DATA SET

This is the energy prediction dataset. The dataset is about usage of appliances with various parameters like temperature, humidity, wind speed etc. The target variable Appliances has been converted to a binary classification variable by taking the median which is 60 and then dividing the values above 60 as Above Average and values below 60 as Below Average. The thresholding is done with the basic assumption that those below 60 will form one half of the distribution and vice versa.

Independent variables I have used from assignment 1:

lights, T1, RH_1, T2, RH_2, T3, RH_3, T4, T6, RH_6, RH_7, T8, RH_8, RH_9, T_out, RH_out, Windspeed

Dependent Variable: Efficiency

The data set is taken, and min-max normalization is done on the selected features. The data set divided into 70 and 30. 70% of the data it's used for training the dataset and the rest is used to test the data and evaluate accuracy of the model.

Algorithm 1: Artificial Neural Networks:

For neural networks, we will use **Keras** API on top of **Tensor flow** as this seems the fastest in terms of computing speed.

For all the experimentations we will be using GridSearchCV to pick the best parameter among a group of parameters. We will be doing the following experimentations:

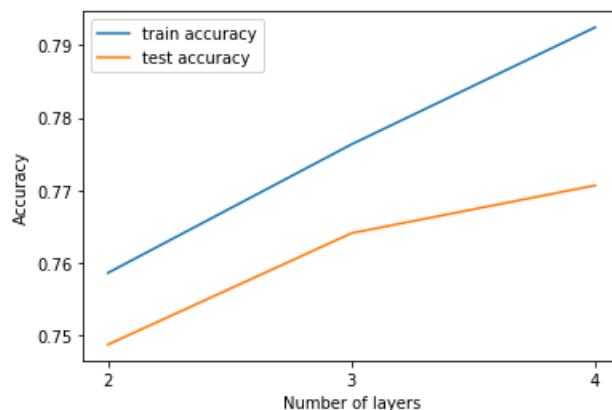
- Number of layers
- Number of nodes
- Activation functions
- Optimizers
- Dropout Rates

Experimentation with number of layers:

Initially, we begin by experimenting with the number of hidden layers from 2 to 4. The number of neurons are selected randomly with the following combinations.

- 2 hidden layers : 50-40
- 3 hidden layers: 50-40-20
- 4 hidden layers: 70-50-40-30

We plot the test and train accuracies to understand the variations in the errors.



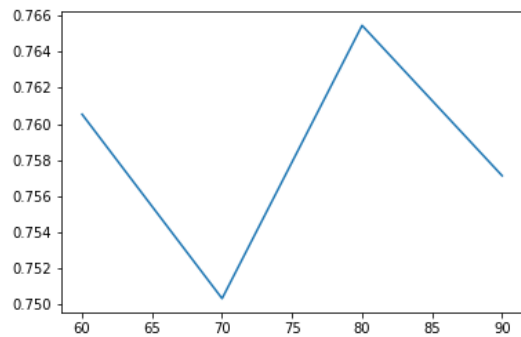
We can see that the training accuracies keep increasing while the test accuracy comes down when the number of hidden layers is 4. Thus we choose 3 hidden layers for this experiment.

Training accuracy: 0.7763

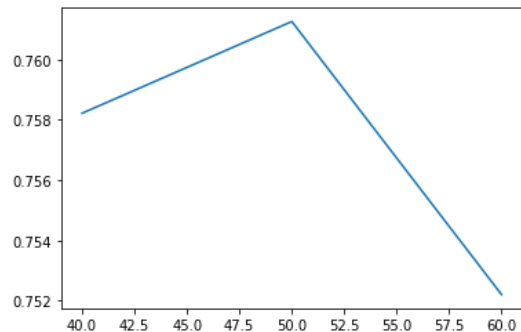
Test accuracy: 0.7641

Experimentation with number of neurons:

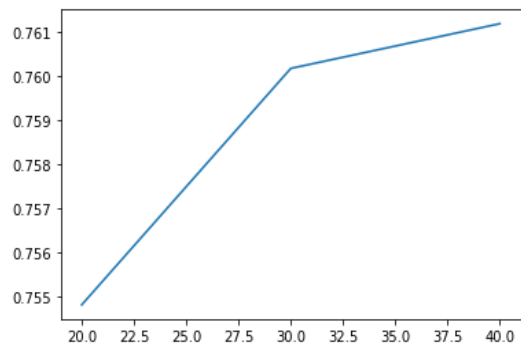
We try different combination of neurons in each of the 3 different hidden layers.



For the first hidden layer we experiment with neurons from 60-90. We see that we get highest accuracy when neurons=80 at 76.545%. Thus we fix neurons=80 for the first layer.



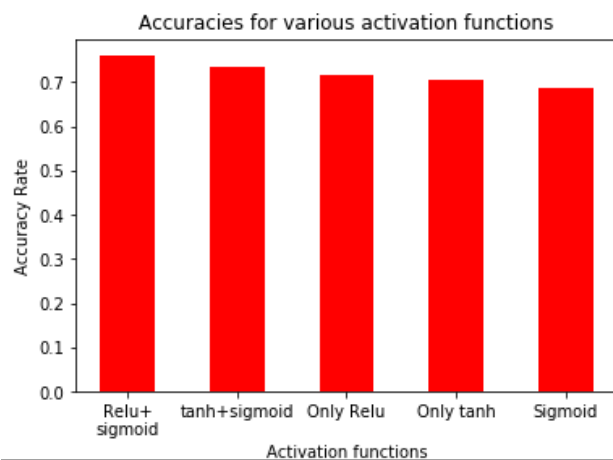
Next we change neurons from 40 to 60 for the second layer. We see that we can achieve maximum accuracy when neurons=50 at 76.126%



Finally we change neurons from 20-40 in the last hidden layer. We achieve maximum accuracy when neurons=40 at 76.12%.

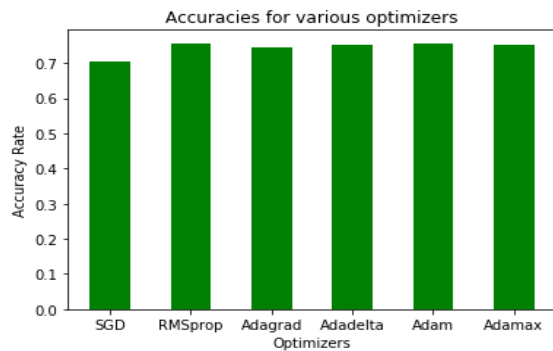
Thus the best combination is **80-50-40**.

Experimentation with activation functions:



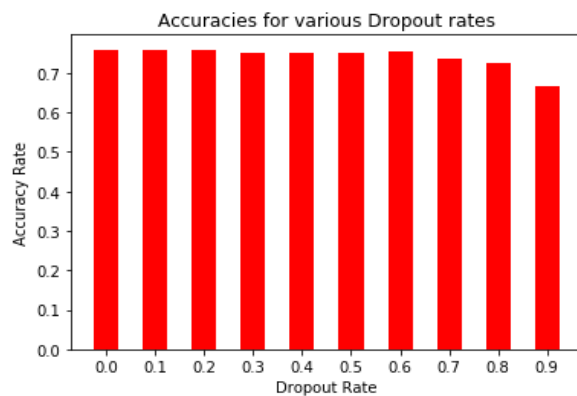
The different activation functions we use are **tanh, sigmoid and relu** as all these work very well with binary classifications. We normally go with sigmoid for the output layer but we also check for the output layer with tanh and relu functions. We can see that **relu for the hidden layers and sigmoid for the output layer** performs best of all with 76% accuracy.

Experimentation with optimization functions:



We experiment with different optimization functions such as SGD, RMSprop, Adagrad, Adadelta, Adam and Adamax. SGD performs worst and **RMSprop** performs the best of all with **75.67%** accuracy.

Experimentation with Dropout rates:



We do not want to over fit the data. Hence dropout is an option to reduce overfitting and hence we try dropout rates from 0.0 to 0.9. We can see that accuracies keep decreasing from 0.0 to 0.9. But the best accuracy is achieved when **dropout=0.1** at **75.91%**

Finally after choosing all the optimum hyper parameters for this model we train the model and test it with our test data. We get the below accuracies.

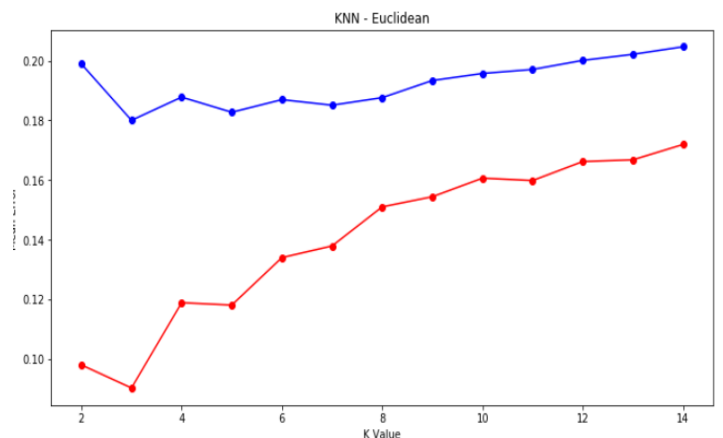
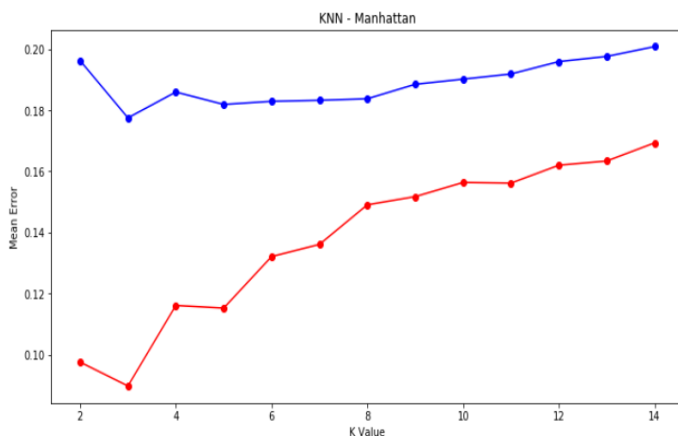
Train accuracy: 77.81%

Test accuracy: 76.74%

Algorithm 2: K Nearest Neighbors (KNN):

We use Scikit package for KNN

Experimentation with distance metrics and number of neighbors:



We are experimenting with Euclidean and Manhattan metrics as these are the widely used metrics. The number of neighbors are experimented from 2 to 15.

Though both give similar performance the Manhattan metric gives slightly better accuracy for all the neighbor values. Also we can see that after K=8, the model seems to generalize well for the test data as well. Till K=7, we can see a clear evidence of over fitting as training error is low and test error is high. Hence we choose **K=8 and Manhattan distance as our best parameters**. We achieve accuracy of **81.62%**.

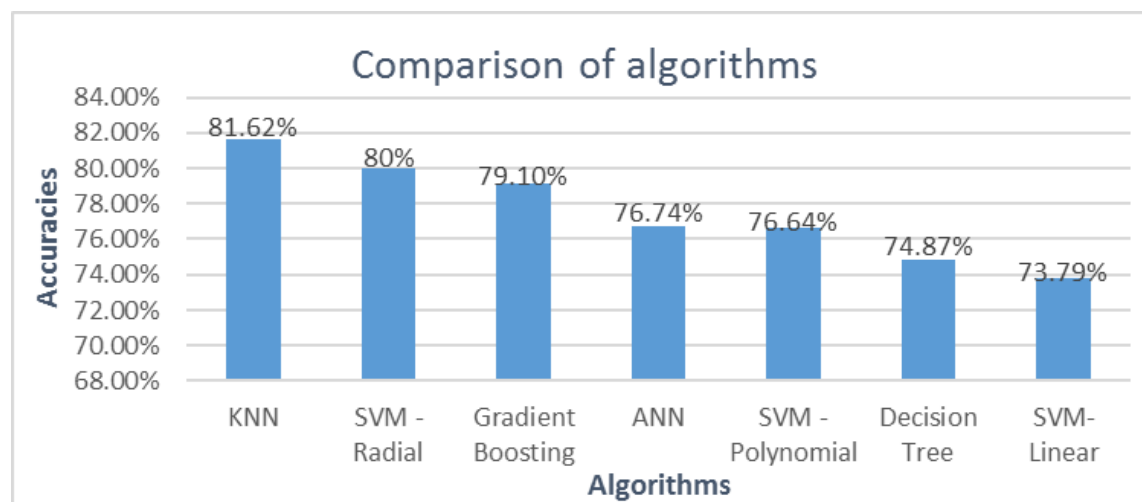
Comparison between the 2 algorithms:

For this dataset, neural networks gives accuracy of 76.74% whereas KNN gives accuracy of 81.62%. Thus KNN performs best for this data. The reason for lower accuracy on neural networks maybe due to more number of neurons and hidden layers, thus increasing the complexity of the model. This may be due to my limited knowledge and experience with choosing the number of neurons and hidden layers. On further research I found out that even a single hidden layer with limited neurons is sufficient in most cases.

Comparison with previous algorithms:

In the previous assignment, the best accuracy was given by SVM Radial with 80%. But KNN beats the performance of SVM. The reason for this is because KNN is not model based. It does not lose any detail and compares every training sample to give the prediction. Hence testing performance will be good with KNN. Thus we conclude that KNN gives the best accuracy on this dataset.

Ranking of Algorithms:



What I could have done?

The performance of neural networks is baffling me as it expected it to do the best of all. I could have used lesser number of neurons like 10 -15 with just one or two hidden layers. Also I could have used Grid search CV with a combination of neurons and hidden layers instead of choosing neurons after fixing hidden layers. Maybe this would have resulted in a better accuracy.

DATASET 2

RAIN IN AUSTRALIA

The objective is to predict whether rain will come tomorrow or not based on the location's weather conditions in Australia. The dataset provides us with various variables such as Wind speed, humidity, temperature, pressure etc. Occurrence of rain is indicated by 1 and non-occurrence is indicated by 0. I have chosen all the continuous variables for my analysis.

Independent Variables: Mintemp, Maxtemp, Rainfall, WindGustSpeed, WindSpeed9am, WindSpeed3pm, Humidity9am, Humidity3pm, Pressure9am, Pressure3pm, Cloud9am, Cloud3pm, Temp9am, Temp3pm

Dependent variable: Rain

I did random sampling of 20000 observations because the original sample contains 1.4 lakhs. 20000 is a good sample for training the model. Min-max normalization is done on the selected features.

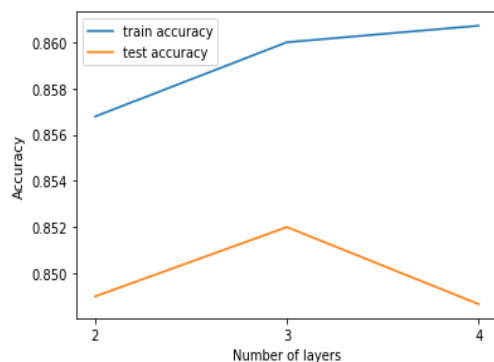
Algorithm 1 - Artificial Neural Networks:

Experimentation with number of layers:

Similar to the previous dataset, initially, we begin by experimenting with the number of hidden layers from 2 to 4. The number of neurons are selected randomly with the following combinations.

- 2 hidden layers : 50-40
- 3 hidden layers: 50-40-20
- 4 hidden layers: 70-50-40-30

We plot the test and train accuracies to understand the variations in the errors.

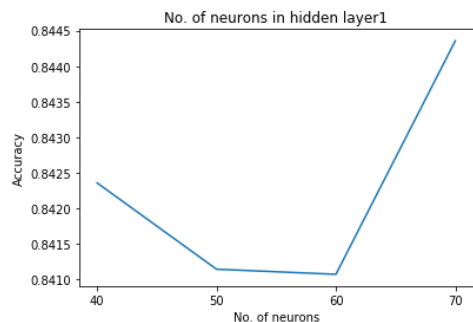


We can see that the test accuracy dips when number of layers is 4. The training accuracy is still high which means the model has over fit the data. Thus **3 hidden layers** seems to be a good bias variance tradeoff and we choose that.

Training accuracy: 86%

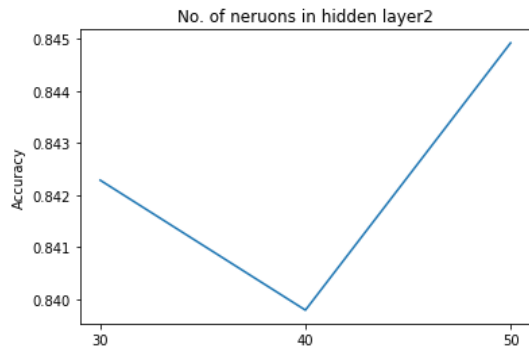
Test accuracy: 85.2%

Experimentation with number of neurons:

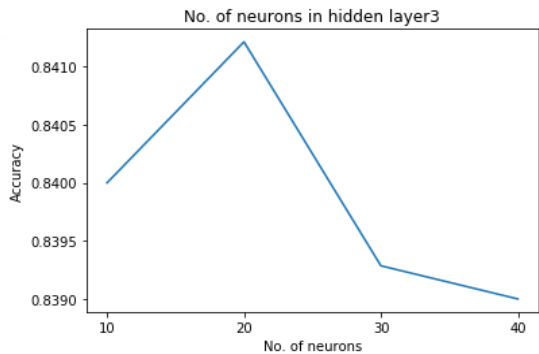


We try different combination of neurons in each of the 3 different hidden layers.

For the first hidden layer we experiment with neurons from 60-90. We see that we get highest accuracy when neurons=70 at 84.44%. Thus we fix neurons=70 for the first layer.



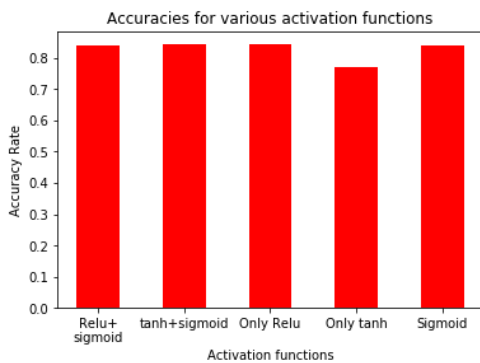
Next we change neurons from 40 to 60 for the second layer. We see that we can achieve maximum accuracy when neurons=50 at 84.49%



Finally we change neurons from 10-40 in the last hidden layer. We achieve maximum accuracy when neurons=20 at 84.12%.

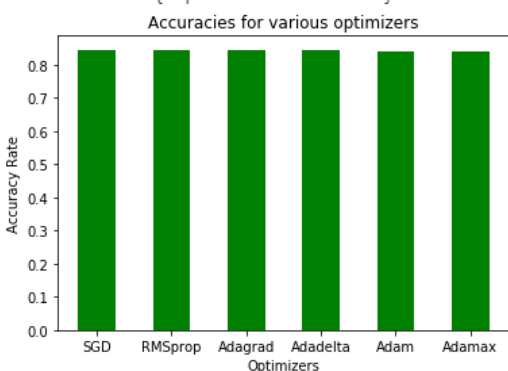
Thus the best combination is **70-50-20**.

Experimentation with activation functions:



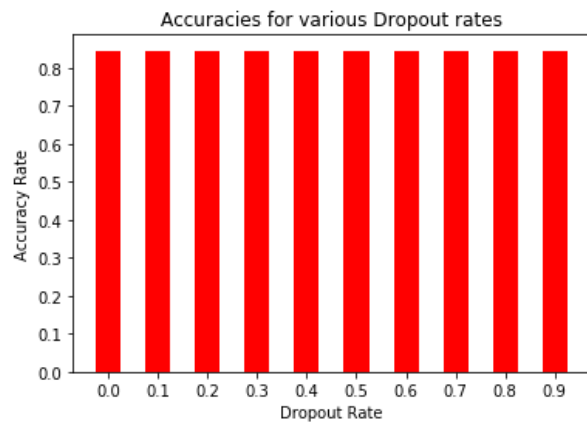
The different activation functions we use are **tanh, sigmoid and relu** as all these work very well with binary classifications. We normally go with sigmoid for the output layer but we also check for the output layer with tanh and relu functions. We can see that **tanh for the hidden layers and sigmoid for the output layer performs best of all with 84.31% accuracy.**

Experimentation with optimization functions:



We experiment with different optimization functions such as **SGD, RMSprop, Adagrad, Adadelata, Adam and Adamax**. All the optimizers perform almost similar. But the highest accuracy is achieved by **Adagrad at 84.48%.**

Experimentation with Dropout rates:



We do not want to over fit the data. Hence dropout is an option to reduce overfitting and hence we try dropout rates from 0.0 to 0.9. We can see that accuracies are almost similar with minute differences. But the best accuracy is achieved when **dropout=0.3 at 84.57%**.

Finally after choosing all the optimum hyper parameters for this model we train the model and test it with our test data. We get the below accuracies.

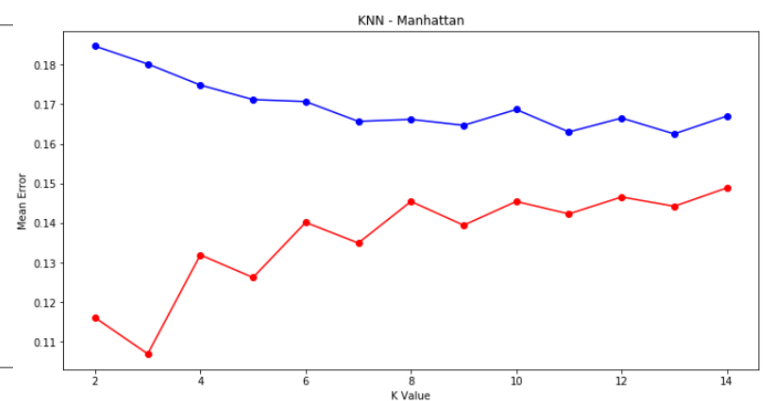
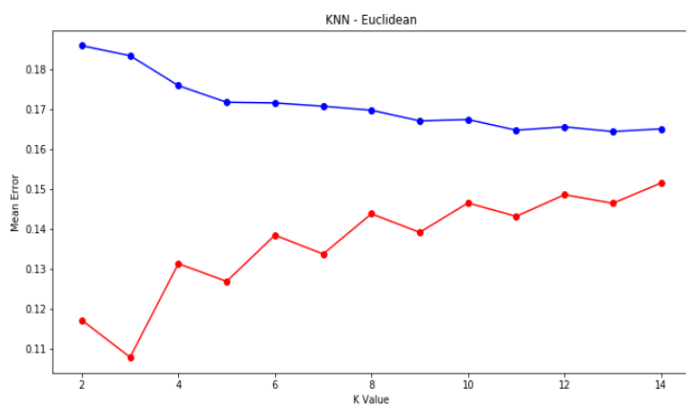
Train accuracy: 84.99%

Test accuracy: 84.95%

Algorithm 2: K Nearest Neighbors (KNN):

Experimentation with distance metrics and number of neighbors:

We are experimenting with Euclidean and Manhattan metrics as these are the widely used metrics. The number of neighbors are experimented from 2 to 15.



Though both give similar performance the Manhattan metric gives slightly better accuracy for all the neighbor values. Also we can see that after K=8, the model seems to generalize well for the test data as well. Till K=7, we can see a clear evidence of over fitting as training error is low and test error is high. Hence we choose **K=8 and Manhattan distance** as our parameter. We achieve maximum accuracy of **83.38%**.

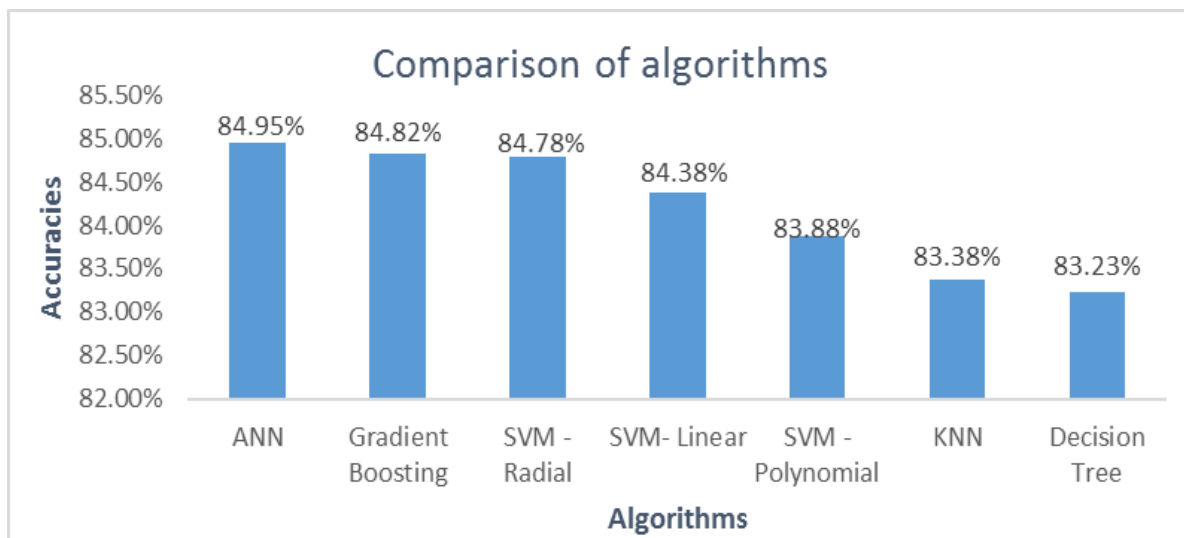
Comparison between the 2 algorithms:

For this dataset, Neural Networks gives 84.95% accuracy and KNN gives 83.38%. Neural networks performs well for this dataset because the choice of neurons and number of hidden layers might be good to handle the data. Also the number of features is 16 which is more than the features we used in the first dataset. Hence ANN is better equipped to handle more number of features than KNN.

Comparison of all algorithms:

For the previous assignment, the highest accuracy was given by Gradient Boosting at 84.82%. **But ANN outperforms it by 0.13% making it the best performer for this dataset.** KNN's poor performance might be due to curse of dimensionality because we use 16 features which is generally high. Hence KNN does not do well.

Ranking of the algorithms:



What I could have done?

Feature selection could have been done to select the 10 most important features as this would have helped to increase the accuracy of KNN to a large extent. Also the accuracy for ANN might have also improved.