# Assignment 3 (10 mark)
## PART 1 - K-Means  - 3 marks

Data Mining in databases  means analyzing  records and related attributes for finding patterns, trends and relationships to form model to be used in business decision making, prediction, and etc.  One of the  data mining method is clustering and a common method of clustering is k-Means. General k-Means algorithm when applies on database  data is  shown below which is from Chapter 28th of Database book by Elmasri:

**Algorithm 28.4.** $k$-Means Clustering Algorithm

**Input:** a database $D$, of $m$ records, $r_1$, ..., $r_m$ and a desired number of clusters $k$

**Output:** set of $k$ clusters that minimizes the squared error criterion

**Begin**
    randomly choose $k$ records as the centroids for the $k$ clusters;
    repeat
    assign each record, $r_i$, to a cluster such that the distance between $r_i$
        and the cluster centroid (mean) is the smallest among the $k$ clusters;
    recalculate the centroid (mean) for each cluster based on the records
        assigned to the cluster;
    until no change;
**End;**

**a)** First use the above algorithm and by using a common similarity metric distance between a record, and also using a value of 3 for K, manually cluster the data of following table. You can assume that the records with RIDs 1, 3, and 5 are used for the initial cluster centroids (means). Try to follow the algorithm  and calculate the centroids in a way that clusters to be optimum. Print the cluster members and centroid information.

| RID | Dimension 1 | Dimension 2 |
|-----|-------------|-------------|
| 1   | 8           | 4           |
| 2   | 5           | 4           |
| 3   | 2           | 4           |
| 4   | 2           | 6           |
| 5   | 2           | 8           |
| 6   | 8           | 6           |

**b)** Now you need to implement the manual k-means clustering you did above, on a software to verify and visualize your results. You can use any software (R or Python) for generating a graphical view of the clusters or a  free trail of SPSS platform  at http://spss.en.softonic.com  (or any platform of your choice) and answer the following questions:

1. Use non hierarchical k-means (k=3) and show clusters' memberships of the above dataset.
2- Use hierarchical clustering with dendrogram which  is a cluster tree for visualization of hierarchical clustering. Print dendogram and explain how many cluster options we can have for a given cutting points.
*There will be bonus mark if you do the above clustering with large dummy data of one of the tables (with several attributes) of the database  that you created in Assignment1.

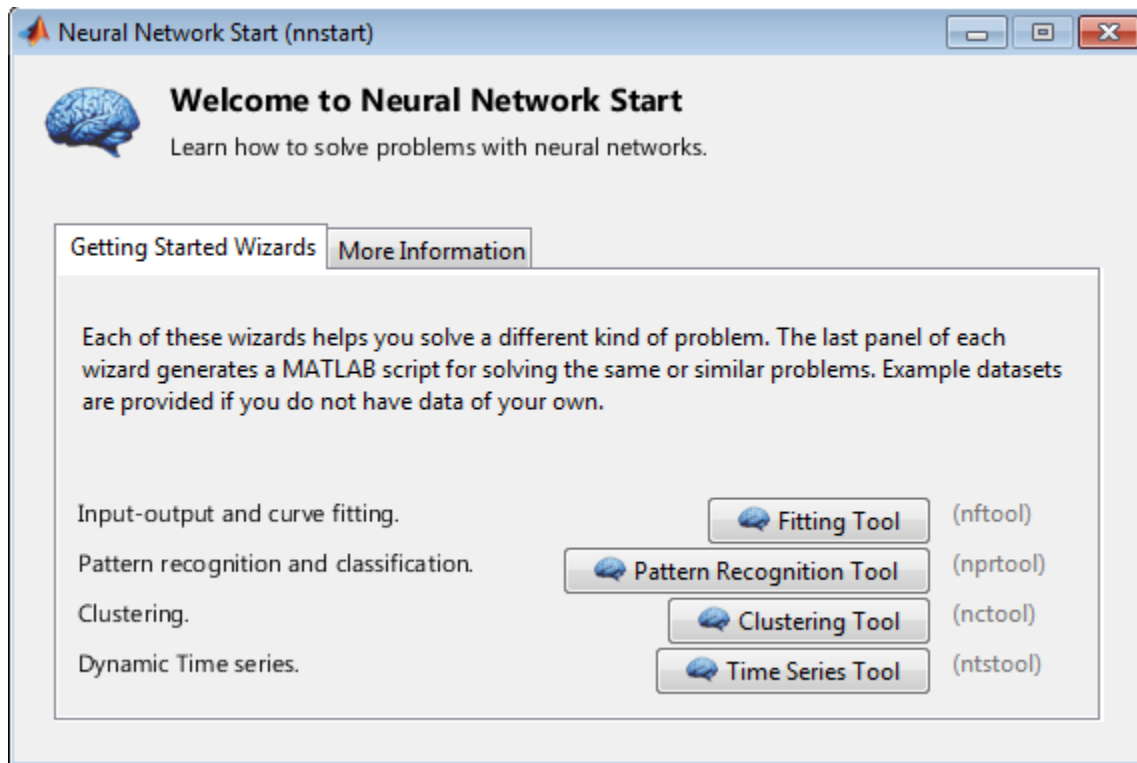# Part 2 – Supervised and Unsupervised ML models- 6 marks

**Section A- Supervised Machine Learning:** In this section, you will use the heart disease data set to train an artificial neural network as explained in next pages compare it with logistic regression analysis ( or any other nonlinear classifier of your choice) and then answer which method is better. First, use two models as the estimator (with the numerical result). Here you need to compare both methods by calculating Errors such as Mean Absolute Error (MAE) and other performance metrics to find which method can do prediction more accurately. Then use two models as the classifier and produce categorical result and confusion matrix.  If it is needed in heart disease dataset change the final result to a class with two labels of severe or non-severe (for example, by considering a threshold).Make sure you separate training set and testing data and there is no overfitting. Explain what is training errors and tests errors in your used methods. Try to use R or Python (the one that you used in Regression analysis) for implementing ANN. However because it is your first assignment of ANN you can use a simple interactive tool in Matlab to build ANN shown as appendix on the next page.. Please note the preference is to implement ANN with R or Python. The following figures in the appendix show different steps of building the model and training of ANN and how you can access the data set..

**Section B- Unsupervised Machine Learning:** In this section use an unsupervised learning method k-means clustering on hear disease data set. Create and visualize the k-means clusters (with k=5)  for the given heart disease dataset. For visualization, draw the scatter plot using the age  and cholesterol features on each group of clusters.  Apply k-means clusters on the heart disease dataset with varying numbers of clusters  from 1 to 10 and compute their corresponding Sum of squared Error (SSE) value. Plot the  graph and estimate the right "k" value . Use SSE vs the number of clusters  to estimate "k" value.  Finally, discuss if applying an unsupervise machine learning method such as clustering on the heart disease dat set can provide any additional information about this data set or not.
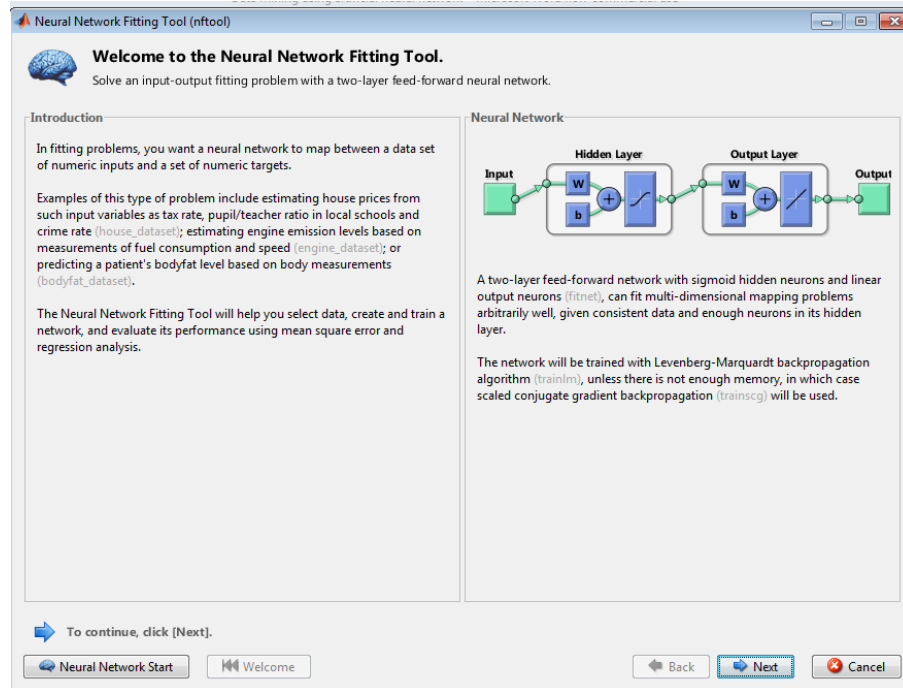
**Appendix related to Part2- Section A: An example shows how to use ANN in Matlab with Heart Disease Data Set**.
ANN model can be first trained using training data set according to following steps:
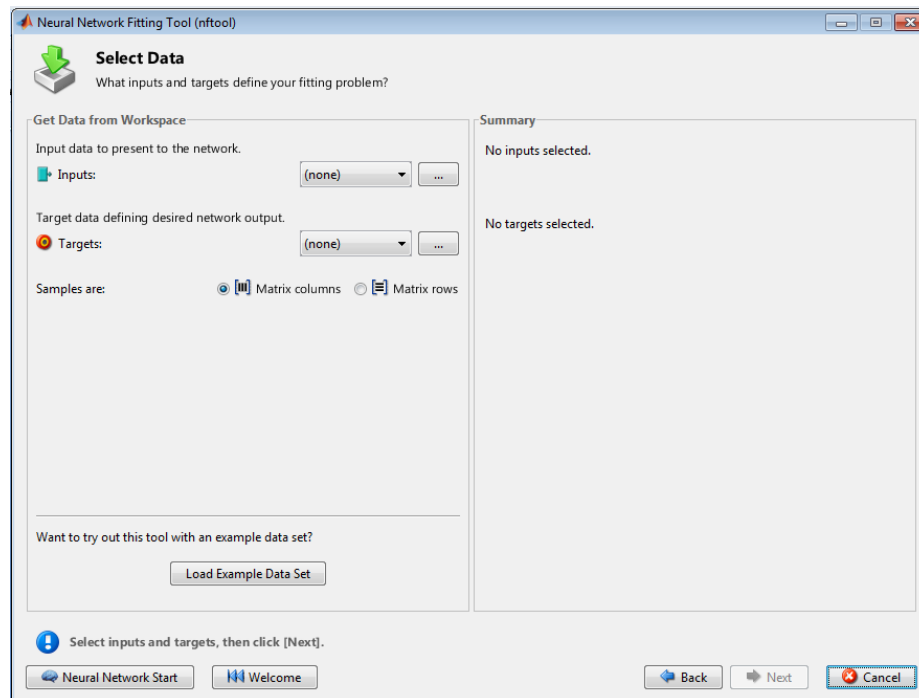
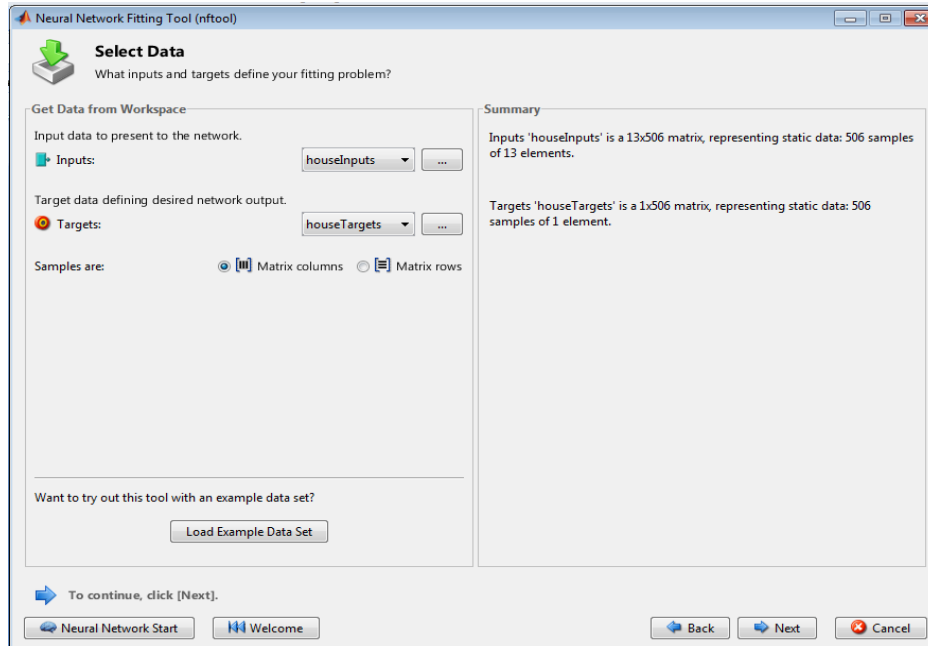1. In Matlab command prompt type **nnstart** and press enter



2. Select any tool which suits your ANN problem. In this example fitting tool was selected.
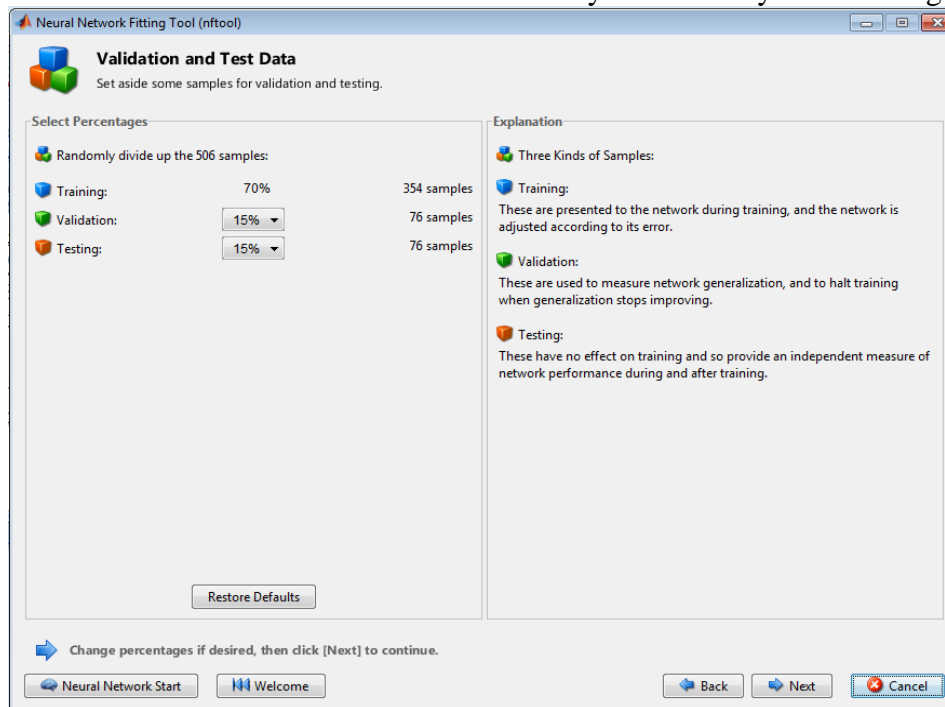
3.   Select data sets, example is loaded by using Load Example Data Set button.
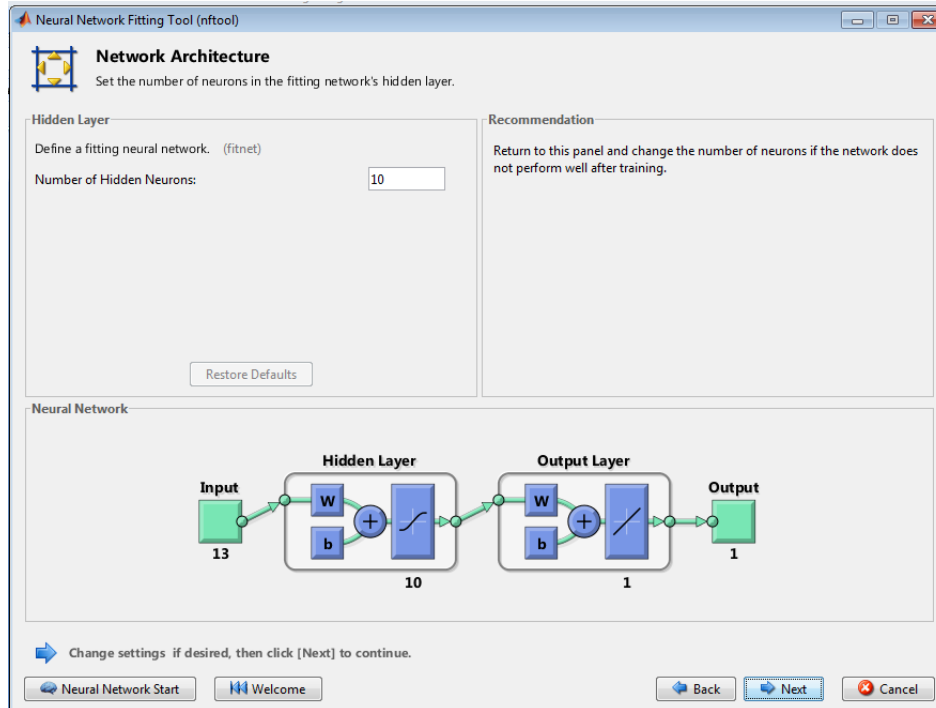


4.   Select Housing Pricing data set. You will notice that Input and target data field are filled with the selected data.
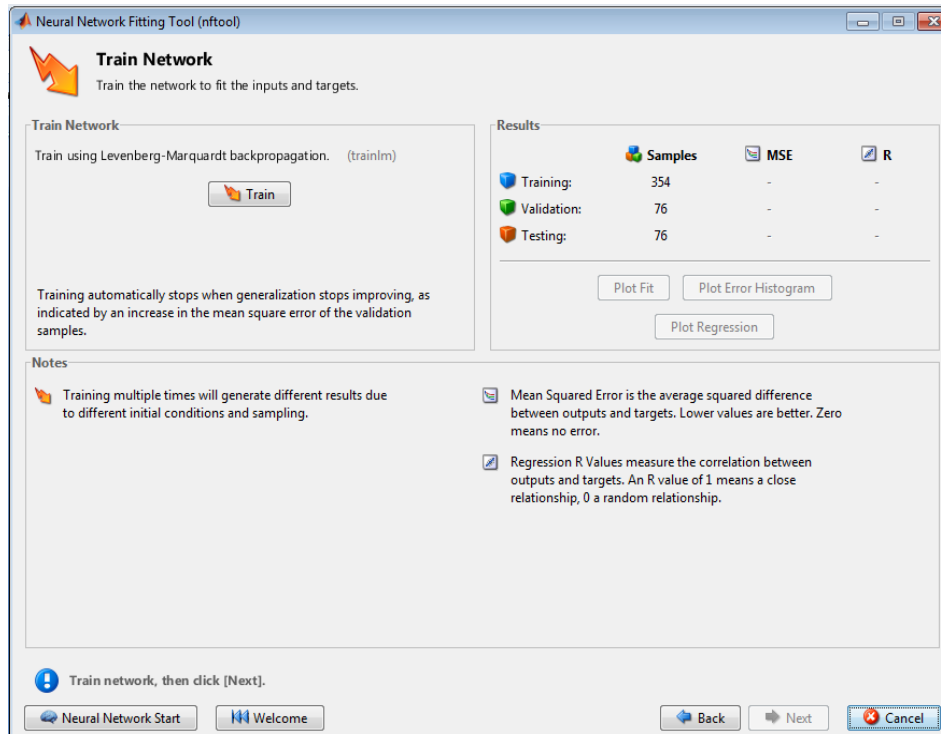
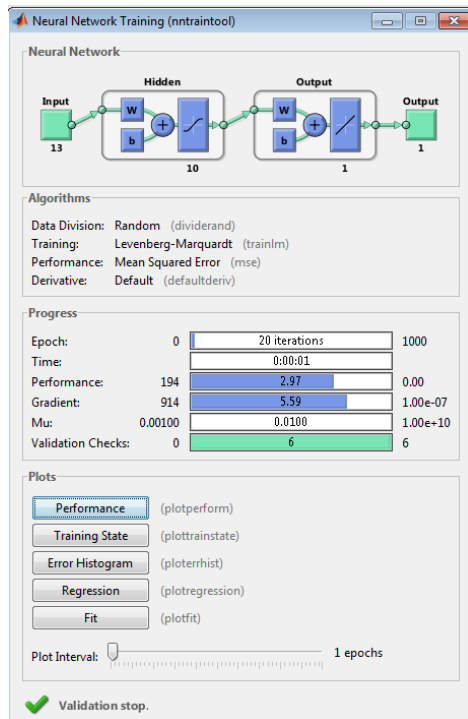5. Data set is divided as 70-15-15 ratio by default but you can change it as you like.

6. This example has 13 inputs and 1 output, therefore ANN model has 13 inputs with 1 output, with 1 hidden layer. Number of neuron in hidden layer can be increase or decrease.
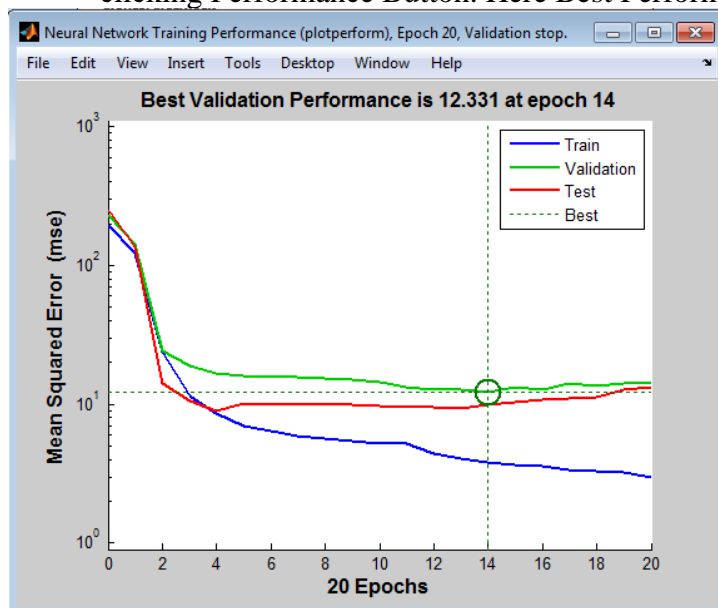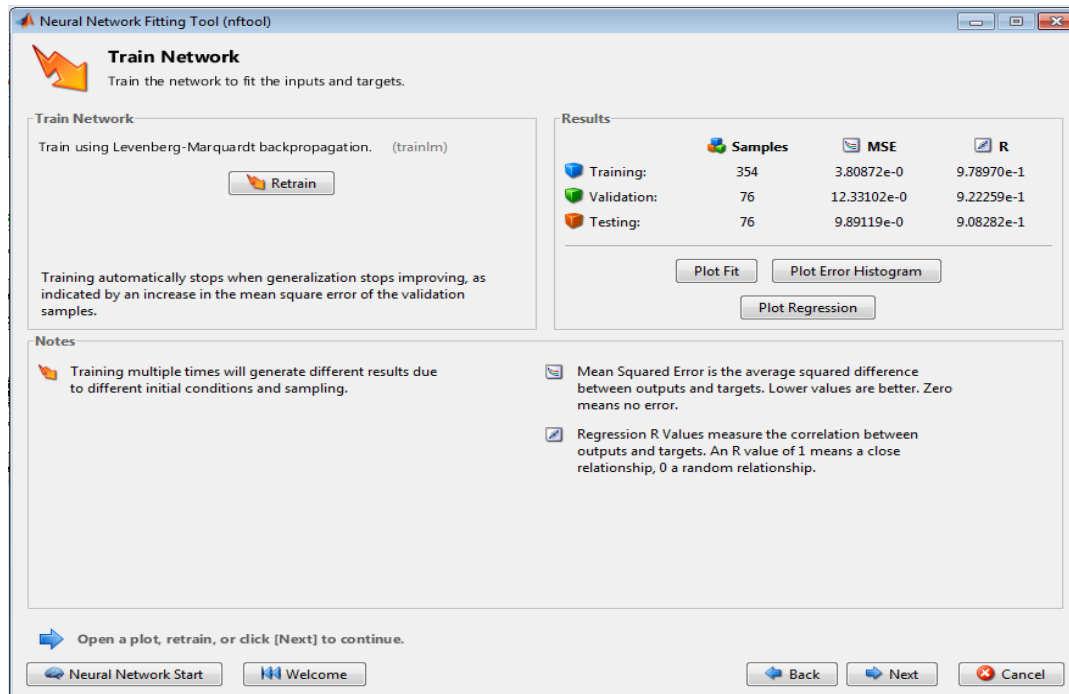


7. To train ANN model click Train button.

8.  Once you start training the ANN network, the following screen shows the progress.
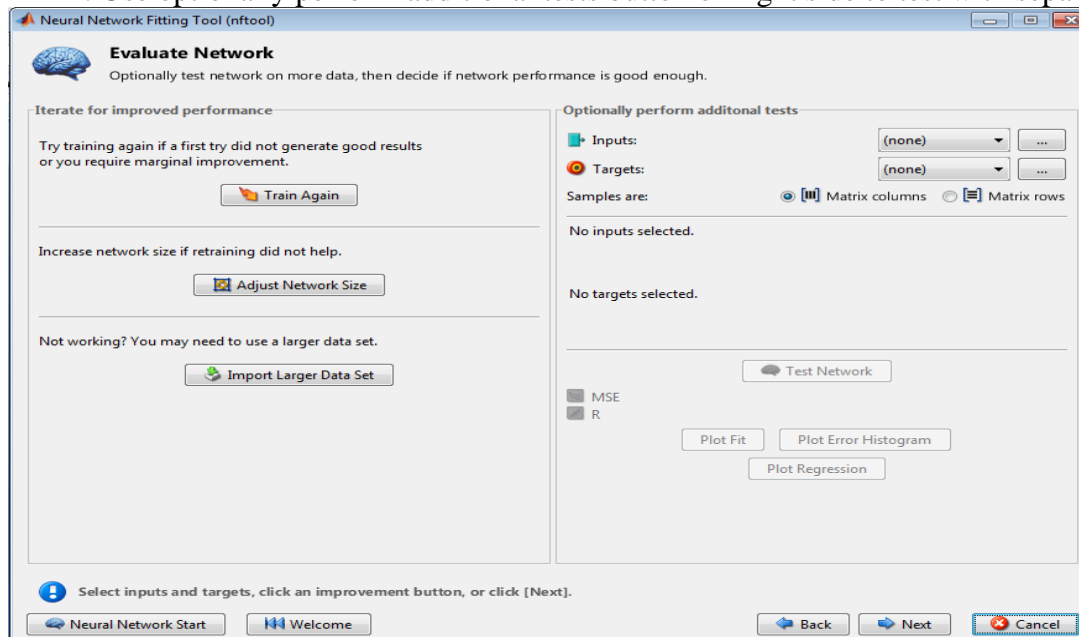


9.  When training is complete. The performance of trained ANN model can be seen by clicking Performance Button. Here Best Performance is achieved at 12.331 at epoch 14.
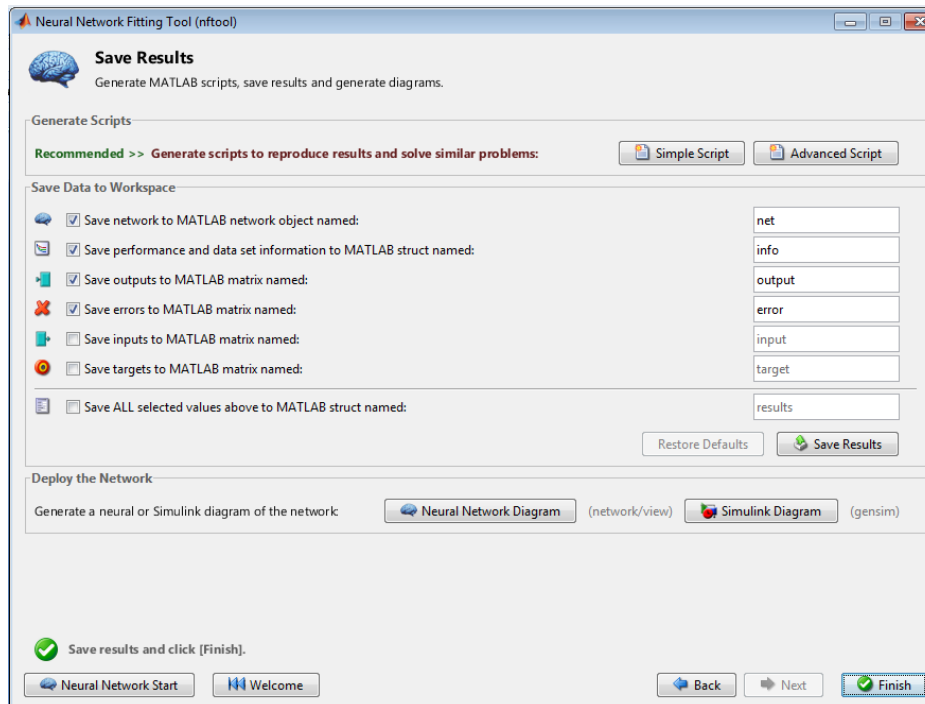
10. Once training is done and if you like the performance then you can test you newly build ANN model with the data which was not included in your training data set. You can also retrain the model. By the way you need to retrain your model multiple times to get best result.



11. Use optionally perform additional tests button on right side to test with separate data set.

12. All the results and the ANN model can be save for further use.



Show your ANN design and write a brief analysis of your ANN result for the above example

**Heart Disease Dataset:** To use ANN in Matlab or R or Python with the heart disease dataset of patients go to the following link:  http://archive.ics.uci.edu/ml/datasets/Heart+Disease

After going to this link you will find two folders: One: Data Folder and two: Dataset description. Data folder that has the dataset.   It is better to use processed cleveland data. In the dataset description folder, you will find the description about the columns' names referring to the14 column of the dataset as the following: The last one attribute (number 14) is the result. You want to use neural network that you created before for training and predicting the result attribute. After designing and training the neural net with an appropriate amount  of  data set (e.g. 70%),  test it with the remaining data and report the accuracy of the  prediction for both data mining techniques by Mean absolute error (MAE) and Root mean squared error (RMSE) and any other performance metrics that you think is useful.