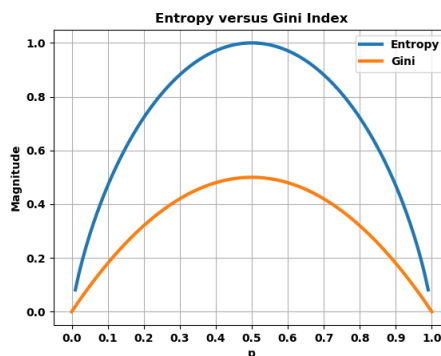


In this LAB, you will learn how to develop Decision Tree classifier manually and how to train the decision tree classifier for real datasets using the pre-pruning and post pruning technique. The LAB is divided into two phases. Set the `random_state = 5805` for all problems. Display numbers with 2-digit decimal precision.

**Phase I:**

1. Plot the entropy and Gini impurity (decision tree classifier) versus probability in one graph. Linewidth = 3. The final figure should be like below figure: [10pts]



2. Without use of python, develop a decision tree model that classifies the Play Tennis target in the following dataset. Graph the final tree and show all your work. Use the [Entropy](#) approach.
  - a. What is the prediction for the test observation [Outlook=sunny, Temp=cool, Humidity=high, Wind=strong]? Show your work. [20pts]
3. Without use of python, develop a decision tree model that classifies the Play Tennis target in the following dataset. Graph the final tree and show all your work. Use the [Gini Impurity](#) approach.
  - a. What is the prediction for the test observation [Outlook=sunny, Temp=cool, Humidity=high, Wind=strong]? Show your work. [20pts]



6. **Post-Pruning**: Using python program, find the optimum alpha in the cost complexity function. Plot the Accuracy score of the train and test set. Display the optimum alfa parameter on the console. With the optimum alpha, develop the pruned decision tree. Display the final tree on the console. Write your comment about the effect of post pruning on the performance of the tree comparing with the pre-pruned tree and no pruned tree.[10pts]
7. **Final classifier selection** Develop a Logistic Regression [LR] classifier for the titanic dataset using the same train and test set created in step 4. Hint: No need to hyperspace search. Use the default parameters for the LR classifier. Display the score of accuracy on the test and train test. [5pts]
8. **Final classifier selection** Create a table and compare the accuracy, confusion matrix, recall, AUC, and ROC [display the curve with the AUC as the legend for three cases: DT pre-pruned, DT post-pruned, and LR] of the pre-pruned, post-pruned tree and logistic regression classifier. Which classifier works better for this dataset? Justify your answer. [15pts]

Upload a formal **report (as a single pdf)** plus **the .py file** through BB by the due date.