

CA03 – Decision Tree Algorithm

1. Data Source and Contents

The dataset is obtained from the Census Bureau and represents salaries of people along with seven demographic variables. The following is a description of our dataset:

- **Number of target classes:** 2 ('>50K' and '<=50K') [Labels: 1, 0]
- **Number of attributes (Columns):** 7
- **Number of instances (Rows):** 48,842

The data is provided in a .csv file. Download the data and other files for this assignment from the following GitHub link. There is a column indicating the rows to be used as “Training Data” and “Testing Data”. You can split the files based on this column value.

<https://github.com/ArinB/MSBA-CA-03-Decision-Trees>

Note that the “continuous” data columns have been “transformed” into certain “data groups” or “data blocks”. This technique is called “binning” and by this process you can transform “continuous” data to “discrete categories”, because for certain applications like this “discrete categories” makes more sense than the exact value in the continuous scale. It’s also called “discretization” of data. **Investigate the data and answer the following question in your assignment** submission. (Probably, you will be able to answer these questions better AFTER you have completed the rest parts of the assignment.)

Q.1.1 Why does it makes sense to discretize columns for this problem?

Q.1.2 What might be the issues (if any) if we DID NOT discretize the columns.

2. Data Quality Analysis (DQA)

Do all of these inside your Notebook:

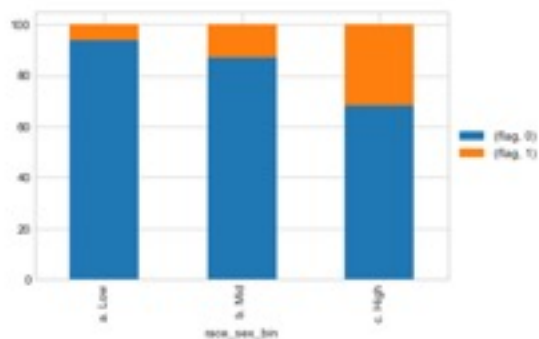
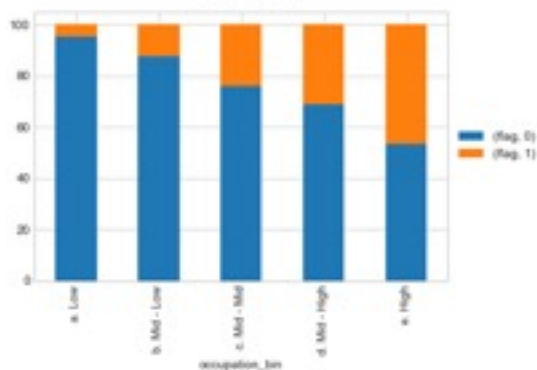
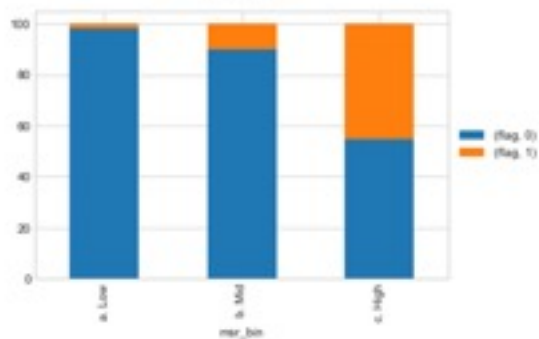
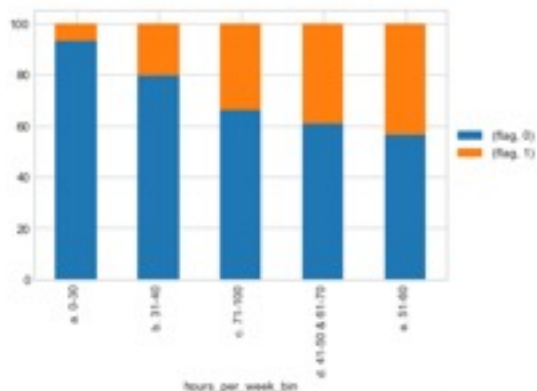
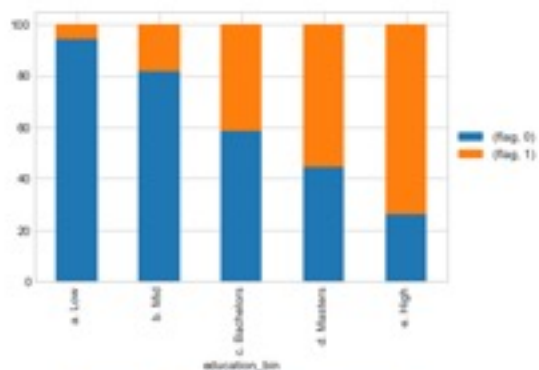
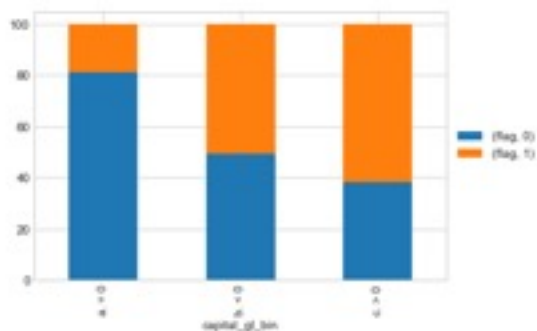
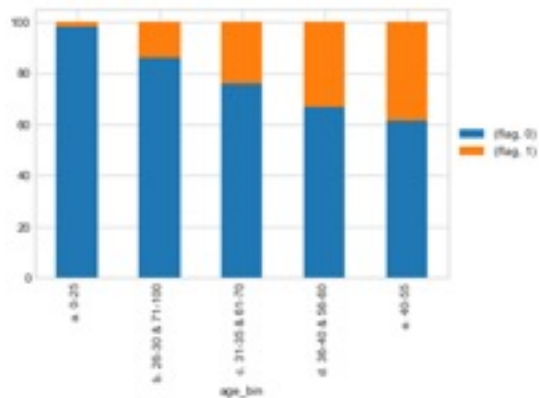
- Perform a Data Quality Analysis to find missing values, outliers, NaNs etc.
- Display descriptive statistics of each column
- Perform necessary data cleansing and transformation based on your observations from the data quality analysis

3. Exploratory Data Analysis (EDA)

Perform EDA of the income group with respect to the seven explanatory variables and display graphical representations as shown below. Do all of these inside your Notebook.

There are 7 explanatory variables:

1. Age (5 bins)
2. Capital Gain / Loss (3 bins)
3. Education (5 bins)
4. Hours per Week (5 bins)
5. Marriage Status and Relationship (3 bins)
6. Occupation (5 bins)
7. Race and Sex (3 bins)



4. Build Decision Tree Classifier Models

Definition: Given a data of attributes together with its classes, a decision tree produces a sequence of rules that can be used to classify the data.

Advantages: Decision Tree is simple to understand and visualise, requires little data preparation, and can handle both numerical and categorical data.

Disadvantages: Decision tree can create complex trees that do not generalize well, and decision trees can be unstable because small variations in the data might result in a completely different tree being generated.

Use “DecisionTreeClassifier” algorithm from scikit learn. Find details of sklearn tree algorithm below. Scikit Learn implements an optimized version of CART algorithm and can be used for binary class as well as multi-class classifications. It can be used for classification, as well as regression. [Study the following link thoroughly, including Section 1.10.5](#) (Tips on Practical Use).

<https://scikit-learn.org/stable/modules/tree.html#tree-algorithms-id3-c4-5-c5-0-and-cart>

Syntax to use the classifier is shown below:

```
from sklearn.tree import DecisionTreeClassifier
dtree = DecisionTreeClassifier(max_depth=10, random_state=101,
                              max_features = None, min_samples_leaf = 15)
dtree.fit(x_train, y_train)
y_pred=dtree.predict(x_test)
```

random_state: This is a number you choose arbitrarily. It’s also called “Random Seed”. If you use a number for this parameter (any number), it ensures that if you run the program multiple times, it will generate the same randomness. Hence, the solution becomes more “reproduceable”.

5. Visualize Your Decision Tree using GraphViz

Get the detail of how to do this from the following link:

<https://medium.com/@rnbrown/creating-and-visualizing-decision-trees-with-python-f8e8fa394176>

Do this inside your Notebook.

6. Evaluate Decision Tree Performance

Calculate and display the following. Do all of these inside your Notebook.

- Confusion Matrix (TP, TN, FP, FN ... etc.)
- Accuracy, Precision, Recall, F1 Score, AUC Value, ROC Curve (graph)

7. Tune Decision Tree Performance

Learn about all hyper-parameters and methods of Scikit Learn DecisionTreeClassifier algorithm at:

<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html#sklearn.tree.DecisionTreeClassifier>

Try varying **FOUR** of the hyperparameters manually, as per the following table, and train / score your model for each set of these hyperparameters. Record your Tree's performance with respect to each of these sets of hyperparameters in the Model Performance section of the following table.

Four Hyperparameters to vary:

- Split Criteria – 'Entropy' or 'Gini Impurity'
- Minimum Sample Split – Minimum number of records required in any node for a further split to be attempted
- Minimum Sample Leaf – Minimum of samples in a leaf node to stop further splitting (becomes a leaf node)
- Maximum Depth – Maximum depth of the tree allowed

Based on the ranges of values you find in your Decision Tree for these hyperparameters, vary them by manually choosing your own reasonable values and manually running the program for each of the Hyper-parameter combinations. An Excel file for this table is provided (Tree Tuning Cases.xlsx). You can enter the values in the given Excel File (Tree Tuning Cases.xlsx) and cut/paste the table in your submission document.

Q.7.1 Decision Tree Hyper-parameter variation vs. performance

Decision Tree Hyperparameter Variations Vs. Tree Performance							
Hyperparameter Variations				Model Performance			
Split Criteria (Entropy or Gini)	Minimum Sample Split	Minimum Sample Leaf	Maximum Depth	Accuracy	Recall	Precision	F1 Score
Entropy	Split Value 1	Leaf Value 1	Depth 1				
	Split Value 1	Leaf Value 1	Depth 2				
	Split Value 1	Leaf Value 2	Depth 1				
	Split Value 2	Leaf Value 1	Depth 1				
Gini Impurity	Split Value 1	Leaf Value 1	Depth 1				
	Split Value 1	Leaf Value 1	Depth 2				
	Split Value 1	Leaf Value 2	Depth 1				
	Split Value 2	Leaf Value 1	Depth 1				

8. Conclusion

Explain your observations from the above performance tuning effort.

Q.8.1 How long was your total run time to train the model?

Q.8.2 Did you find the BEST TREE?

Q.8.3 Draw the Graph of the BEST TREE Using GraphViz

Q.8.4 What makes it the best tree?

9. Automation of Performance Tuning

Can you “automate” generation of the above “Hyperparameter Vs. Model Performance Table” in your code itself and display the table in your Notebook?

(You can enter the “Hyperparameter Variation” part of the table manually in an Excel file, save it as CSV with header and read it in a dataframe in Python to repeatedly train/score the model for each of the lines of hyperparameters dataframe.)

10. Prediction using your “trained” Decision Tree Model

Based on the Performance Tuning effort in the previous section, **pick your BEST PERFORMING TREE**. Now **make prediction of a “new” individual’s** Income Category (<=50K, or >50K) with the following information. Do this in your Notebook.

- Hours Worked per Week = 48
- Occupation Category = Mid - Low
- Marriage Status & Relationships = High

- Capital Gain = Yes
- Race-Sex Group = Mid
- Number of Years of Education = 12
- Education Category = High
- Work Class = Income
- Age = 58

Q.10.1 What is the probability that your prediction for this person is accurate?

11. Deliverables

Your assignment outputs will have the following components:

- (1) Fully functional Notebook, Data, Readme file in a single folder at GitHub
- (2) Word / PDF document with answers to all questions marked Q.1, Q.2 etc.