

# CS 584-04: Machine Learning

Autumn 2019 Assignment 3

---

## Question 1 (20 points)

Please provide information about your Data Partition step.

- a) (5 points). Please provide the frequency table (i.e., counts and proportions) of the target variable in the Training partition?

Answer:

Frequency tables are a basic tool you can use to explore data and get an idea of the relationships between variables. A frequency table is just a data table that shows the counts of one or more categorical variables. One of the most useful aspects of frequency tables is that they allow you to extract the proportion of the data that belongs to each category. With a one-way table, you can do this by dividing each table value by the total number of records in the table:

CAR_USE			
Commercial	2652	2652	36.777146
Private	4559	4559	63.222854
Total	7211	7211	100.000000

- b) (5 points). Please provide the frequency table (i.e., counts and proportions) of the target variable in the Test partition?

Answer:

CAR_USE			
Commercial	2652	2652	36.777146
Private	4559	4559	63.222854
Total	7211	7211	100.000000

- c) (5 points). What is the probability that an observation is in the Training partition given that  $CAR\_USE = Commercial$ ?

Answer:

$P\_Train = 0.7$ ;  $P\_Test = 0.3$

$P\_Comm\_train = y\_train[y\_train == 'Commercial'].count() / y\_train.shape[0] = 0.3677714602690334$

$P\_Comm\_test = y\_test[y\_test == 'Commercial'].count() / y\_test.shape[0] = 0.3678421222905209$

$Prob\_Training\_Commercial = P\_Comm\_train * P\_Train / ((P\_Comm\_train * P\_Train) + (P\_Comm\_test * P\_Test))$

Probability (Training|Commercial): 0.6999596538317057

- d) (5 points). What is the probability that an observation is in the Test partition given that CAR\_USE = Private?

Answer:

$P_{\text{Pri\_train}} = y_{\text{train}}[y_{\text{train}} == \text{'Private'}].\text{count}() / y_{\text{train}}.\text{shape}[0] = 0.6322285397309666$

$P_{\text{Pri\_test}} = y_{\text{test}}[y_{\text{test}} == \text{'Private'}].\text{count}() / y_{\text{test}}.\text{shape}[0] = 0.6321578777094792$

$\text{Prob\_Test\_Pri} = P_{\text{Pri\_test}} * P_{\text{Test}} / ((P_{\text{Pri\_train}} * P_{\text{Train}}) + (P_{\text{Pri\_test}} * P_{\text{Test}}))$

Probability (Test|Private): 0.29997652823125087

## Question 2 (40 points)

Please provide information about your decision tree.

- a) (5 points). What is the entropy value of the root node?

Answer: As we learnt in class, entropy is a measure of a split. Higher the entropy means the distribution is impure.

Entropy for root node is given as 0.9489621493401781

- b) (5 points). What is the split criterion (i.e., predictor name and values in the two branches) of the first layer?

Answer: The split criteria for Level 0 is on variable Occupation: with categories ('Blue Collar', 'Unknown', 'Student') on left node and remaining categories on right node.

We have chosen to split on variable Occupation because it has lowest entropy.

- c) (10 points). What is the entropy of the split of the first layer?

Answer: Entropy of split of the first layer is 0.7148805225259208

- d) (5 points). How many leaves?

Answer: There are 4 leaf nodes on Level 2 of the tree.

- e) (15 points). Describe all your leaves. Please include the decision rules and the counts of the target values.

Answer:

Decision rules are as follows:

Confidence of occurrence of Event:'Commercial' at leaf nodes

0.24647887323943662

0.8504761904761905

0.006151953245155337

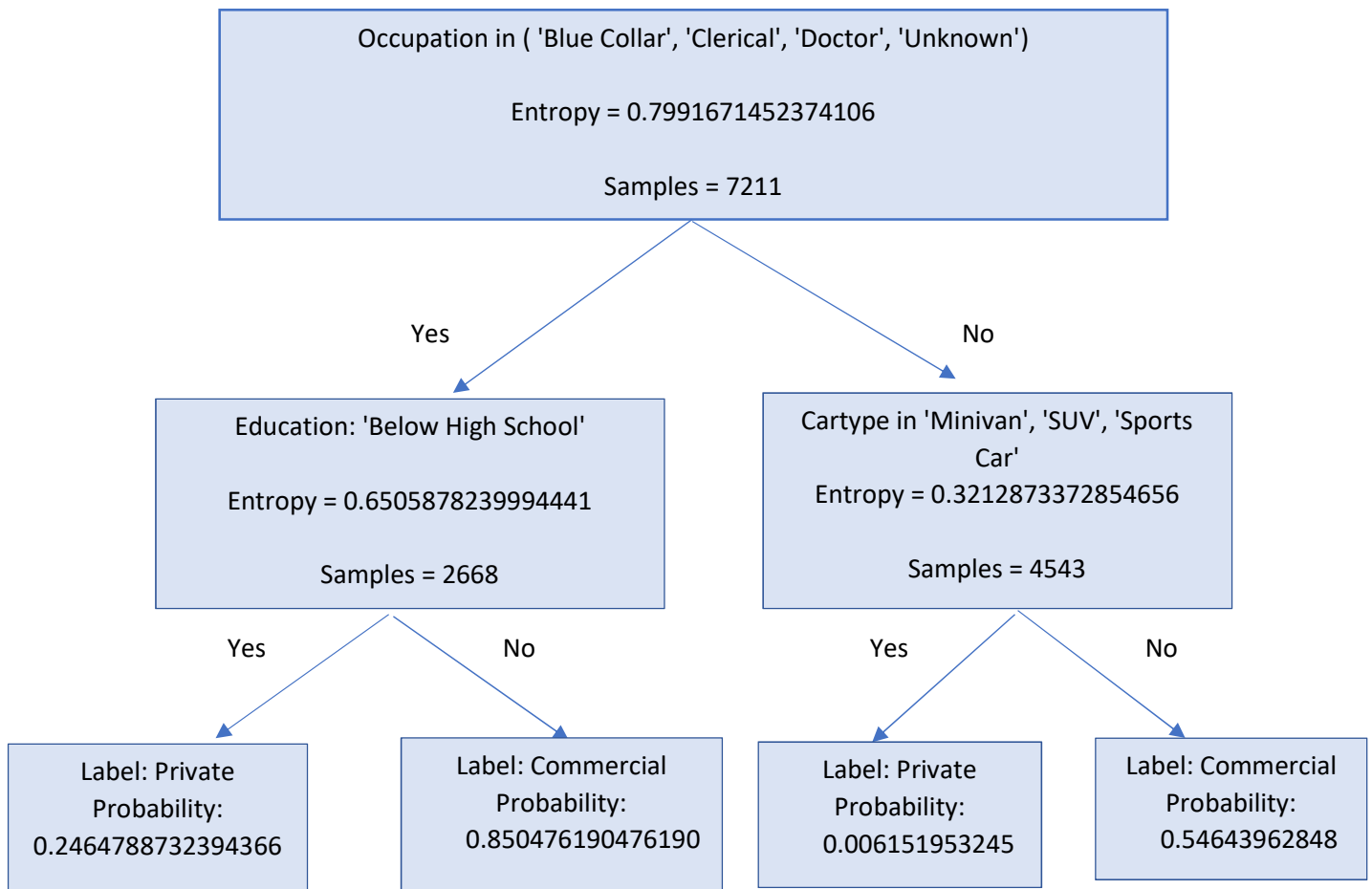
0.5464396284829721

Count of Target values on all the leaf nodes are as follows:

140  
1786  
20  
706

Threshold: 0.3677714602690334

Decision trees being formed is as follows:



Working of the tree:

Tree has 3 levels: Level 0, Level 1, and Level2

The Root node is split on Variable occupation.

On Level 1 there are two splits on variables: Education and Cartype.

Finally, the we calculate the decision rules and apply it on testing dataset. After this we get suitable predicted probabilities as per the rules. This dataset is then assigned suitable labels once the predicted probabilities are compared with the threshold.

### Question 3 (40 points)

Please apply your decision tree to the Test partition and then provide the following information.

- a) (10 points). Use the proportion of target Event value in the training partition as the threshold, what is the Misclassification Rate in the Test partition?

Answer:

An observation is misclassified if the observed target value is not equal to the predicted target value.

The misclassification rate formula is:  $(FP+FN)/(TP+TN+FP+FN)$

In our case misclassification rate is: 0.1708185053380783

- b) (10 points). What is the Root Average Squared Error in the Test partition?

Answer:

RASE in our case is: 0.3300251195213117

- c) (10 points). What is the Area Under Curve in the Test partition?

Answer:

Area Under Curve is: Area Under the curve 0.9114708659772841

- d) (10 points). Generate the Receiver Operating Characteristic curve for the Test partition. The axes must be properly labeled. Also, don't forget the diagonal reference line.

The curve is as follows:

