

Tut 7: Decision Tree Models

May/June 2022

Classification Tree

1. Use **gain ratio** to determine which split is better:

Split 1: Leaf $A = [20+, 15-]$; Leaf $B = [5+, 20-]$

Split 2: Leaf $A = [10+, 2-]$; Leaf $B = [15+, 33-]$

Remark: The larger “information gain” and “gain ratio”, the better.

2. (Jan 2022 Final Q4(b)) A classification tree is being constructed to predict whether the credit card application approval is positive. Consider the two splits below:

- **Split 1:** The left node has 178 observations with 68 positives and the right node has 295 observations with 144 positives.
- **Split 2:** The left node has 136 observations with 83 positives and the right node has 337 observations with 129 positives.

By calculating the information gains, determine which split is better.

(7 marks)

3. (May 2020 Final Q4(b)(ii)) In trying to build a model that is able to predict whether or not an email message is spam based on the following predictors:

- to_multiple: Indicator for whether the email was addressed to more than one recipient;
- image: Indicates whether any images were attached;
- attach: Indicates whether any files were attached;
- dollar: Indicates whether a dollar sign or the word ‘dollar’ or ‘ringgit’ appeared in the email;
- winner: Indicates whether “winner” appeared in the email;
- num_char: The number of characters in the email, in thousands;

- format: Indicates whether the email was written using HTML (e.g. may have included bolding or active links) or plaintext;
- re_subj: Indicates whether the subject started with “Re:”, “RE:”, “re:”, or “rE:”;
- number: Factor variable saying whether there was no number, a small number (under 1 million), or a big number.

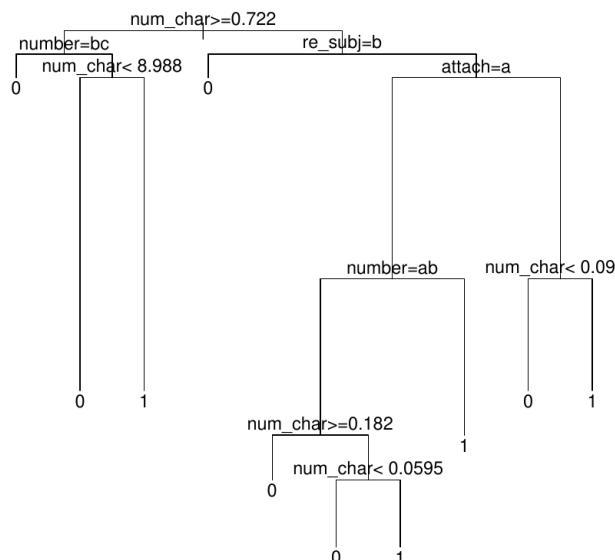
Note that “spam” is denoted with the value 1 while “non-spam” is denoted with the value 0. The trained logistic regression model has the parameters given in Figure 4.2.

Table 4.2: Coefficients of Logistic Regression

| Coefficients: | | | | | |
|----------------|-----------|------------|---------|----------|------|
| | Estimate | Std. Error | z value | Pr(> z) | |
| (Intercept) | -1.468478 | 0.181285 | -8.100 | 5.48e-16 | *** |
| to_multipleyes | -2.152057 | 0.349538 | -6.157 | 7.42e-10 | *** |
| imageyes | -1.467843 | 0.797895 | -1.840 | 0.065820 | . |
| attachyes | 0.957716 | 0.281455 | 3.403 | 0.000667 | *** |
| num_char | -0.014651 | 0.007199 | -2.035 | 0.041849 | * |
| dollaryes | 0.453477 | 0.197009 | 2.302 | 0.021346 | * |
| winneryes | 1.994563 | 0.392252 | 5.085 | 3.68e-07 | *** |
| numbersmall | -1.227981 | 0.186300 | -6.591 | 4.36e-11 | *** |
| numberbig | -0.561313 | 0.263563 | -2.130 | 0.033195 | * |
| formatPlain | 1.032511 | 0.171915 | 6.006 | 1.90e-09 | *** |
| re_subjyes | -2.447223 | 0.398309 | -6.144 | 8.05e-10 | *** |
| --- | | | | | |
| Signif. : | 0 | ‘***’ | 0.001 | ‘**’ | 0.01 |
| | | | ‘*’ | 0.05 | ‘.’ |
| | | | | 0.1 | ‘ ’ |
| | | | | | 1 |

If an email does not address to multiple, has no image, no attached file(s), no “dollar” sign, does not have the word “winner”, has 20.133×10^3 number of characters and is in HTML format, has no subject starting with “Re:” and has a small number in the email. **Determine** whether the email is a spam using the trained logistic regression model and using the decision tree model (you will need to interpret the decision tree model based on your knowledge of “rpart” algorithm) given in Figure 4.3.

Figure 4.3: The trained decision tree model.



(4.5 marks)

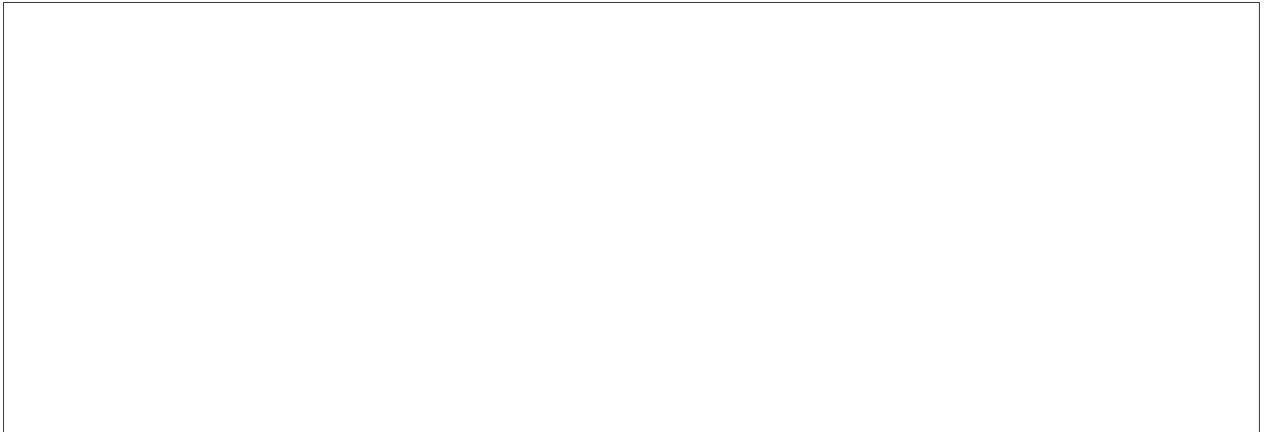


4. (Jan 2021 Final Q2(a)) The dataset in Table 2.1 is used to build a classification tree which predicts if a student pass predictive modelling (Pass or Fail, P, F for short), based on their previous GPA (High, Medium, or Low, H, M, L for short) and whether they have or have not (Y or N in short) put in significant efforts in their study.

Table 2.1: Training dataset for classification problem.

| GPA | Studied | Pass |
|-----|---------|------|
| L | N | F |
| L | Y | P |
| M | N | F |
| M | Y | P |
| H | N | P |
| H | Y | P |

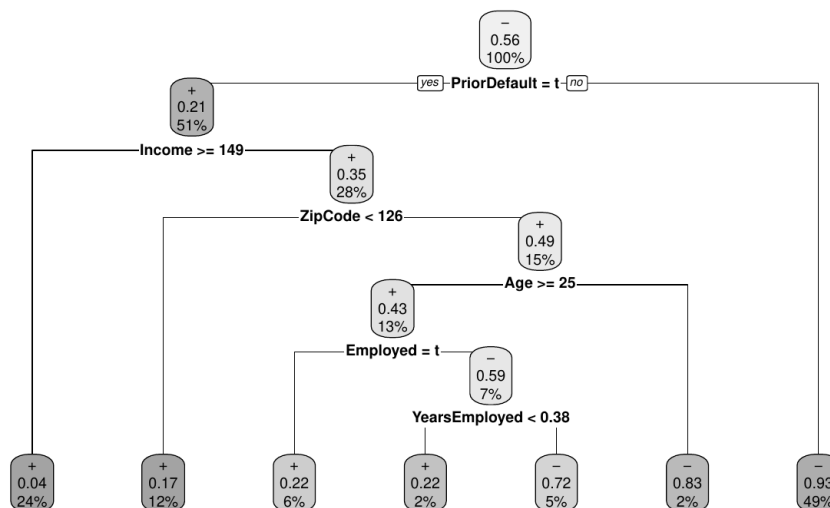
Construct and plot the ID3 classification tree (using information gain) with appropriate labels. You must show all the calculation steps. (5 marks)





5. (Jan 2022 Final Q2(b)) For the same training data (as Tutorial 4 Q1, i.e. Jan 2022 Final Q2(b)), use the CART tree in Figure 2.1 to predict the the credit card application being approved (positive or negative) for a male of age 22.08 with a debt of 0.83 unit who has been employed for 2.165 years with no prior default and is currently unemployed, has a credit score 0 and a zip code 128 with income 0.

Figure 2.1: CART tree for credit card application approval data.



You need to show your workings by explaining the steps to move left or right in the tree traversal to reach the prediction. (4 marks)

6. (Jan 2022 Final Q2(c)) Compare the ability of the logistic regression model and the C4.5 tree model in the handling missing values and the prediction of highly nonlinear data. (4 marks)