

ISLR | Chapter 8 Exercises

Marshall McQuillen

10/10/2018

Conceptual

1

(sketch on following page)

2

Consider the regression setting; when each tree has depth one (meaning there is only one split), each tree will only consist of one split on one attribute, and the predicted response for those two groups will be the mean response (of the training data) of those two groups. Therefore, the base model will take the form $f_j(x_j) = \beta_j I(x_{ji} \leq t) + \beta'_j I(x_{ji} > t)$, where t is the value in the range of x_j that dictates the split (in order to minimize the RSS). Knowing this, one can walk through algorithm 8.2, as shown below:

$$f(\hat{x}) = 0 \tag{1}$$

$$f_1(x_1) = \beta_1 I(x_{1i} \leq t) + \beta'_1 I(x_{1i} > t) \tag{2}$$

$$f(\hat{x}) = f(\hat{x}) + \lambda f_1(x_1) \tag{3}$$

$$r_i = r_i - \lambda f_1(x_i) \tag{4}$$

$$f_2(x_2) = \beta_2 I(x_{2i} \leq t) + \beta'_2 I(x_{2i} > t) \tag{5}$$

$$f(\hat{x}) = f(\hat{x}) + \lambda f_2(x_2) \tag{6}$$

$$r_i = r_i - \lambda f_2(x_i) \tag{7}$$

$$\vdots \tag{8}$$

$$f_j(x_j) = \beta_j I(x_{ji} \leq t) + \beta'_j I(x_{ji} > t) \tag{9}$$

$$f(\hat{x}) = f(\hat{x}) + \lambda f_j(x_j) \tag{10}$$

$$r_i = r_i - \lambda f_j(x_i) \tag{11}$$

$$\vdots \tag{12}$$

$$f(x) = \sum_{j=1}^P f_j(X_j) \quad \text{where} \quad f_j(X_j) = \beta_j I(x_{ji} \leq t) + \beta'_j I(x_{ji} > t) \tag{13}$$

$$\tag{14}$$

Note that β_j is the mean of the response for all values of x_j that are less than or equal to t , and β'_j is the mean response for all values of x_j that are greater than t .

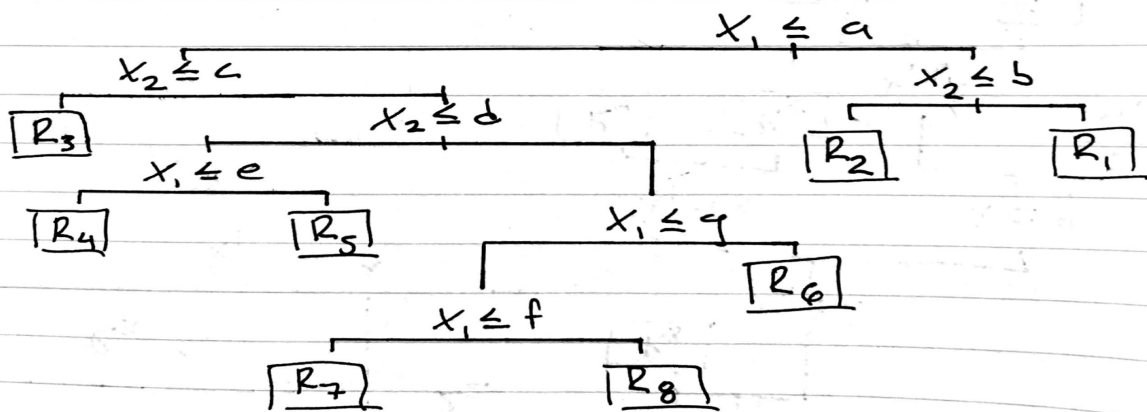
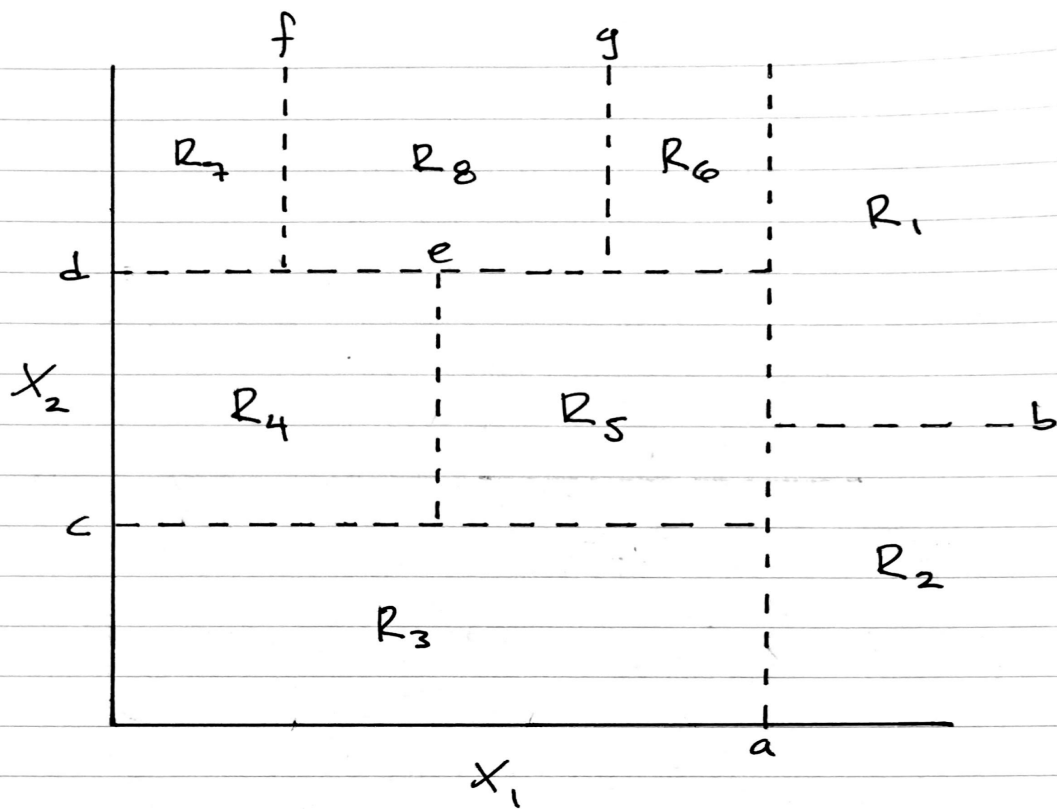


Figure 1: "Conceptual Exercise 1"

```

suppressPackageStartupMessages(library(ggplot2))
pm1 <- seq(0.01, 0.99, 0.01)

# defining functions
two_class_gini <- function(prob) {
  return((prob * (1-prob)) + (prob * (1-prob)))
}

two_class_entropy <- function(prob) {
  return(-1*(prob*log(prob) + (1 - prob)*log(1 - prob)))
}

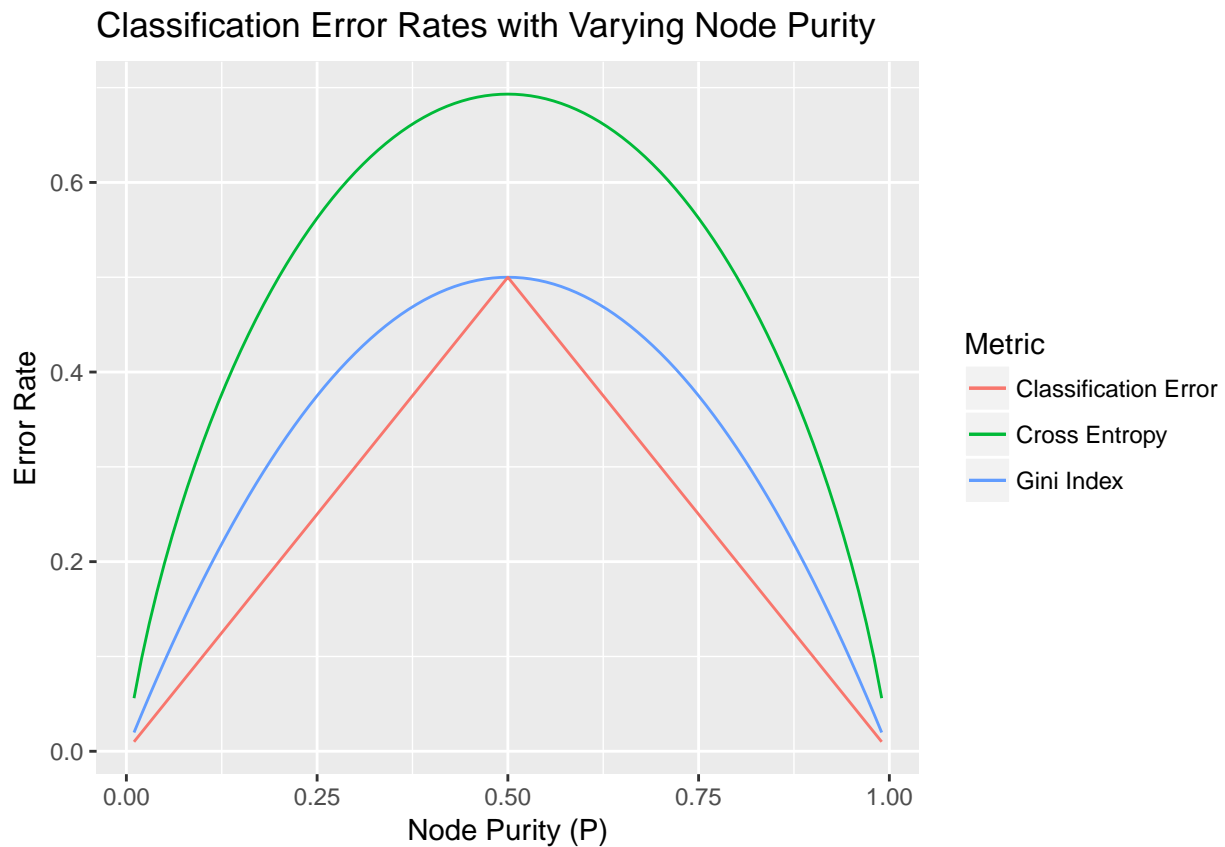
class_error <- function(prob) {
  return(1 - pmax(prob, (1 - prob)))
}

# error calculation
classification_error <- class_error(pm1)
gini_error <- two_class_gini(pm1)
entropy_error <- two_class_entropy(pm1)

# bring all data together for plotting
error_df <- data.frame(x = pm1,
                      gini = gini_error,
                      class = classification_error,
                      entropy = entropy_error)

# error plot
ggplot(error_df, aes(x = x, y = gini)) +
  geom_line(aes(color = 'Gini Index')) +
  geom_line(aes(y = entropy, color = 'Cross Entropy')) +
  geom_line(aes(y = class, color = 'Classification Error')) +
  guides(color = guide_legend(title = "Metric")) +
  ggtitle("Classification Error Rates with Varying Node Purity") +
  ylab("Error Rate") +
  xlab("Node Purity (P)")

```



4

(sketch on following page)

5

```
probs <- c(0.1, 0.15, 0.2, 0.2, 0.55, 0.6, 0.6, 0.65, 0.7, 0.75)
```

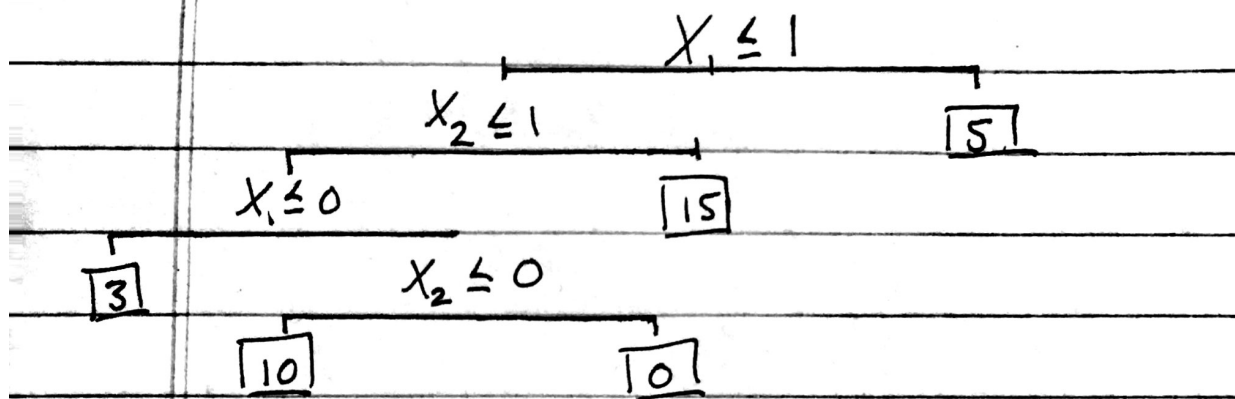
```
# majority vote chooses red
majority_vote <- ifelse(probs >= 0.5, 'red', 'green')
table(majority_vote)
```

```
## majority_vote
## green  red
##      4   6
```

```
# average chooses green
ifelse(mean(probs) >= 0.5, 'red', 'green')
```

```
## [1] "green"
```

4) a.



4) b.

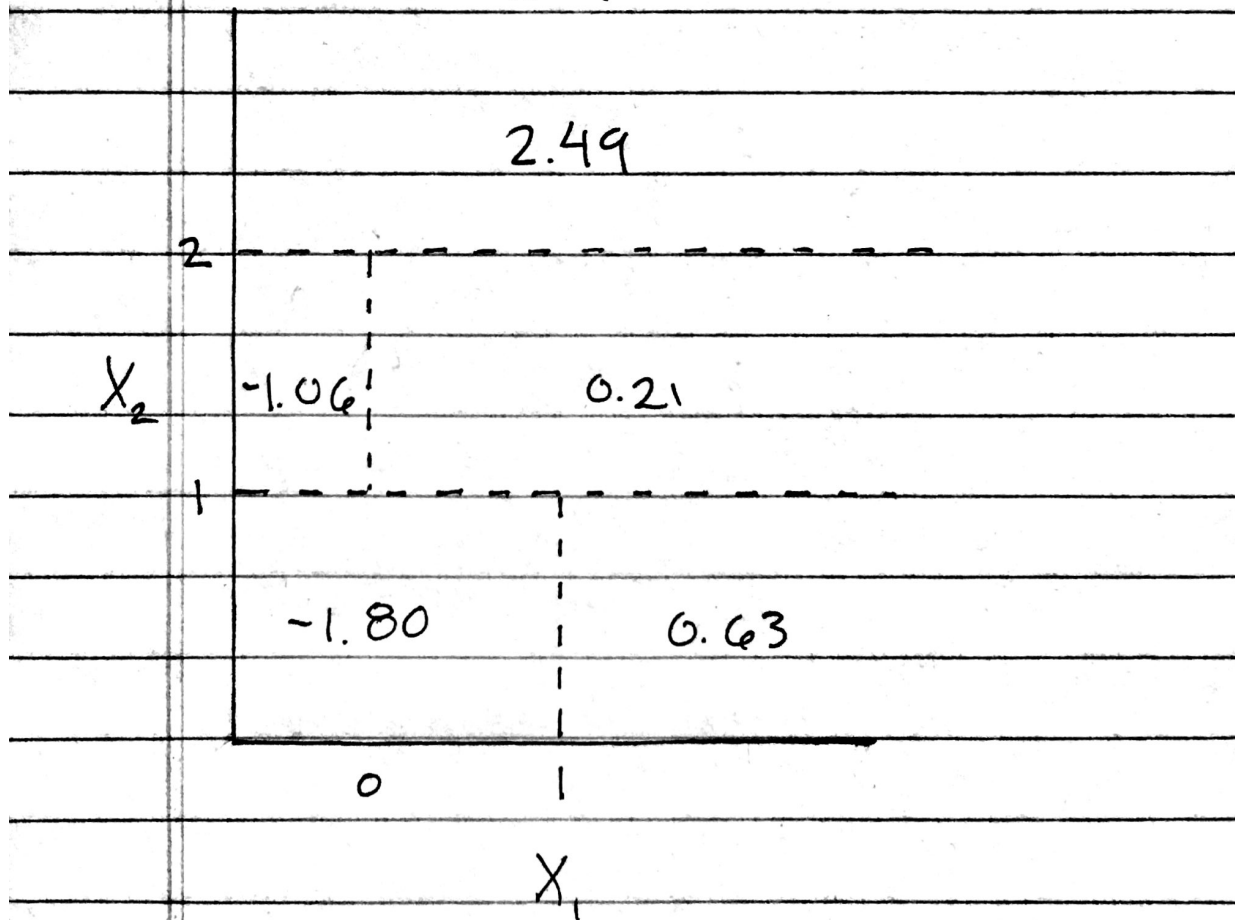


Figure 2: "Conceptual Exercise 4"

6

Broad Strokes - Use recursive binary splitting to continually split the data until there are fewer than x training observations in each terminal node, where x is predefined. The prediction for each terminal node will be the average of the response values of the training observations in that terminal node.

Detail - Starting with all the observations in the training set in one group, choose **the** value that, when split on, produces the lowest RSS among the two groups dictated by that split value.

- When the data set is split into two groups, the **single prediction** for the observations that are in group A will be the mean of the response values of the observations that are in group A, and the **single prediction** for the observations that are in group B will be the mean of the response values of the observations that are in group B.
- To calculate the (training) RSS of a particular split, sum the squared differences between the predicted value and the actual value of the response for all observations.

Now the data is partitioned into two groups. Repeat the same process again, treating each of those two groups as the *parent* group, which will be split into two *child* groups, leading to four total groups/divisions of the data set. Repeating again would lead to eight total division of the data set. Repeat this process until there are less than x training observations in each terminal node. Note that if only one of the *parent* groups has more than x training observations in it, but the other has less than x , only the group with more than x would be split further (in other words, the tree doesn't have to split evenly all the way down).