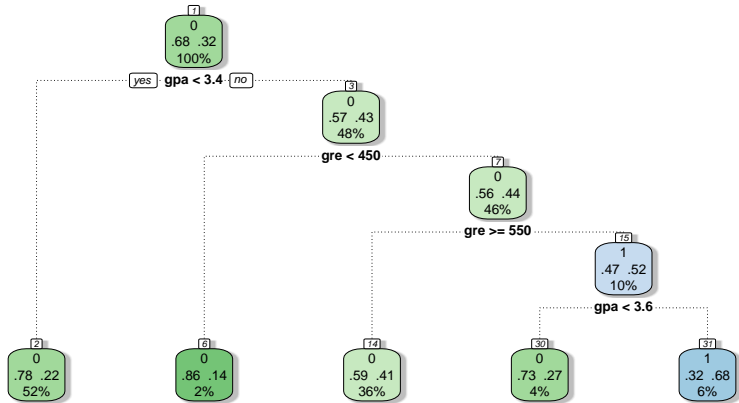


# Decision Trees

# Decision Trees

- Build a model to predict admit

admit	gre	gpa
0	380	3.61
1	660	3.67
1	800	4.00
1	640	3.19
0	520	2.93
1	760	3.00



# Making Prediction

- ▶ Predict the outcome of an applicant with 700 GRE score and 3.5 GPA



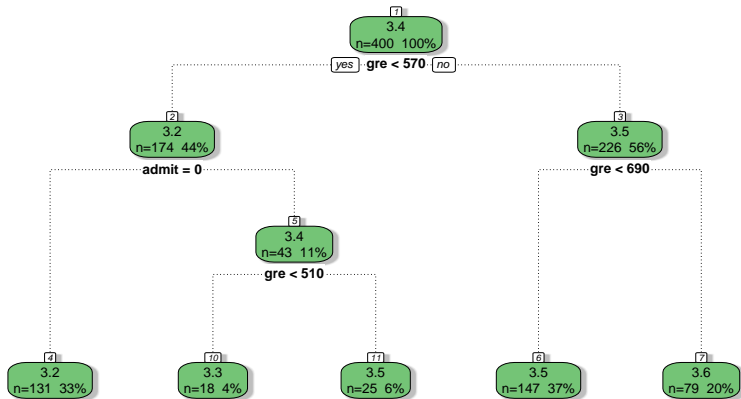
## Decision Trees on the coordinate

- ▶ Plot the previous tree on the coordinate of GRE and GPA

# Regression Trees

- Build a model to predict gpa

admit	gre	gpa
0	380	3.61
1	660	3.67
1	800	4.00
1	640	3.19
0	520	2.93
1	760	3.00





## Example

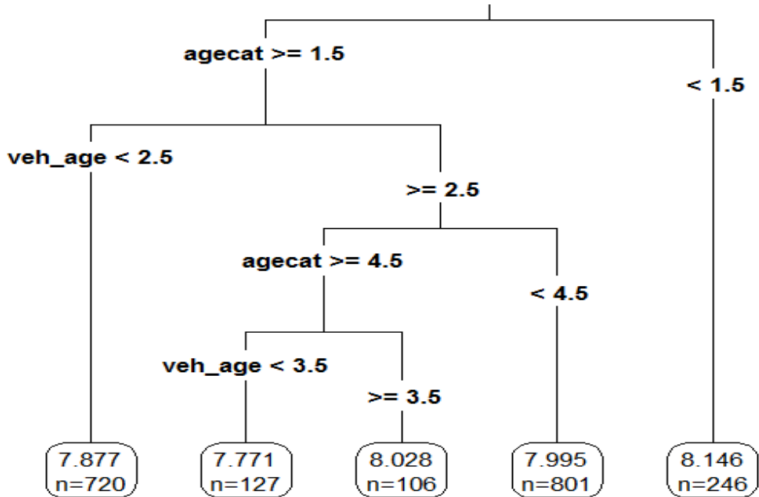
The regression tree shown below was produced from a dataset of auto claim payments. Age Category (agecat: 1, 2, 3, 4, 5, 6) and Vehicle Age (veh\_age: 1, 2, 3, 4) are both predictor variables, and log of claim amount (LCA) is the dependent variable. Rank the estimated LCA of Autos I, II, and III.

I: An Auto in Age Category 1 and Vehicle Age 4

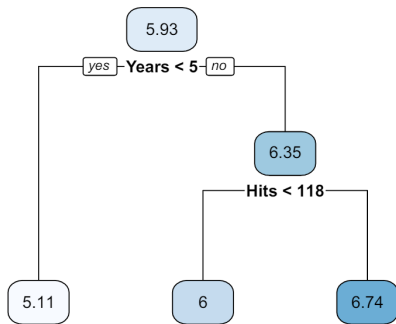
II: An Auto in Age Category 5 and Vehicle Age 5

III: An Auto in Age Category 5 and Vehicle Age 3

Calculate the estimated LCA of Autos I, II, and III.

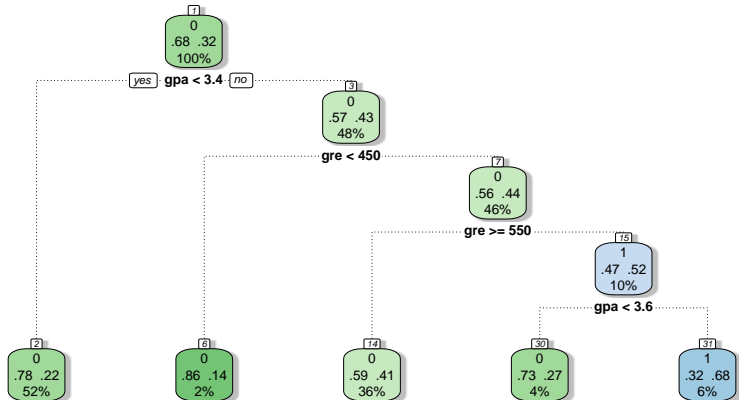


Consider the following regression tree to predict the log salary of a baseball player.



- Calculate the predicted salary for a player who has played 10 years and has 100 hits.

# Growing a Tree



- ▶ A tree is a combination of a sequence of splits
- ▶ Every split of a tree is defined by
  - ▶ The predictor it split at
  - ▶ The cutoff value
- ▶ For example, a split at gpa (predictor) of 3.5 (cutoff value)
- ▶ How to determine the first split, the second split, and so on?

- ▶ For the first split: Consider all possible split (a combination of all possible predictors and cutoffs values) then choose the best one
- ▶ For the second split: Consider all possible split after the first split then choose the best one
- ▶ And so on

# Classification Tree



## Impurity of a Node

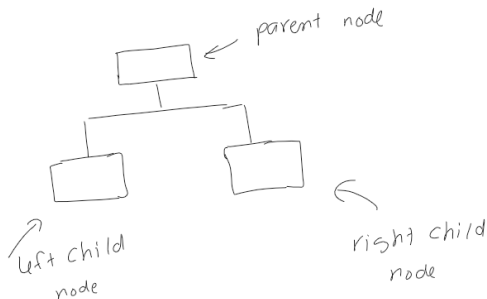
- ▶ Impurity can be measured by: classification error, Gini Index, and Entropy.
- ▶ Let  $p_0$  and  $p_1$  be the proportion of class 0 and class 1 in a node.

By Classification Error:  $I = \min\{p_0, p_1\}$

By Gini Index:  $I = 1 - p_0^2 - p_1^2$

By Entropy:  $I = -p_0 \log_2(p_0) - p_1 \log_2(p_1)$

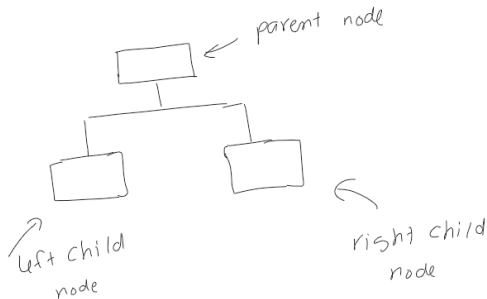
## Children Impurity of a Split



$$I_{children} = \frac{N_{left}}{N} I_{left} + \frac{N_{right}}{N} I_{right}$$

- ▶  $N_{left}$  and  $N_{right}$  are the number of points in the left child node and right child node, respectively.
- ▶  $N_{left} + N_{right} = N$

## Impurity Gain of a Split



$$IG = I_{parent} - I_{children}$$

► IG is Impurity Gain of the split

- ▶ The split with the highest impurity gain is the best
- ▶ Coming from the same parent node, the split with the lowest total children impurity is the best

## Example 1

admit	gre	gpa
0	380	3.61
1	660	3.67
1	800	4.00
1	640	3.19
0	520	2.93
1	760	3.00

- ▶ Split 1: Split at  $gre \geq 700$
- ▶ Split 2: Split at  $gre \geq 3.5$

Which split is the best?



## Example 2

A classification tree is being constructed to predict if an insurance policy will lapse. A random sample of 100 policies contains 30 that lapsed. You are considering two splits:

- ▶ Split 1: One node has 20 observations with 12 lapses and one node has 80 observations with 18 lapses.
- ▶ Split 2: One node has 10 observations with 8 lapses and one node has 90 observations with 22 lapses.

The total Gini index after a split is the weighted average of the Gini index at each node, with the weights proportional to the number of observations in each node.

The total entropy after a split is the weighted average of the entropy at each node, with the weights proportional to the number of observations in each node.

Determine the best split based on the total entropy.





# Regression Trees

# Regression Trees

- ▶ The tree will search for all combination of predictors and cutoff value to decide the best split
- ▶ In Regression tree, the best split is the split that minimizes

$$\underbrace{\sum_{i:\mathbf{x}_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2}_{\text{RSS of obs. in left branch}} + \underbrace{\sum_{i:\mathbf{x}_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2}_{\text{RSS of obs. in right branch}}$$

- ▶  $\hat{y}_{R_1}$  and  $\hat{y}_{R_2}$  are the means of the responses falling in to the left branch and right branch, respectively.

## Example

$X_1$	$X_2$	$Y$
1	0	1.2
2	1	2.1
3	2	1.5
4	1	3.0
2	2	2.0
1	1	1.6

Using the RSS to decide the best split among

- ▶ Split 1: Region 1  $X_1 < 4$ , Region 2  $X_1 \geq 4$
- ▶ Split 2: Region 1  $X_2 < 2$ , Region 2  $X_2 \geq 2$



