

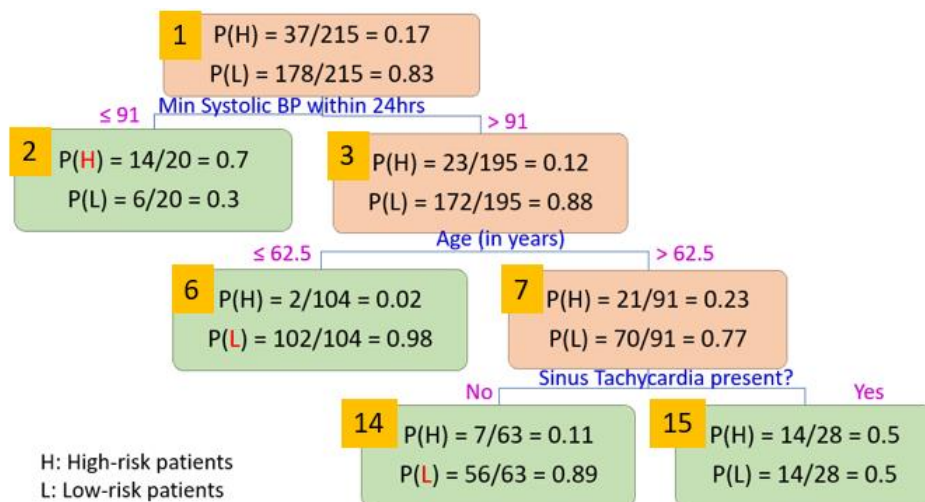
## Exercise 8.1 CART (Part 1)

### Understanding the Search for Best Split:

1. In default10.xlsx, the training set consists of 10 cases with 3 input variables (Home Owner, Marital Status, and Annual Income). The outcome variable is Defaulted Borrower. Let  $y = 0$  for no default, and  $y = 1$  for defaulted. Using Excel for all calculations,
  - (a) What are the proportions of cases before any splitting?
  - (b) Calculate the Entropy, Gini index and Misclassification error of the two descendant nodes if the splitting criteria is:
    - (i) Annual Income  $\leq 100K$
    - (ii) Marital Status = Single
  - (c) Breiman et al (1984) preferred the Gini index. Which of the above 2 split would you choose if the Gini index is used?
  - (d) Misclassification error is not a good formula for selecting the best split, in general. Why?

*Note: Here, we had only compared 2 possible splits. CART will have to consider and compare all possible splits from all the X variables and variable values in order to find the best split.*

### Understanding the CART Model Structure:



2. What is the decision rule, model prediction and misclassification error at node 6?
3. What is the decision rule, model prediction and misclassification error at node 7?

*An alternative perspective on Pruning by considering the portion of the tree to be pruned away instead of considering the tree that remains after pruning. This perspective directly defines the strength of the link to descendants.*

Define  $T_t$  to be a subtree with root at node  $t$ , and define  $R_\alpha(t) = R(t) + \alpha$

If  $R_\alpha(T_t) < R_\alpha(t)$ , then the contribution to the cost complexity of the subtree  $T_t$  is less than that for node  $t$  (if it becomes a terminal node). Do not prune.

As  $\alpha$  increases, equality is achieved when

$$\alpha^* = \frac{R(t) - R(T_t)}{|T_t| - 1}$$

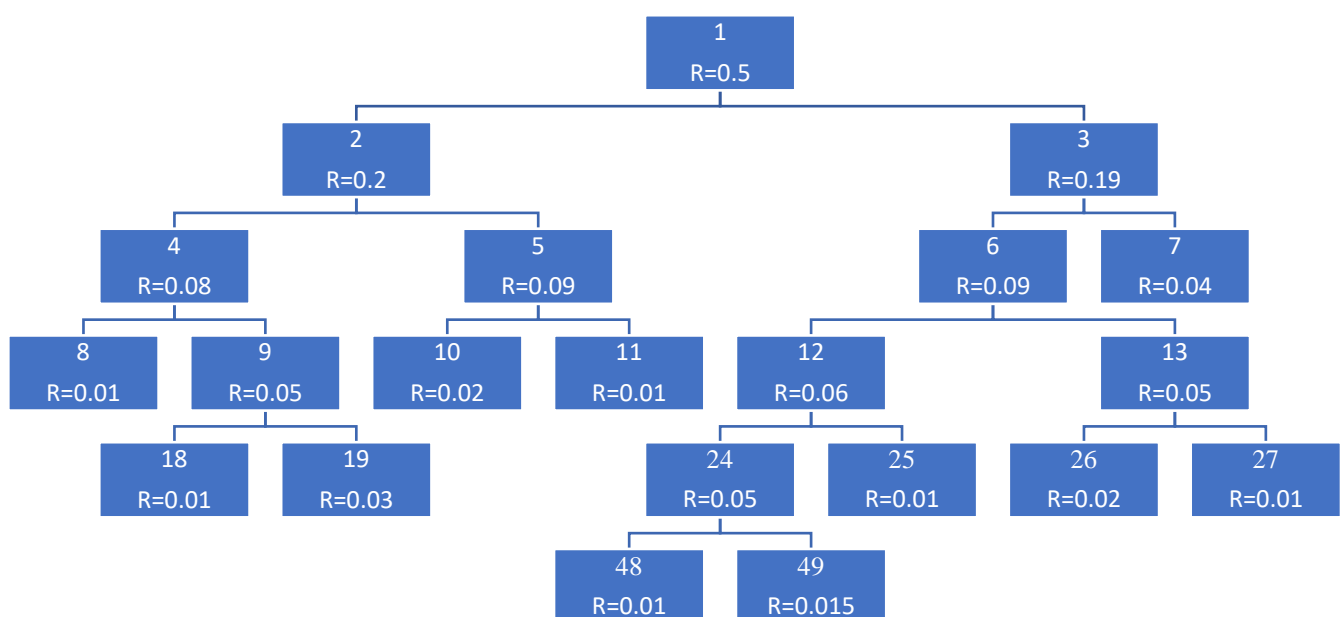
and termination of the tree at  $t$  is preferred if  $\alpha \geq \alpha^*$  i.e. prune all descendant nodes at node  $t$ .

Thus, we can define

$$g(t) = \frac{R(t) - R(T_t)}{|T_t| - 1}$$

as the strength of the link from node  $t$  to its descendants.

4. A fully grown tree is shown below with  $R(t)$  calculated at each node. Calculate and write down  $g(t)$  in the table below to get the weakest link for pruning.



$t$	$R(t)$	$R(T_t)$	$ T_t  - 1$	$g(t)$
1	0.5			
2	0.2	$R(8)+R(18)+R(19)+R(10)+R(11) =$		
3				
4	0.08			
5	0.09			
6	0.09			
7*	0.04			
8*	0.01			
9	0.05			
10*	0.02			
11*	0.01			
12	0.06			
13	0.05			
18*	0.01			
19*	0.03			
24	0.05			
25*	0.01			
26*	0.02			
27*	0.01			
48*	0.01			
49*	0.015			

## The CART Pruning Sequence

After growing the tree using the Gini impurity index, and obtaining  $g(t)$  for all internal nodes and  $R(t)$  for all nodes. The CART pruning sequence is as follows:

1. Find the node  $t$  with the lowest  $g(t)$ . i.e. node  $t_x$  where  $g(t_x) \leq g(t_i)$ , for all nodes  $t_i \neq t_x$   
(There maybe more than one node with the minimum  $g(t)$  value.)
2. Prune at nodes  $t_x$ .
3. Re-compute  $g(t)$  at all ancestor (and internal) nodes of  $t_x$
4. Repeat steps 1 to 3 until only the root node remains.