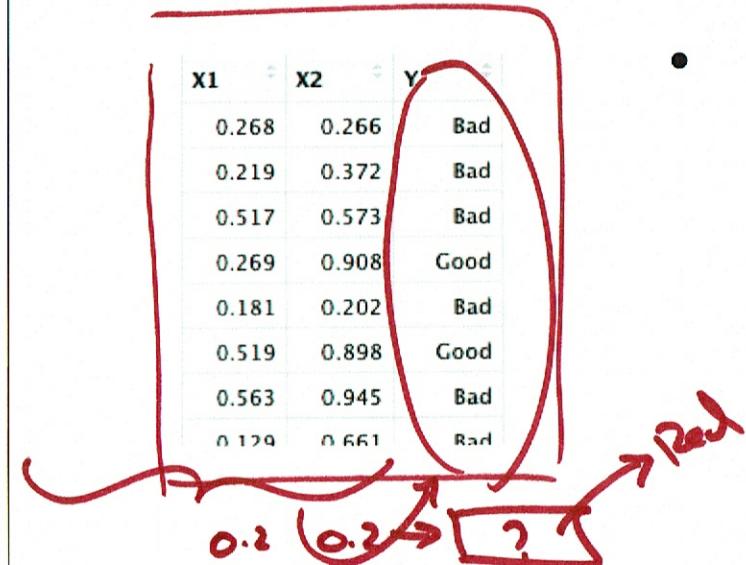


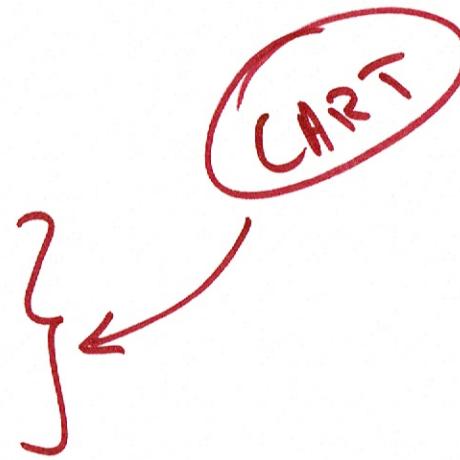
Classification and Regression *Learning for Life*

- Decision Trees can be used for both



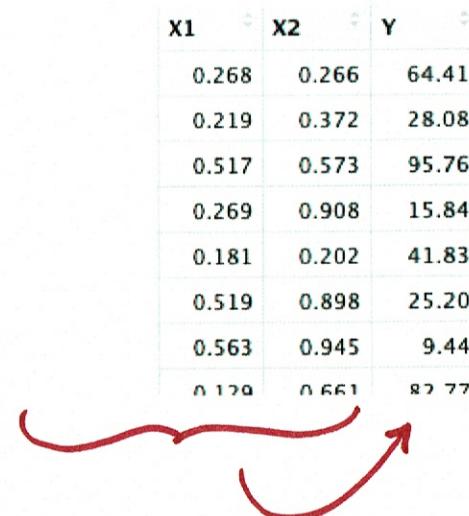
- Classification

- Spam / not Spam
- Admit to ICU /not
- Lend money / deny
- Intrusion detections

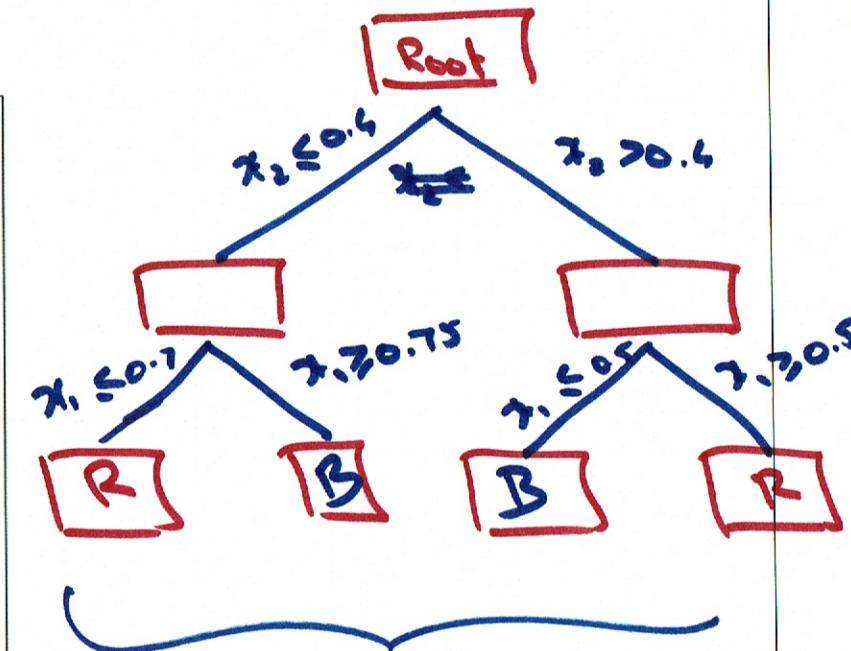
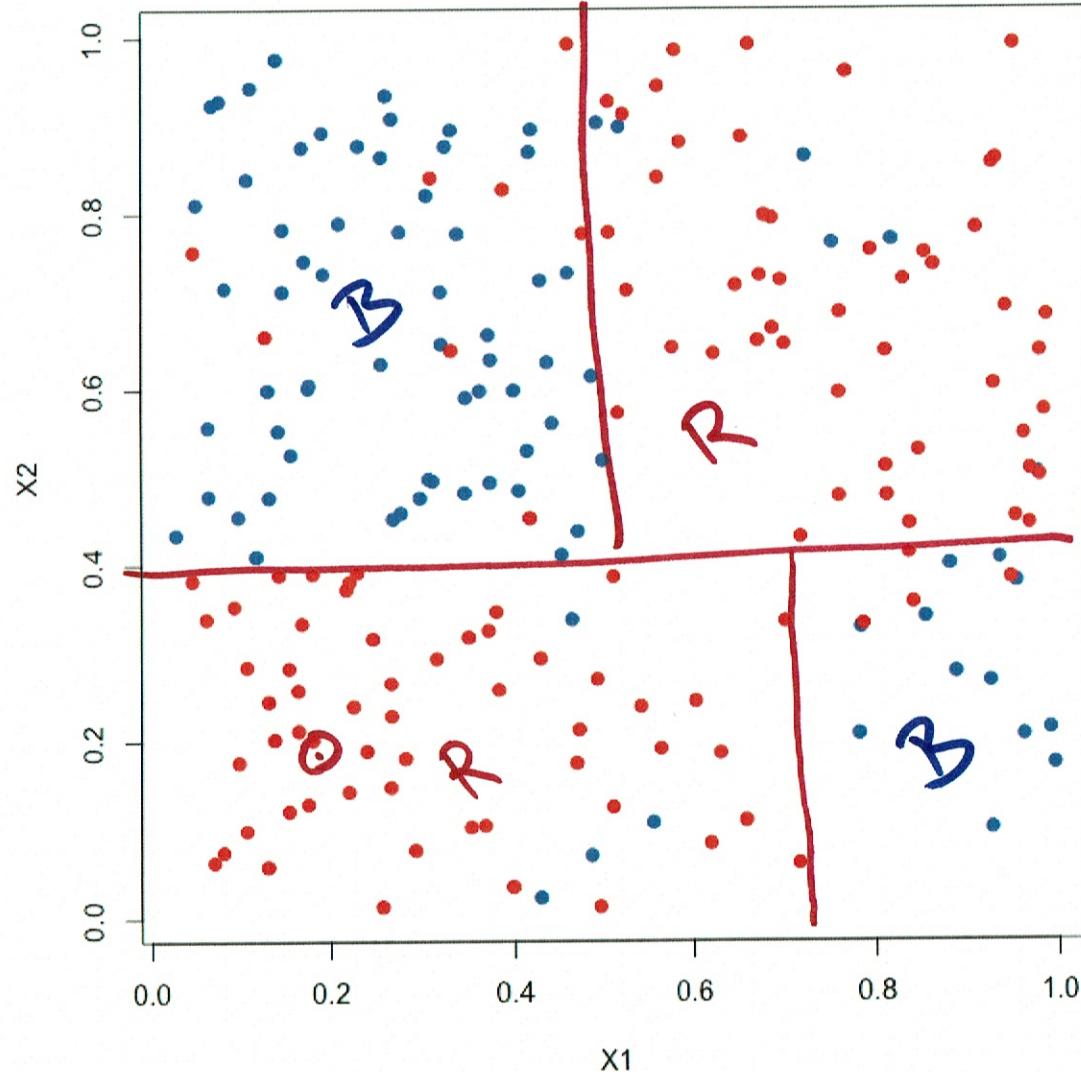


- Regression

- Predict stock returns
- Pricing a house or a car
- Weather predictions (temp, rain fall etc)
- Economic growth predictions
- Predicting sports scores

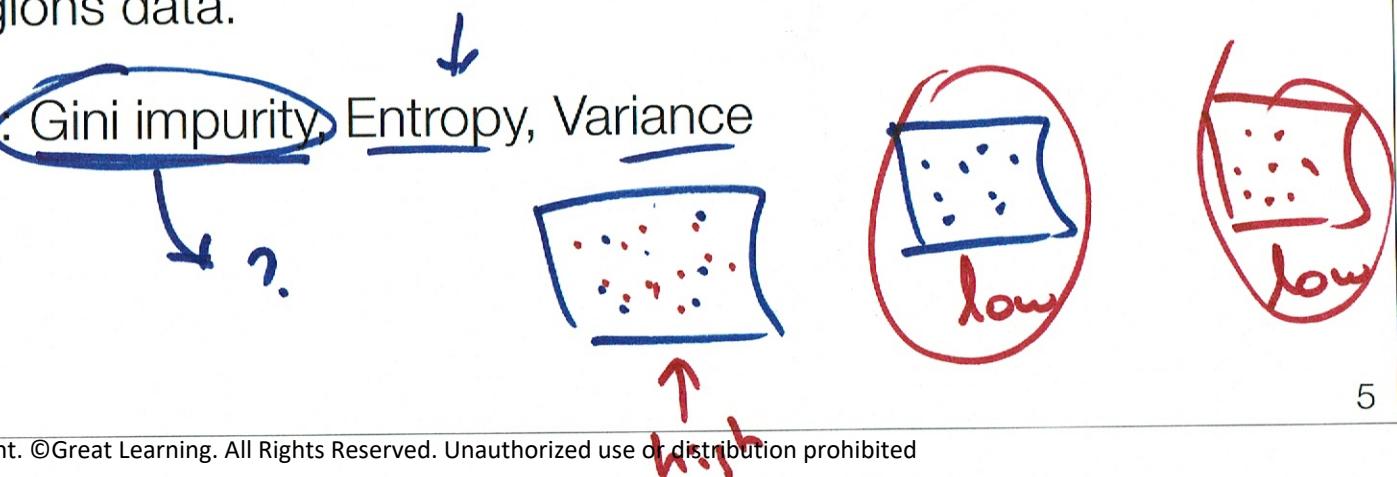


Visualizing Classification as a Tree



Metrics

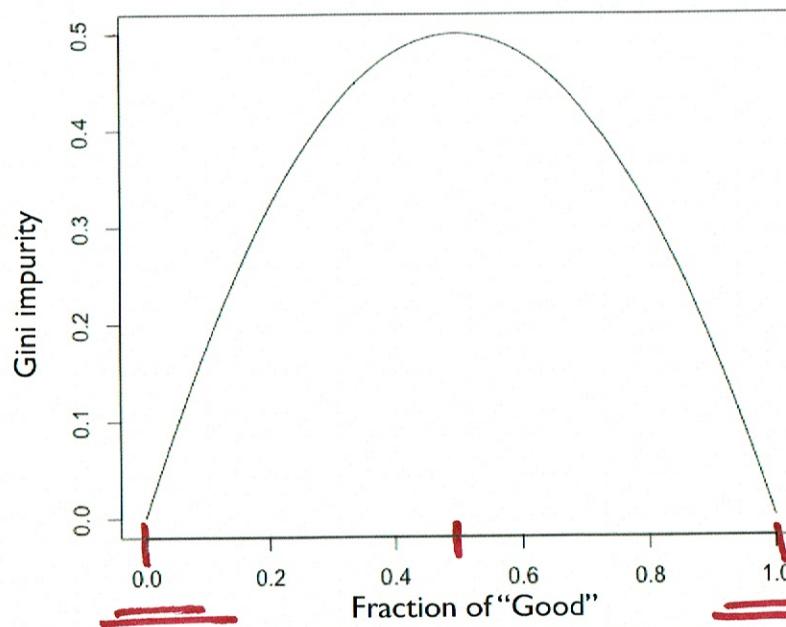
- Algorithms for constructing decision trees usually work top-down, by choosing a variable at each step that best splits the set of items.
- Different algorithms use different metrics for measuring “best”
- These metrics measure how similar a region or a node is. They are said to measure the impurity of a region.
- Larger these impurity metrics the larger the “dissimilarity” of a nodes/regions data.
- Examples: Gini impurity, Entropy, Variance

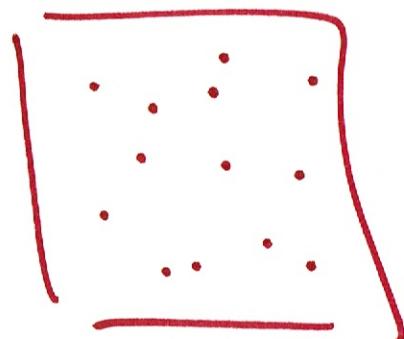
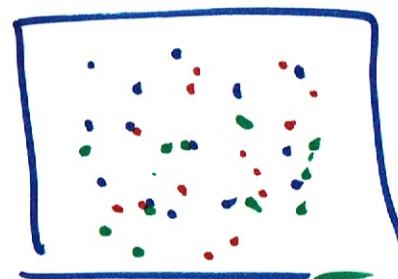


Gini impurity

- Used by the CART
- Is a measure of how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset.
- Can be computed by summing the probability of an item with label i being chosen (p_i), times the probability of a mistake ($1 - p_i$) in categorizing that item.
- Simplifying gives, the Gini impurity of a set:

$$1 - \sum_{i=1}^J p_i^2$$



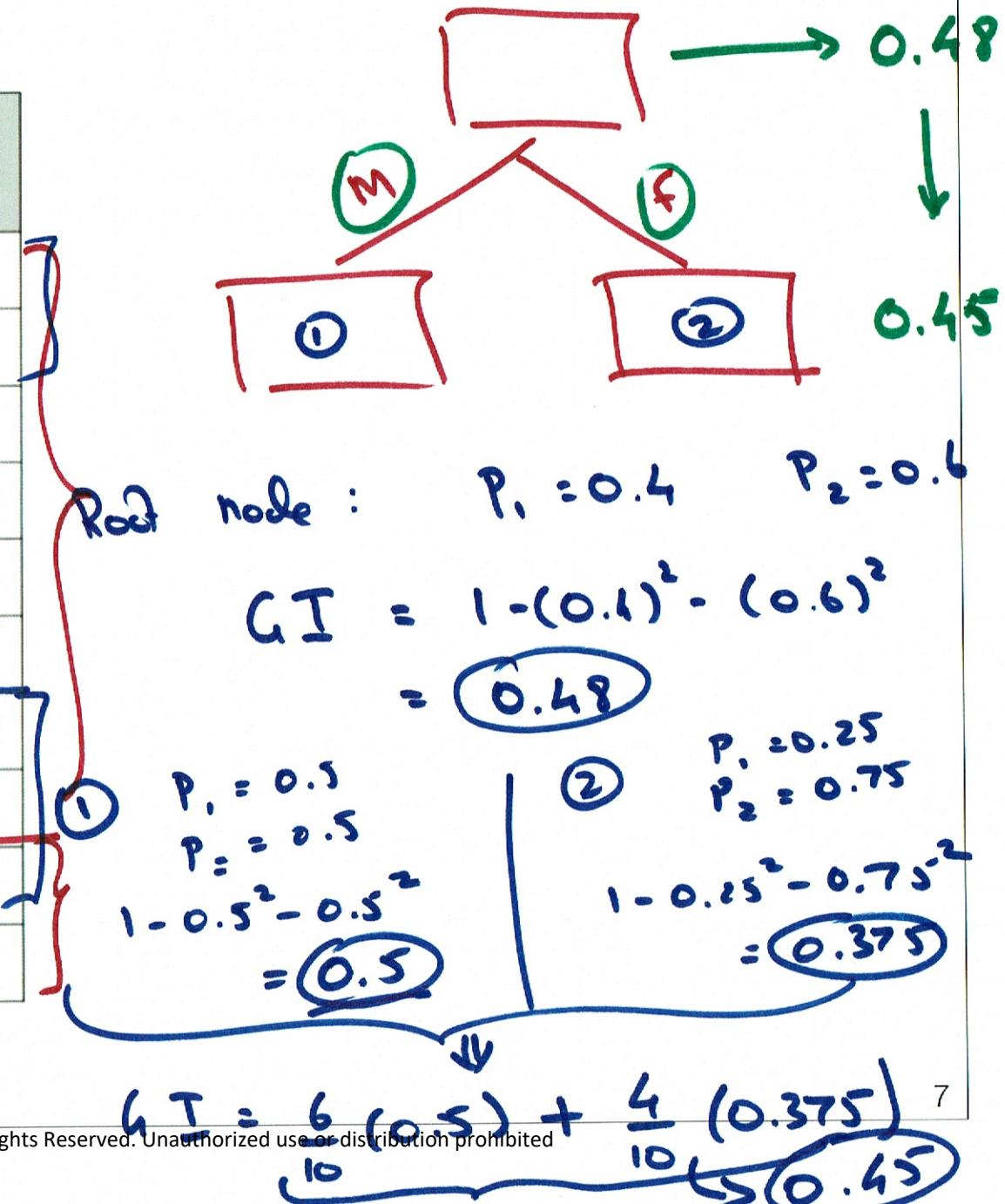


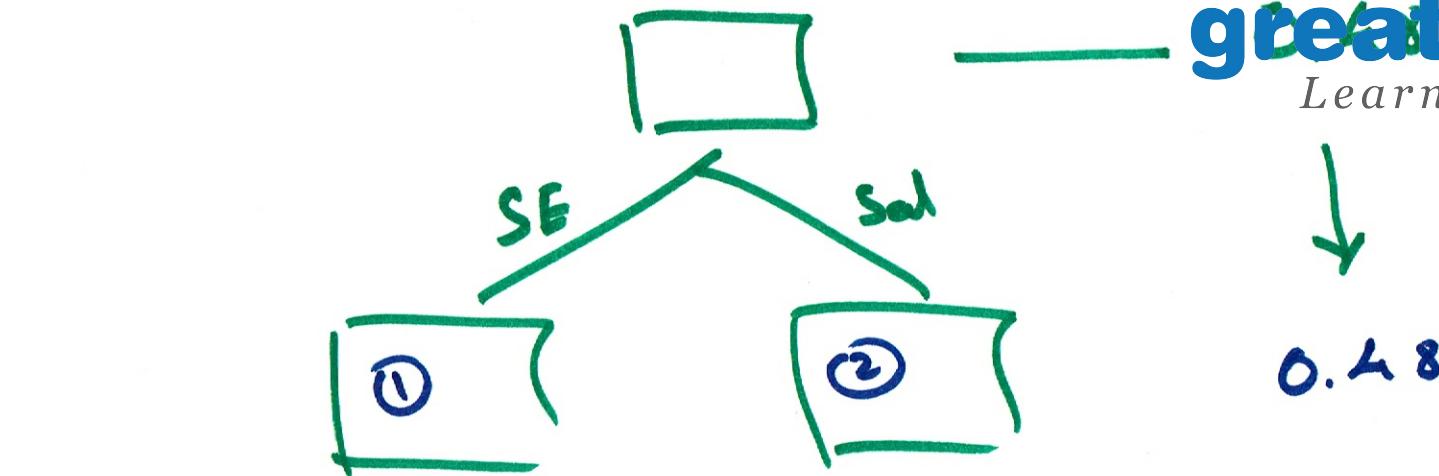
$$P_1 | P_2 | P_3$$

$$\begin{aligned}
 & \xrightarrow{P_1} ① \Rightarrow P_1(1 - P_1) && \leftarrow P_1 P_2 + P_1 P_3 \\
 & \xrightarrow{P_2} ② \Rightarrow P_2(1 - P_2) && \leftarrow P_2 P_3 + P_2 P_1 \\
 & \xrightarrow{P_3} ③ \Rightarrow P_3(1 - P_3) && \leftarrow \underbrace{P_3 P_1 + P_3 P_2} \\
 & \downarrow \\
 & \sum P_i(1 - P_i) \\
 & \sum P_i - \sum P_i^2 \Rightarrow 1 - \sum P_i^2
 \end{aligned}$$

CART: An Example

Cust_ID	Gender	Occupation	Age	Target
1	M	Sal	22	1
2	M	Sal	22	0
3	M	Self-Emp	23	1
4	M	Self-Emp	23	0
5	M	Self-Emp	24	1
6	M	Self-Emp	24	0
7	F	Sal	25	1
8	F	Sal	25	0
9	F	Sal	26	0
10	F	Self-Emp	26	0

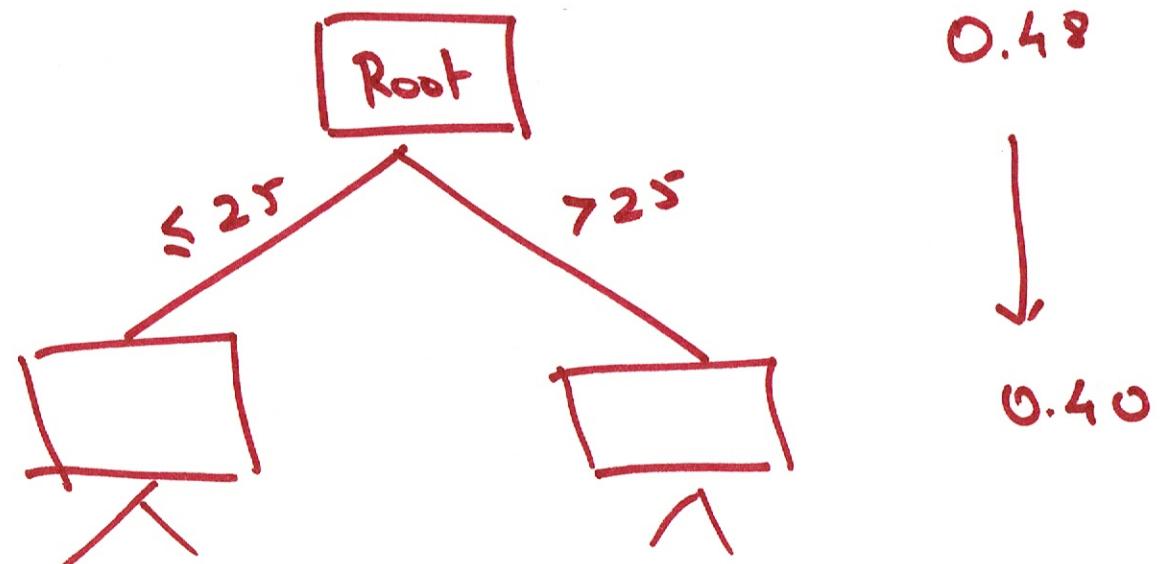




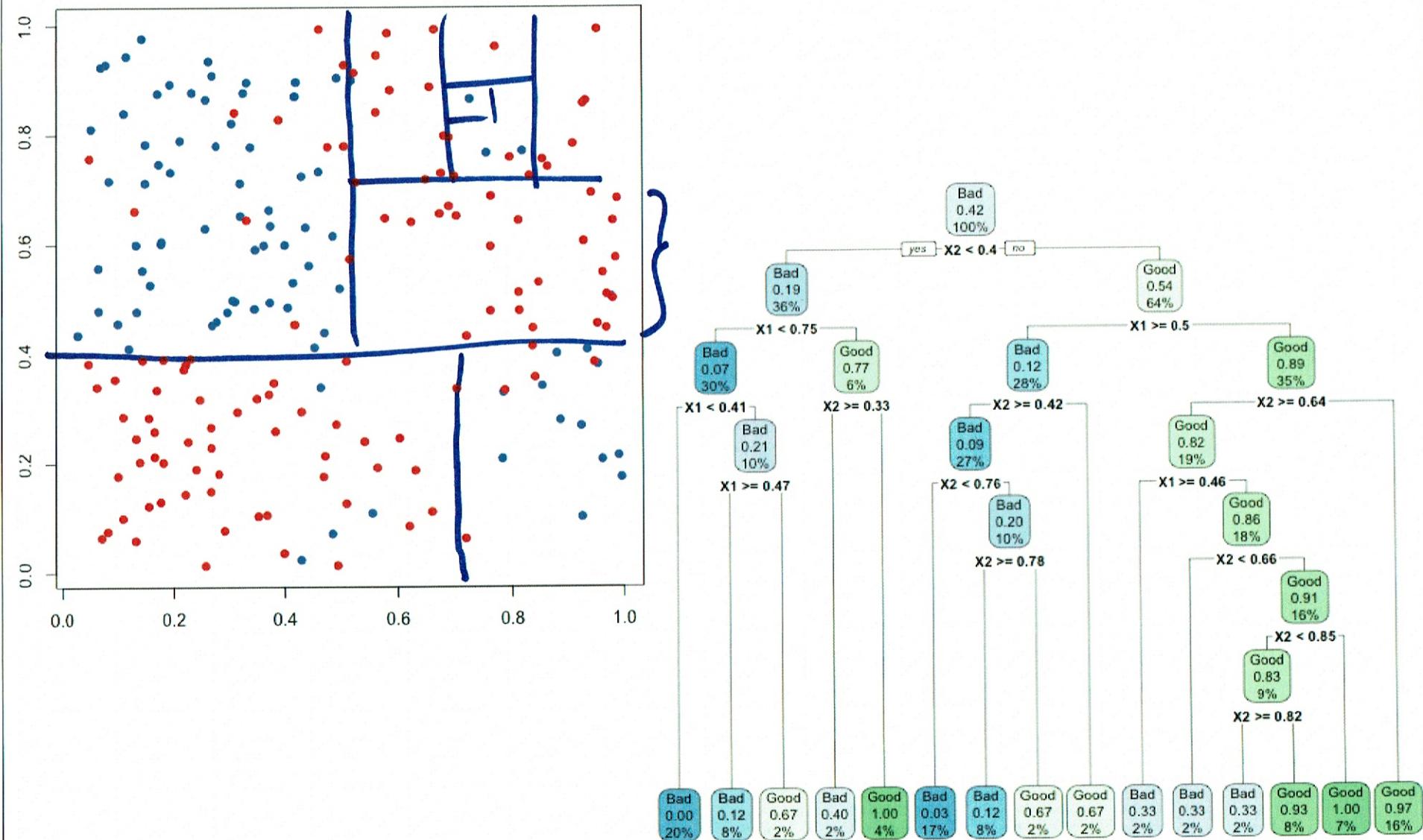
$$\begin{array}{c}
 \textcircled{1} \quad G.I = 1 - 0.4^2 - 0.6^2 \\
 \qquad\qquad\qquad = 0.48 \\
 \textcircled{2} \quad G.I = 1 - 0.4^2 - 0.6^2 \\
 \qquad\qquad\qquad = 0.48 \\
 \\
 G.I = \frac{\sum}{10} (0.48) + \frac{\sum}{10} (0.48) = 0.48
 \end{array}$$

	Left	Right	Gini Split
$\leq 22, > 22$	0.5	0.47	0.48
$\leq 23, > 23$	0.5	0.44	0.47
$\leq 24, > 24$	0.5	0.38	0.45
$\leq 25, > 25$	0.5	0	0.40

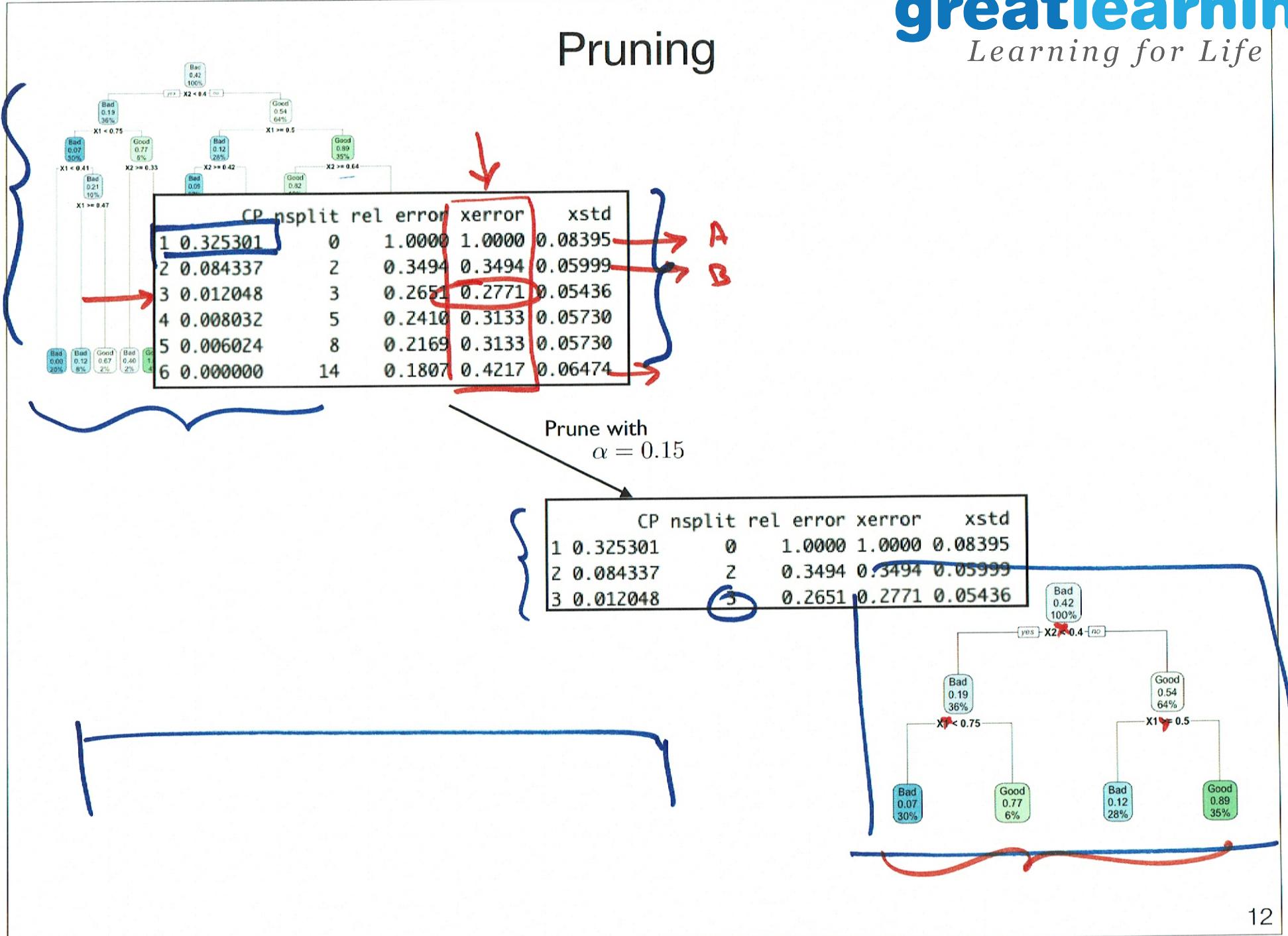
Gain $\Rightarrow 0.08$

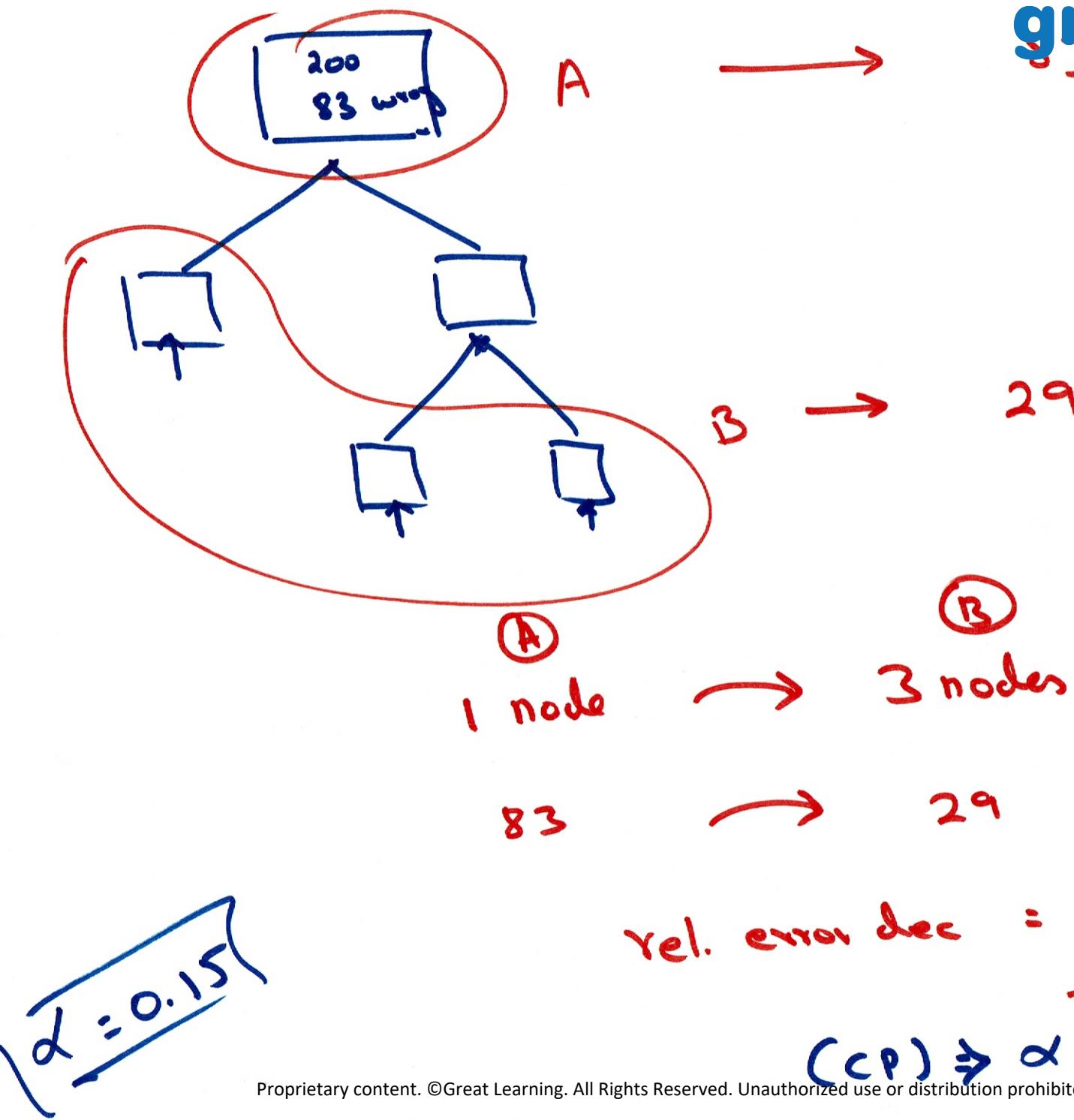


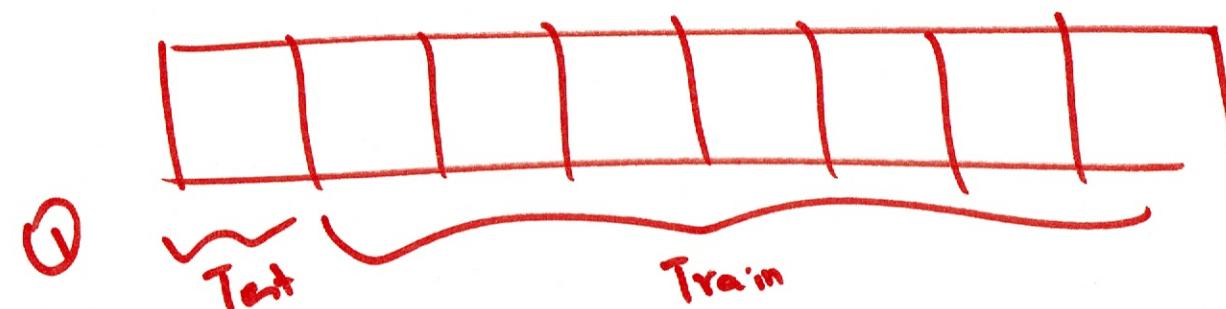
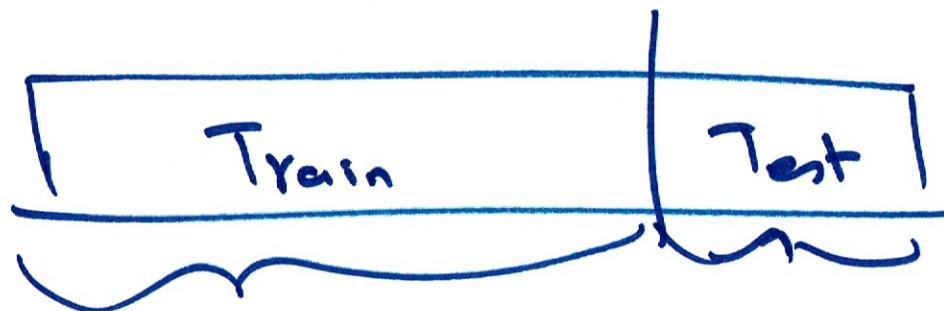
Overfitting in Decision Trees



Pruning







Regression Trees

	Price	Country	Reliability	Mileage	Type
Acura Integra 4	11950	Japan	Much better	NA	Small
Dodge Colt 4	6851	Japan	NA	NA	Small
Dodge Omni 4	6995	USA	Much worse	NA	Small
Eagle Summit 4	8895	USA	better	33	Small
Ford Escort 4	7402	USA	worse	33	Small
Ford Festiva 4	6319	Korea	better	37	Small
GEO Metro 3	6695	Japan	NA	NA	Small
GEO Prizm 4	10125	Japan/USA	Much better	NA	Small
Honda Civic 4	6635	Japan/USA	Much better	32	Small
Hyundai Excel 4	5899	Korea	worse	NA	Small
Mazda Protege 4	6599	Japan	Much better	32	Small

