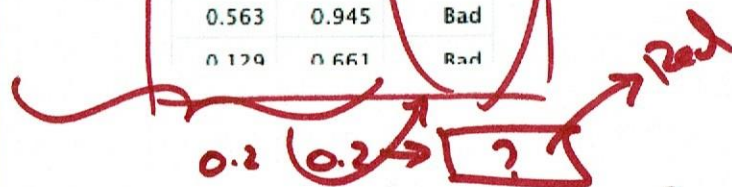


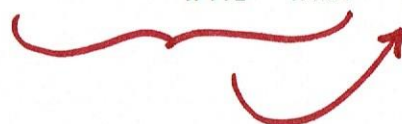
# Classification and Regression *Learning for Life*

- Decision Trees can be used for both

x1	x2	y
0.268	0.266	Bad
0.219	0.372	Bad
0.517	0.573	Bad
0.269	0.908	Good
0.181	0.202	Bad
0.519	0.898	Good
0.563	0.945	Bad
0.129	0.661	Bad

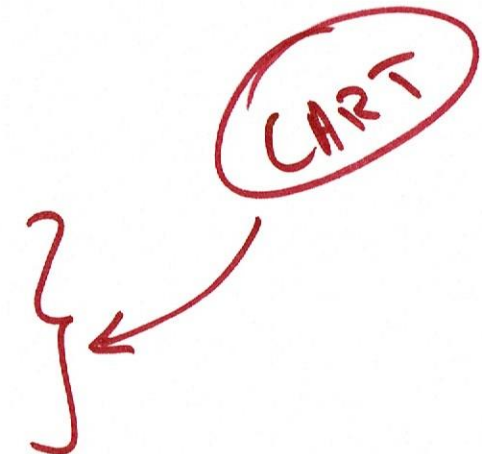


x1	x2	Y
0.268	0.266	64.41
0.219	0.372	28.08
0.517	0.573	95.76
0.269	0.908	15.84
0.181	0.202	41.83
0.519	0.898	25.20
0.563	0.945	9.44
0.129	0.661	82.77



## Classification

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

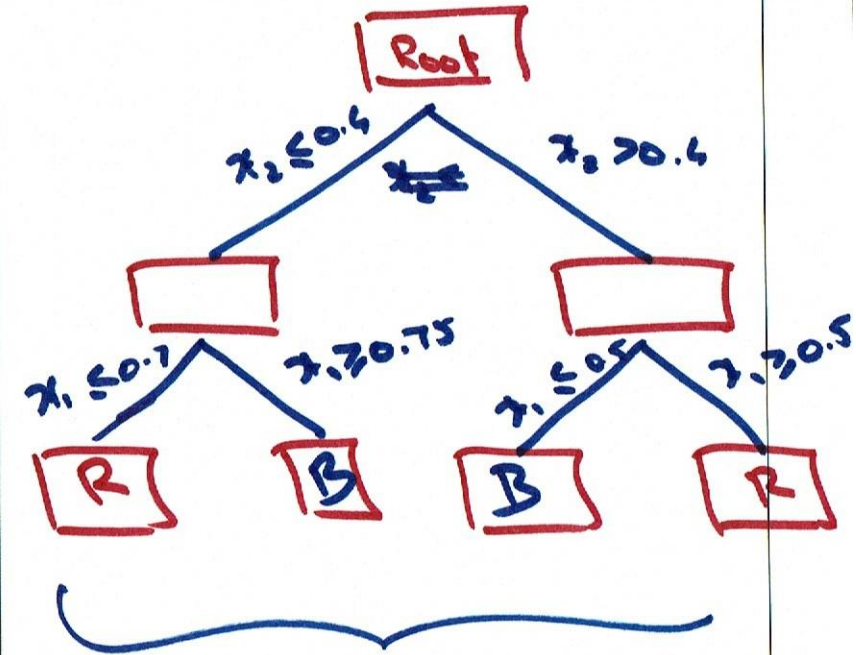
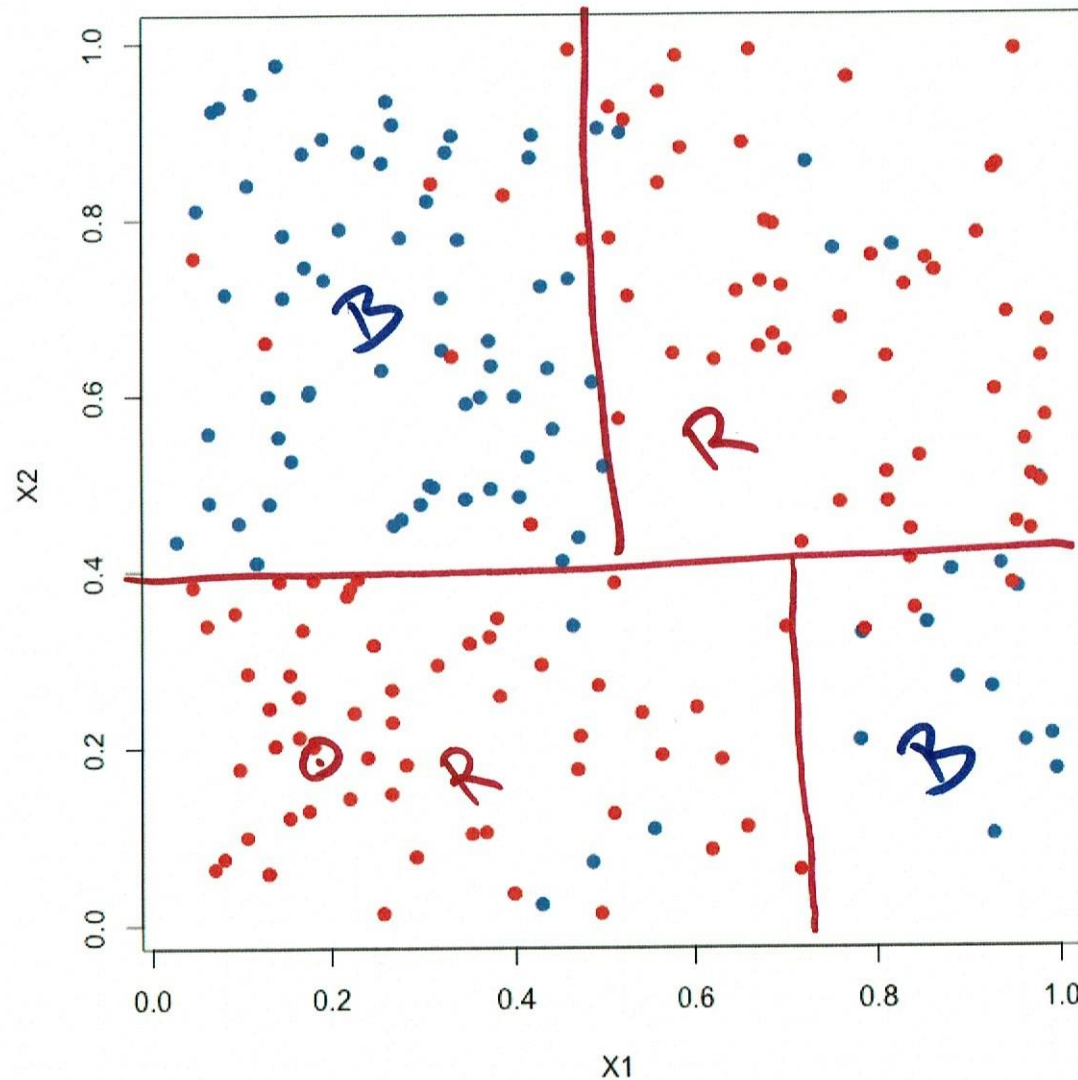


## Regression

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

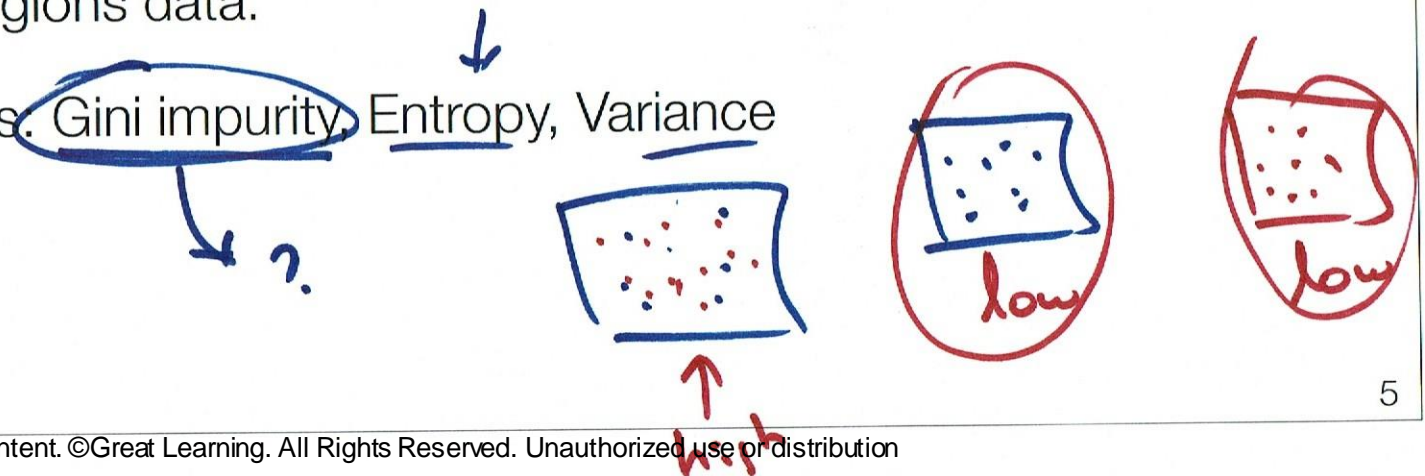


# Visualizing Classification as a Tree



# 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

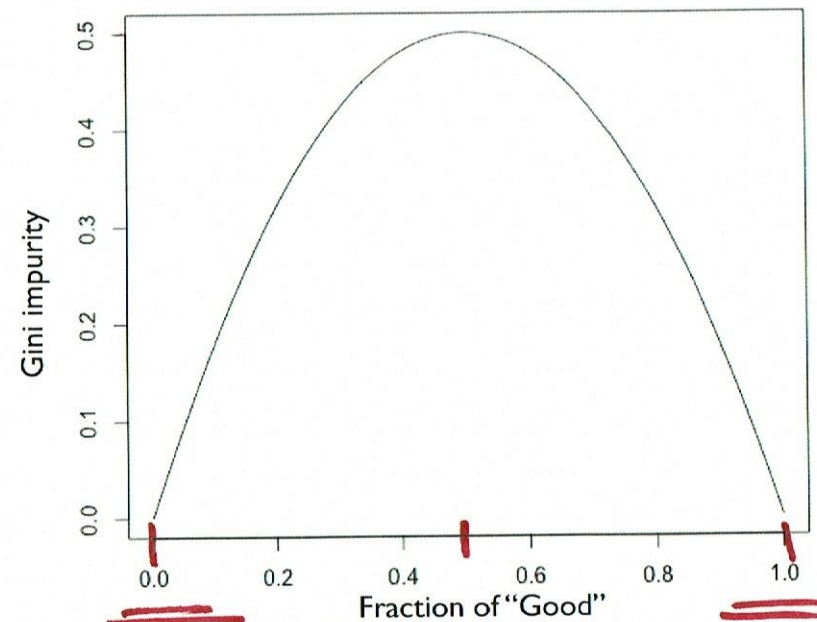
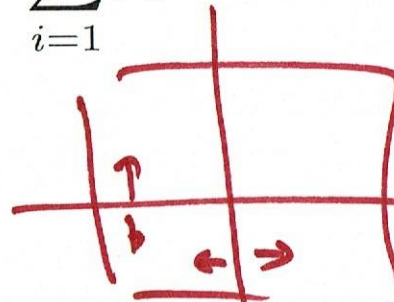




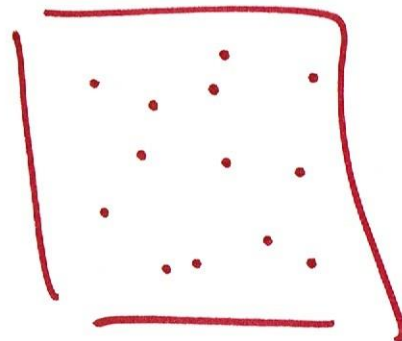
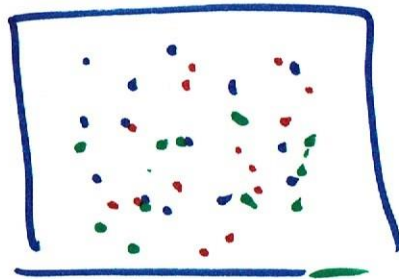
# Gini impurity

- 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:

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



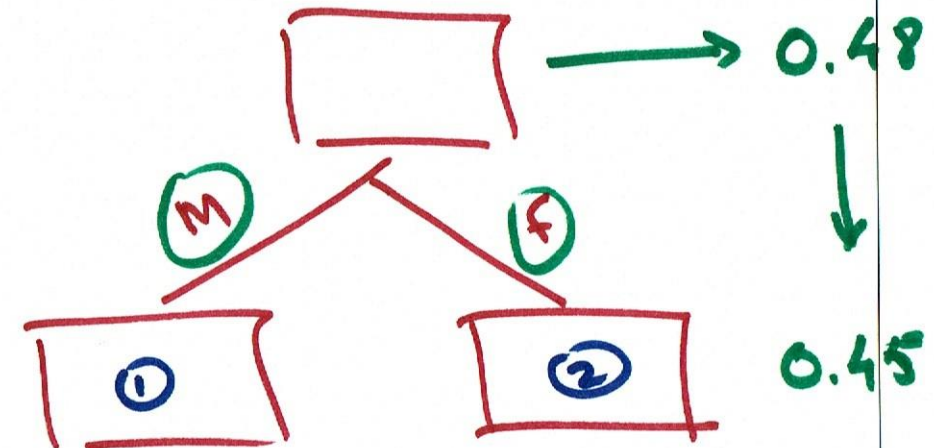
$$\boxed{p_1 \mid p_2 \mid p_3}$$



$$\begin{aligned}
 & \begin{array}{l} p_1 \rightarrow \textcircled{1} \Rightarrow p_1(1-p_1) \\ p_2 \rightarrow \textcircled{2} \Rightarrow p_2(1-p_2) \\ p_3 \rightarrow \textcircled{3} \Rightarrow p_3(1-p_3) \end{array} \quad \begin{array}{l} \leftarrow p_1 p_2 + p_1 p_3 \\ \leftarrow p_2 p_3 + p_2 p_1 \\ \leftarrow p_3 p_1 + p_3 p_2 \end{array} \\
 & \quad \quad \quad \Downarrow \\
 & \quad \quad \quad \Sigma p_i(1-p_i) \\
 & \quad \quad \quad \rightarrow \left( \Sigma p_i - \Sigma p_i^2 \right) \Rightarrow \left( 1 - \Sigma p_i^2 \right)
 \end{aligned}$$

# CART: An Example

Cust_ID	Gender	Occupati on	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



Root node :  $P_1 = 0.4$   $P_2 = 0.6$

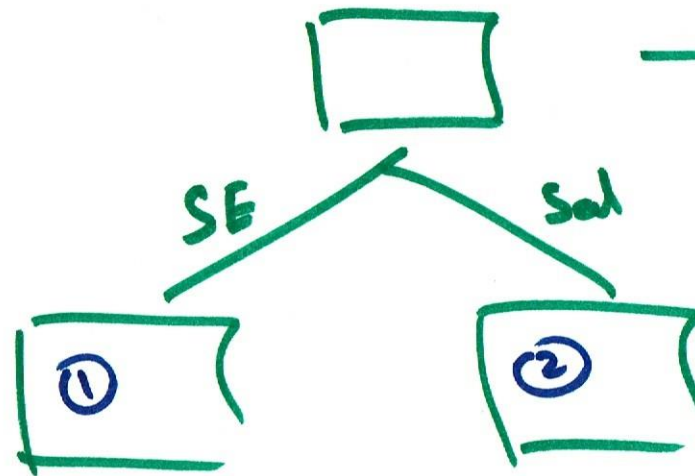
$$GI = 1 - (0.4)^2 - (0.6)^2 = 0.48$$

①  $P_1 = 0.5$   
 $P_2 = 0.5$   
 $1 - 0.5^2 - 0.5^2 = 0.5$

②  $P_1 = 0.25$   
 $P_2 = 0.75$   
 $1 - 0.25^2 - 0.75^2 = 0.375$

$$GI = \frac{6}{10} (0.5) + \frac{4}{10} (0.375) = 0.45$$





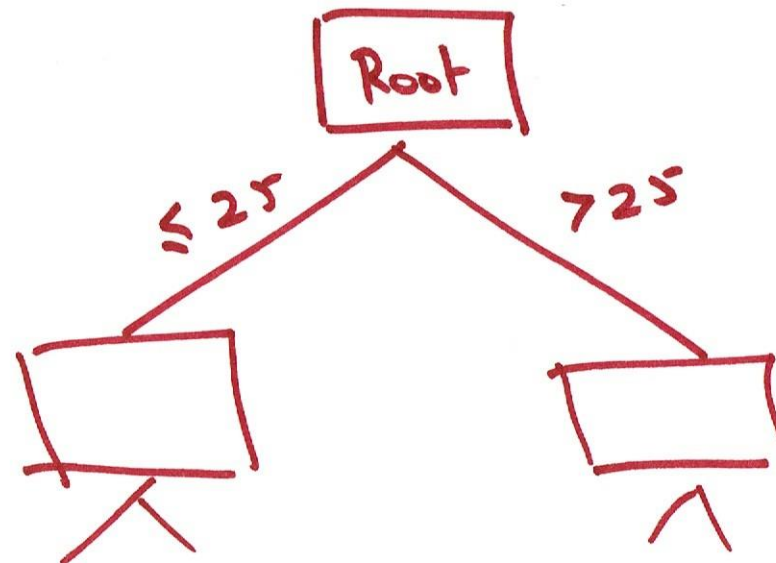
0.48

$$\textcircled{1} \quad G.I = 1 - 0.4^2 - 0.6^2 = 0.48 \quad \bigg| \quad \textcircled{2} \quad G.I = 1 - 0.4^2 - 0.6^2 = 0.48$$

$$G.I = \frac{5}{10} (0.48) + \frac{5}{10} (0.48) = 0.48$$

	Left	Right	Gain 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



0.48  
↓  
0.40