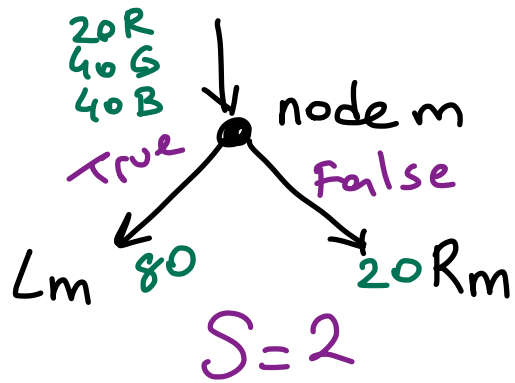


# Univariate Decision Trees



$$L_m = \{x \mid x_j > w_{m0}\}$$

threshold  
feature index

$$R_m = \{x \mid x_j \leq w_{m0}\}$$

$$P_{m1} = \frac{20}{100}$$

$$P_{m2} = \frac{40}{100}$$

$$P_{m3} = \frac{40}{100}$$

$100 \Leftarrow N_m = \# \text{ of data points that reach node } m$

$(20, 40, 40) N_{mc} = \# \text{ of data points that reach node } m \text{ from class } c$

$(80, 20) N_{m,s} = \# \text{ of data points that reach node } m \text{ and take split } s$

$3 \Leftarrow k = \# \text{ of classes}$

$$P_{mc} = \hat{\Pr}(y=c \mid \mathcal{X}_m) = \frac{N_{mc}}{N_m}$$

impurity of a node  $\Leftarrow I_m = - \sum_{c=1}^k P_{mc} \log_2(P_{mc}) \Leftarrow \text{Entropy}$

impurity of a split  $\Leftarrow I_m' = \sum_{s=1}^S \left[ \underbrace{\frac{N_{m,s}}{N_m}}_{\text{weight of a child node}} \underbrace{\left[ - \sum_{c=1}^k P_{m,s,c} \log_2(P_{m,s,c}) \right]}_{\text{impurity of a node}} \right]$

$$\frac{80}{100} I_m(L_m) + \frac{20}{100} I_m(R_m)$$

Entropy:

$$-p \cdot \log_2(p) - (1-p) \log_2(1-p)$$

$$\frac{N_+}{N_+ + N_-}$$

$$\frac{N_-}{N_+ + N_-}$$

$p$  = ratio of positively labeled data points  
 $1-p$  = ratio of negatively labeled data points  
 $0 \log_2(0) \triangleq 0$

Gmi Index:

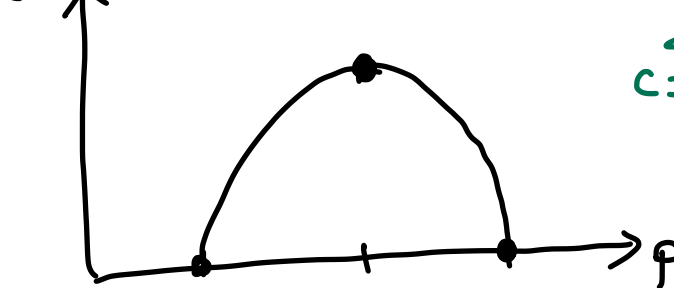
$$2 \cdot p(1-p)$$

$$p=0 \Rightarrow 2 \cdot 0 \cdot (1-0)$$

$$p=0.5 \Rightarrow 2 \cdot \frac{1}{2} \cdot \frac{1}{2}$$

$$p=1.0 \Rightarrow 2 \cdot 1 \cdot (1-1)$$

Gmi index



$$\sum_{c=1}^K p_c(1-p_c)$$

0% 100% 0%  
 $\downarrow \downarrow \downarrow$   
 0 0 0

all data points are from (-) class

all data points are from (+) class

Misclassification Error:

$$1 - \max(p, 1-p)$$

60% 40%

classification accuracy

OR

$$\min(p, 1-p)$$

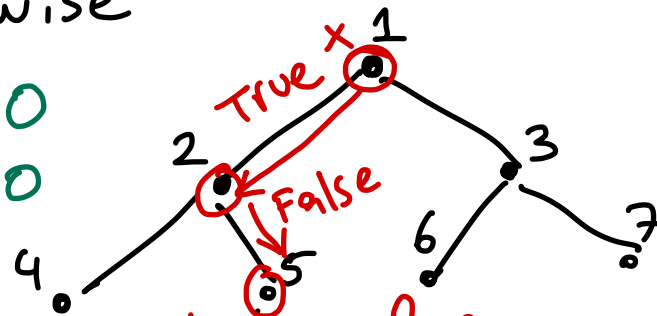
misclassification rate

multiclass classification  $\Rightarrow 1 - \max(p_1, p_2, \dots, p_K)$   
 classification accuracy when majority label is used as prediction

# Regression Trees

$$b_m(x) = \begin{cases} 1 & \text{if } x \in \mathcal{X}_m \text{ (x reaches node m)} \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{array}{lll} b_1(x) = 1 & b_2(x) = 1 & b_3(x) = 0 \\ b_4(x) = 0 & b_5(x) = 1 & b_6(x) = 0 \\ & b_7(x) = 0 & \end{array}$$



path  $\Rightarrow 1-2-5$

error of  
a node

$$E_m = \frac{1}{N_m} \sum_{i=1}^N \left[ \overset{\text{observed / true value}}{y_i} - \overset{\text{predicted value at node m}}{g_m} \right]^2 \cdot b_m(x_i)$$

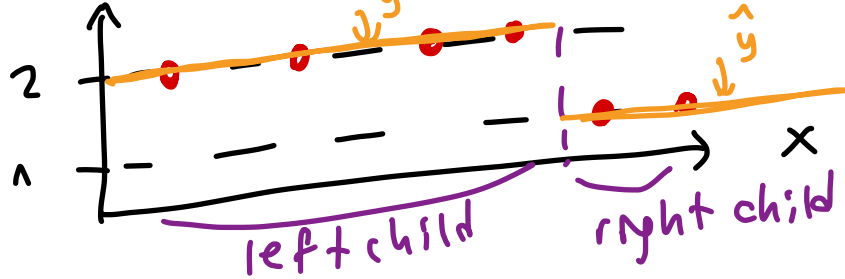
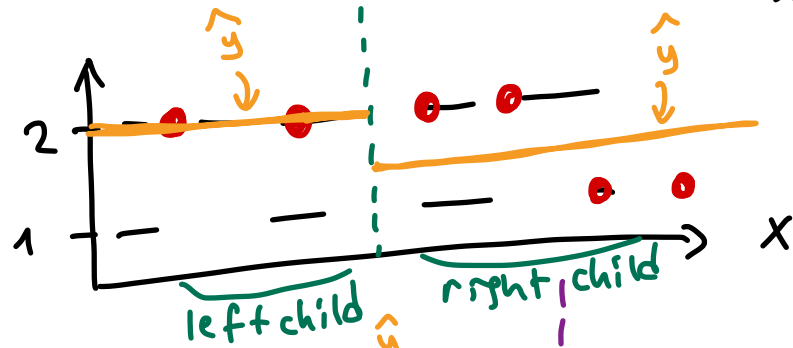
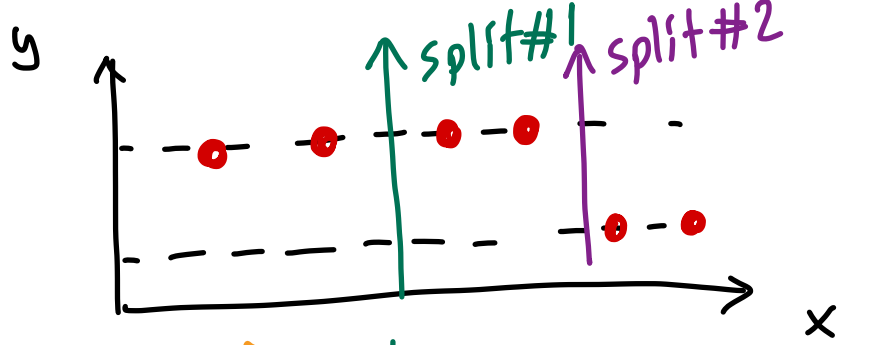
$\rightarrow$  # of data points that reach node m

$$g_m = \frac{\sum_{i=1}^N [y_i b_m(x_i)]}{\sum_{i=1}^N b_m(x_i)} \quad \left. \vphantom{\frac{\sum_{i=1}^N [y_i b_m(x_i)]}{\sum_{i=1}^N b_m(x_i)}} \right\} \text{average response (sample mean)}$$

error of  
a split

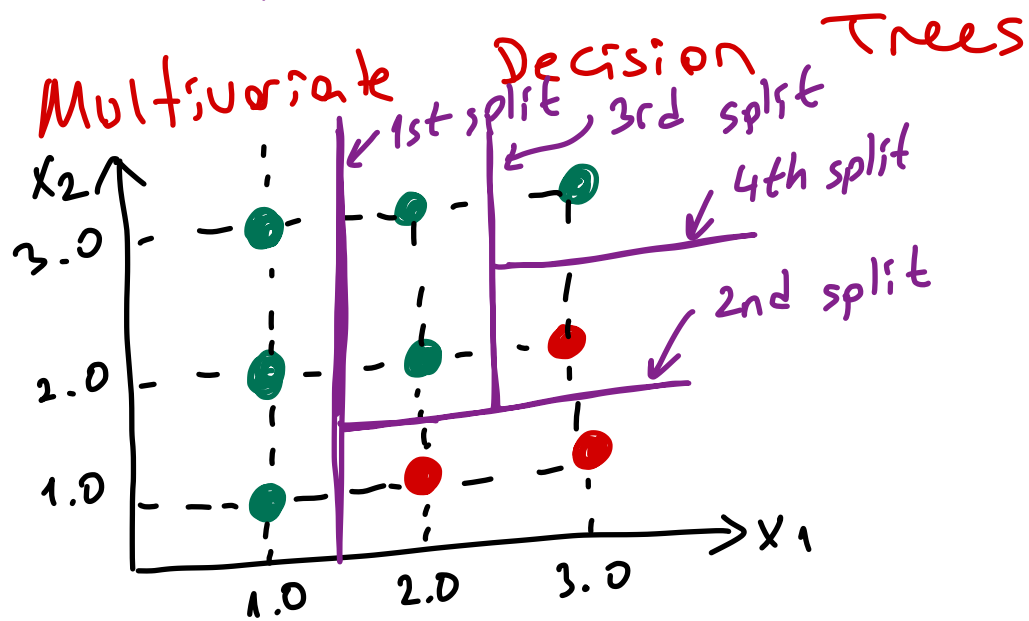
$$\bar{E}_m = \frac{1}{N_m} \sum_{s=1}^S \sum_{i=1}^N \left[ (y_i - g_{ms})^2 b_{ms}(x_i) \right]$$

$$\sum_{s=1}^S \frac{\cancel{N_{ms}}}{N_m} \sum_{i=1}^N \frac{1}{\cancel{N_{ms}}} \left[ (y_i - g_{ms})^2 b_{ms}(x_i) \right]$$

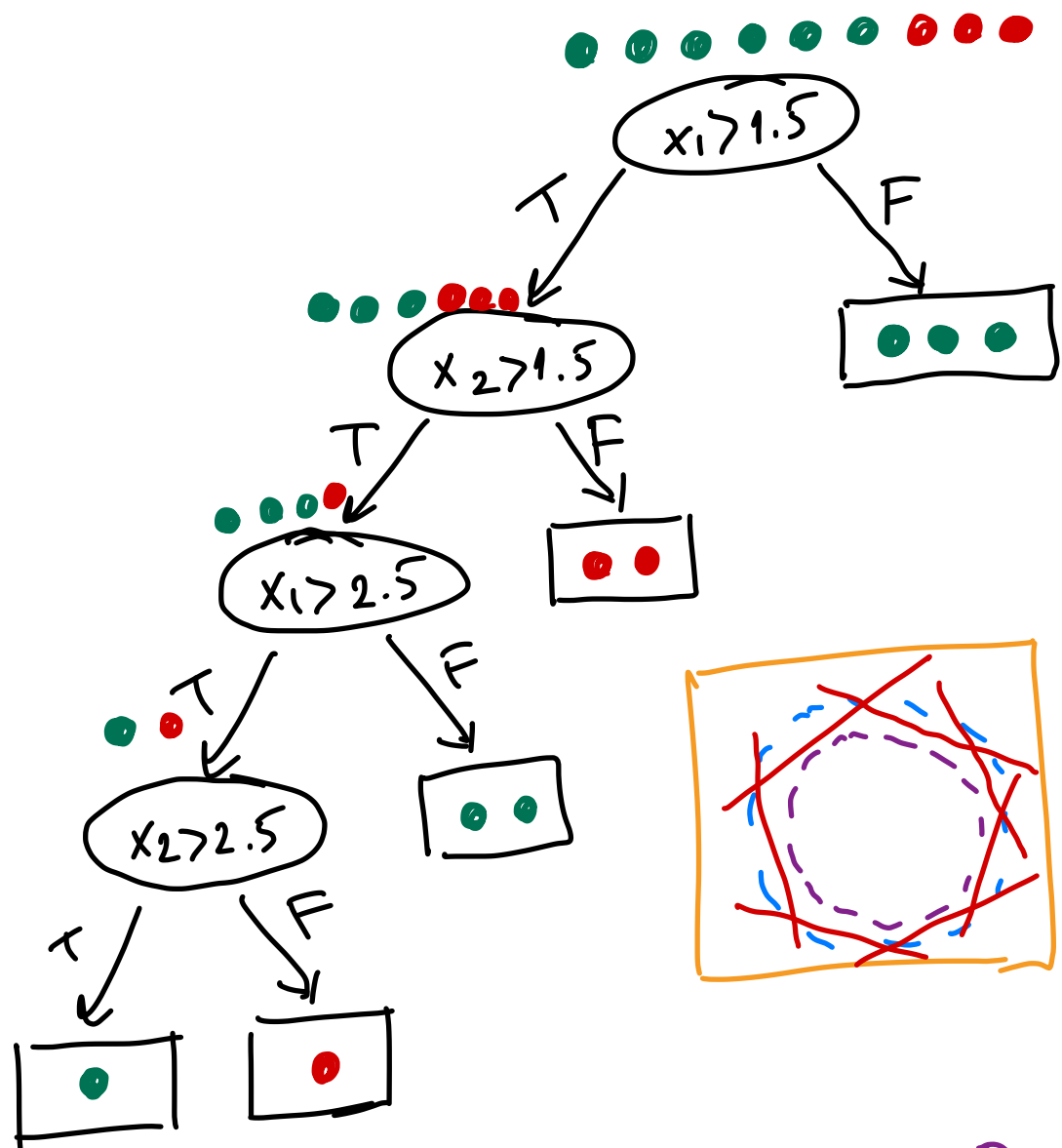


$$E(S_1) = \frac{1}{6} \left[ \begin{matrix} (2-2)^2 + (2-2)^2 + (2-1.5)^2 + \\ (2-1.5)^2 + (1-1.5)^2 + (1-1.5)^2 \end{matrix} \right] = \frac{1}{6}$$

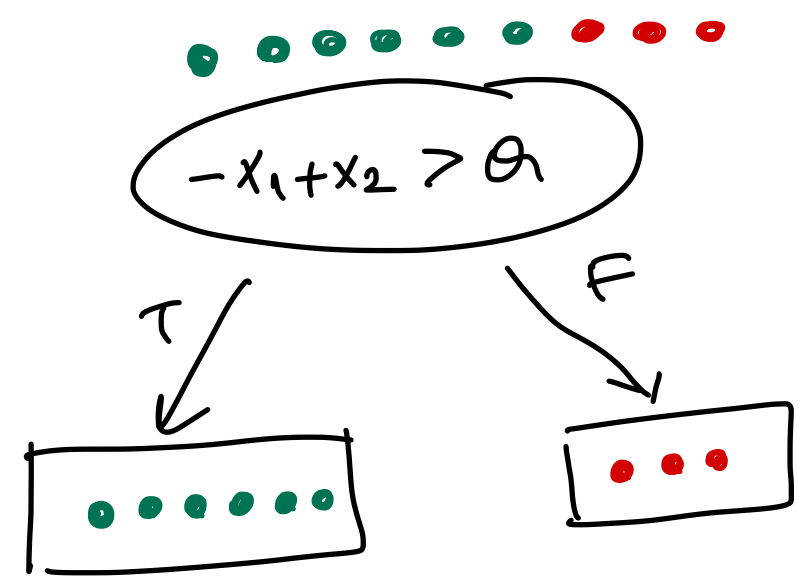
$$E(S_2) = \frac{1}{6} \left[ \begin{matrix} (2-2)^2 + (2-2)^2 + (2-2)^2 + \\ (2-2)^2 + (1-1)^2 + (1-1)^2 \end{matrix} \right] = 0$$



1st split  $x_1 > 1.5$   
 2nd split  $x_2 > 1.5$   
 3rd split  $x_1 > 2.5$   
 4th split  $x_2 > 2.5$



$$\underbrace{w_1}_{-1} x_1 + \underbrace{w_2}_{+1} x_2 + \underbrace{w_0}_{-9} > 0$$

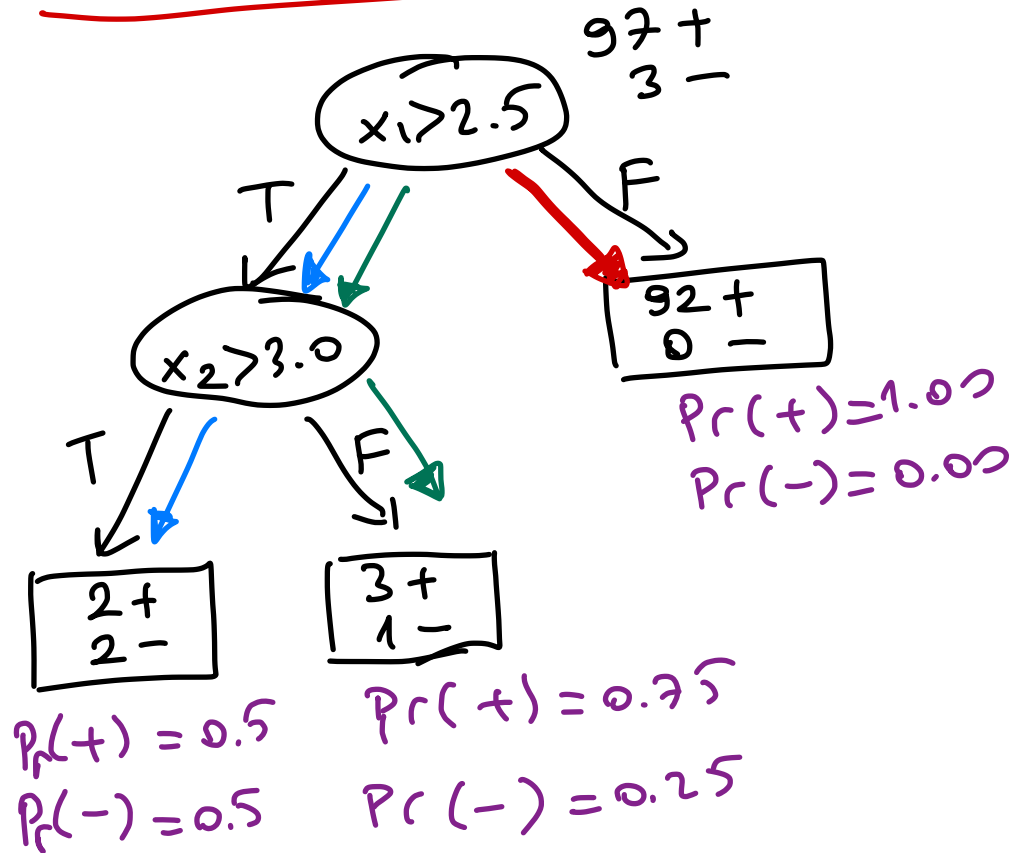


univariate  $f_m(\mathbf{x})$  :  $x_j > w_{m0} \Rightarrow x_j - w_{m0} > 0$

multivariate  $f_m(\mathbf{x})$  :  $\mathbf{w}_m^T \cdot \mathbf{x} + w_{m0} > 0$

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} - w_{m0} > 0$$

# RULE EXTRACTION

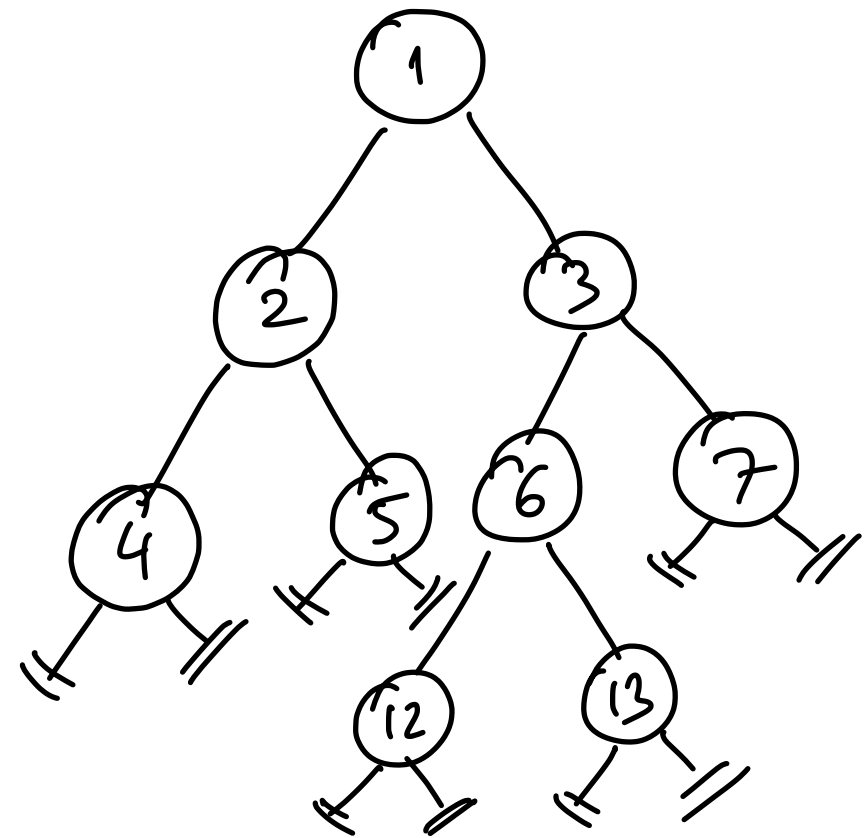


# of rules = # of leaves

Path 1  $x_1 > 2.5 \wedge x_2 > 3.0$  ?

Path 2  $x_1 > 2.5 \wedge x_2 \leq 3.0$  +

Path 3  $x_1 \leq 2.5$  +



left child =  $2 * \text{parent}$

right child =  $2 * \text{parent} + 1$

parent =  $\lfloor \text{child} / 2 \rfloor$

↑  
floor function

