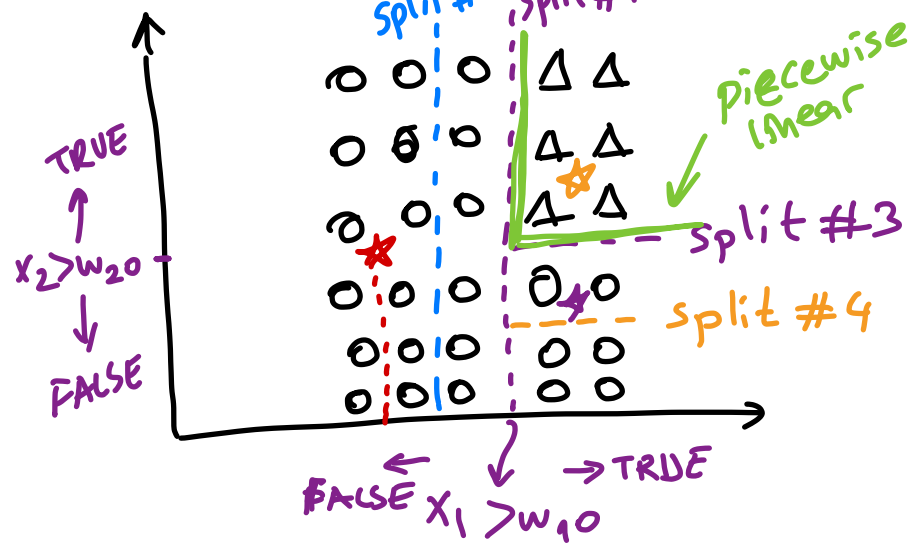
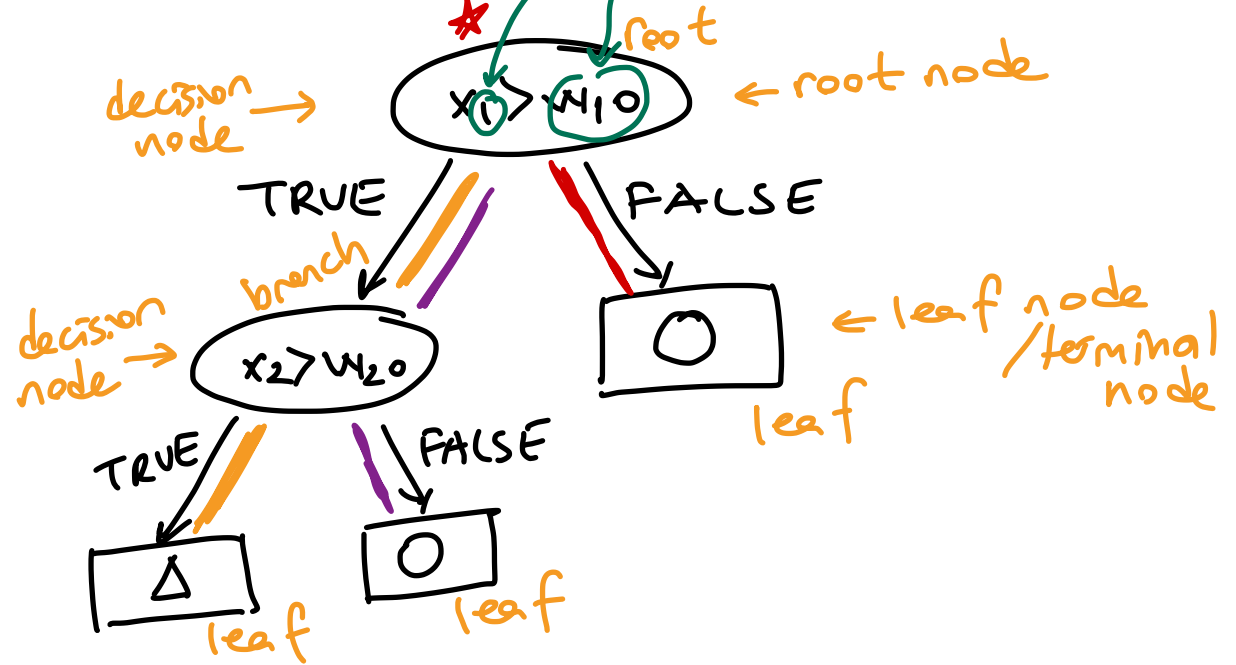


# Decision Trees



1st class: 0  
2nd class: Δ

which feature to split  
where to split



★ predicted label = ?

$x_1 > w_{10}$

FALSE  
↓

★ predicted label = ?

$x_1 > w_{10}$

TRUE  
↓

TRUE  $x_2 > w_{20}$   
↓  
Δ

★ predicted label = ?

$x_1 > w_{10}$

TRUE  
↓

$x_2 > w_{20}$   
↓  
FALSE  
↓  
0

if  $x_1 \leq w_{10} \Rightarrow \hat{y} = 0$

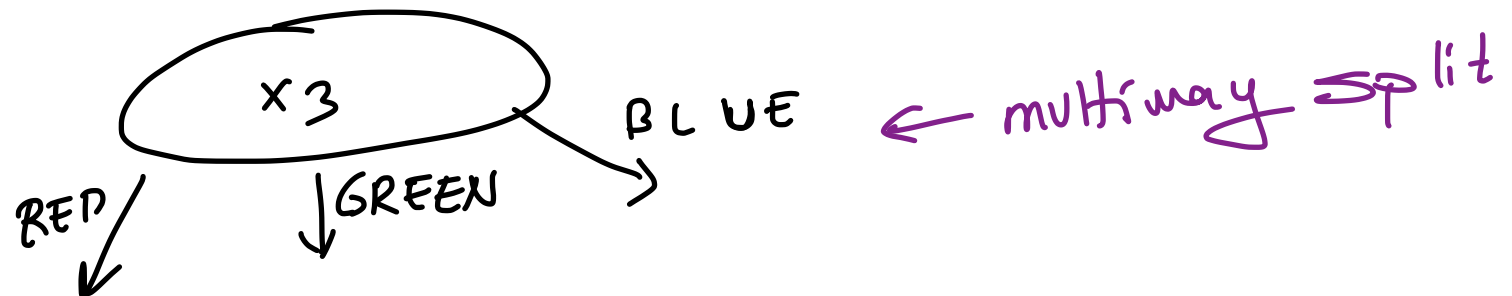
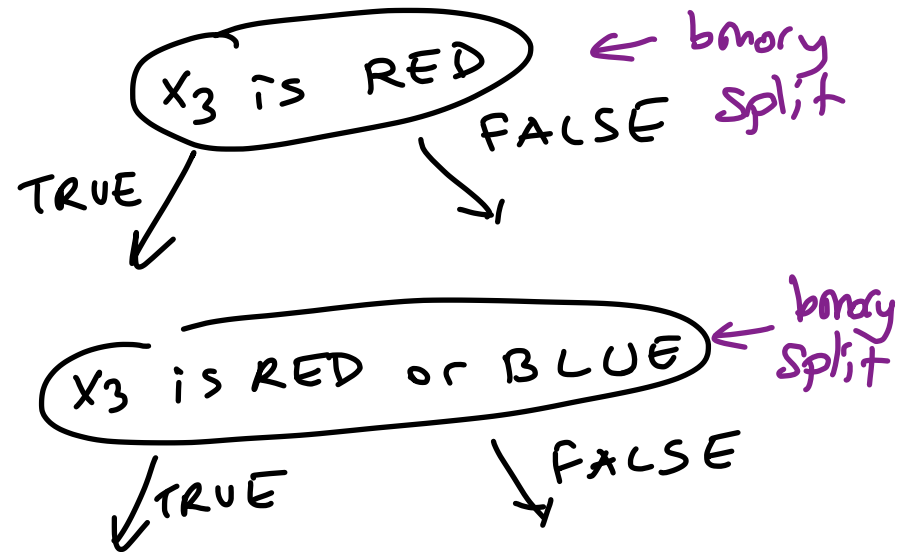
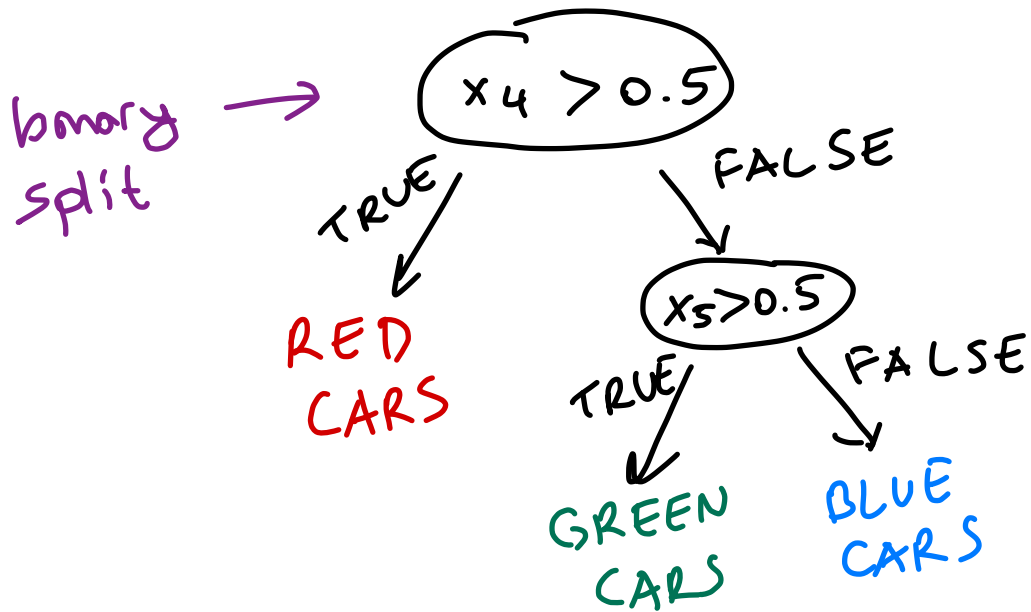
if  $x_1 > w_{10} \wedge x_2 \leq w_{20} \Rightarrow \hat{y} = 0$

if  $x_1 > w_{10} \wedge x_2 > w_{20} \Rightarrow \hat{y} = \Delta$

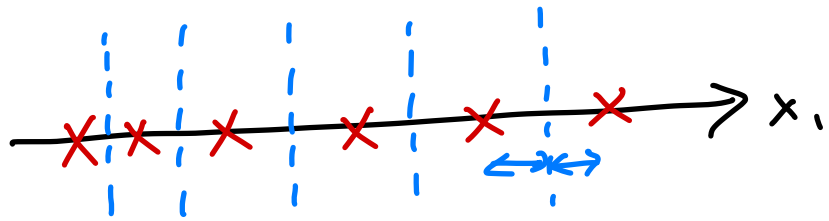
$X_3 = \begin{cases} \text{RED} \\ \text{GREEN} \\ \text{BLUE} \end{cases}$

$$\Rightarrow$$

	$x_4$	$x_5$	$x_6$
R	1	0	0
G	0	1	0
B	0	0	1



How can we learn on which feature and where to split?



$N$  data points  $\Rightarrow (N-1)$  possible splits  
 $D$  features  $\Rightarrow D(N-1)$  possible splits in total.

## Univariate Trees

Each decision node uses only one feature.

$$f_m(x): \quad x_j \stackrel{?}{>} w_{m0} \quad [x_j = w_{m0}]$$

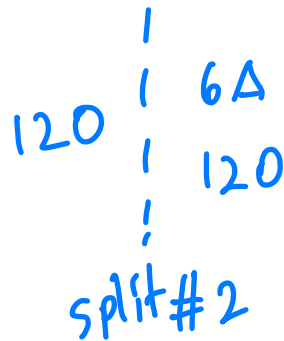
TRUE

FALSE

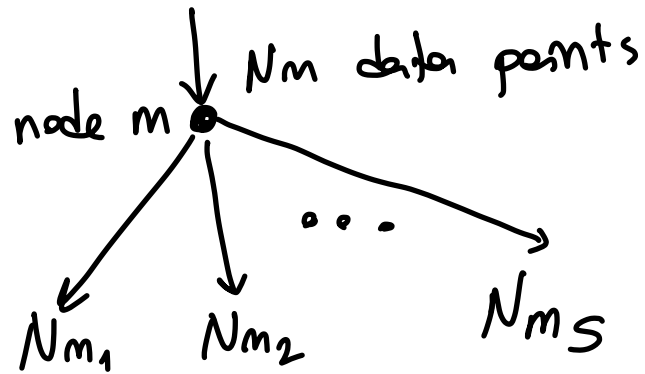
Left child

$$L_m = \{x \mid x_j > w_{m0}\} \quad [x_j = w_{m0}] \leftarrow$$

$$R_m = \{x \mid \underbrace{x_j \leq w_{m0}}_{[x_j \neq w_{m0}]}\} \quad \text{Right child}$$



Goodness of a split  $\Rightarrow$  Is Split #1 better than Split #2?



$S = \#$  of splits (branches)  
 $N_m = \#$  of data point that reach node  $m$   
 $K = \#$  of classes.

$$N_m = N_{m,1} + N_{m,2} + \dots + N_{m,S}$$

$$N_m = N_{m,1} + N_{m,2} + \dots + N_{m,K}$$

$$N_m = \sum_{s=1}^S N_{m,s} \quad (\text{splits})$$

$$N_m = \sum_{c=1}^K N_{m,c} \quad (\text{classes})$$

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

$$0 \cdot \log_2(0) \triangleq 0$$

$$I_m = - \sum_{c=1}^K P_{mc} \log_2(P_{mc})$$

Impurity of a node

$$I_m = - \left[ \frac{18}{18} \cdot \log_2\left(\frac{18}{18}\right) + \frac{0}{18} \cdot \log_2\left(\frac{0}{18}\right) \right] = 0$$

$$I_m = - \left[ \frac{6}{12} \cdot \log_2\left(\frac{6}{12}\right) + \frac{6}{12} \cdot \log_2\left(\frac{6}{12}\right) \right] = 1$$

$$18 \left\{ \begin{matrix} 18 & 0 \\ 0 & \Delta \end{matrix} \right.$$

$$12 \left\{ \begin{matrix} 6 & 0 \\ 6 & \Delta \end{matrix} \right.$$

12 0  
0 8  
12 0  
6 4

$$I_m = - \left[ \frac{12}{12} \log_2 \left( \frac{12}{12} \right) + \frac{0}{12} \cdot \log_2 \left( \frac{0}{12} \right) \right] = 0$$

$$I_m = - \left[ \frac{12}{18} \cdot \log_2 \left( \frac{12}{18} \right) + \frac{6}{18} \cdot \log_2 \left( \frac{6}{18} \right) \right] \approx 0.918$$

$I_m' = \sum_{s=1}^S \left( \frac{N_{m,s}}{N_m} \right) \left[ - \sum_{c=1}^K p_{m,s,c} \log_2(p_{m,s,c}) \right]$

*Impurity of a split*

*weights*

*impurity of a child node*

*impurity of the split*

*class index*  
*split index*  
*node index*

$$I_m'(\text{split \#1}) = \left[ \frac{18}{30} \cdot [0] + \frac{12}{30} \cdot [1] \right] = 0.4$$

*minimum is better*

$$I_m'(\text{split \#2}) = \left[ \frac{12}{30} \cdot [0] + \frac{18}{30} \cdot [0.918] \right] = 0.55$$

Split #1 is better than Split #2.

- ⇒ at each decision (internal) node
- ⇒ - for all features
    - for all possible splits
    - calculate impurity
    - pick the best split among all possible splits
- ⇒ stop when all terminal nodes are "pure"
- POSSIBLE PROBLEM ⇒ OVERFITTING  
(Training accuracy is 100%)

## PRUNING

### ① Prepruning

- ① [ - fix maximum depth  
- if you reach this depth, stop ]
- ② [ - you will not split if your node  
has a specified amount of your  
data set ]

### ② Postpruning

- grow your tree until it is completely pure
- prune your tree step by step until your misclassification error starts increasing on a validation data set.