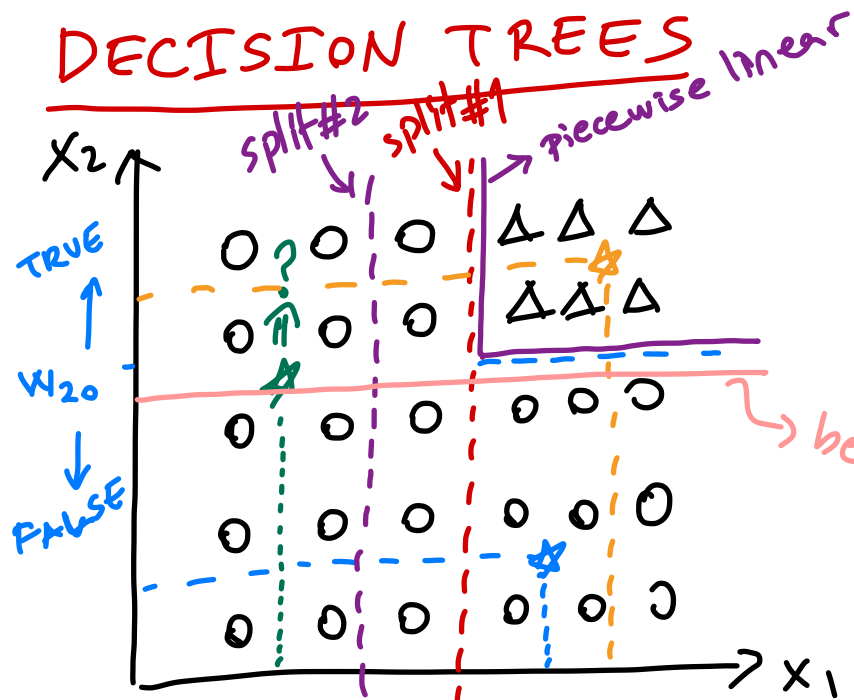


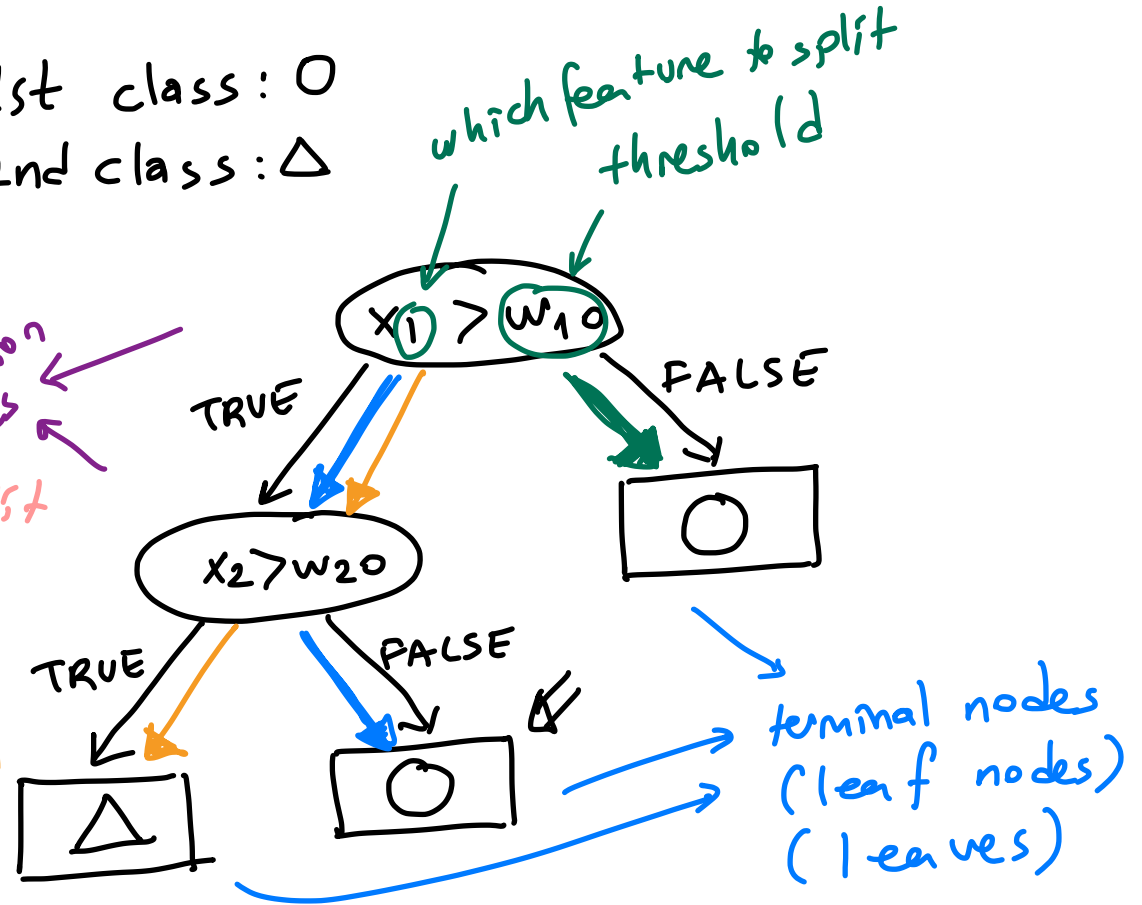
# DECISION TREES



FALSE  $\leftarrow w_{10}$  → TRUE

1st class: 0  
2nd class: Δ

decision nodes



★ (prediction = 0)

$x_1 > w_{10}$   
FALSE  
0

★ (prediction = 0)

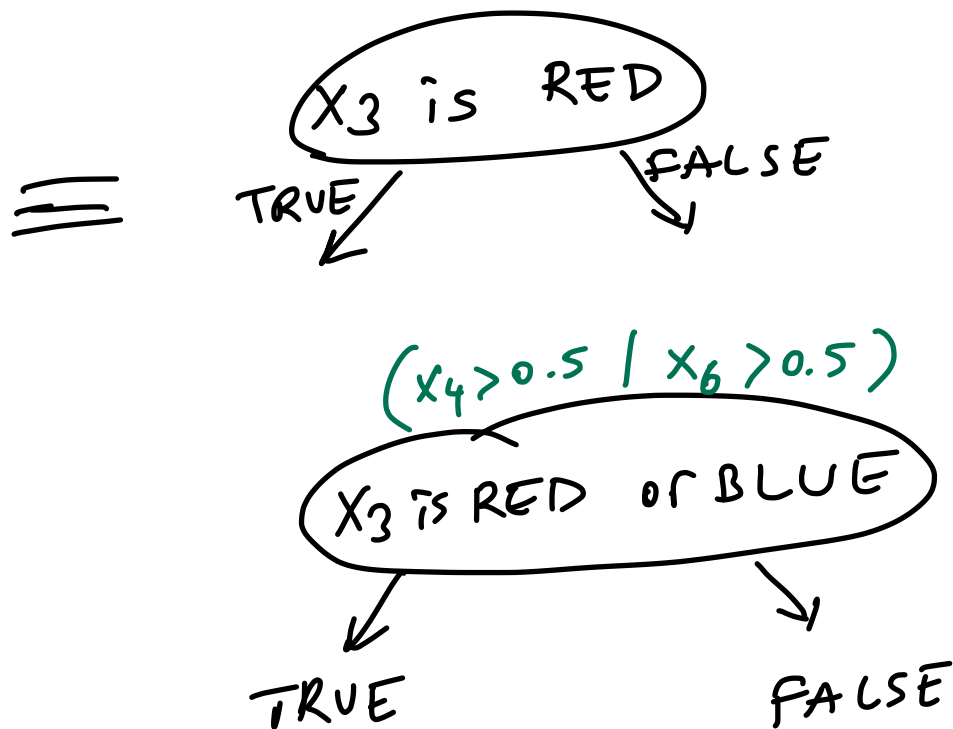
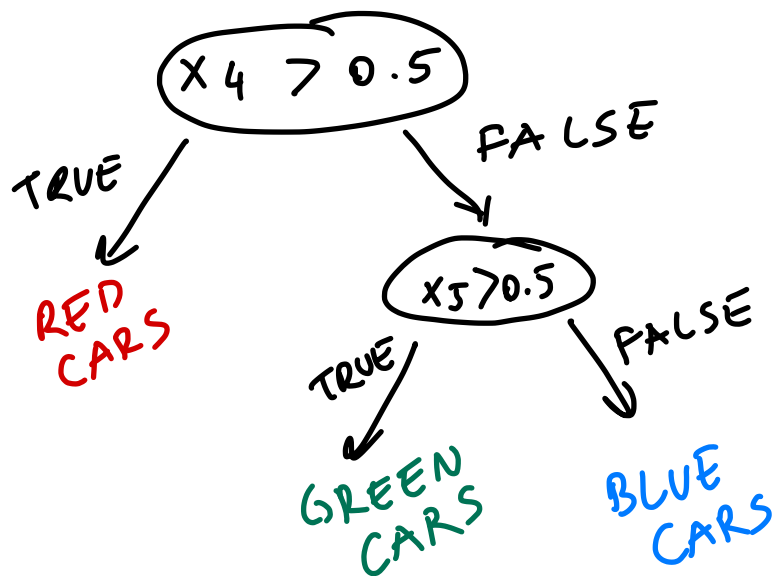
$x_1 > w_{10}$   
TRUE  
 $x_2 > w_{20}$   
FALSE  
0

★ (prediction = Δ)

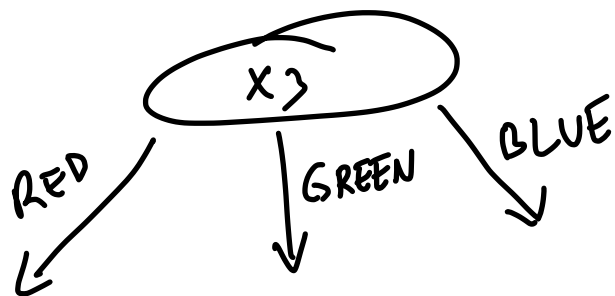
$x_1 > w_{10}$   
TRUE  
 $x_2 > w_{20}$   
TRUE  
Δ

- ① if  $x_1 \leq w_{10} \Rightarrow \hat{y} = 0$
- ② if  $x_1 > w_{10} \wedge x_2 \leq w_{20} \Rightarrow \hat{y} = 0$
- ③ if  $x_1 > w_{10} \wedge x_2 > w_{20} \Rightarrow \hat{y} = \Delta$

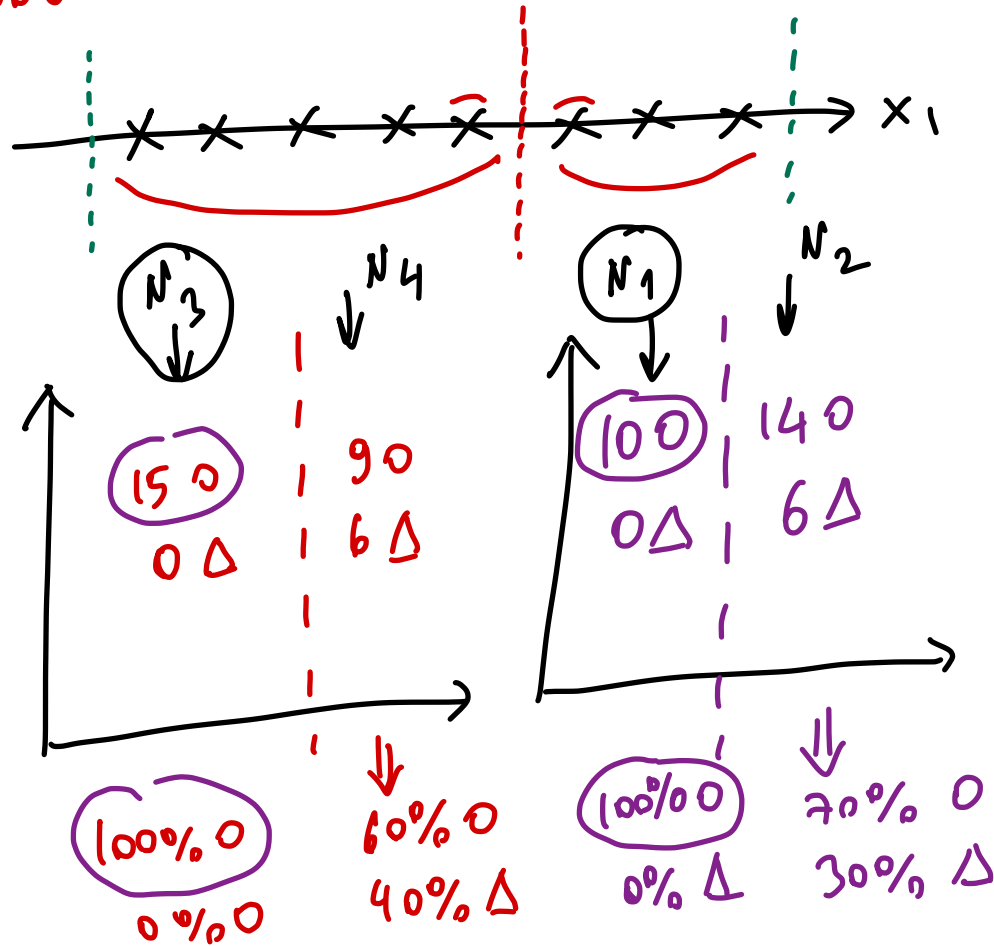
$$x_3 = \begin{cases} \text{RED} \\ \text{GREEN} \\ \text{BLUE} \end{cases} \Rightarrow \begin{array}{ccc} & R & G & B \\ x_4 & x_5 & x_6 & \\ 1 & 0 & 0 & \\ 0 & 1 & 0 & \\ 0 & 0 & 1 & \end{array}$$



multiway split



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

Univariate Trees

Each decision node uses only one feature.

$$f_m(x) : x_j > w_{mo} [x_j = w_{mo}]$$

TRUE

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

Left child

$$x_j = w_{mo}$$

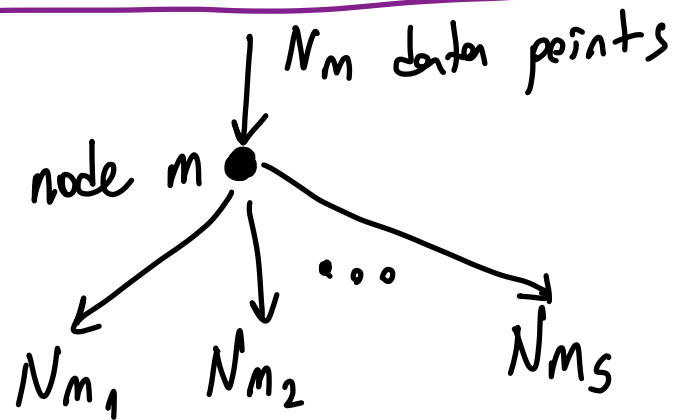
FALSE

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

Right child

$$w_j \neq w_{mo}$$

# Goodness of a split



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

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

$$P_{mc} = \hat{P}(y=c | x_m) = \frac{N_{m,c}}{N_m}$$

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

$$\begin{matrix} 10 & 0 \\ 0 & 0 \end{matrix}$$

$$N_1 \Rightarrow$$

$$I_m = - \left[ \frac{10}{10} \cdot \log_2\left(\frac{10}{10}\right) + \frac{0}{10} \cdot \log_2\left(\frac{0}{10}\right) \right] = 0$$

$$\begin{matrix} 14 & 0 \\ 6 & 6 \end{matrix}$$

$$N_2 \Rightarrow$$

$$I_m = - \left[ \frac{14}{20} \log_2\left(\frac{14}{20}\right) + \frac{6}{20} \log_2\left(\frac{6}{20}\right) \right] = 0.8813$$

Is split #1 is better than Split #2?

$S$  = # of splits (branches)

$N_m$  = # of data points that reach node m

$K$  = # of classes

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

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

$$0 \cdot \log_2(0) \equiv 0$$

$${}^{15}_0\Delta N_3 \Rightarrow I_m = - \left[ \frac{15}{15} \log_2 \left( \frac{15}{15} \right) + \frac{0}{15} \cdot \log_2 \left( \frac{0}{15} \right) \right] = 0$$

$${}^{30}_6\Delta N_4 \Rightarrow I_m = - \left[ \frac{9}{15} \log_2 \left( \frac{9}{15} \right) + \frac{6}{15} \log_2 \left( \frac{6}{15} \right) \right] = 0.9710$$

$$\underset{\text{impurity}}{\overset{I'_m}{\circ}} = \sum_{s=1}^S \underset{\substack{\text{weights.} \\ \text{impurity of a child node}}}{\overset{\frac{N_{m,s}}{N_m}}{\circ}} \left[ - \sum_{c=1}^k P_{msc} \log_2 (P_{msc}) \right]$$

$\swarrow$  class index  
 $\swarrow$  split index  
 $\swarrow$  node index

impurity of the split

$$I'_m(\text{split \#1}) = \left[ \frac{15}{30} \cdot [0] + \frac{15}{30} \cdot [0.9710] \right] = \underset{\text{minimum is better.}}{\circ} 0.4855$$

$$I'_m(\text{split \#2}) = \left[ \frac{10}{30} [0] + \frac{20}{30} [0.8813] \right] = 0.5875$$

Split #1 is better than Split #2.

⇒ at each internal (decision) node

- ⇒ [
- for all features
  - for all possible splits
  - calculate impurity
  - pick the best split among all possible splits (the one with minimum impurity)
- ]

⇒ Stop when all terminal nodes are "pure"

POSSIBLE PROBLEM

⇒ OVERFITTING

(Training accuracy is 100%)

## PRUNING

### ① Prepruning

- ① [
- fix the 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 Post pruning

- 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.