

# Decision Trees

## Purity

purity of a region  $R_j$

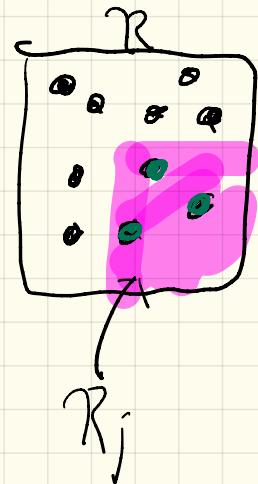
fraction of the mixture of classes of pts in  $R_j$   
by majority label

let  $D_j$  be the pts in  $D$  w/in  $R_j$

$$n_j = |D_j|$$

$n_{ji}$  = # of pts in  
 $D_j$  w/ label  $i$

$$\text{purity}(D_j) = \max_i \left\{ \frac{n_{ji}}{n_j} \right\}$$



# Decision Tree Algo

Given dataset  $D$

leaf size  $\eta \leftarrow (\text{eta})$

purity threshold  $\pi$

Recursively:

pick "best" split point  $(X_j, v)$

partition data in  $D_Y, D_N$

recurse( $D_Y$ )

recurse( $D_N$ )

Stop:

- # of pts in partition is below  $\eta$  (avoid overfitting)

or - purity is above  $\pi$

---

## Algorithm 19.1: Decision Tree Algorithm

---

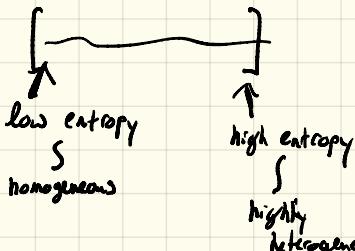
DECISIONTREE ( $D, \eta, \pi$ ):

```
1  $n \leftarrow |\mathbf{D}|$  // partition size
2  $n_i \leftarrow |\{x_j | x_j \in \mathbf{D}, y_j = c_i\}|$  // size of class  $c_i$ 
3  $\text{purity}(\mathbf{D}) \leftarrow \max_i \left\{ \frac{n_i}{n} \right\}$ 
4 if  $n \leq \eta$  or  $\text{purity}(\mathbf{D}) \geq \pi$  then // stopping condition
5    $c^* \leftarrow \arg \max_{c_i} \left\{ \frac{n_i}{n} \right\}$  // majority class
6   create leaf node, and label it with class  $c^*$ 
7   return
8 ( $\text{split point}^*, \text{score}^*$ )  $\leftarrow (\emptyset, 0)$  // initialize best split point
9 foreach (attribute  $X_j$ ) do
10   if ( $X_j$  is numeric) then
11      $(v, \text{score}) \leftarrow \text{EVALUATE-NUMERIC-ATTRIBUTE}(\mathbf{D}, X_j)$ 
12     if  $\text{score} > \text{score}^*$  then  $(\text{split point}^*, \text{score}^*) \leftarrow (X_j \leq v, \text{score})$ 
13   else if ( $X_j$  is categorical) then
14      $(V, \text{score}) \leftarrow \text{EVALUATE-CATEGORICAL-ATTRIBUTE}(\mathbf{D}, X_j)$ 
15     if  $\text{score} > \text{score}^*$  then  $(\text{split point}^*, \text{score}^*) \leftarrow (X_j \in V, \text{score})$ 
16 // partition  $\mathbf{D}$  into  $\mathbf{D}_Y$  and  $\mathbf{D}_N$  using  $\text{split point}^*$ , and call
    recursively
17  $\mathbf{D}_Y \leftarrow \{x^T | x \in \mathbf{D} \text{ satisfies } \text{split point}^*\}$ 
18  $\mathbf{D}_N \leftarrow \{x^T | x \in \mathbf{D} \text{ does not satisfy } \text{split point}^*\}$ 
19 create internal node  $\text{split point}^*$ , with two child nodes,  $\mathbf{D}_Y$  and  $\mathbf{D}_N$ 
20 DECISIONTREE( $\mathbf{D}_Y$ ); DECISIONTREE( $\mathbf{D}_N$ )
```

---

## Split point evaluation measures:

**Entropy** — amount of disorder in a system



$$H(D) = - \sum_{i=1}^k P(c_i | D) \log P(c_i | D)$$

# of class labels in D

prob of a class  $C_i$  in dataset D

$$\frac{\text{# of records w/ class label } C_i}{\text{# of records in } D}$$

$$P(c_e | D) = 1$$

$$\Rightarrow \forall j \neq e \quad P(c_j | D) = 0$$

$$\Rightarrow H(D) = - [$$

$$\begin{array}{c} 0 \cdot \log 1 \\ -0 \cdot \log 1 \\ \vdots \end{array}$$

$$\begin{array}{c} i=1 \\ i=2 \\ \vdots \\ i=e \end{array}$$

$$1 \cdot \log(1)$$

$$i=e$$

$$\begin{array}{c} 0 + \frac{1}{e} \\ \hline i \\ = 0 \end{array}$$

—

$$0$$

$$\Rightarrow \text{no entropy}$$

$$\begin{aligned} P(c_i | D) &= \frac{1}{k} \quad \forall i \\ \Rightarrow H(D) &= - \sum_{i=1}^k \frac{1}{k} \log(\frac{1}{k}) \end{aligned}$$

$$\begin{aligned} &= -\frac{1}{k} \sum_{i=1}^k \log(\frac{1}{k}) \\ &= -\frac{1}{k} \log \frac{k}{k} \end{aligned}$$

$$= -\frac{1}{k} \log k$$

$$= (\frac{1}{k}) \cancel{k} \log k$$

$$= \log(k) \leq \max_{\text{entropy}}$$

Let  $D_y, D_N$  be partitions of  $D$

$$n = |D| \quad n_y = |D_y| \quad n_N = |D_N|$$

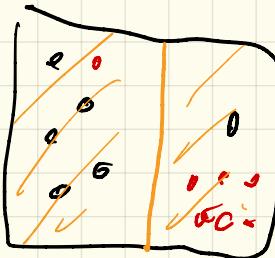
split entropy: weighted entropy of each partition

$$H(D_y, D_N) = \frac{n_y}{n} H(D_y) + \frac{n_N}{n} H(D_N)$$

information gain at a split point as decrease in entropy

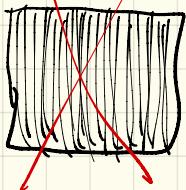
$$\text{Gain}(D, D_y, D_N) = H(D) - H(D_y, D_N)$$

$\Rightarrow$  high information gain means better split point  
pick split point w/ highest information gain

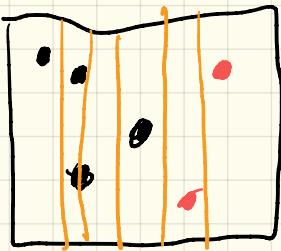


Eval split pts

Alg 0 for fingering split pts  
compute  $D_y$ ,  $D_N$  and  $P(c_i | D_y)$ ,  $P(c_i, D_N)$ , Gain  $(D, D_y, D_N)$   
try all



Alg 1



$n$  pts  
 $\Rightarrow n-1$  places to consider splitting