

# Split Points

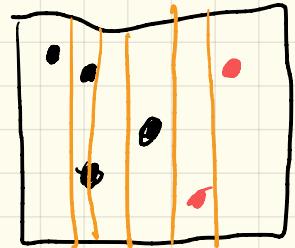
Let  $U = \{u_1, \dots, u_m\}$  be our set of midpoints

$$\exists u_1 < u_2 < \dots < u_m$$

We need to estimate

$$\hat{P}(c_i | D_y) = \hat{P}(c_i | X \leq u)$$

$$\hat{P}(c_i | D_N) = \hat{P}(c_i | X > u)$$



Let  $I$  indicator random var

$1$  when arg is true  
 $0$  when arg is false

E.g. Let  $B$  be a boolean statement

$$I(B) = \begin{cases} 0 & \text{if } B \text{ is false} \\ 1 & \text{if } B \text{ is true} \end{cases}$$

Let  $\vec{x}_i$  be labeled  $y_i = c_3$

$$y_i = c_1$$

$$I(y_i = c_1) = 0$$

$$I(y_i = c_2) = 0$$

$$I(y_i = c_3) = 1$$

$$I(y_i = c_4) = 0$$

using Bayes thm

$$\hat{P}(c_i | X \leq u) = \frac{\hat{P}(X \leq u | c_i) \hat{P}(c_i)}{\hat{P}(X \leq u)} = \frac{\hat{P}(X \leq u | c_i) \hat{P}(c_i)}{\sum_{j=1}^k \hat{P}(X \leq u | c_j) \hat{P}(c_j)}$$

note  $\sum_{i=1}^k \hat{P}(X \leq u | c_i)$

note 2  $P(X \leq u | c_j) = \frac{P(X \leq u \text{ and } c_j)}{P(c_j)} \Leftrightarrow P(X \leq u \text{ and } c_j) = P(X \leq u | c_j) \cdot P(c_j)$

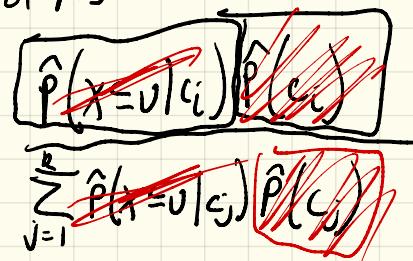
$$\hat{P}(c_i | x \leq v) = \frac{\hat{P}(x \leq v | c_i) \hat{P}(c_i)}{\sum_{j=1}^k \hat{P}(x \leq v | c_j) \hat{P}(c_j)}$$

# of pts w/ label i

$$\hat{P}(c_i) = \frac{1}{n} \sum_{j=1}^n I(y_j = c_i) = \frac{n_i}{n}$$

# of pts

$y_i$  is the  
label of pt  $\hat{x}_i$



Next  $\hat{P}(x \leq v | c_i) = \frac{N_{vi}}{n_i}$

Let  $N_{vi}$  be the number of pts w/ class  $c_i$   
and val less than our split

$$N_{vi} = \sum_{j=1}^n I(x_j \leq v \text{ and } y_j = c_i)$$

$$\hat{P}(x \leq v | c_i) = \frac{\hat{P}(x \leq v \text{ and } c_i)}{\hat{P}(c_i)} = \frac{\frac{1}{n} N_{vi}}{\left(\frac{n_i}{n}\right)} = \frac{N_{vi}}{n_i}$$

$$\hat{P}(c_i | \mathbf{D}_y) = \frac{\hat{P}(X \leq v | c_i) \hat{P}(c_i)}{\sum_{j=1}^k \hat{P}(X \leq v | c_j) \hat{P}(c_j)} = \frac{\frac{N_{vi}}{n}}{\sum_{j=1}^k \frac{N_{vji}}{n}} = \frac{\frac{N_{vi}}{n}}{\frac{1}{n} \sum_{j=1}^k N_{vji}}$$

Simplify with:

$$\hat{P}(X \leq v | c_i) \hat{P}(c_i) = \left( \frac{N_{vi}}{n_i} \right) \left( \frac{n_i}{n} \right) = \frac{N_{vi}}{n}$$

$$\hat{P}(c_i | \mathbf{D}_W) = \hat{P}(c_i | X > v)$$

using same idea as before:

$$\hat{P}(c_i | \mathbf{D}_N) = \hat{P}(c_i | X > v) = \frac{\hat{P}(X > v | c_i) \hat{P}(c_i)}{\sum_{j=1}^k \hat{P}(X > v | c_j) \hat{P}(c_j)}$$

$$\hat{P}(X > v | c_i) = 1 - \hat{P}(X \leq v | c_i) = 1 - \frac{N_{vi}}{n_i} = \frac{n_i - N_{vi}}{n_i}$$

and plug back in and simplify

$$\hat{P}(c_i | \mathbf{D}_N) = \hat{P}(c_i | X > v) = \frac{\hat{P}(X > v | c_i) \hat{P}(c_i)}{\sum_{j=1}^k \hat{P}(X > v | c_j) \hat{P}(c_j)} = \frac{n_i - N_{vi}}{\sum_{j=1}^k (n_j - N_{vj})}$$

Now we can evaluate split points

compute midpts  
sort midpts  
 $O(n \log n)$

Compute  $N_{vi}$ 's  
 $O(nk)$

$O(nk)$

$O(k)$

---

### Algorithm 19.2: Evaluate Numeric Attribute (Using Gain)

---

#### EVALUATE-NUMERIC-ATTRIBUTE ( $\mathbf{D}, X$ ):

```
1 sort  $\mathbf{D}$  on attribute  $X$ , so that  $x_j \leq x_{j+1}, \forall j = 1, \dots, n - 1$ 
2  $\mathcal{M} \leftarrow \emptyset$  // set of midpoints
3 for  $i = 1, \dots, k$  do  $n_i \leftarrow 0$ 
4 for  $j = 1, \dots, n - 1$  do
    5 if  $y_j = c_i$  then  $n_i \leftarrow n_i + 1$ 
    // running count for class  $c_i$ 
    6 if  $x_{j+1} \neq x_j$  then
        7  $v \leftarrow \frac{x_{j+1} + x_j}{2}; \mathcal{M} \leftarrow \mathcal{M} \cup \{v\}$  // midpoints
        8 for  $i = 1, \dots, k$  do
            9  $N_{vi} \leftarrow n_i$  // Number of points such that  $x_j \leq v$  and
             $y_j = c_i$ 
10 if  $y_n = c_i$  then  $n_i \leftarrow n_i + 1$ 
// evaluate split points of the form  $X \leq v$ 
11  $v^* \leftarrow \emptyset; score^* \leftarrow 0$  // initialize best split point
12 forall  $v \in \mathcal{M}$  do
13 for  $i = 1, \dots, k$  do
14  $\hat{P}(c_i | \mathbf{D}_Y) \leftarrow \frac{N_{vi}}{\sum_{j=1}^k N_{vj}}$   $\leftarrow O(k)$ 
15  $\hat{P}(c_i | \mathbf{D}_N) \leftarrow \frac{n_i - N_{vi}}{\sum_{j=1}^k n_j - N_{vj}}$ 
16  $score(X \leq v) \leftarrow Gain(\mathbf{D}, \mathbf{D}_Y, \mathbf{D}_N)$  // use Eq. (19.5)
17 if  $score(X \leq v) > score^*$  then
18  $v^* \leftarrow v; score^* \leftarrow score(X \leq v)$ 
19 return  $(v^*, score^*)$ 
```

---

Note: DO NOT COMPUTE  
MIDPOINTS LIKE THIS  
CAN CAUSE OVERFLOW  
 $x_j + \frac{x_{j+1} - x_j}{2}$

↓  
total time

$O(n \lg n + nk)$  as  $k$  is usually small (e.g. # labels)  
we treat  $k$  as a const

$O(n \lg n)$  work to find best split in <sup>each</sup> 1 dim

$n = |\mathbf{D}|$

assume # of labels  $k$  is const

$\Rightarrow$  complexity of computing split point in a dim is  $O(n \log n)$

---

#### Algorithm 19.1: Decision Tree Algorithm

---

**DECISIONTREE ( $\mathbf{D}, \eta, \pi$ ):**

```

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  $purity(\mathbf{D}) \leftarrow \max_i \left\{ \frac{n_i}{n} \right\}$ 
4 if  $n \leq \eta$  or  $purity(\mathbf{D}) \geq \pi$  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  $(split point^*, score^*) \leftarrow (\emptyset, 0)$  // initialize best split point
9 foreach (attribute  $X_j$ ) do
10  if ( $X_j$  is numeric) then
11     $(v, score) \leftarrow EVALUATE-NUMERIC-ATTRIBUTE(\mathbf{D}, X_j)$ 
12    if  $score > score^*$  then  $(split point^*, score^*) \leftarrow (X_j \leq v, score)$ 
13  else if ( $X_j$  is categorical) then
14     $(V, score) \leftarrow EVALUATE-CATEGORICAL-ATTRIBUTE(\mathbf{D}, X_j)$ 
15    if  $score > score^*$  then  $(split point^*, score^*) \leftarrow (X_j \in V, score)$ 
16 // partition  $\mathbf{D}$  into  $\mathbf{D}_Y$  and  $\mathbf{D}_N$  using  $split point^*$ , and call
17 // recursively
18  $\mathbf{D}_Y \leftarrow \{x^T | x \in \mathbf{D} \text{ satisfies } split point^*\}$ 
19  $\mathbf{D}_N \leftarrow \{x^T | x \in \mathbf{D} \text{ does not satisfy } split point^*\}$ 
18 create internal node  $split point^*$ , with two child nodes,  $\mathbf{D}_Y$  and  $\mathbf{D}_N$ 
19 DECISIONTREE( $\mathbf{D}_Y$ ); DECISIONTREE( $\mathbf{D}_N$ )

```

---

find good split point  
for each dim d  
run  $O(n \log n)$   
time alg  
to find a split  
point in dim



$O(dn \log n)$   
for each  
recursive call

How many times do we recurse?  
 $\Rightarrow$  How deep can our tree be  
at worst  $O(n)$

$\Rightarrow$  routine that is  
 $O(dn \log n)$   
 $O(n)$  times

$\Rightarrow O(n \cdot d \cdot n \log n)$  time =  $O(dn^2 \log n)$

