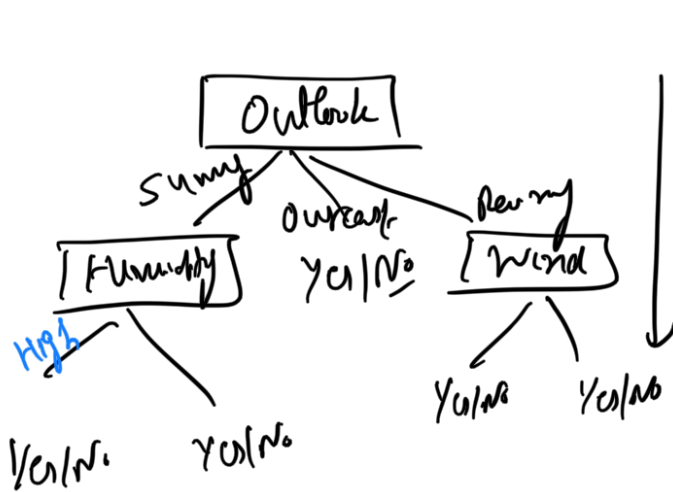


Coh 7 & 4  
Machine Learning  
Nov 27, 2020

Logistics  
Make Extra  
Class:-  
Sunday Nov 29  
8am - 9am

Last Class:- Decision Trees:-



Boolean =  $2^L$   
 ↳ Axis || Rectangles.  
 Learning Decision Trees from  
 Data  $\{x^i, y^i\}_{i=1}^m$

Inductive Bias:  
 Grow trees which are  
 "Small" in size

Decision Tree Learning Algorithm:-

Grow Tree (D)  $\Sigma$   
 // Leaf creation  
 if  $\forall i, (x^i, y^i) \in D, y^i = 0$   
 Return createLeaf(0);  
 if  $\forall i, (x^i, y^i) \in D, y^i = 1$   
 Return createLeaf(1);  
 // split  
 $x_d \leftarrow \text{ChooseBestAttrToSplit}(D);$   
 $n_d = \text{NewNode}(x_d, \text{NULL});$   
 $\forall \ell \in \{1 \dots k\} \Sigma$  # of attribute  
 # value  
 $x_d \in \{v_1 \dots v_k\}$   
 $D_\ell = \{x^i, y^i\} \in D \text{ s.t. } x_d^i = v_\ell$   
 $n_d.add(\text{GrowTree}(D_\ell, v_\ell));$

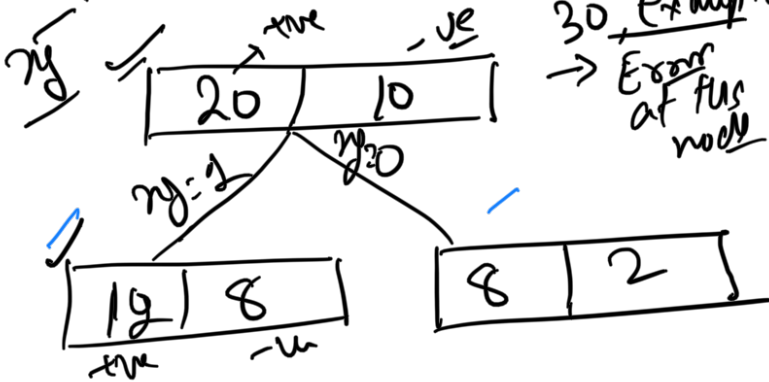
Data points  
 such that  
 value of  $x_d$   
 attribute is  $v_\ell$

3

return  $n_{ji}$

3

How do we choose the Attribute to split on?



$$E_{x_2} = 10$$

$$E_{x_2=v_1} = 8$$

$$E_{x_2=v_2} = 2$$

$$E_{x_2 \text{ split}} = 8 + 2 = 10$$

No Reduction on Error "Progress"

Entropy :- Measure Randomness in the system

Consider a Random Variable

$$Y \in \{1, \dots, K\}$$

$$P(Y=K) = P_K$$

Entropy

$$H(Y) = - \sum_{K=1}^K P_K \log P_K$$

$$= -E[\log P(Y=K)]$$

Assume a simple Boolean Var.  $n \leq 2$

$$[p_0, p_1] = 1$$

$$H(Y) = -[p_0 \log p_0 + p_1 \log p_1]$$

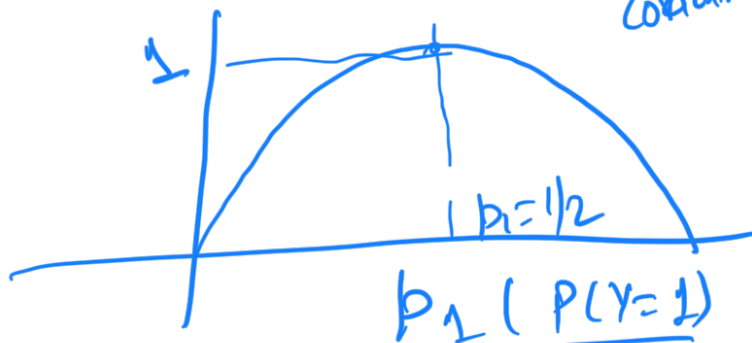
$$\text{so } \rightarrow = -[p \log p + (1-p) \log(1-p)]$$

$$p \log p :- \text{ as } p \rightarrow 0 \quad p \log p \rightarrow 0$$

$$\hookrightarrow \text{Using } \lim_{p \rightarrow 0} \frac{\log p}{1/p} = \frac{d(\log p)}{d(1/p)} = \frac{1/p}{-1/p^2} = -p \rightarrow 0$$

$$p \log p :- \text{ as } p \rightarrow 1 \quad p \log p \rightarrow 0$$

$\Rightarrow$  The <sup>(max)</sup> value of  $-[p \log p]$  is ~~maximum~~ obtained at  $\underline{0.5}$  concave function



$H(Y)$

$Y \in \{0, 1\}$

Deep connection with Information Theory:-  
Entropy, "bits" required to transmit a message is maximum

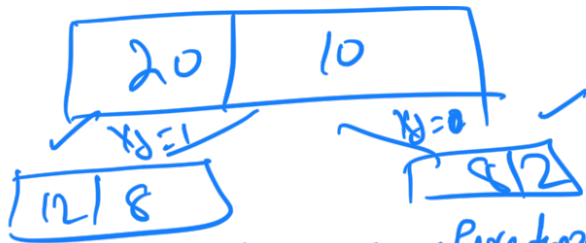
$Y \in \{1, \dots, r\}$

$$\hookrightarrow H(Y) = \sum_{k=1}^r -p_k \log p_k$$

when  $p_k = 1/r \quad \forall k$

$\hookrightarrow$  Amount of randomness in the system

Our goal:- To reduce the entropy of class label.



$$H(Y) = \frac{2}{3} \log \frac{3}{2} + \frac{1}{3} \log \frac{3}{1}$$

How do we characterize a split?

$$H(Y|X) = \sum_x P(X=x) H(Y|X=x)$$

we'll define  
Entropy of the distribution  
 $P(Y|X=x)$

Which attribute to split on?

$\#j$ : split on the attribute which  
minimizes  $H(Y|X_j)$

$$\#j^* = \underset{j}{\operatorname{argmin}} H(Y|X_j)$$

↓ Attribute to split on

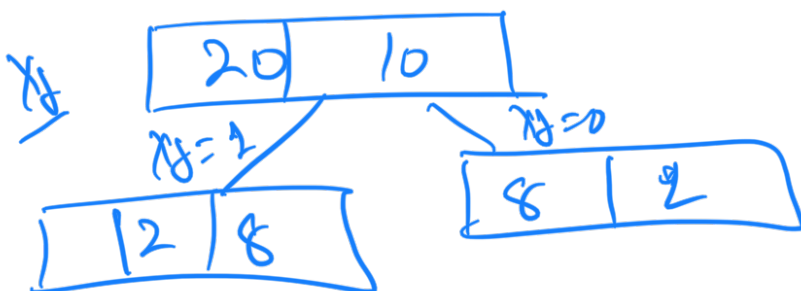
$$\boxed{MI(Y, X) = H(Y) - H(Y|X)}$$

Symmetry property

$$= MI(X, Y)$$

$$= H(X) - H(X|Y)$$

$$\Rightarrow \boxed{\begin{aligned} \#j^* &= \underset{j}{\operatorname{argmax}} MI(Y, X_j) \\ &= \underset{j}{\operatorname{argmax}} H(Y) - H(Y|X_j) \end{aligned}}$$



$$H(Y) = \frac{2}{3} \log \frac{3}{2} + \frac{1}{3} \log \frac{3}{1}$$

= 1

$$H(Y|X_j=1) = \frac{3}{5} \log \frac{5}{3} + \frac{2}{5} \log \frac{5}{2} \quad \text{--- (2)}$$

$$H(Y|X_j=0) = \frac{4}{5} \log \frac{5}{4} + \frac{1}{5} \log \frac{5}{1} \quad \text{--- (3)}$$

$$MI(Y, X_j) = H(Y) - \underline{H(Y|X_j)}$$

compute these values (Hint)

$$= \cancel{\frac{2}{3}} H(Y) - \left[ H(Y|X_j=1) \frac{2}{3} + H(Y|X_j=0) \frac{1}{3} \right]$$

$$= \underline{\underline{??}} \quad \underline{\underline{> 0}}$$

"Progress"

⇒ A better metric than using plain cross classification / count

Entropy - Most popular metric for computing the attribute to split on in decision trees.