

### PROBLEM 3.

Given an arbitrary decision tree, it might have repeated queries splits (feature  $f$ , threshold  $t$ ) on some paths from root to leaf. Prove that there exists an equivalent binary tree only with distinct splits on each path.

Ans. let  $a_1, \dots, a_n$  denote the attributes characterizing the data under consideration

let  $D_1, \dots, D_n$  denote the corresponding domains (i.e.  $D_i$  represents the set of values for attribute  $a_i$ )

let  $c_1, \dots, c_m$  represent the decision classes associated with dataset.

Now, two sets of rules could result in the equivalent decision tree.

This is done assigning a

There are decisions which solely depends on one attribute; equivalent tree could be designed by assigning a 'don't care' value to rest of the attributes.

For example,

we have attributes  $a_1, a_2$  and  $a_3$  with the same domain containing  $v_1, v_2$  and  $v_3$  as possible values.

Assuming that the following rules correspond to class  $C_1$ :

rule 1:  $C_1 \leftarrow a_1 = v_1 \wedge a_2 = v_2$

rule 2:  $C_1 \leftarrow a_1 = v_3$

Equivalent representations of these two rules will be following

rule 1:  $C_1 \leftarrow a_1 = v_1 \wedge a_2 = v_2 \wedge a_3 = \text{don't care}$

rule 2:  $C_1 \leftarrow a_1 = v_3 \wedge a_2 = \text{don't care} \wedge a_3 = \text{don't care}$

Considering real life example,

$C_1 = \text{Movie}$        $a_1 = \text{friends visiting}$  :  $v_1 = \text{yes}$

$c_2 = \text{Tennis}$        $v_2 = \text{no}$ .

$c_3 = \text{Shop}$

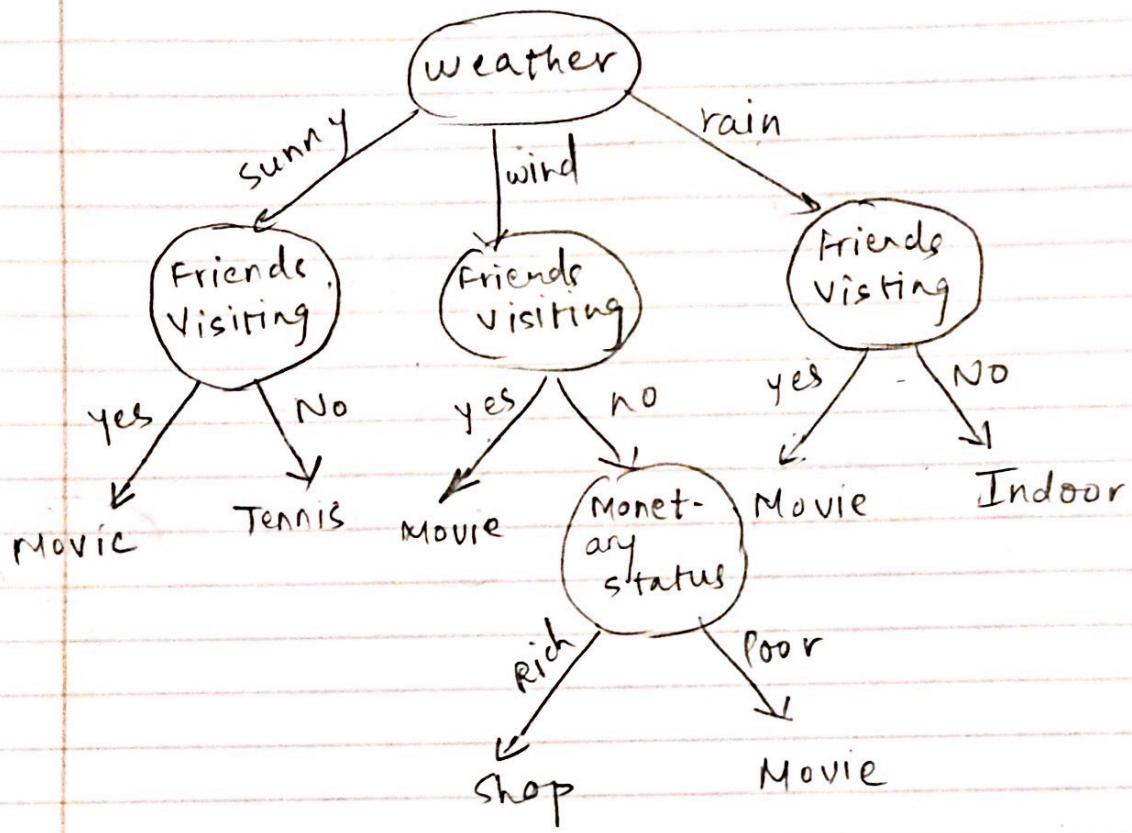
$c_4 = \text{Indoors}$        $a_2 = \text{weather}$  :  $v_1 = \text{sunny}$

$v_2 = \text{wind}$

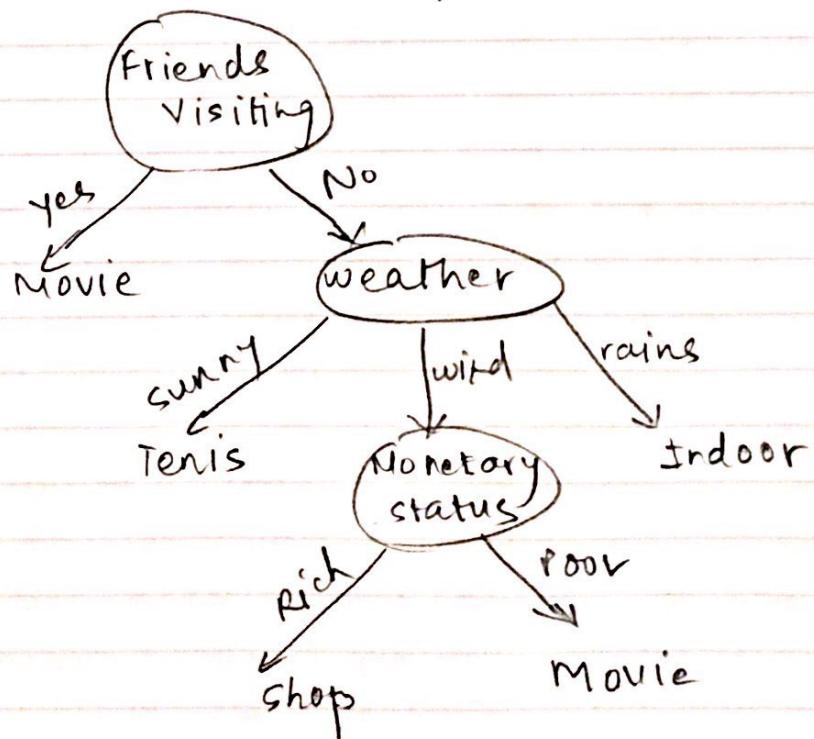
$v_3 = \text{rain}$

$a_3 = \text{Monetary status}$  :  $v_1 = \text{rich}$

$v_2 = \text{poor}$ .



This tree would have an equivalent tree where the attribute won't be repeated as follows.

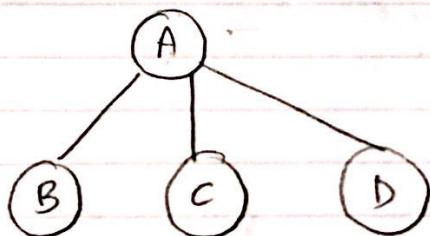


### PROBLEM 3

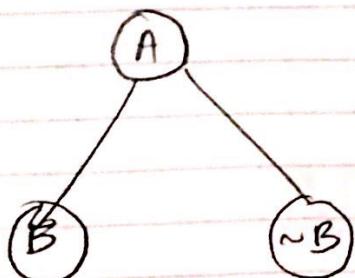
- (a) Prove that for any arbitrary tree, with possible unequal branching ratios throughout, there exists a binary tree that implements the same classification functionality.

Consider a Node A,

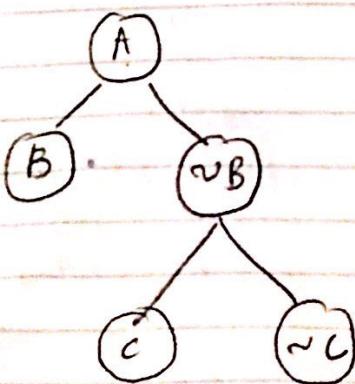
Now, there is a choice of traversing to nodes B, C, & D from A based on the requirement



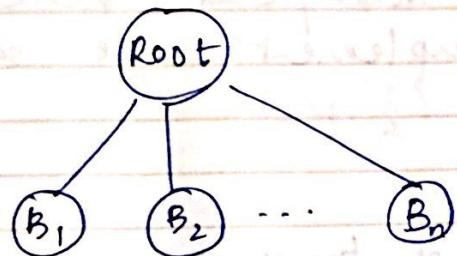
since they are mutually exclusive, this tree could be split into binary decision i.e B or not B



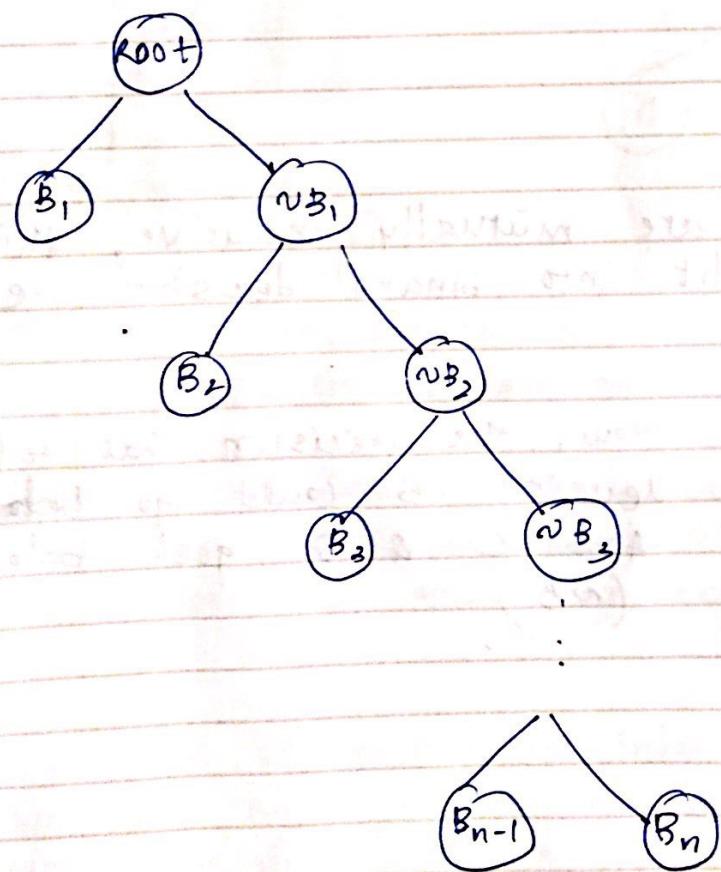
Now, the decision that follows the decision B could go below it, And C & D goes below not B ( $\sim B$ )



Thus, considering a root node and  $B_n$  leaf nodes in a level-2 decision tree;



there will be one extra node for each of the node starting from  $B_3$  till  $B_n$ .



(b) Consider a tree with just two levels, - a root node connected to  $B$  leaf nodes ( $B \geq 2$ ) . What are then upper and lower limits on the number of levels of functionally equivalent binary tree , as a function of  $B$  ?

Ans. let levels of existing  $k$ -ary tree =  $2 = L$   
Number of branches =  $B$

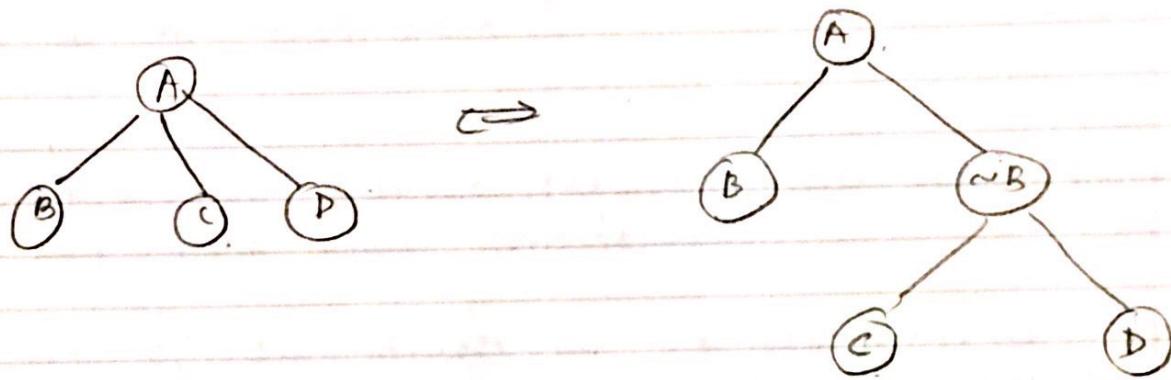
$$\text{Best case (lower limit)} = B + L - 1 = B + 2 - 1 = B + 1$$

$$\text{Upperlimit} = (L-1) * B + 1$$

$$\text{Worstcase} = (2-1) * B + 1 = B + 1$$

Here, since the level is 2 , the Upperlimit & lower limit for the equivalent binary tree is same .

(c) As in (b), what are the upper and lower limits on number of nodes in functionally equivalent binary tree?



lower limit : 1 extra node for every node for  $(B-2)$  such Nodes.

If  $k$ -ary tree has a root node connected to  $B$  leaf nodes ( $B \geq 2$ ) then,

$$\begin{aligned} \text{number of nodes in binary tree} &= (B-2) + B \\ &= 2B - 2 + \text{Root Node} \end{aligned}$$

Upper limit : When each node has an extra ~Node corresponding to it.

Thus, if  $k$ -ary tree has a root node connected to  $B$  leaf nodes ( $B \geq 2$ ) then,

$$\text{number of nodes in binary tree} = 2B + \text{Root node}$$

#### PROBLEM 4.

(a) Prove that the decrease in entropy by a split on a binary yes/no feature can never be greater than 1 bit.

Let number of items =  $N$

Let these items fall into 2 categories i.e (Binary yes/no feature).

Number of items that has Label 1 =  $n$

Number of items that has Label 2 =  $N - n = m$

To get our data more ordered, we group them by labels.

Calculate the ratio,

$$p = \frac{n}{N} \quad \text{and} \quad q = \frac{m}{N} = 1 - p$$

The entropy of our set is given by the following equation:

$$E = -p \log_2(p) - q \log_2(q)$$

Now, consider a split when there is no item with label 1 in the set i.e  $p=0$   
OR set is full of items with Label 1 i.e  $p=1$ .

$$E = -p \log_2(p) - q \log_2(q)$$

$$= -p \log_2(p) - (1-p)[\log_2(1-p)]$$

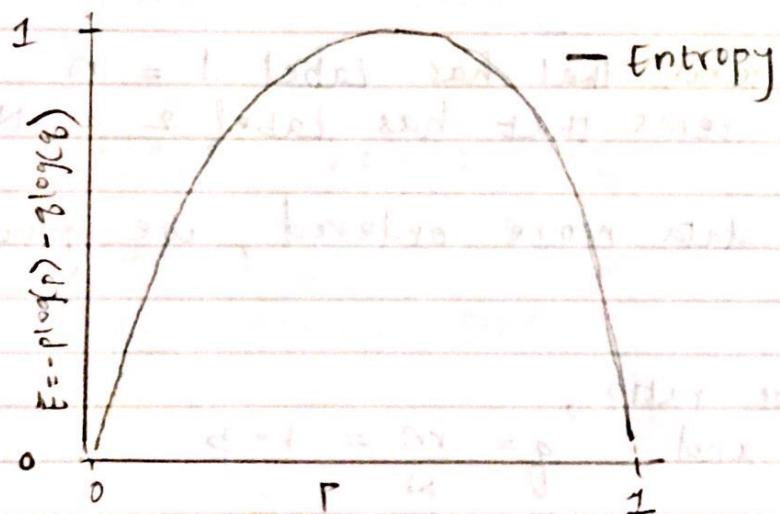
$$E = 0$$

Now, consider a split with set that has half Label 1 and half Label 2 i.e.  $p = \frac{1}{2}$

$$E = -p \log(p) - (1-p) \lceil \log_2(1-p) \rceil$$

$$= -\frac{1}{2} \log\left(\frac{1}{2}\right) - \left(\frac{1}{2}\right) \lceil \log_2\left(\frac{1}{2}\right) \rceil$$

$$= 1.$$



The figure represents the Entropy for a binary-state variable with probabilities  $p$  &  $q$  where  $p+q=1$ .

Thus entropy can only vary between 0 and 1 as probability varies between 0 and 1. Hence, decrease in entropy by split can never be greater than 1.

From the perspective of information theory, it specifies the minimum number of bits of information needed to encode the classification of an arbitrary member of set  $S$ . (i.e. an item drawn from set at random with uniform probability)

for example, if  $p=1$ , the receiver knows that the drawn sample will be positive, so no message needs to be sent and the entropy is zero.

Whereas, if  $p=0.5$ , 1 bit is required to indicate whether the sample is +ve or negative.

- (b) Generalize this result to the case of arbitrary branching  $B > 1$

Ans. Now, consider a case where the target can take up  $b$  possible (different) values. Then entropy,

$$E = \sum_{i=1}^b -p_i \log_2 p_i$$

where  $p_i$  is the probability of the  $i^{th}$  class.

Now, if the target attribute can take up  $b$  possible values, the Entropy could be as large as

$$\log_2 b$$

## PROBLEM 5.

Derive explicit formulas for normal equations solution presented in class for the case of one input dimension.

Ans. The data is  $(x_i, y_i)$  where  $i = 1, 2, \dots, m$  and we are looking at  $h(x) = ax + b$  that realizes the minimum mean square error.

The equation  $h(x_i) = ax_i + b$  can be put into matrix form as :

$$\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

So the normal equations are :

$$\cancel{\begin{bmatrix} x^T & x \end{bmatrix} \vec{\theta}} \quad x^T x \vec{\theta} = x^T y$$

This could be further simplified as :

$$\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_m \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ -x_1 & -x_2 & \dots & -x_m \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

which could be simplified as :

$$\underbrace{\begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}}_{x^T x} \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{\theta} = \underbrace{\begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}}_{x^T y}$$

Now solving the  $2 \times 2$  linear system,

$$\begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

$$= \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{bmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{bmatrix} \begin{bmatrix} \sum x_i y_i \\ \sum x_i y_i \end{bmatrix}$$

$$= \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{bmatrix} (\sum x_i^2)(\sum y_i) - (\sum x_i)(\sum x_i y_i) \\ -(\sum x_i)(\sum y_i) + n(\sum x_i y_i) \end{bmatrix}$$

$$= \frac{1}{\frac{1}{n} \sum x_i^2 - \left( \frac{1}{n} \sum x_i \right)^2} \begin{bmatrix} \left( \frac{1}{n} \sum x_i^2 \right) \left( \frac{1}{n} \sum y_i \right) - \left( \frac{1}{n} \sum x_i \right) \left( \frac{1}{n} \sum x_i y_i \right) \\ - \left( \frac{1}{n} \sum x_i \right) \left( \frac{1}{n} \sum y_i \right) + \left( \frac{1}{n} \sum x_i y_i \right) \end{bmatrix}$$

### PROBLEM 7

The convex hull of a set of vectors  $x_i$ , where  $i=1, 2, \dots, n$  is the set of all vectors of the form

$$x = \sum_{i=1}^n x_i x_i \quad \longrightarrow \quad (1)$$

where the coefficients  $x_i$  are non-negative and sum to one. Given two sets of vectors, show that either they are linearly separable or their convex hulls intersect.

Ans. Given:  $x_i \geq 0$ ;  $\sum_{i=1}^n x_i = 1$

Now, consider 2<sup>nd</sup> set of points corresponding to 2<sup>nd</sup> hull.  $z_j$ , where  $j=1, 2, \dots, m$

The two set of points will be linearly separable if there exists a vector  $\hat{w}$  and a scalar  $w_0$  such that

$$\hat{w}^T x_i + w_0 > 0 \text{ for all } x_i, \text{ and}$$

$$\hat{w}^T z_j + w_0 < 0 \text{ for all } z_j.$$

Now, calculate the linear discriminant for the points belonging to the two convex hulls.

Linear discriminant for points in set I,

$$y(x) = \hat{w}^T x^* + w_0 \quad \longrightarrow \quad (2)$$

substituting (1) in (2), we get :

$$y(x) = \hat{w}^T \left( \sum_{i=1}^n \alpha_i x_i \right) + w_0 \quad \text{--- (3)}$$

Bringing the  $\alpha$  outside,

$$\begin{aligned} y(x) &= \sum_{i=1}^n \alpha_i (\hat{w}^T x_i) + w_0 \\ &= \sum_{i=1}^n \alpha_i (\hat{w}^T x_i + w_0) \quad \text{--- (4)} \end{aligned}$$

Given :  $\sum_{i=1}^n \alpha_i = 1$ .

Similarly, we find linear discriminant for points belonging to 2<sup>nd</sup> set

$$y(z) = \sum_{i=1}^m \beta_i (\hat{w}^T z_i + w_0) \quad \text{--- (5)}$$

Given :  $\beta_i \geq 0$  &  $\sum_{i=1}^m \beta_i = 1$ .

consider the scenario where convex hull intersects.

This means there must be atleast one point in common between  $\{x\}$  &  $\{z\}$ .

let the common point. = a

since a belong to both convex hulls, there must be set of  $\alpha$  &  $\beta$  that gives rise to a.

linear discriminant for a from equation (4) & (5),

$$\begin{aligned}y(a) &= \sum_{i=1}^n \alpha_{+i} (\hat{w}^T x_i + w_0) \\&= \sum_{i=1}^m \beta_i (\hat{w}^T z_i + w_0) \quad (6)\end{aligned}$$

For linear separability, we must have

$$\begin{aligned}y(x) &= \hat{w}^T x + w_0 > 0 \quad (4) \\y(z) &= \hat{w}^T z + w_0 < 0 \quad (7)\end{aligned}$$

From the non-negativity and simplex constraints on  $\alpha$  &  $\beta$  and equations (6) & (7); we have a contradiction.

The linear discriminant  $y(a)$  has to be simultaneously greater than and less than zero which is impossible.

Now, consider the scenario where patterns are linearly separable,

$$\begin{aligned}y(x) &= \hat{w}^T x + w_0 > 0 \\y(z) &= \hat{w}^T z + w_0 < 0 \quad (8)\end{aligned}$$

Assuming that there is a point  $a$  lying in the intersection of the convex hulls.

From equation (6), we have

$$y(a) = \sum_{i=1}^n \alpha_i (\hat{w}^T x_i + w_0) = \sum_{i=1}^m \beta_i (\hat{w}^T z_i + w_0)$$

Equation is not possible since from equation (8), we have the fact that contradicts the above equation.