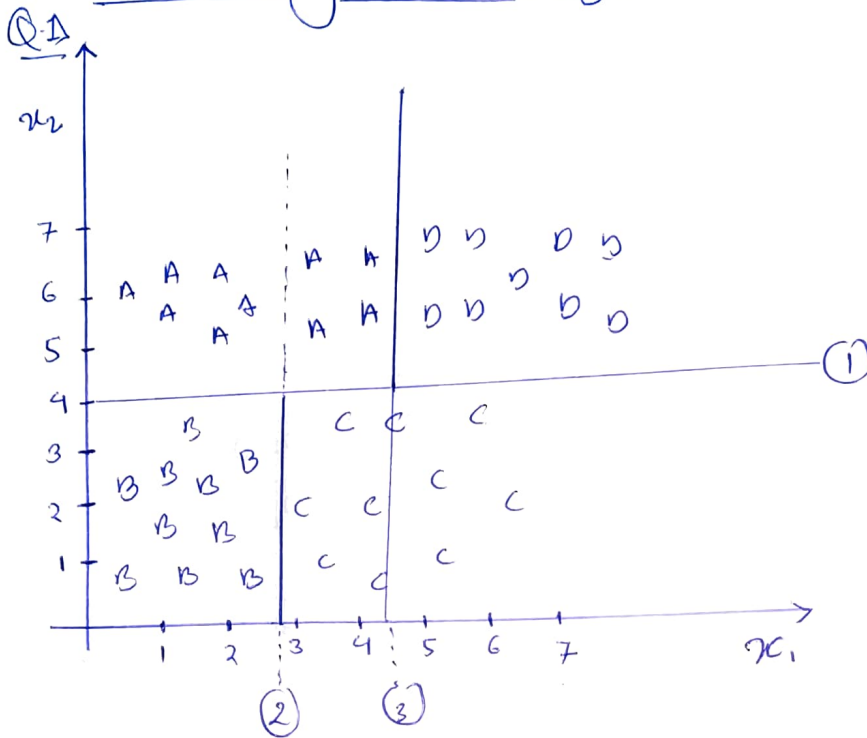


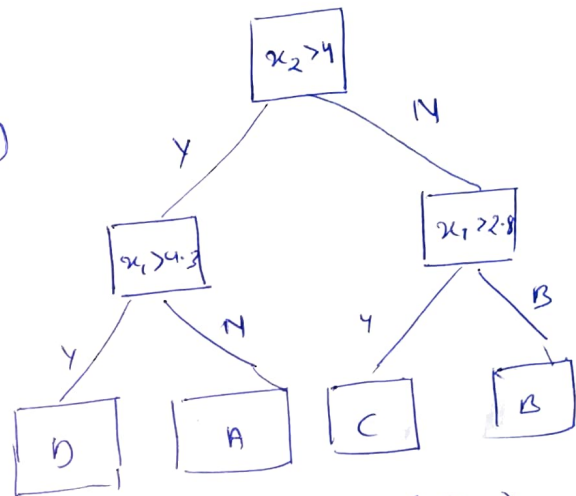
# SMA1- Quiz (on Decision Trees)

Name: Aditya Kumar Singh

ID: 2021701010

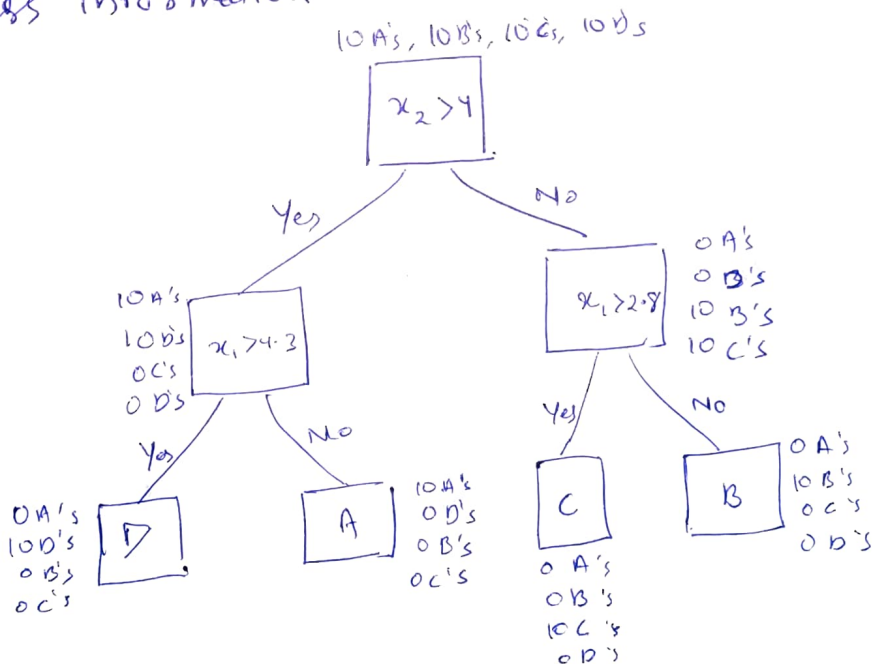


Given decision Tree structure :



Calculate Information Gain,  $\Delta i$ , for each split (i.e., 1, 2 & 3).

Ans let's first denote the no. of examples at each node with all the class information.



So, the Entropy at root node (or the node with  $x_2 > 4$  decision) =  $-\sum_{i=1}^4 p_i \log(p_i)$

$p_i = i^{th}$  class probability

$H(\text{root}) = H(x_2 > 4) := \text{Entropy for decision } x_2 > 4$

$$= - \sum_{i=1}^4 P_i \log P_i$$

$$= - \sum_{i=1}^4 \frac{1}{4} \log_2 \frac{1}{4} \quad (\text{as } P_A = P_B = P_C = P_D = \frac{10}{40} = \frac{1}{4})$$

$$= - 4 \cdot \frac{1}{4} (\log_2 1 - \log_2 2^2) = 2 \quad \text{--- (1)}$$

Similarly:

$$H(x_1 > 4.3) = H(\textcircled{3}) = - \sum_{i=1}^4 P_i \log P_i$$

$$= - \left[ \frac{10}{20} \log_2 \frac{10}{20} + \frac{10}{20} \log_2 \frac{10}{20} \right]$$

$$= 1 \quad \text{--- (2)} \quad (\text{where } P_A = P_D = \frac{10}{20+0+0} = \frac{1}{2})$$

$$H(x_1 > 2.8) = H(\textcircled{2})$$

$$= - \left[ \cancel{P_A \log P_A} + \cancel{P_B \log P_B} + P_C \log P_C + P_D \log P_D \right]$$

$$= - \left[ 0 + 0 + \frac{10}{20+0} \log_2 \frac{10}{20+0} + \right.$$

$$\left. \frac{10}{20+0} \log_2 \frac{10}{20+0} \right] = - 2 \cdot \frac{1}{2} \log_2 \left( \frac{1}{2} \right) = 1 \quad \text{--- (3)}$$

$$H(D) = - \left[ \cancel{P_A \log P_A} + \cancel{P_B \log P_B} + \cancel{P_C \log P_C} + P_D \log P_D \right]$$

$$= - \left[ 0 + 0 + 0 + \frac{10}{10+0} \log_2 \frac{10}{10+0} \right] = 0 \quad \text{--- (4)}$$

Similarly:  $H(A) = H(C) = H(B) = 0$  (as all they contain only data-points of single class i.e., all are pure nodes)  $\rightarrow$  Hence no mixing of classes  $\Rightarrow$  entropy (confusion) is 0.

Let's calculate information gain  $\Delta_i$  at each node:

$$\Delta(x_2 > 4) = \mathcal{IG}(x_2 > 4) = H(x_2 > 4) - \left[ \overset{x_2 > 2.8}{P(D_2)} H(x_1 > 2.8) + \overset{x_2 > 4.3}{P(D_3)} H(x_1 > 4.3) \right]$$

$$= 2 - \left[ \frac{20}{40} \cdot 1 + \frac{20}{40} \cdot 1 \right] \quad (\text{from (1), (2) \& (3)})$$

8

$$\Delta(x_2 > 4) = 2 - 1 = 1 \quad \left( \begin{array}{l} \text{where } \frac{20}{40} = \text{probability of datasets} \\ \text{at each child nodes of} \\ \text{root node.} \end{array} \right)$$

$\downarrow$  Before split      After split      (5)

Similarly

$$\checkmark \Delta(8) = \Delta(x_1 > 4.3) = \text{Entropy before split at } (D_3 = \text{Decision 3})$$

$$= \Delta(x_1 > 4.3) - \text{Entropy after split at } (D_2 = \text{Decision 2} = x_1)$$

Both D & A node)

$$= 1 - \left[ \underbrace{\frac{10}{20} \cdot 0}_{P(\text{datasets at D})} + \underbrace{\frac{10}{20} \cdot 0}_{P(\text{datasets at A})} \right] = 1 \quad \text{--- (6)}$$

$\downarrow$   $H(D)$        $\downarrow$   $H(A)$

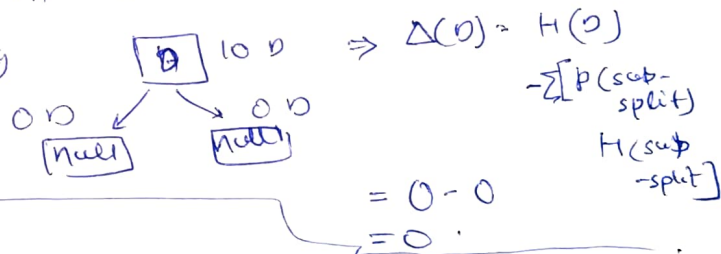
$$\checkmark \Delta(2) = \Delta(x_1 > 2.8) = \underbrace{H(2)}_{\text{Entropy before split}} - \underbrace{\left[ \frac{10}{20} \cdot H(C) + \frac{10}{20} \cdot H(B) \right]}_{\text{Entropy after split}}$$

$$= 1 - \left[ 0 \cdot \frac{10}{20} + 0 \cdot \frac{10}{20} \right] = 1 \quad \text{--- (7)}$$

$$\checkmark \Delta(D) = \Delta(A) = \Delta(C) = \Delta(B) = 0 \quad \text{--- (8)}$$

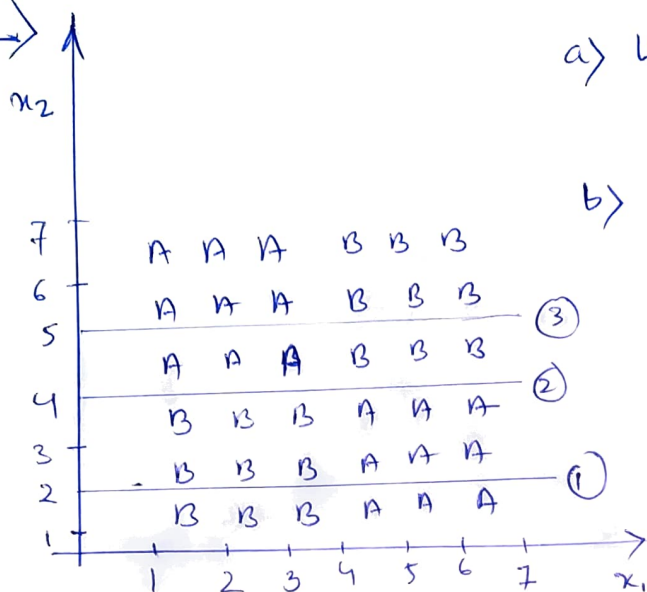
$\because H(D) = H(A) = H(C) = H(B) = 0$  & all they have are "null"

as child nodes were anyway  $p_A = p_B = p_C = p_D = 0$  as no split is happening at these leaf nodes. e.g.



Nodes	H = Entropy	$\Delta = IG = \text{Information Gain}$
① = $x_2 > 4$	2	1
③ = $x_1 > 4.3$	1	1
② = $x_1 > 2.8$	1	1
D	0	0
A	0	0
C	0	0
B	0	0

Q-2 →



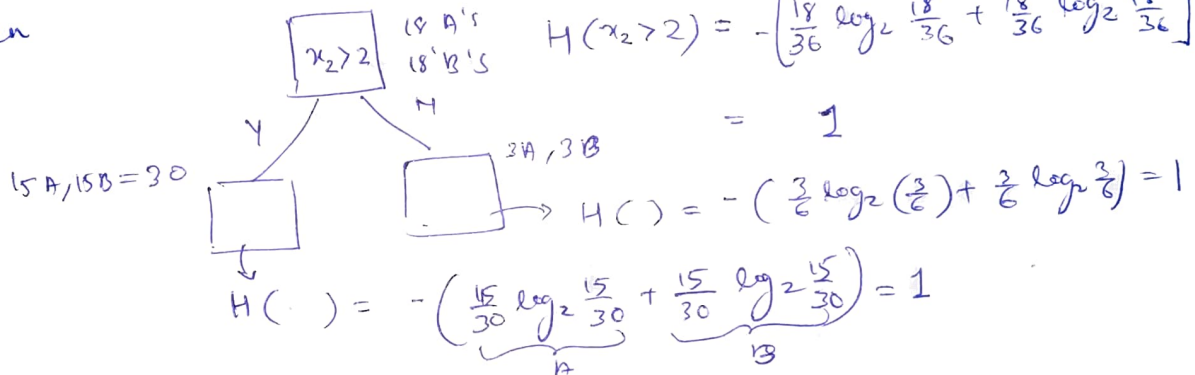
a) What is the problem with this data.

b) How can we overcome this?

Ans

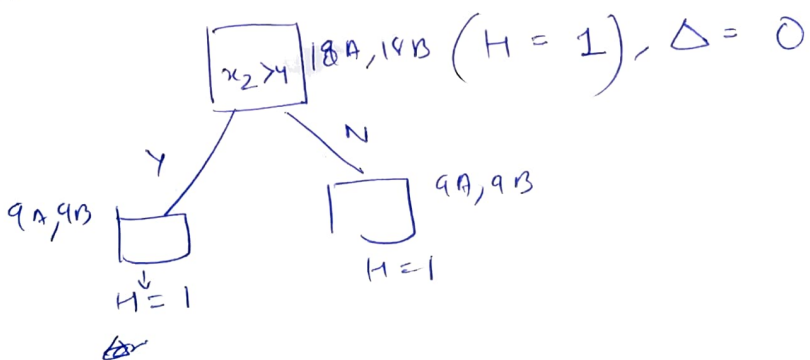
① Say we start our root node with  $x_2 > 2$  (i.e. classifier ①)

then



$$\Delta(x_2) \Rightarrow \Delta(x_2 > 2) = 1 - \left( \frac{30}{36} \cdot 1 + \frac{6}{36} \cdot 1 \right) = 0.$$

② Similarly, for if we have taken decision ( $x_2 > 4$ ), then



& so on.

So what we observed, that for "any" horizontal (not supposedly passing through datapoints rather in betn them) split, we're getting "even-split" i.e., each class have equal number of candidates at child nodes.



→ And since it's two class-classification, this even-split results in '0' Entropy (i.e., for any node/Decision having data with containing equal no. of labels from each class (ofcourse, <sup>for</sup> binary/two class) ~~file~~ → results in 'zero' entropy) → e.g. at child node of  $x_2 > 4$ ,  $q_1, q_2 \Rightarrow$  even splits of data at this node into two labels  $\Rightarrow H() = 0$  & so on.

→ Not only horizontal, for all sorts of vertical (parallel to  $x_2$ ) decision boundaries would produce the same result. This is due to the fact that:

1) We're constrained to choose decision boundary parallel to axes.

2) Dataset are also well-aligned with axes.

3) Further they are "diagonally" fused/mixed up → which ensures, results will be same for both vertical

Δ horizontal decision boundaries to start with.

→ Hence, the root node decision splits the data into two halves <sup>no matter what decision boundary (parallel to axes)</sup> for which the "proportion" of classes in both remain the same → which accounts for "Entropy after split" = 1

→ Now since, we have equal no. of both class to begin with, the entropy of root node, will always be = 1.

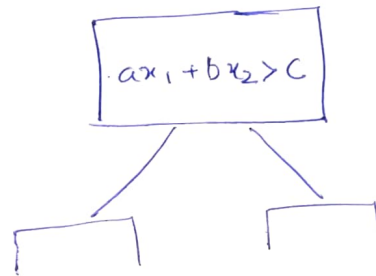
⇒ Information Gain = 0 (always for any vertical/horizontal decision) → further splitting/construction of trees would STOP. [∵ info-gain varies bet  $[0, 1]$ , with 0 nothing can be done].

How to overcome this ? (b).

★ Now the whole problem lies with the fact that both our data ~~are~~ & decision boundaries to be constructed are axes-aligned.

★ If somehow, we start with a "tilted" decision boundaries

at root node (i.e. combination of function of combination of both features) then the ~~case~~ chance of getting equal "proportion" of labels in both the split is zero. In other words, 'no' even split can take place which implies the  $IG =$  ~~Information~~ Entropy won't be 1  $\Rightarrow IG \neq 0$  (won't be 0).



$$\begin{aligned} a &\neq 0 \\ b &\neq 0 \\ c &\in \mathbb{R} \end{aligned}$$

$\Rightarrow$  Now with above structure & with possible optimisation of  $a, b, c$ , we can obtain such a decision boundary that outputs ~~max~~ min entropy (= less fusion among classes = less confusion) which is indirectly max' information gain.

$\Rightarrow$  Once that set, we can start of our journey for constructing

DT further from root node.

(e.g. C4.5 DT algo involves such possibilities).