

Solution

XOR can be implemented by making use of AND and OR.

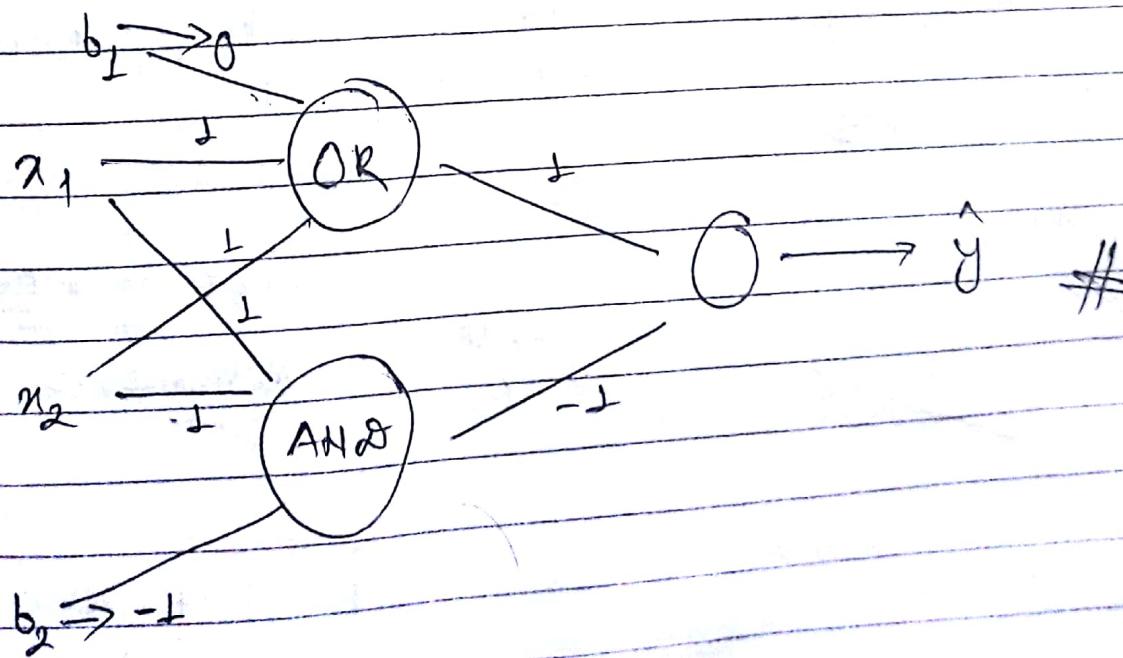
x_1	x_2	OR	AND	XOR
0	0	0	0	0
0	1	1	0	1
1	0	1	0	1
1	1	1	1	0

XOR is the same as OR
except for the case when

$$x_1 = x_2 = 1$$

\Rightarrow subtracting AND from OR
gives XOR.

So, our network becomes



b. Solution

Gaussian Naive Bayes classifier gives us

$$P(Y=1|X) = \frac{1}{1 + e^{w_0 + \sum_{i=1}^n w_i x_i}}$$

From Tom Mitchell's Book Chapter 3 eq. 22

where weights $w_1 - w_n$ are given by

$$w_i = \frac{\mu_i D - \mu_D}{\sigma_i^2}$$

$$w_0 = \ln \frac{1-\pi}{\pi} + \sum_i \frac{\mu_D - \mu_i}{2 \sigma_i^2}$$

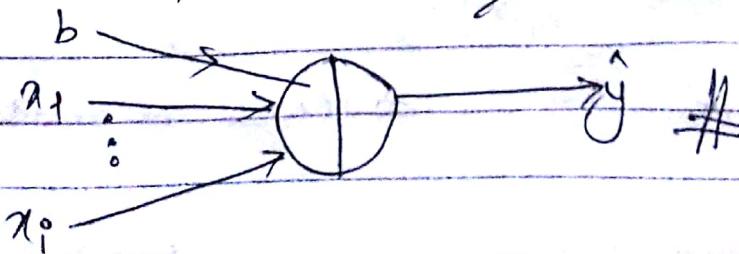
Also, we have the ^{parametric} probabilistic approach of Logistic Regression assumed is

$$P(Y=1|X) = \frac{1}{1 + e^{w_0 + \sum_{i=1}^n w_i x_i}}$$

from eq. 16

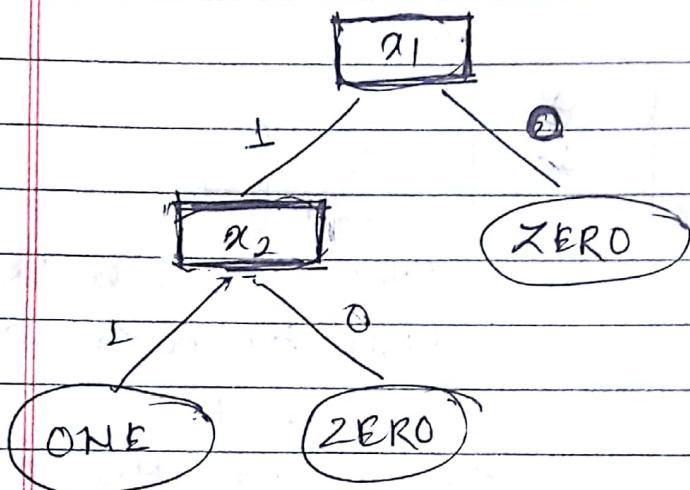
As observed, both the classifiers have an identical form.

Since a logistic regression unit is just a single layer Neural Network with 1 sigmoid unit, this implies that any Naive Bayes classifier with Boolean features can be implemented using a Neural Network.



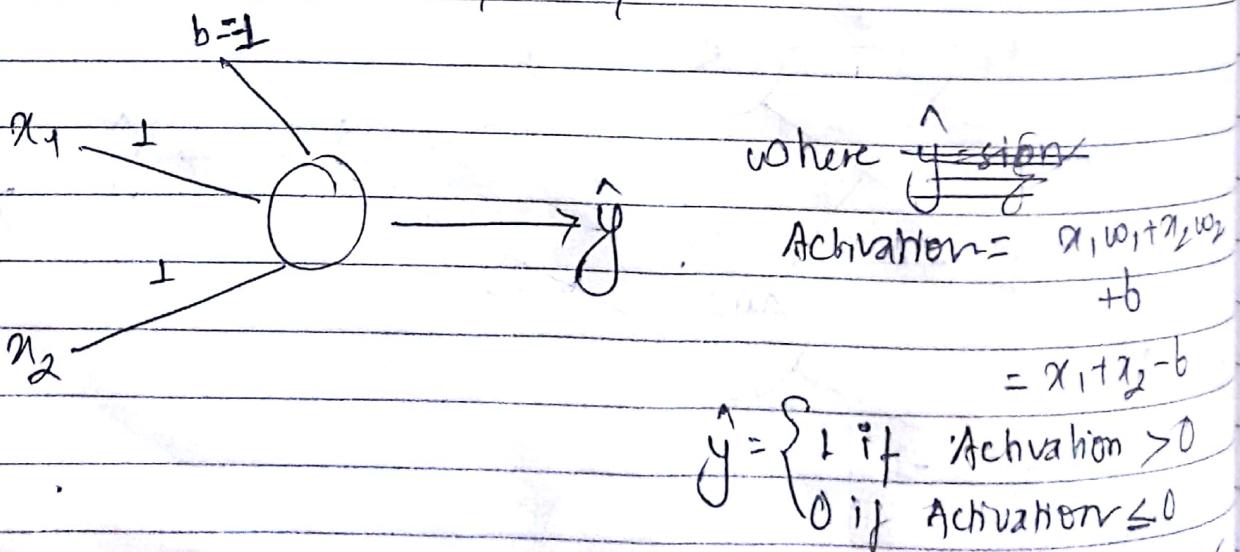
c). Solution:

We know that a decision tree can represent any boolean-valued function. For example, it can represent AND as $x_1 \text{ AND } x_2$.



We also know that a 2-layer NN can represent any Boolean function.

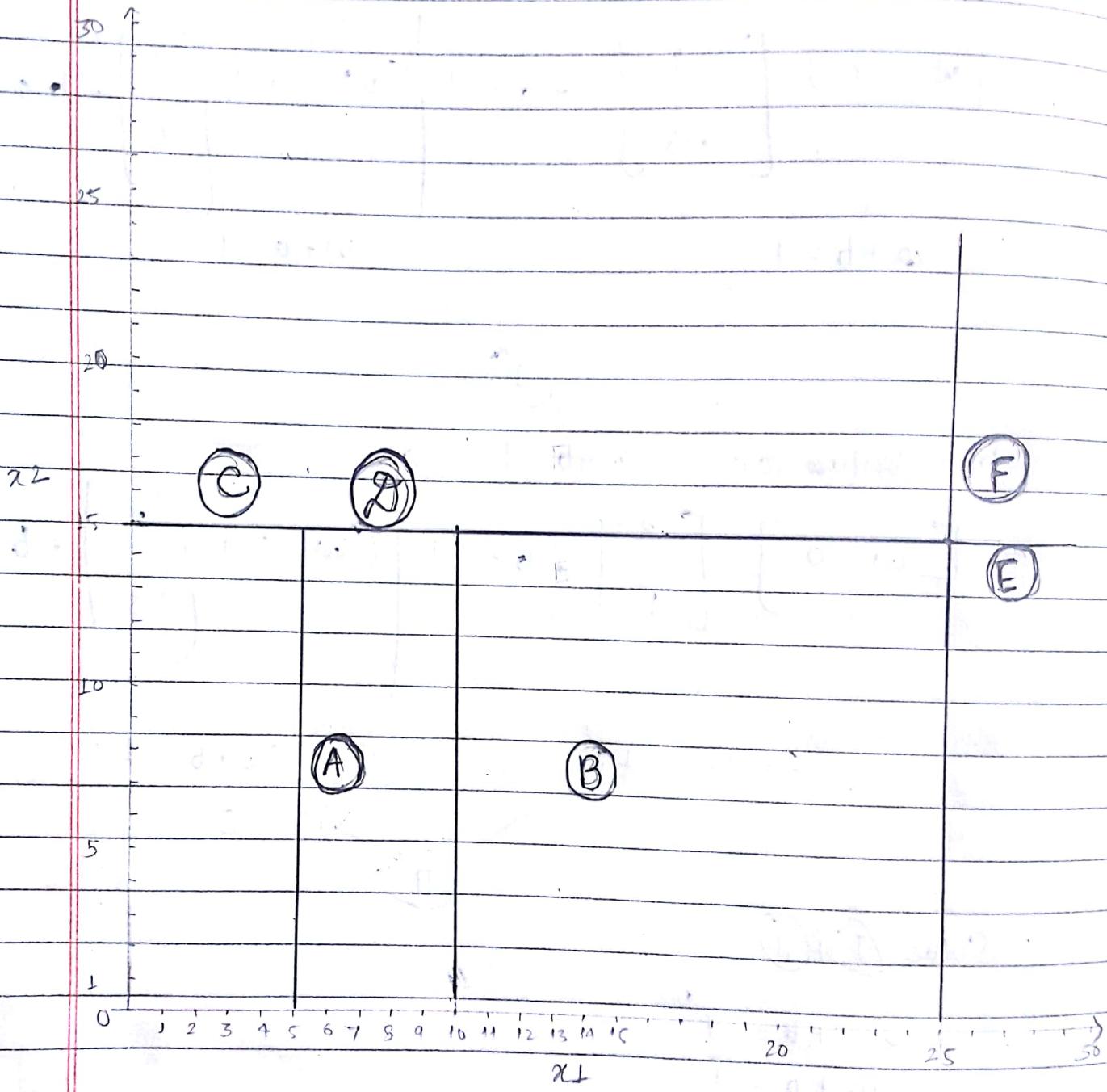
For example, $x_1 \text{ AND } x_2$ can be represented using only 1-layer NN



Thus, a NN can be constructed to implement a decision tree classifier with boolean features.

d.

a) solution



b) Solution

$$x_2 < 15$$

T

F

$$x_2 < 25$$

T

F

T

F

$$x_1 < 10$$

E

$$x_1 > 25$$

A

B



$$x_1 < 5$$

Ans: Another decision tree for the same question.

In my view, this does not pose any problem. The only difference is the length of the decision tree. This may make a difference in computational time but both the trees reach the same final solution. So, the rate of learning may vary but the convergence is identical.

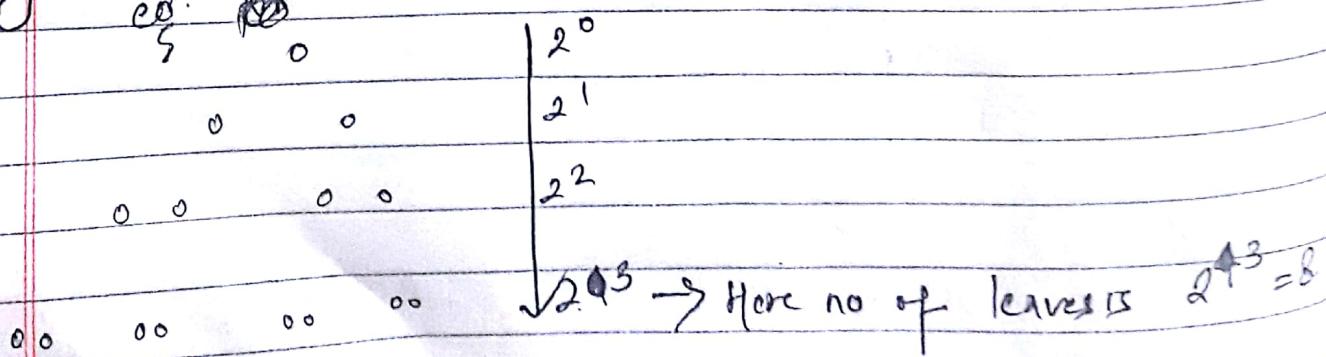
3 & Solution

Since we are choosing the nodes at random without looking at the gain, I think the number of leaf nodes is going to be It is because we never stop choosing a node i.e., we remove the fact that we choose the node with maximum ^{highest} gain so, we don't run out of the features and it just

3a Solution

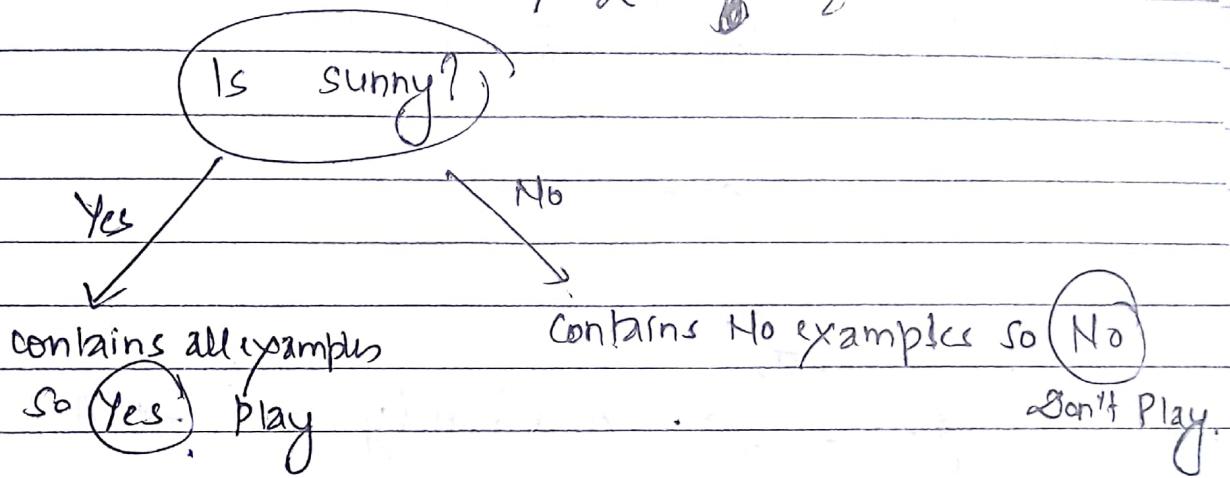
In such a scenario, the maximum number of nodes is possible if only we take the worst feature that reduces the task of classification by the least amount i.e., at each internal node, we choose a feature that gives us the test least value of information gain. This happens when the dataset is divided such that the no of examples in each label is equal and also, the number of examples for each possible value of a feature x_i is ^{same} for every possible value for every feature. In such a case, the classification decreases logarithmically and the [no of leaf nodes maximum possible] given that the ~~splits~~ splits are binary is 2^m .

e.g. ~~10~~



b) Solution

If we use the original mutual information version, we get the best possible solution. So, it is always going to be at least as small as the answer in the first case but never bigger. In the best case, it could contain only 2 leaves e.g.



c) Solution

The accuracy of the decision tree produced does not change because a problem can have multiple optimal decision trees. The ordering of the features that we choose doesn't affect the fact that a node will always classify a label of "Play Tennis" as "Play Tennis" no matter at what point of time it is classified. The fact that we consider all the features no matter the order in the decision tree makes sure we get an accurate tree. The only problem with the random choice is that we get a longer complex decision tree.

4. Solution

We have 3 features.

for all 3, we find out Entropy H.

$$\cancel{H_A} =$$

$$\text{At the root } H(y=1) = -\frac{3}{6} \times \log_2 \frac{3}{6} - \frac{3}{6} \times \log_2 \frac{3}{6}$$

$$= \frac{1}{2} + \frac{1}{2}$$

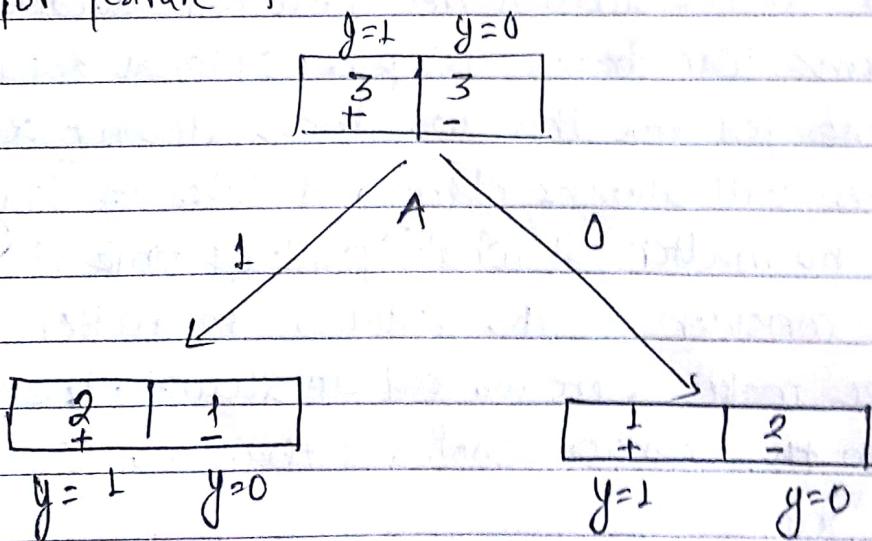
$$= 1.$$

For feature A,

~~$$H(y|A) = p(A=0) H(y|A \neq 0) + p(A=1) H(y|A=1)$$~~

$$= \frac{1}{2} \times$$

For feature A



$$\text{Left Branch } H(y|x=\text{true}) = -\frac{2}{3} \times \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \times \log_2 \left(\frac{1}{3}\right)$$

$$= \frac{0.584 \times 2}{3} + \frac{1.584}{3}$$

$$= 0.9133$$

$$\text{Right Branch } H(y|x=\text{false}) = -\frac{1}{3} \times 10^0 \log_2(1/3) - \frac{2}{3} \times 10^0 \log_2(2/3)$$

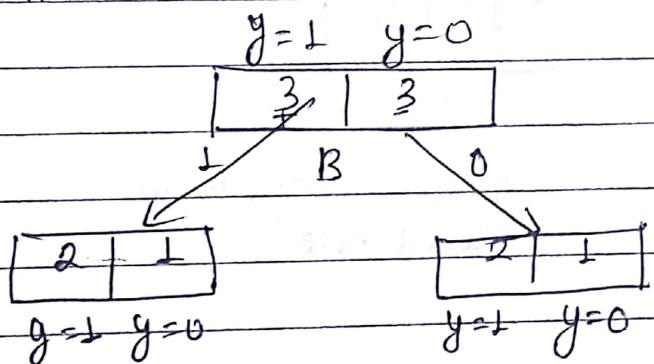
$$= 0.9133$$

\therefore Combined uncertainty is the weighted entropy of all the branches

$$= \frac{3}{6} \times 0.9133 + \frac{3}{6} \times 0.9133$$

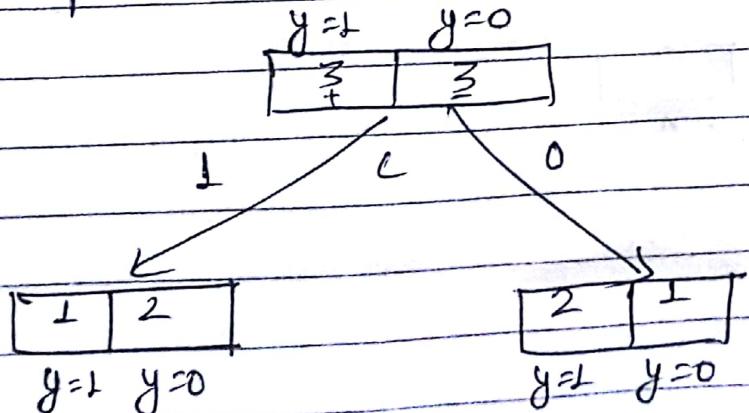
$$= 0.9133. \Rightarrow H(y|A) = 0.9133$$

for feature B



Since it is identical to feature A, $H(y|B) = 0.9133$

for feature C

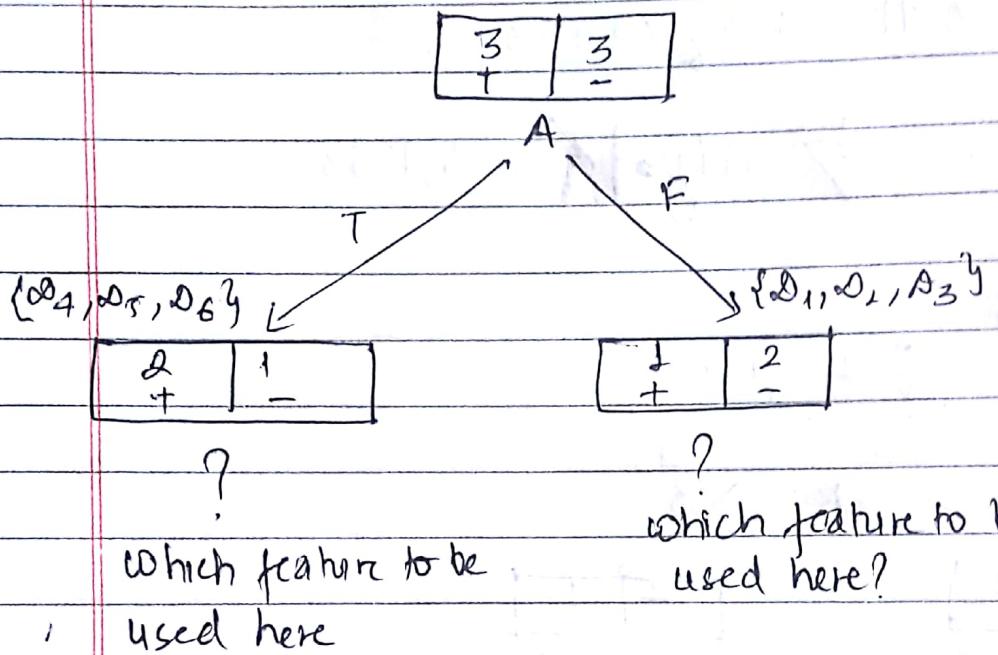


Identical $\Rightarrow H(y|C) = 0.9133$

So, Information gain, $I(x,y) = H(y) - H(y|x)$

is going to be same for all 3 features.

Thus, we can use any of the 3 features in the root node and proceed.

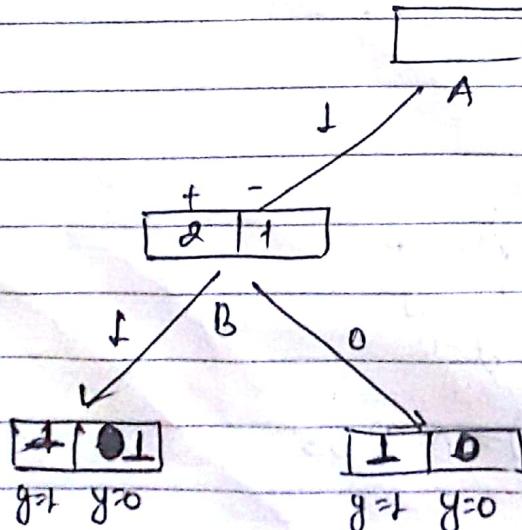


which feature to be used here

which feature to be used here?

$$H(y) \text{ at left branch} = -\frac{2}{3} \times \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \times \log_2 \left(\frac{1}{3}\right) \\ = 0.9133.$$

If we put feature B,



If we put B there

$$\therefore \text{left branch } H(y|B=\text{true})$$

$$= -\frac{1}{2} \times \log_2(1/2) - \frac{1}{2} \times \log_2(1/2)$$

$$= \frac{1}{2} + \frac{1}{2} \\ = 1.$$

$$\text{right branch } H(y|B=\text{false})$$

$$= -\frac{1}{2} \times \log_2(1) - 0$$

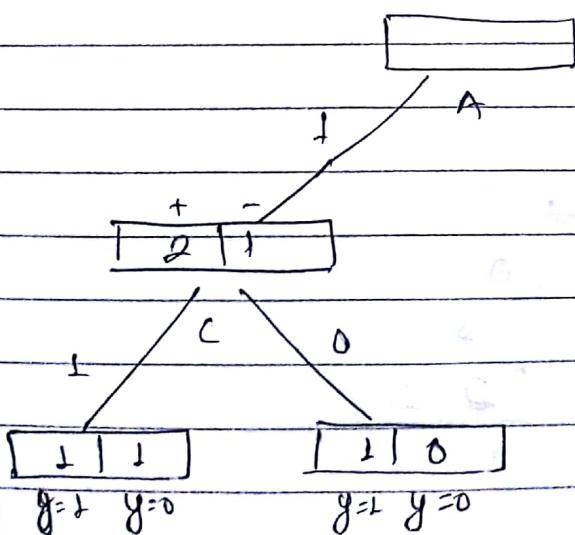
$$= 0$$

\Rightarrow combined uncertainty

$$= \frac{2}{3} \times 1 + 0$$

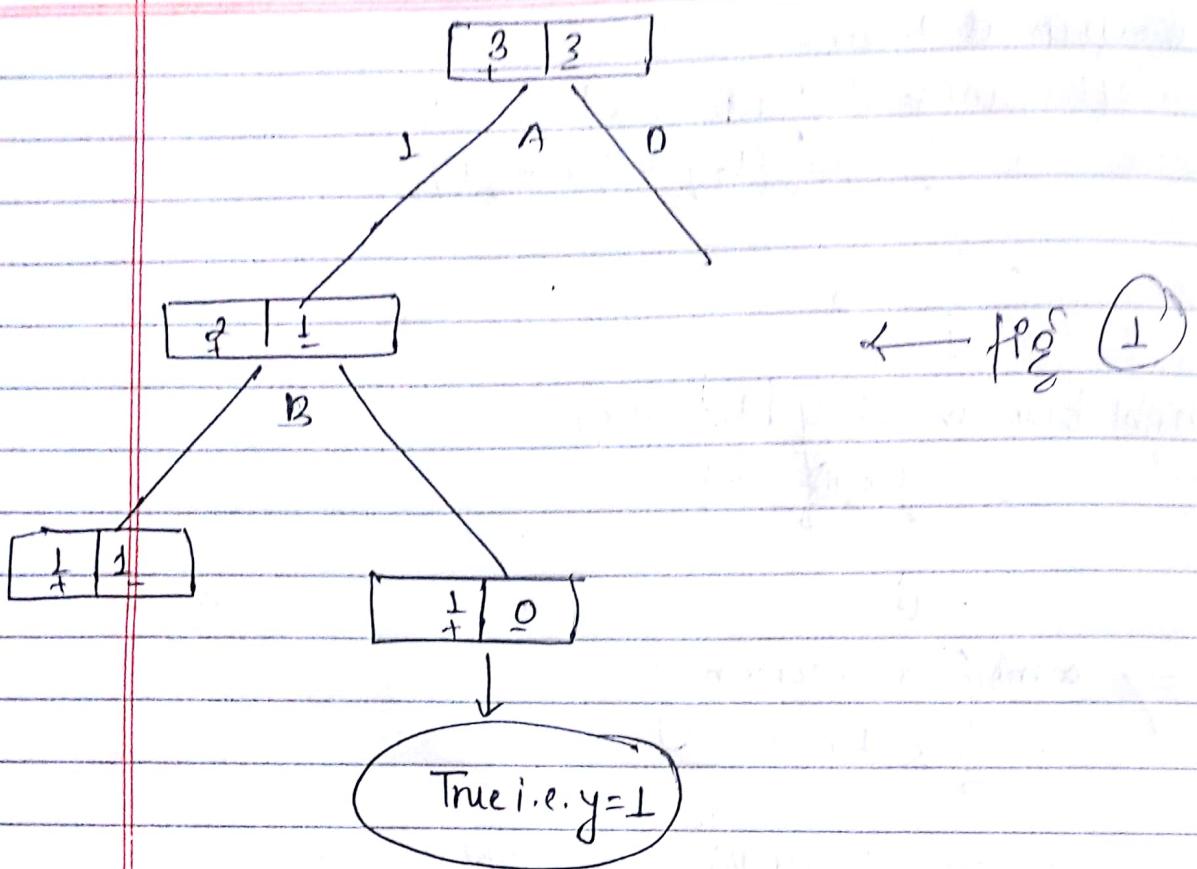
$$= \frac{2}{3} \quad H(y|B) = 2/3 = 0.66..$$

If we put feature C,

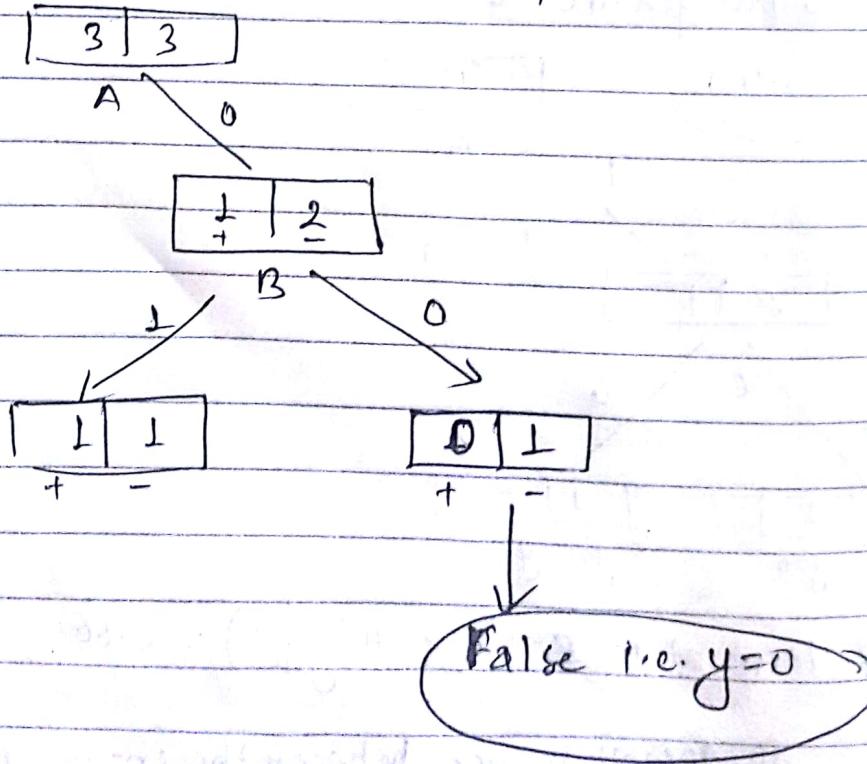


$$\text{Identical to B } \Rightarrow H(y|C) = 0.66.$$

So, doesn't matter between the choice of B & C here.



For the right branch of A, let's say we use B.



~~we used~~

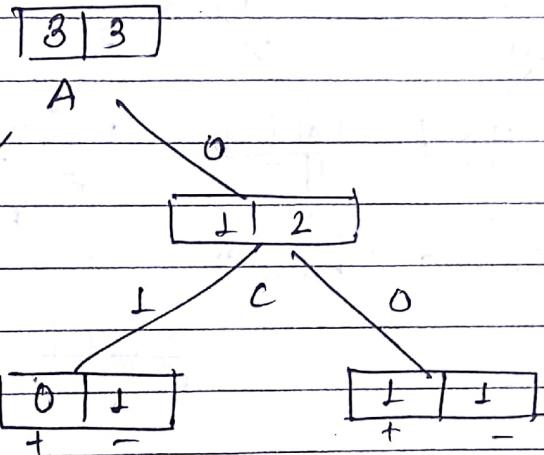
$$H(y|B=\text{true}) = -\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right)$$

$$= 1$$

$$H(y|B=\text{false}) = 0$$

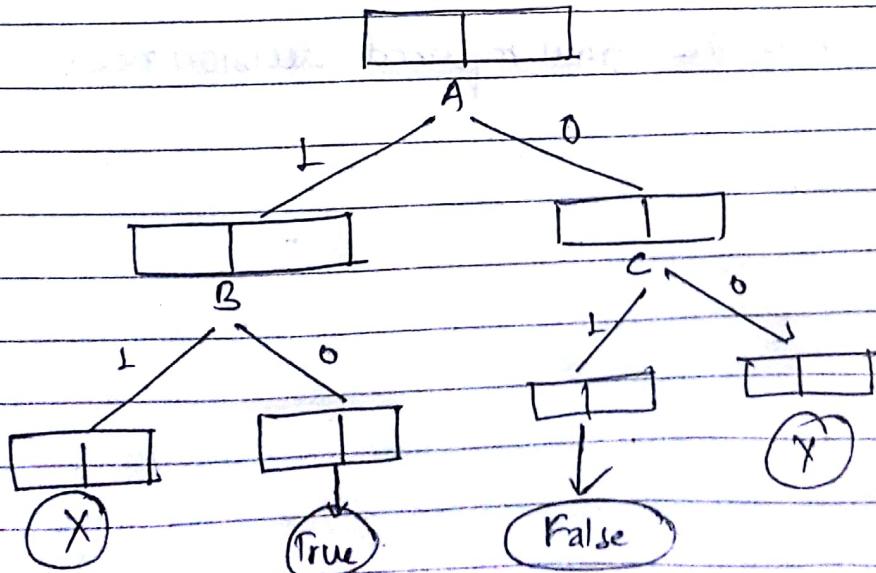
$$\text{combined } H(y|B) = 0.66.$$

If we put C instead,

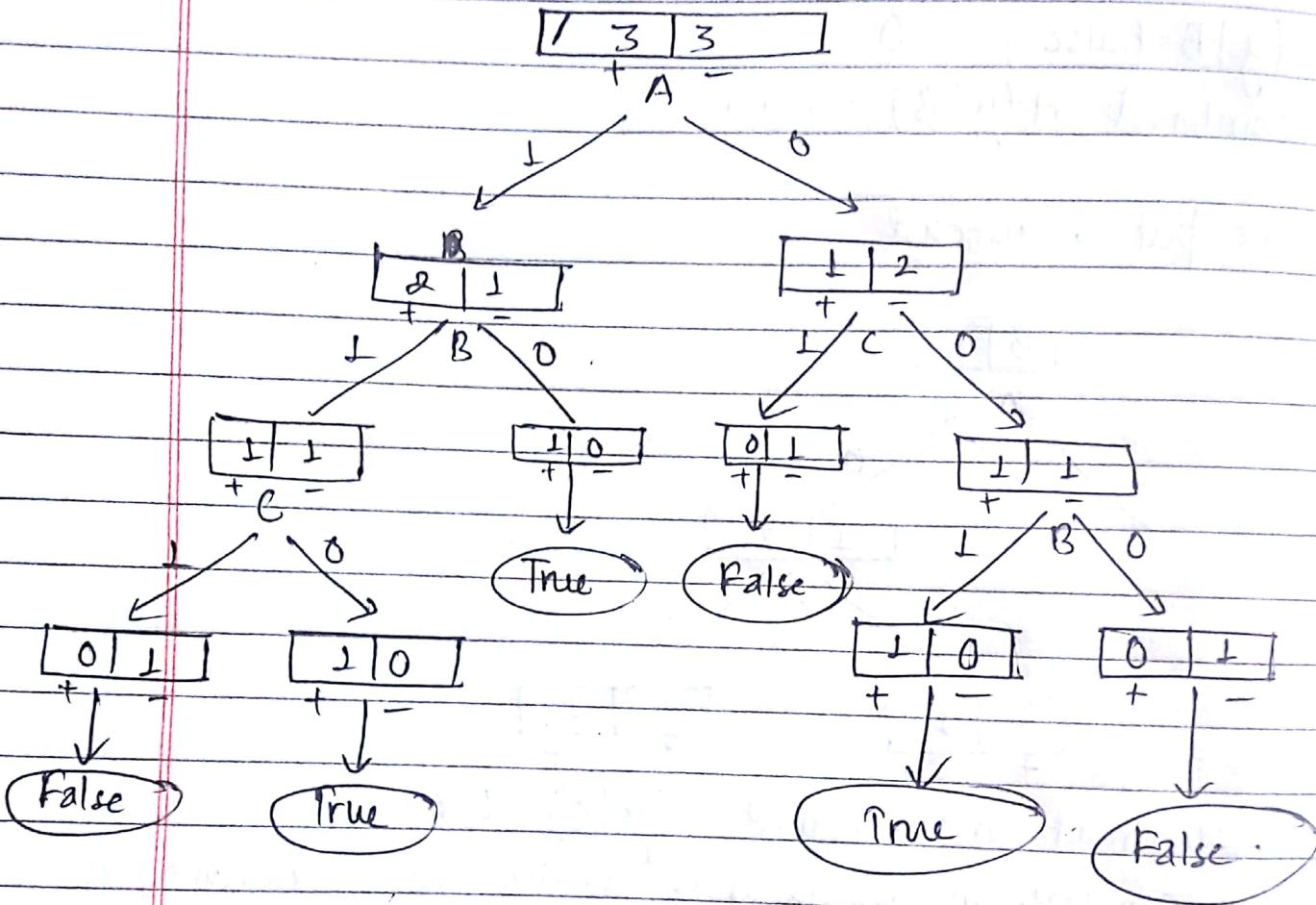


Identical so, combined $H(y|C) = 0.66$.

We can use any of B or C as the right branch of A.



Now at branches X & Y, we can use features C & ~~A~~ B respectively because those are the only features left to use there.



which is the final required decision tree.

Q. 5. Solution

$$\sum_{j=1}^N \omega_{i+1}(j) I(h_i(x_j) \neq y_j) = \sum_{\substack{h_i(x_j) \neq y_j}} \omega_i(j) e^{\alpha_i} \quad (1)$$

Normalizing the sum of $\omega_i(j) e^{\alpha_i}$ to make sure it is a distribution that adds to a probability of 1.

Note: $\omega_i(j)$ gives a probability of seeing this example in the next iteration.

$$\Rightarrow \frac{1}{Z_i} \sum_{\substack{h_i(x_j) \neq y_j}} \omega_i(j) e^{\alpha_i} = p(h_i(x_j) \neq y_j)$$

$$\text{Also, } \frac{1}{Z_i} \sum_{\substack{h_i(x_j) \neq y_j}} \omega_i(j) e^{\alpha_i} = e^{\alpha_i} \cdot \epsilon_i \quad \text{where } \epsilon_i = \text{error}(h_i, S, \omega_i) + \alpha_i = \frac{1}{2} \ln \left(\frac{1-\epsilon_i}{\epsilon_i} \right)$$

Now, eqn (1) becomes

$$\sum_{j=1}^N \omega_{i+1}(j) I(h_i(x_j) \neq y_j) = \frac{\epsilon_i}{Z_i} \sqrt{\frac{1-\epsilon_i}{\epsilon_i}} \quad (2)$$

$$\sum_{j=1}^N \omega_{i+1}(j) I(h_i(x_j) = y_j) = \frac{\sqrt{(1-\epsilon_i)\epsilon_i}}{Z_i} \quad (3)$$

However we know that the sum of probabilities is 1

$$\text{So, } \frac{1}{Z_i} \left[\sum_{\substack{h_i(x_j) \neq y_j}} \omega_i(j) + \sum_{\substack{h_i(x_j) = y_j}} \omega_i(j) \right] = 1$$

But from eqns (2) & (3) the two terms are equal.
Hence,

$$\sum_{j=1}^N \alpha_{i+1}(j) I(h_i(x_j) \neq y_j) = 50\%.$$

Hence, proved. $\#$

Q.6. Solution

From AdaBoost Algorithm, the probability of seeing an example in the next iteration is given by

$$\alpha_{l+1}(i) = \alpha_l(i) \times \begin{cases} e^{-\alpha_l} & \text{if } h_l(x_i) \neq y_i \\ e^{\alpha_l} & \text{if } h_l(x_i) = y_i \end{cases}$$

Combining in a single form gives

$$\alpha_{l+1}(i) = \alpha_l(i) e^{-y_i \alpha_l h_l(x_i)}$$

However, the recursive formulation for α_l can be made as

$$\begin{aligned} \alpha_l(i) &= \alpha_{l-1}(i) e^{-y_i \alpha_{l-1} h_{l-1}(x_i)} \\ &= \alpha_{l-2}(i) e^{-y_i \alpha_{l-2} h_{l-2}(x_i)} \times e^{-y_i \alpha_{l-1} h_{l-1}(x_i)} \\ &= \alpha_l(i) e^{-y_i \alpha_1 h_1(x_i)} \times e^{-y_i \alpha_2 h_2(x_i)} \times \dots \end{aligned}$$

$$= \alpha_l(i) e^{-y_i \sum_{t=1}^{l-1} \alpha_t h_t(x_i)} \quad (1)$$

$$\text{From question, } w_l^i = e^{-y_i \sum_{t=1}^{l-1} \alpha_t h_t(x_i)} \quad (2)$$

By comparison of (1) + (2), $\boxed{\alpha_l(i) \propto w_l^i}$, proved. $\#$