

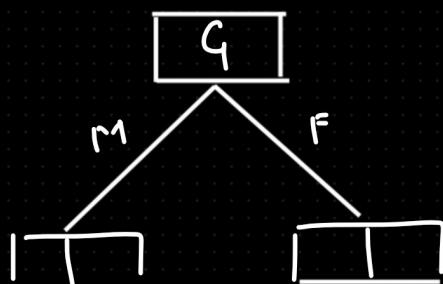
- 1 Decision tree regressor
- 2 Hyperparameter tuning
- 3 Cross validation
- 4 Task | few interview question
- 5 Pruning (Pre, Post)

Categorical

↓
Gender | Heart Disease

Gender	Heart Disease
M	Y
F	N
M	Y
F	Y
M	N

Gender [categorical]
Dependent [categorical]



(Numerical)

↓
Weight | Heart Disease

Weight	Heart Disease
220	Y
180	Y
225	Y
150	N
125	N

=
Categorical | Regression
Classification

- ① Sort the value
- ② take avg of adjustant value
- ③ wrt every avg value i need to find out
[Gini impurity | entropy Eq]

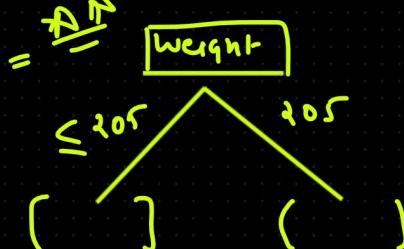
Weight	Heart-Disease
220	Y
180	Y
215	Y
190	N
155	N

Wherever we will get a low impurity will take that threshold only

Weight	HD
167.5	
185	i55 Y -
205	80 N -
225	90 N -
220	Y -
155	Y -

Info gain always will be high than only feature we will select-

\Rightarrow Gini impurity will only have lowest one



(3) \Rightarrow

$\frac{31/2N}{Weight}$

≤ 167.5 > 167.5

$\frac{0.1/N}{155}$ $\frac{3.1/1N}{215}$

Gini impurity

$$1 - \sum_{i=1}^n p_i^2$$

$$\text{Gini impurity}[L] = 0$$

$$\begin{aligned} \text{Gini imp}[R] &= 1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right] \\ &= 1 - \left(\frac{9}{49} + \frac{16}{49} \right) \end{aligned}$$

$$= 1 - \left[\frac{10}{16} \right]$$

$$= \frac{16-10}{16} = \frac{6}{16} = \underline{\underline{0.37}}$$

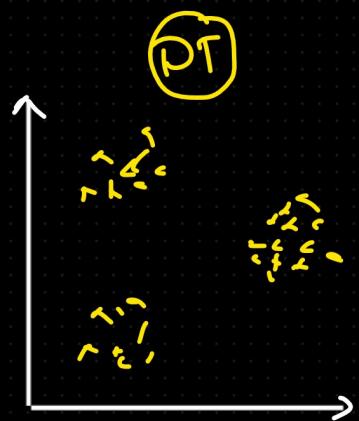
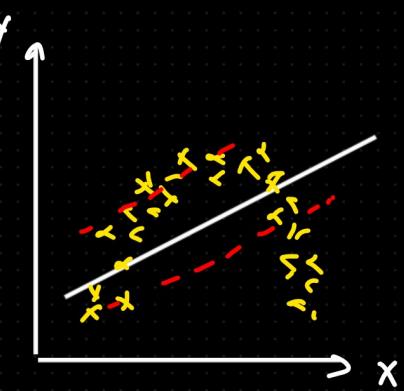
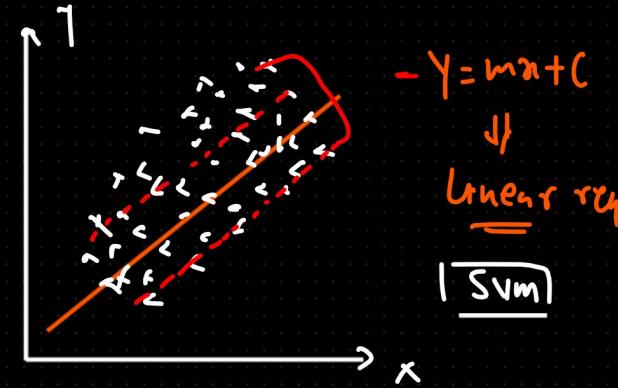
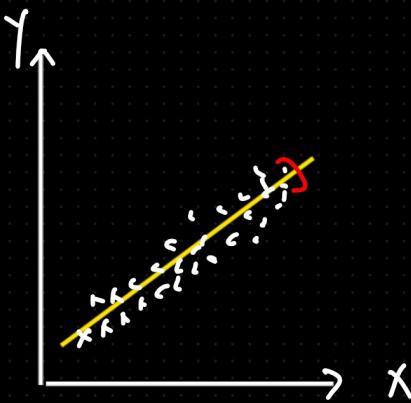
Weight Gini impurity

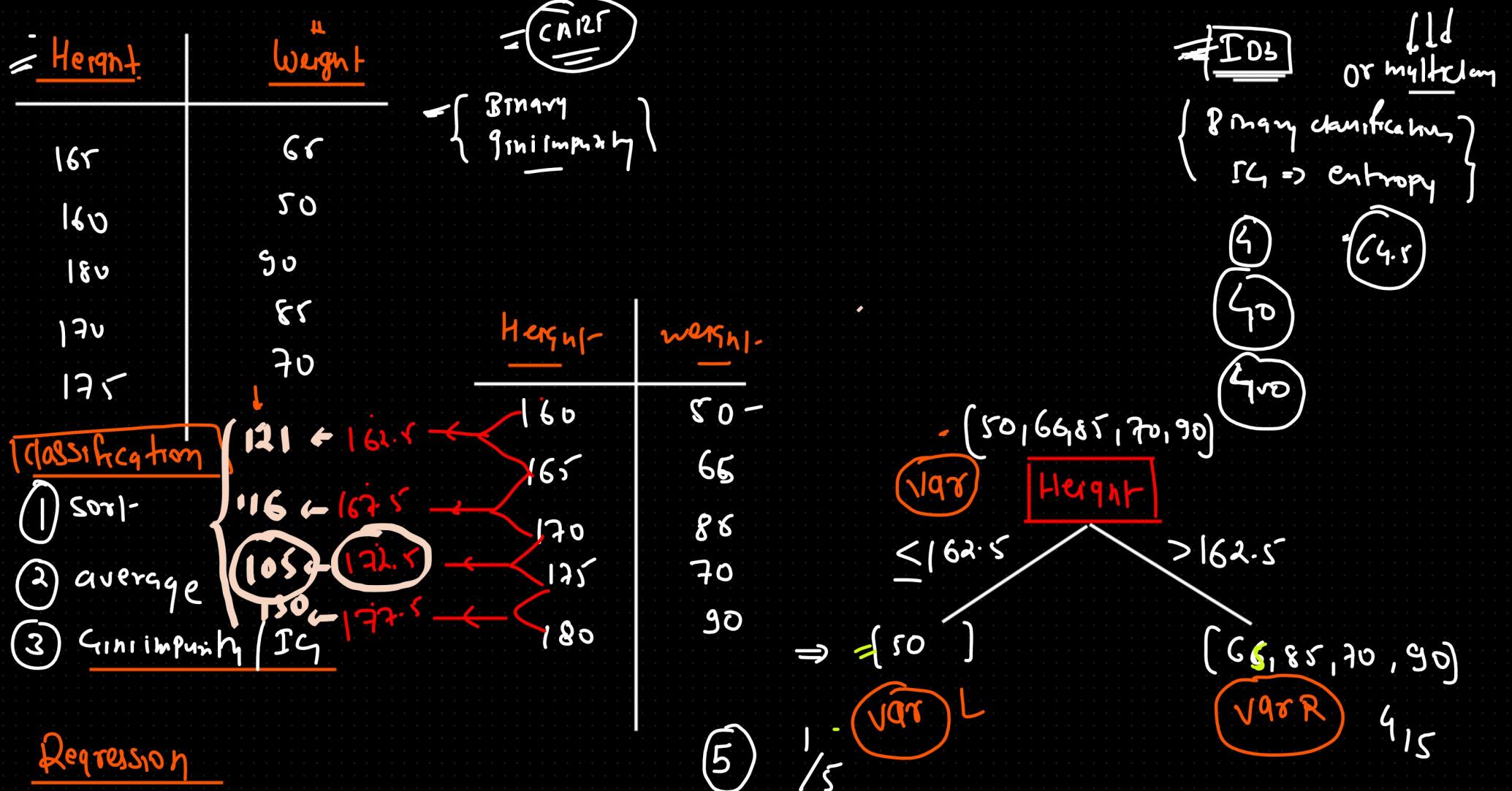
$$\rightarrow (L) \quad d \quad (R) \\ = \frac{1}{5} \times 0 + \left(\frac{4}{5} \times 0.37 \right)$$

$$= 0 + 0.296 = \boxed{0.296}$$

Decision-tree-regressor

<u>Height</u>	<u>dependent weight</u>	<u>I</u>	<u>Gender</u>	<u>Dependent- Wt</u>	<u>Weight-</u>
165	65	M		66	
160	50	F		50	
180	90	m		90	
170	85	F		85	
175	70	m		70	





$$\text{Var}_{\text{MSE}} = [50, 65, 85, 70, 90]$$

$$\text{Var}_{\text{MSE}} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2$$

$$\sqrt{Var} = ? \quad \underline{\text{Std dev (RMS E)}}$$

$$Var(L) = 0$$

$$Var(R) = [65, 85, 70, 90]$$

$$Avg = \frac{65 + 85 + 70 + 90}{4} \\ = 77.5$$

$$Var(R) = \frac{(77.5 - 65)^2 + (77.5 - 85)^2 + (77.5 - 70)^2 + (77.5 - 90)^2}{4}$$

$$Var(R) = 106.25$$

$$Var(R_{out}) = 206$$

$$Var(R) = 106.25$$

$$Var(L) = 0$$

$$\text{Mean} = \frac{65 + 85 + 70 + 90}{4} \\ = 72$$

$$Var(R_{out}) = \frac{(72 - 65)^2 + (72 - 85)^2 + (72 - 70)^2 + (72 - 90)^2}{4} \\ = \frac{(22)^2 + (7)^2 + (13)^2 + (2)^2 + (18)^2}{5}$$

$$\boxed{Var(R_{out}) = 206}$$

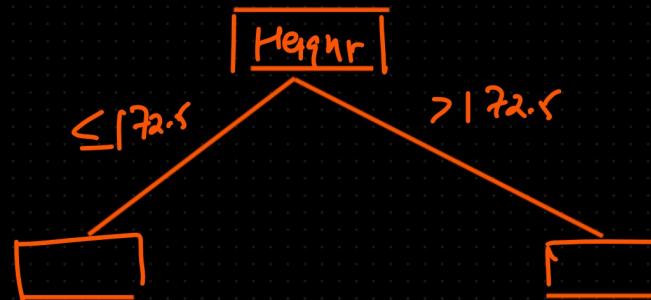
↓

$$\text{Reduction in Var} = \frac{Var(R_{out}) - \sum_{i=1}^h w_i \times Var(\text{Child})}{h}$$

$$= 206 - \left[\frac{1}{5} \times 0 + \frac{4}{5} \times 106.25 \right]$$

$$= 206 - [0 + 85] = \boxed{121}$$

lowest Reduction in var



Classification

D.P

Weight	Gender	0 NODb
50	m	0
60	f	0
70	m	No
80	f	0
90	m	No

= Weight | Gender | height

D.P

Pruning

PrePruning

Why?

Oversfitting

Post Pruning

Cutting

Pre Pruning \Rightarrow while i am building Decision

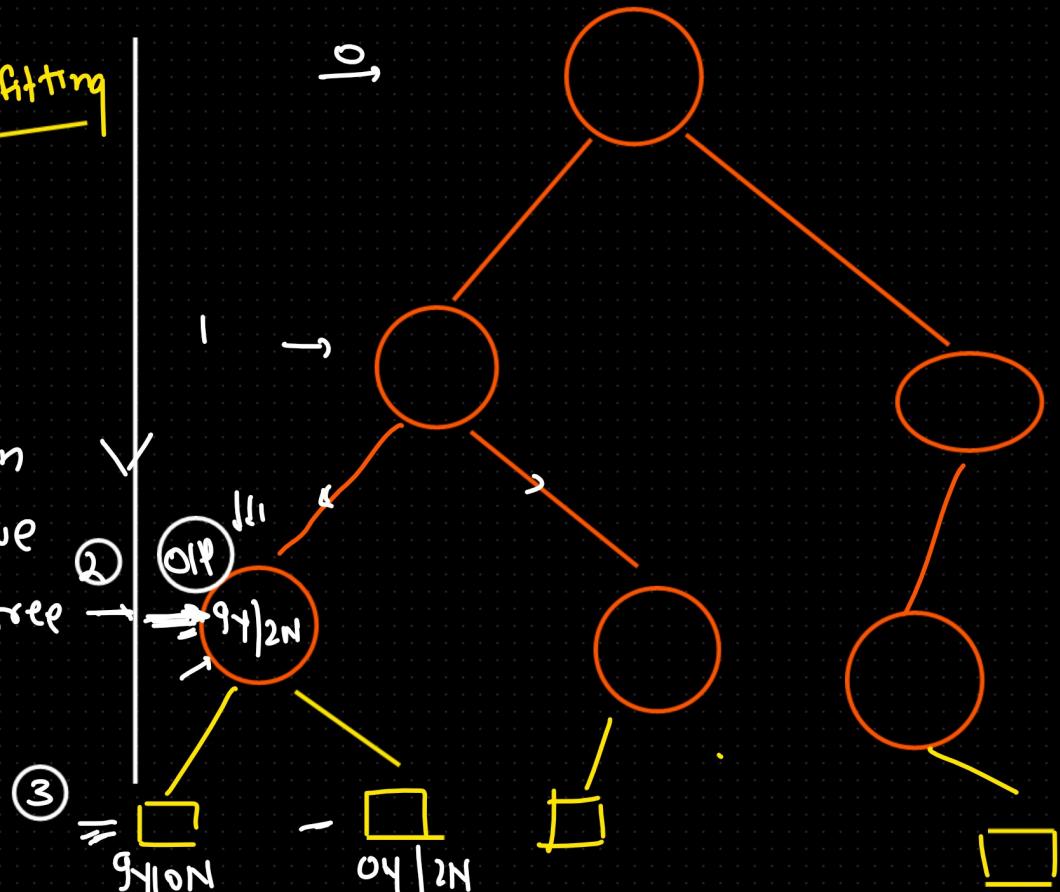
tree at that time only we
decide the depth of the tree

(Structure)

1 max_depth = 2

- 2 min-sample-Split
- 3 min-Sample-leaf
- 4 max-feature

} hyperparameter



Oversfitting \Rightarrow training accuracy will be high
test accuracy will low

ROLE

training accuracy
and

test accuracy
should be near to each other
there won't be much diff

Underfitting \Rightarrow training acc is low
test acc is also low

hyperparameter \Rightarrow $y = mx + c$

Non
Parametric
Model



Hyperparameter Tuning

- ① Grid Search $[CV]$ \Rightarrow cross validation
- ② Random Search $[CV]$

$$\underline{\text{Max-Depth}} = [2, 3, 4]$$

$$\underline{\text{Min-Sample-Split}} = [10, 15, 16]$$

max_Depth		
min_Sample_Split	10	15
10	✓	✓
15	✓	✓
16	✓	✓

$\Rightarrow 9$

How many times
g times

max_Depth, min_Sample_Split, leaf



more complex

= Grid search \Rightarrow In Grid Search we take every combination

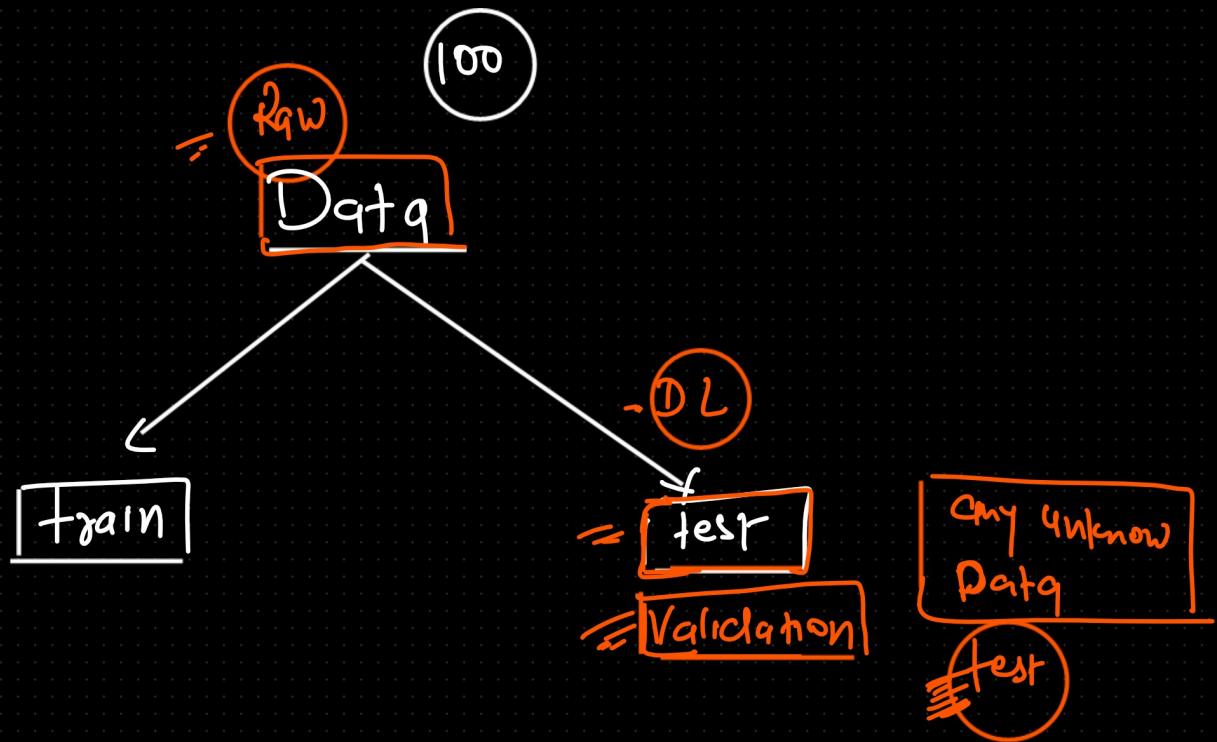
Random Search \Rightarrow We take - random combination

less complex

Cross Validation

Diffr \Rightarrow Random [split] [train-test-split]

- = 1 87 accuracy
- = 2 88 accuracy
- = 3 85 accuracy
- = 4 80 accuracy



\downarrow
 = Cross-Validation (testing)
CV

\downarrow
 K-fold CV

10 Data Point

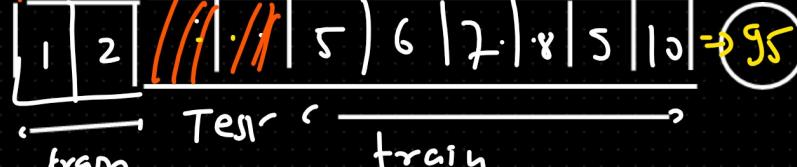
$k=5$ fold

$$\frac{10}{5} = 2 \Rightarrow \text{test Data}$$

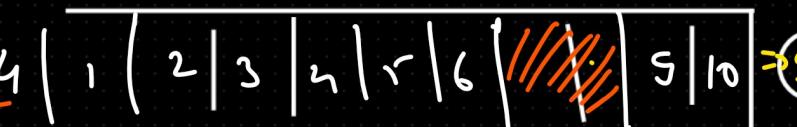
$$\frac{10}{2} = 5$$

$$10 - 2 = 8 \Rightarrow \text{training data}$$

- $k=1$ 
 $1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 = 90$

- $k=2$ 
 $1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 \Rightarrow 95$

- $k=3$ 
 $1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 \Rightarrow 91$

- $k=4$ 
 $1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 \Rightarrow 90$

- $k=5$ 
 $1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 \Rightarrow 85$

Assumption

Data

5/6

Grid Search CV

max-Depth = 3

min-sample-split = 3 \Rightarrow $3 \times 3 \times 5$

K-fold CV $K = 5$ ~~5~~ \Rightarrow 45

Instead of Random Split \Rightarrow Cross Validation

If Cross testing \Rightarrow iterate on each Data Point

disadvantage:

overfitting

if there any complex pattern of the data then it won't perform
it is little sensitive missing value
for the large dataset it won't be good

how to use decision tree for the feature selection?

if my data is imbalanced then should I use the DT and it is robust to DT?

first handle the imbalance data and then go for the building up this DT.

10==> 8 yes and 2 no

if you are giving such data to decision tree it might lead to wrong prediction

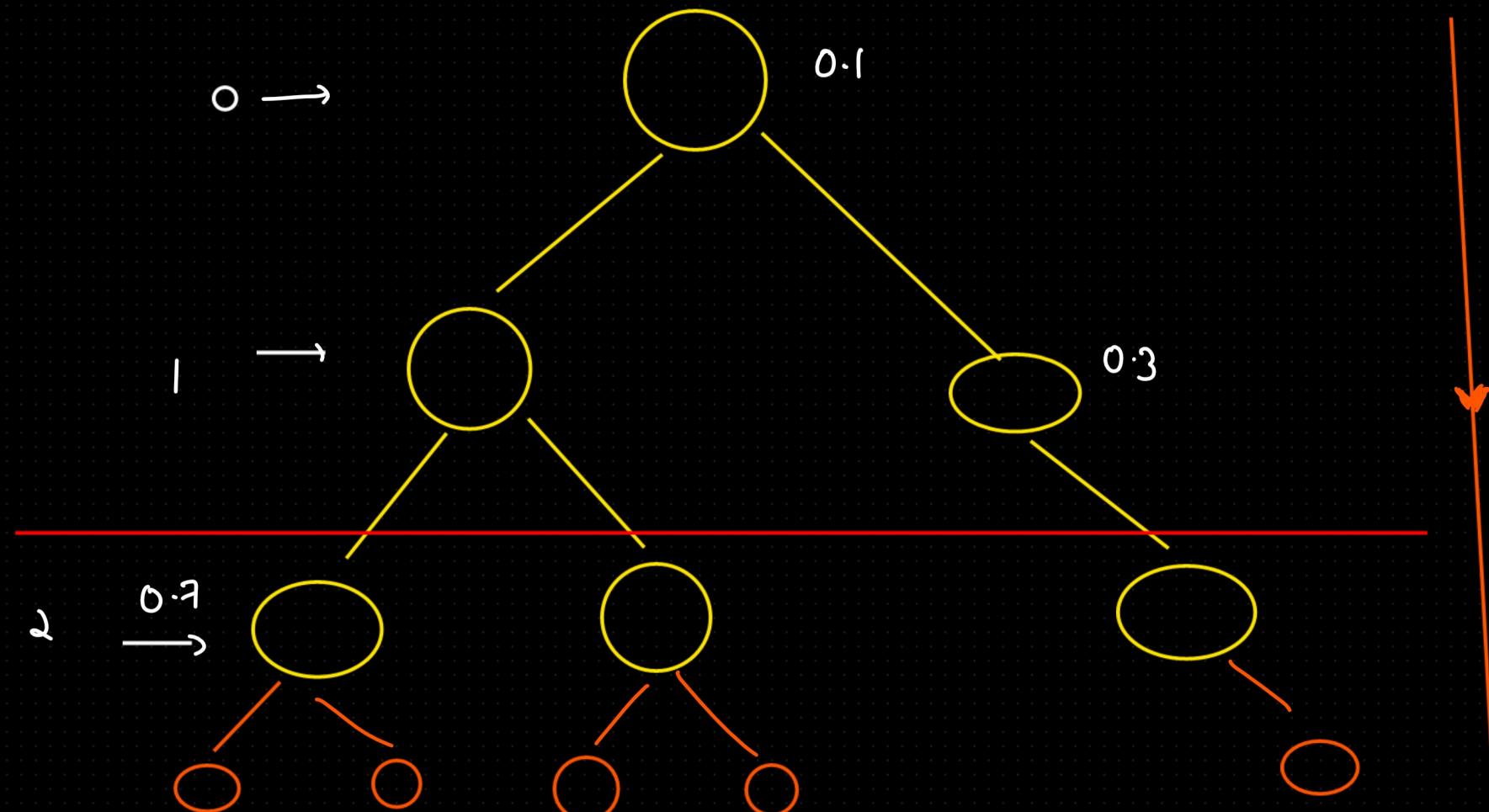
q10) can you give me those scenarios where I should not use DT?

can you tell them disadvantage

Post Pruning \Rightarrow first we building the entire tree till complete Depth.

and then we are going to prune it.

cutting



How it IS Possible \Rightarrow Build up our DT

With respect to every Node we have impurity value

Setup the threshold \rightarrow ~~(0.5)~~ $\text{CCP_alpha}_q = [0.5, 0.9, 0.8, 0.6]$

0.9 0.8 0.6

PrePruning \Rightarrow Decide the Depth of the tree while building up.

~~PostPruning~~ \Rightarrow Cut the branches after building up the complete tree

CP-Alpha - [-----]