

Classification  $\rightarrow$

Ex Mail  $\rightarrow$  Spam / Not Spam

Online Transaction  $\rightarrow$  Fraudulent / Non Fraudulent

Tumor  $\rightarrow$  Malignant / Benign

+ve -ve

where  $y \in \{0, 1\}$

0  $\Rightarrow$  -ve class

1  $\Rightarrow$  +ve class

Here  $y$  has Discrete Value.

In above case OIP has only 2 class  $\{0, 1\}$  } Binary classification.

To check Weather.

Possible OIP

Windy

Sunny

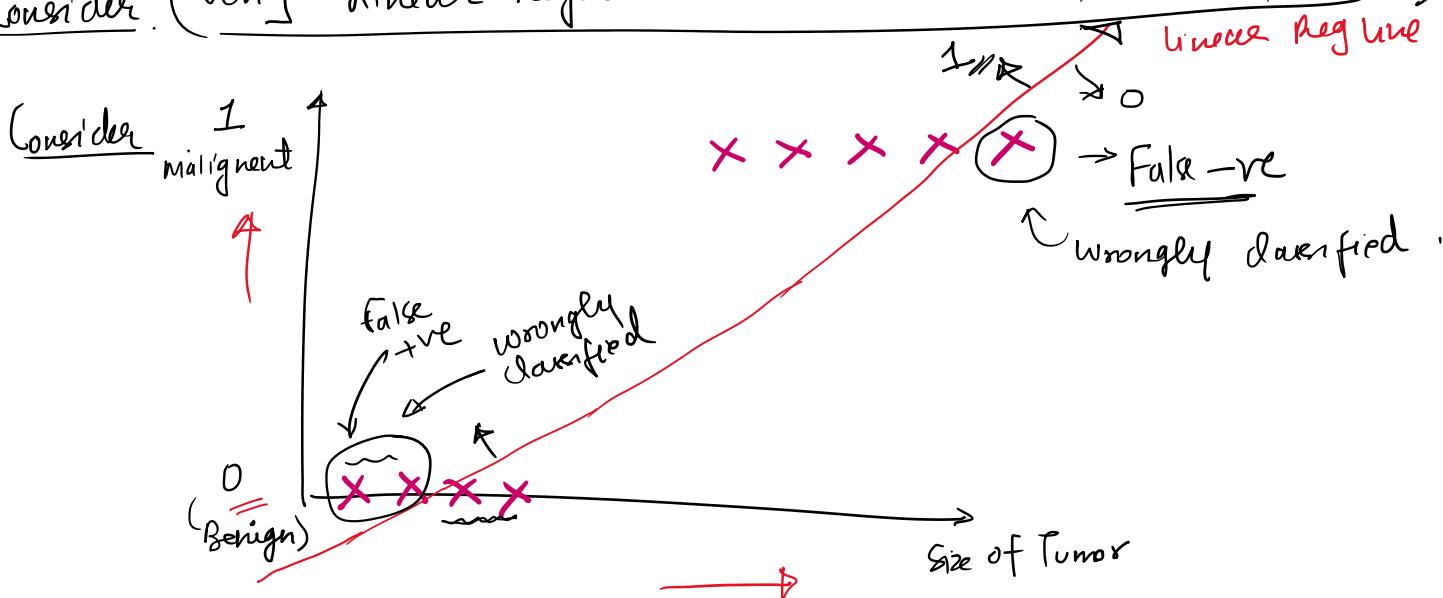
Cloudy

Rainy

OIP has more than one class [ multi class classification ]

\* Classification bothers about label and not the Exact value.

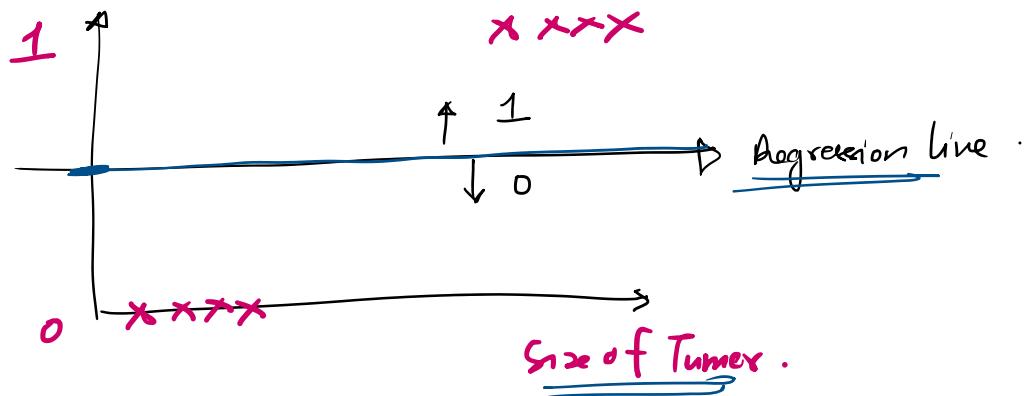
Consider (Why linear regression is not used for classification)



Suppose

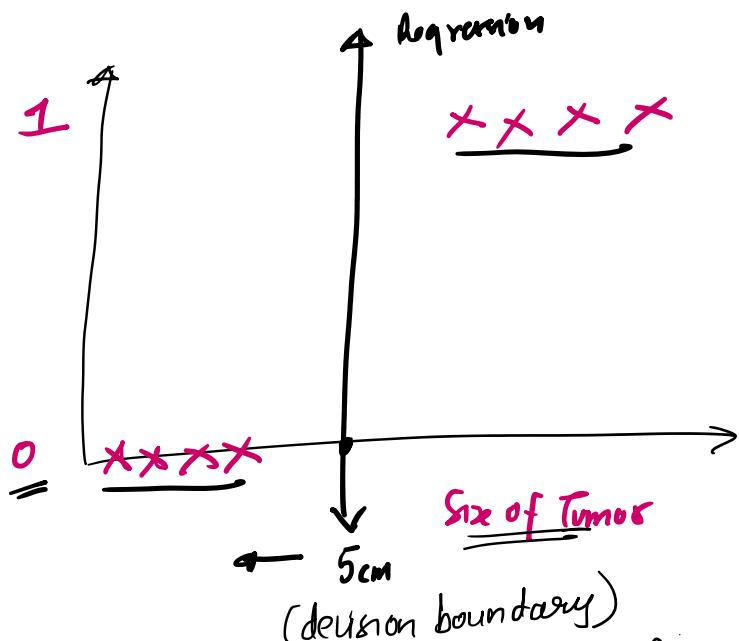
Malignancy:

Not Correct



Consider

Malignancy:



here we can observe for above regression line

If tumor size  $\leq 5\text{cm}$   $\rightarrow$  Yes (1) Malignant

If tumor size  $\leq 5\text{cm}$   $\rightarrow$  Yes (1) Malignant  
 tumor size  $> 5\text{cm}$   $\rightarrow$  No (0) Benign

We can say let 'P' denote Probability that  $y=1$  when  $\underline{\underline{X=x}}$ .

$$P = \underline{\underline{P}}(y=1 | \underline{\underline{X=x}}) = \frac{\beta_0 + \beta_1 x}{\text{In linear Reg}}$$

$P$  = probability lies bet<sup>n</sup>  $\underline{\underline{0 \text{ to } 1}}$

But linear function are unbounded.

and Expected o/p here is 0 or 1

So we cannot use regression to build classifier

$\therefore$  linear regression is not suitable for classification.

For classification we will use logistic regression.

\* In logistic regression we get probability score.

\* It predicts the probability of occurrence of event

$$\text{Odd} = \frac{\text{No of time the Event happens}}{\text{No of time the Event will not happen}}$$

Odd =  
Represents  
chances that  
the event  
will occur

Ex. If the odd of India winning against W.I. is  $\underline{4:1} = \frac{\text{No of India win}}{\text{No of India not win}}$   
 $= \frac{4}{1}$

Best Case  $\Rightarrow$  The odd of India winning against W.I. =  $\underline{\infty}$

If odd of W.I. winning against India is 1:4

$$= \frac{\text{No of W.I. win}}{\text{No of W.I. not win}} = \frac{1}{4}$$

Worst Case  $\Rightarrow$  W.I. is winning 0 match = Odd =  $\underline{0}$

Range of value that odd can take = 0 to  $\infty$

Relationship between Odd & Probability = Odd =  $\frac{P}{1-P}$

$$\therefore \text{odd} = \frac{P}{1-P}$$

Here we know that odd has range 0 to  $\infty$

$\rightarrow$  There is no upper bound for odd.

$\rightarrow$  But odd has lower bound.

$\rightarrow$  To remove lower bound, to have symmetrical analysis

Ex  $\text{odd} = 1:6 \quad 1/6 = 0.167 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{not symmetrical}$

$$\text{odd} = 6:1 \quad 6/1 = 6 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{not symmetrical}$$

But  $\ln(1/6) = -1.79 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{symmetrical}$

$$\ln(6/1) = 1.79 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{symmetrical}$$

By taking log we have overcomed the lower bound

also

$\log(\text{odd})$   $\nearrow$  will not have upper bound  
 $\searrow$  will not have lower bound.

$$\underline{\underline{\log(\text{odd})}} = y = \boxed{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n}$$

$$\log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n.$$

$$\text{let } z = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n.$$

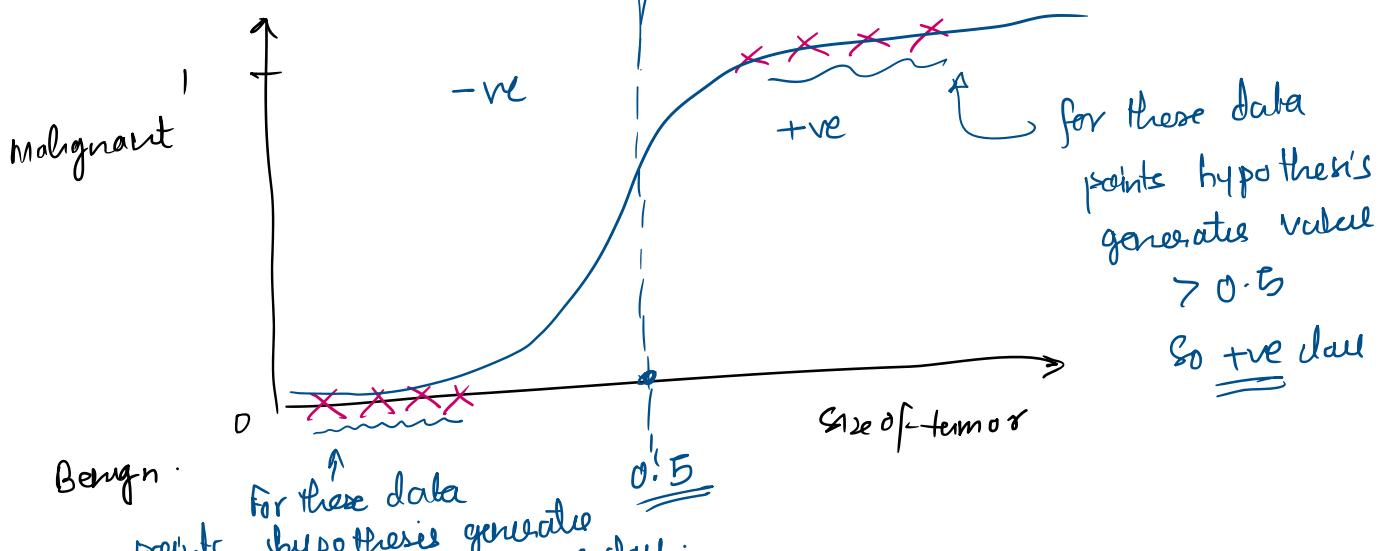
$$\therefore \log_e \left( \frac{P}{1-P} \right) = z$$

$$\therefore \frac{P}{1-P} = e^z$$

$$P = -e^z + e^z$$

$$P(1+e^z) = e^z$$

$$P = \frac{e^z}{1+e^z} = \frac{1}{1+e^{-z}} = \frac{1}{1 + \frac{1}{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}} = \text{Sigmoid function}$$



logistic regression Model

in linear reg  $h_{\theta}(x) = \theta^T x$  → range  $-\infty$  to  $+\infty$   
 ↑ feature vector  
 ↑ parameter vector  
 hypothesis (prediction).

In logistic regression

$$0 \leq h_{\theta}(x) \leq 1$$

↑ predicted value

$$h_{\theta}(x) = g(\theta^T x)$$

↑ Sigmoid

sigmoid

Let  $z = \mathbf{Q}^T \mathbf{x}$

$$h_{\theta}(x) = g(z) = \frac{1}{1 + e^{-z}}$$

↑  
prediction.

Estimated probability that  $y=1$  on given input  $x$ .

If  $\underline{h_{\theta}(x) = 0.7}$

This means There is 70% chance that tumor is Malignant.

Since  $h_{\theta}(x) = 0.7 > 0.5$

So  $\boxed{y=1}$

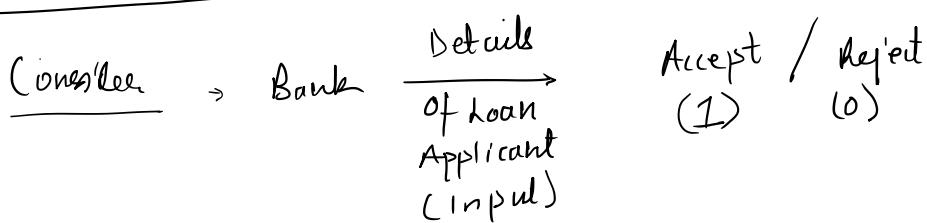
If  $\underline{h_{\theta}(x) = 0.3}$

This means There is 30% chance of tumor being malignant

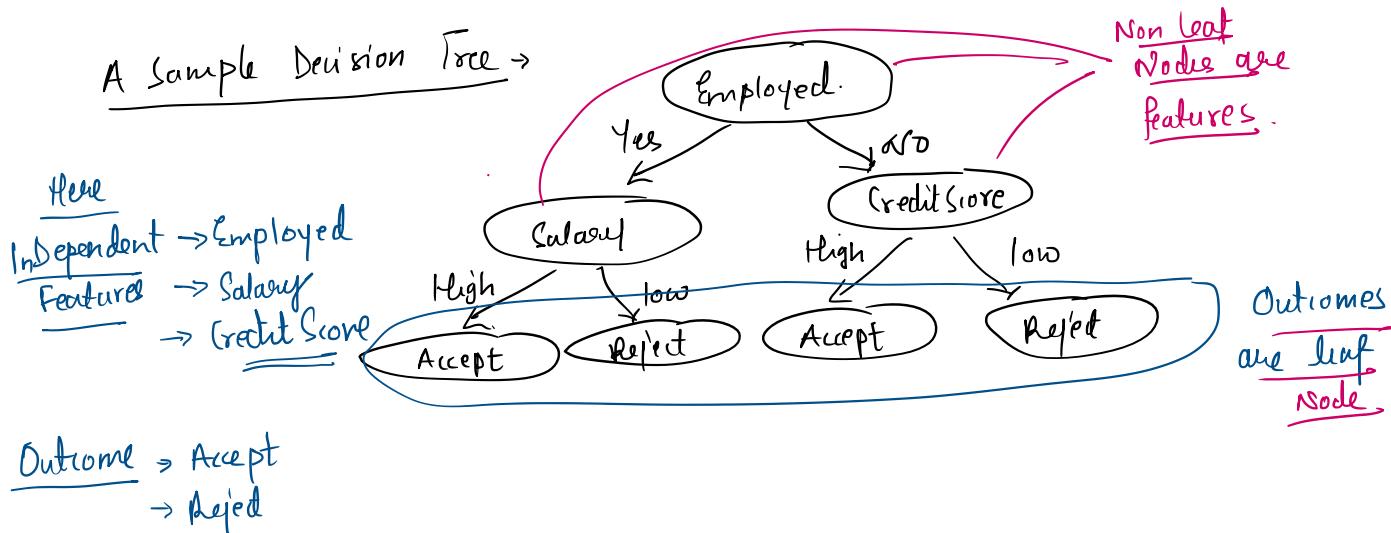
Since  $h_{\theta}(x) = 0.3 < 0.5$

So  $\boxed{y=0}$

## Decision Tree :



### A Sample Decision Tree



Q) Create Decision Tree for following using Gini Index  
(Classification & Regression Tree) → CART. ✓

Weekend	Weather	Parent	Money	Decision	(y)	P
w1	Sunny	Yes	Rich	Cinema	-	
w2	Sunny	No	Rich	Tennis	-	
w3	Windy	Yes	Rich	Cinema	-	
w4	Rainy	Yes	Poor	Cinema	-	
w5	Rainy	No	Rich	Stay In	-	
w6	Rainy	Yes	Poor	Cinema	-	
w7	Windy	No	Poor	Cinema	-	
w8	Windy	No	Rich	Shopping	-	
w9	Windy	Yes	Rich	Cinema	-	
w10	Sunny	No	Rich	Tennis	-	

Solution → Independent features: Weather  
Parent  
Money

Decision/Outcome  
Cinema  
Tennis  
Stay In  
Shopping.

↳ Shopping.

Step 1

We will calculate Gini Index for Overall collection of Outcomes of Training Examples.

These are 4 possible outcomes for decision

Cinema — 6 instances  
Tennis — 2 instances  
StayIn — 1 instance  
Shopping — 1 instance.

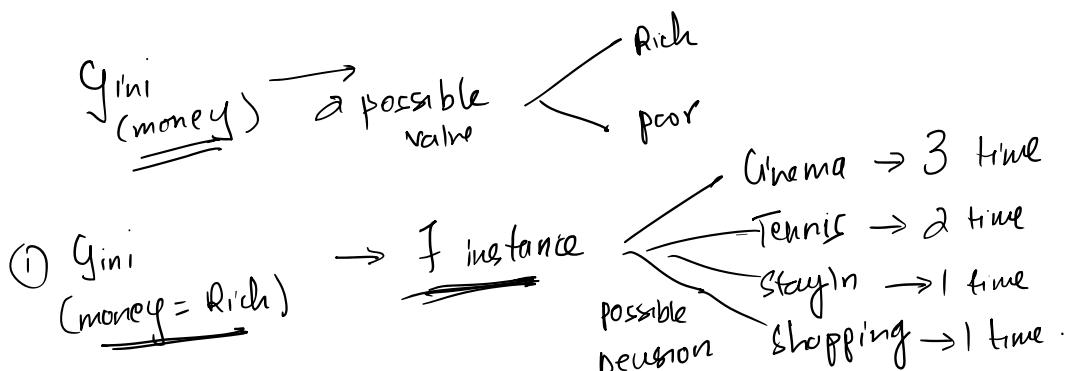
$$G_{ini} = 1 - \left( \left(\frac{6}{10}\right)^2 + \left(\frac{2}{10}\right)^2 + \left(\frac{1}{10}\right)^2 + \left(\frac{1}{10}\right)^2 \right)$$

$$(deusion) = 1 - \left( \frac{42}{100} \right) = 0.58$$

Note → In Machine Learning, Gini index/coefficient is utilized as an Impurity measure in decision tree for Classification.

$$G_{ini} = 1 - \sum_{i=1}^n (P_i)^2 \text{ where } P_i \text{ probability of outcome of specific data}$$

Step 2 To find Gini Index for Money



$$G_{ini} (\underline{\text{money}} = \text{Rich}) = 1 - \left( \left(\frac{3}{7}\right)^2 + \left(\frac{2}{7}\right)^2 + \left(\frac{1}{7}\right)^2 + \left(\frac{1}{7}\right)^2 \right)$$

$$= 0.694$$

(2)  $G_{ini}$  → n class → Cinema

(2)  $Gini_{(money=poor)}$   $\rightarrow$  3 instance  $\xrightarrow{\text{Decision}}$  Cinema.

$$= 1 - ((3/3)^2) = 0_{//}$$

Weighted Average  $Gini_{(money)} = (Gini_{(money=Rich)} * \text{proportion of Rich}) + (Gini_{(money=poor)} * \text{proportion of poor})$

$$= (0.694 * 7/10) + (0 * 3/10)$$

$Gini_{(money)}$   
 $= 0.485$

### Step 3 Gini Index on Parent

For Parent feature  $\xrightarrow[\text{values}]{\text{possible values}}$  Yes  
No

$Gini_{(parent=Yes)} = 5 \text{ instances} \xrightarrow[\text{possible decision}]{\text{Cinema}}$

$$= 1 - ((5/5)^2) = 0_{//}$$

$Gini_{(parent=No)} = 5 \text{ instances} \xrightarrow[\text{possible decision}]{\text{Tennis} \rightarrow 2 \text{ times}, \text{StayIn} \rightarrow 1 \text{ times}, \text{Shopping} \rightarrow 1 \text{ time}, \text{Cinema} \rightarrow 1 \text{ time}}$

$$= 1 - ((2/5)^2 + (1/5)^2 + (1/5)^2 + (1/5)^2)$$

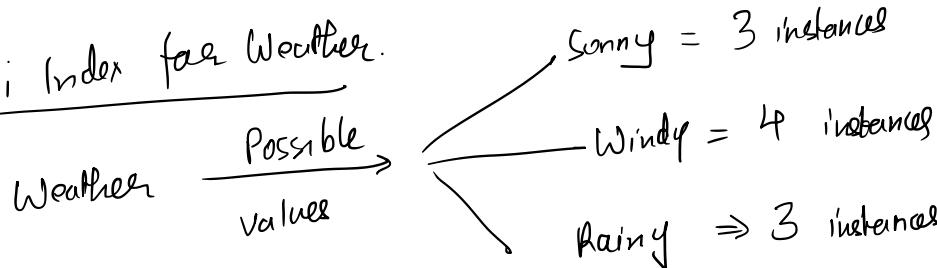
$$= 1 - (9/25) = 0.72$$

Weighted Average of  $Gini_{(parent)} = (0 * 5/10) + (0.72 * 5/10) = 0.36$

$Gini_{(parent)} = 0.36$

$$\boxed{Gini_{(parent)} = 0.36}$$

Step 4) Gini Index for Weather.



$Gini_{(\text{weather} = \text{Sunny})} \Rightarrow$

3 instances

possible

outcomes

$$= 1 - \left( \left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right) = \underline{\underline{0.444}}$$

3 time Cinema

1 time Tennis

$Gini_{(\text{weather} = \text{Windy})} \rightarrow$

4 instances

possible

outcome

$$= 1 - \left( \left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right) = \underline{\underline{0.375}}$$

$Gini_{(\text{weather} = \text{Rainy})} \rightarrow$

3 instances

2 cinema  
1 stay In  
possible outcomes

$$= 1 - \left( \left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right) = \underline{\underline{0.444}}$$

Weighted Average  $Gini_{(\text{Weather})} = \left(0.444 \times \frac{3}{10}\right) + \left(0.375 \times \frac{4}{10}\right) + \left(0.444 \times \frac{3}{10}\right)$

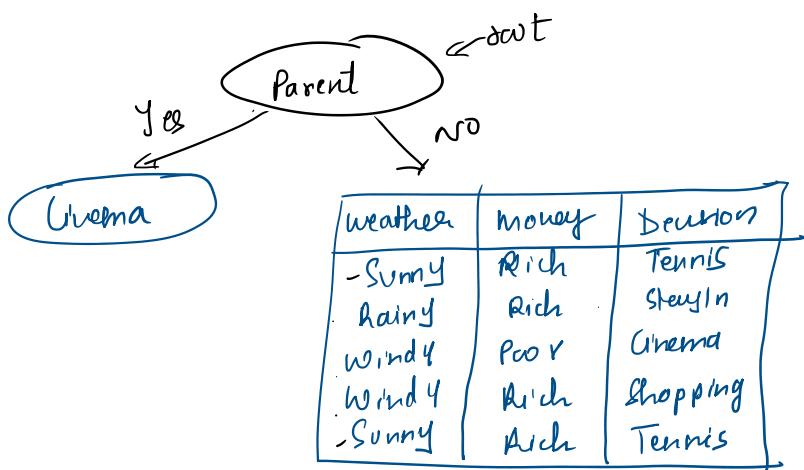
$$= \underline{\underline{0.414}}$$

$$\boxed{1 - r_{avg} = 0.486}$$

$$\boxed{\begin{aligned} Gini(\text{money}) &= 0.486 \\ Gini(\text{parent}) &= 0.36 \\ Gini(\text{weather}) &= 0.416 \end{aligned}}$$

Minimum Gini  $\rightarrow$  Minimum Impurity in Decision.  
Here Minimum Gini Value =  $Gini(\text{parent}) = 0.36$

So the root of Decision is Parent



We need to find  $Gini(\text{parent}=\text{No} \text{ and } \text{weather})$

Also  $Gini(\text{parent}=\text{No} \text{ and } \text{money})$

$\xrightarrow{\text{steps}} Gini(\text{parent}=\text{No} \text{ and } \text{weather})$   $\xrightarrow{\text{possible values}} \begin{cases} \text{Sunny - 2 times} \\ \text{Windy - 2 times} \\ \text{Rainy - 1 time} \end{cases}$   
5 instances

$$\begin{aligned} Gini(\text{parent}=\text{No} \text{ and } \text{weather}=\text{Sunny}) &\xrightarrow{\text{1 instance}} \text{Tennis} = 1 - \left( \left(\frac{1}{2}\right)^2 \right) = 0 // \\ Gini(\text{parent}=\text{No} \text{ and } \text{weather}=\text{Windy}) &\xrightarrow{\text{2 instances}} \begin{aligned} \text{Cinema} &= 1 - \left( \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right) = 0.5 \\ \text{Shopping} &= \text{possible outcome} \end{aligned} \end{aligned}$$

$$\text{Gini}'_{(\text{parent} = \text{No} \text{ and } \text{Weather} = \text{Rainy})} \xrightarrow[\text{1 instance}]{\text{possible outcome}} \text{StayIn} = 1 - ((Y_1)^2) = 0$$

$$\boxed{\text{Weighted Average Gini}_{(\text{parent} = \text{No} \text{ and } \text{Weather})} = 0.5 \times \frac{4}{5} = \underline{\underline{0.2}}}.$$

Step 6)  $\text{Gini}'_{(\text{parent} = \text{No} \text{ and Money})}$

Rich = 4 time  
poor = 1 time  
possible values.

5 instance

$$\text{Gini}'_{(\text{parent} = \text{No} \text{ and Money} = \text{Rich})} \xrightarrow[\text{(4P)}]{\text{possible outcome}} \begin{array}{l} \text{Tennis} - 2 \\ \text{StayIn} - 1 \\ \text{Shopping} - 1 \end{array}$$

$$= 1 - \left( \left( \frac{2}{4} \right)^2 + \left( \frac{1}{4} \right)^2 + \left( \frac{1}{4} \right)^2 \right) = \underline{\underline{0.625}}$$

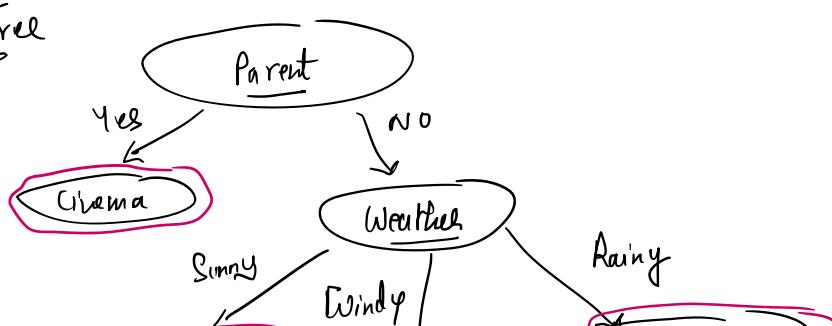
$$\text{Gini}'_{(\text{parent} = \text{No} \text{ and Money} = \text{Poor})} \xrightarrow{\text{possible outcome}} \text{Cinema} = 1 - ((Y_1)^2) = 0$$

$$\boxed{\text{Weighted Average Gini}'_{(\text{parent} = \text{No} \text{ and Money})} = 0.625 \times \frac{4}{5} = 0.5}$$

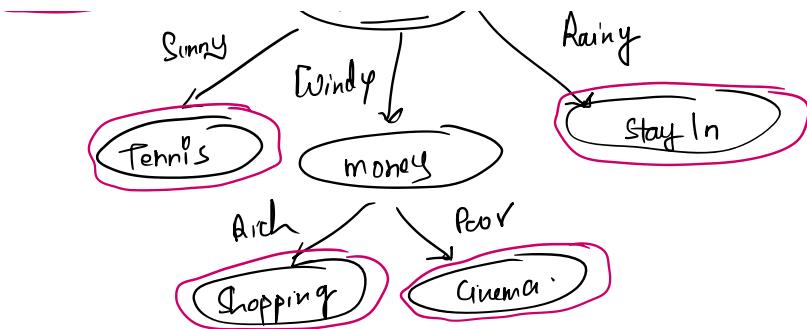
Here  $\text{Gini}'_{(\text{parent} = \text{No} \text{ and Weather})}$  has smallest value

So now Next Node = Weather

\* Updated Decision Tree



Aho



HW

Construct an optimal Decision Tree for following

outlook	Temperature	Humidity	windy	Play (Decision)
Sunny	Hot	High	False	NO
Sunny	Hot	High	True	NO
Overcast	Hot	High	F	Y
Rainy	Mild	High	F	Y
Rainy	Cool	Normal	F	Y
Rainy	Cool	Normal	T	NO
overcast	Cool	Normal	T	Y
Sunny	Mild	High	F	Y
Sunny	Cool	Normal	F	Y
Rainy	Mild	Normal	F	Y
Sunny	Mild	Normal	T	Y
Overcast	Mild	High	T	Y
Overcast	Hot	Normal	F	Y
Rainy	Mild	High	T	NO