

- 1 Data ingested
- 2 EDA
- 3 Processing or FE
- 4 Model building
- 5 Evaluation

**Core step**

## Model building

- 1 Supervised ML
- 2 Unsupervised ML

Regression

Classification

- multiclassification

= 1 Linear regression  $\Rightarrow$  Regression

$L_1, L_2, L_1+L_2$   
 $\uparrow$  lasso     $\uparrow$  ridge     $\downarrow$  elasticnet

= 2 Logistic regression  $\Rightarrow$  Classification

= 3 SVM  $\Rightarrow$  Regression  
Classification

~~Tree based approach (condition based)~~

① DT  $\rightarrow$  DTC  
 $\searrow$  DTR

② RF

③ AB  
 $\searrow$  XBG

④ GBM

## Decisiontree

- ① DTC
- ② DTR

- = ① DTC
- ② Entropy
- ③ Gini impurity | Gini coeff | Gini index
- ④ IG
- ⑤ Pruning = Pre → Post-
- ⑥ DTQ

ID<sub>3</sub>, CART, C<sub>4.5</sub>

! 1.  
age = 15      2.  
age = 17      3.  
age = 23

if (age ≤ 15) :

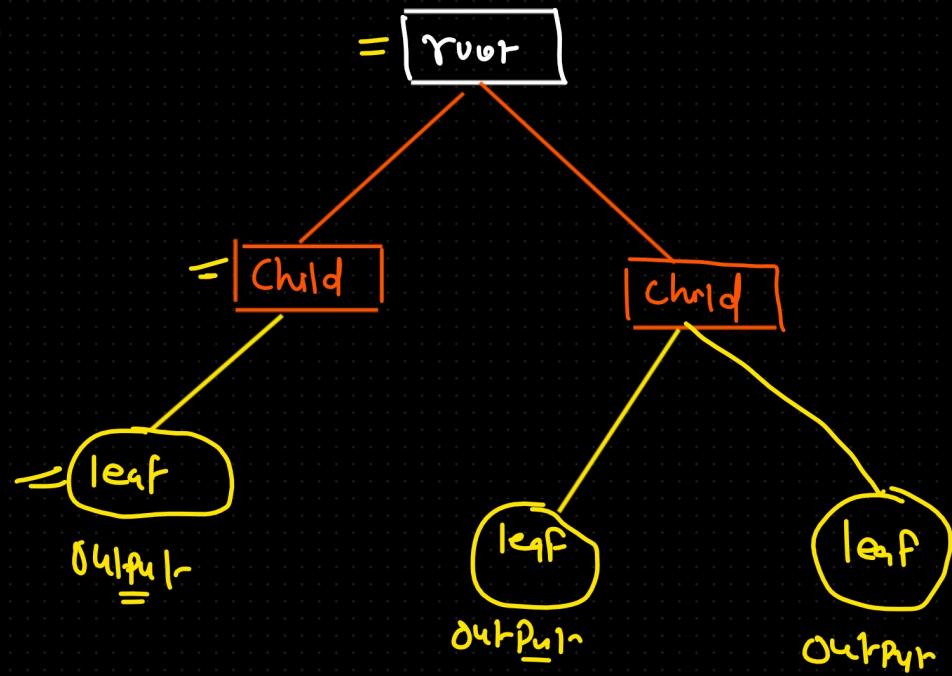
    Print ("schw1")

elif (age > 15 and age ≤ 21) :

    Print ("College")

else :

    Print ("Wrong")



Decision-tree - ① ID3 Iterative Dichotomiser 3  $\Rightarrow$  C4.5  
- ② CART Classification and regression tree

### ID3

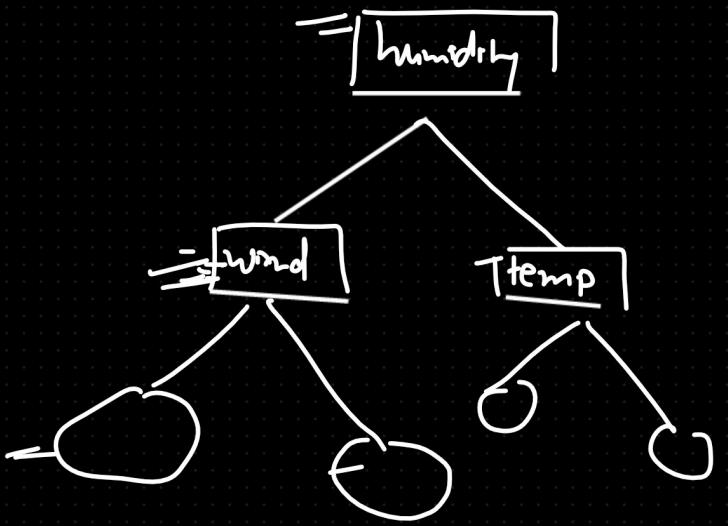
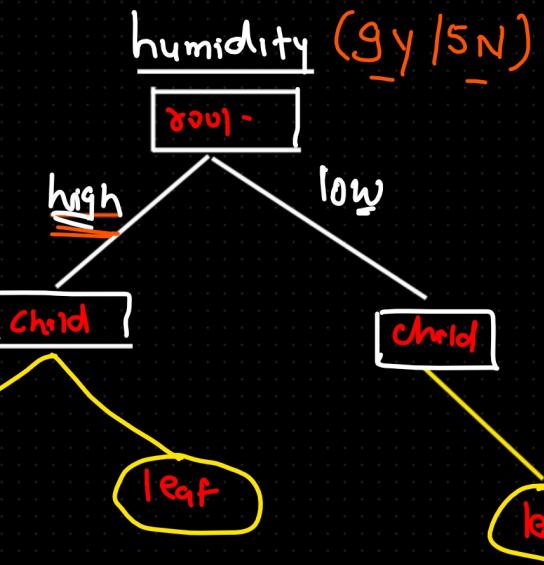
- ① Entropy  $\Rightarrow$  Information gain ( $I_q$ )
- ② Categorical (Classification) -  
it <sup>may</sup> divide the root feature  
more than 2 category
- ③

### CART

- ① Gini impurity
- ② Classification and regression
- ③ It always divide root-  
feature into two category.

	$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	$f_6$	$f_7$
Outlook	Temperature	Humidity	Wind	Played football(yes/no)			
Sunny	Hot	High	Weak	No	+tzyt		
Sunny	Hot	High	Strong	No			
Overcast	Hot	High	Weak	Yes			
Rain	Mild	High	Weak	Yes			
Rain	Cool	Normal	Weak	Yes			
Rain	Cool	Normal	Strong	No			
Overcast	Cool	Normal	Strong	Yes			
Sunny	Mild	High	Weak	No			
Sunny	Cool	Normal	Weak	Yes			
Rain	Mild	Normal	Weak	Yes			
Sunny	Mild	Normal	Strong	Yes			
Overcast	Mild	High	Strong	Yes			
Overcast	Hot	Normal	Weak	Yes			
Rain	Mild	High	Strong	No			

classification



Stage 1

$\Rightarrow -\text{high}$

$-[3y | 5N]$   
leaf node

$-[\text{Normal}]$

$-[gy | N]$

$gy | 5N (14)$

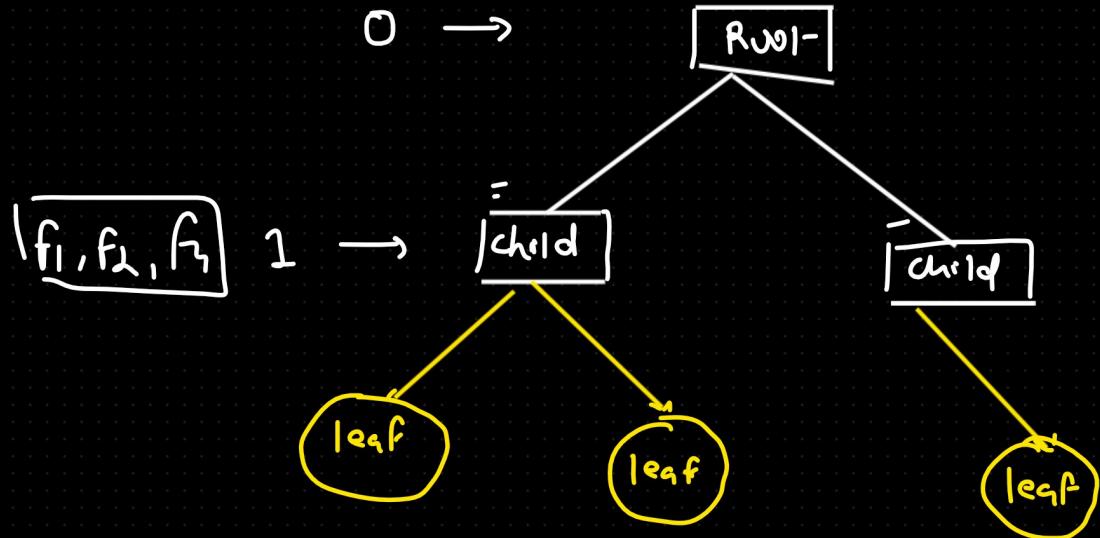
$gy | 5N$

humidity

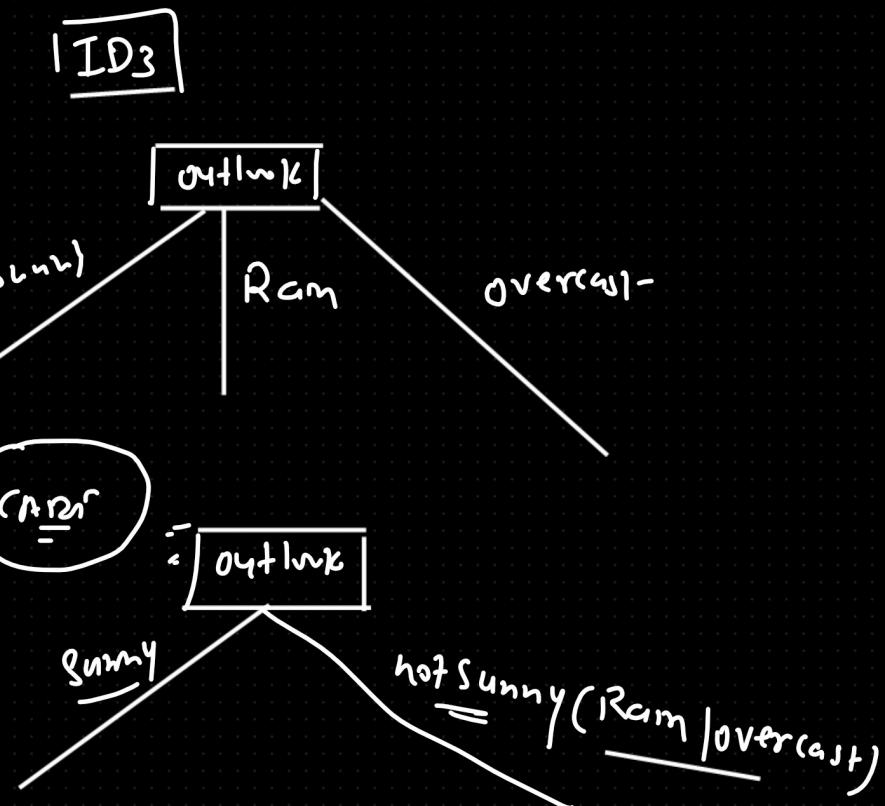
① DT  $\xrightarrow{DTR} =$   
 $\xrightarrow{DTC}$

Datq  
-  $f_1$   $F_2$   $f_3$   $f_4$  OIP

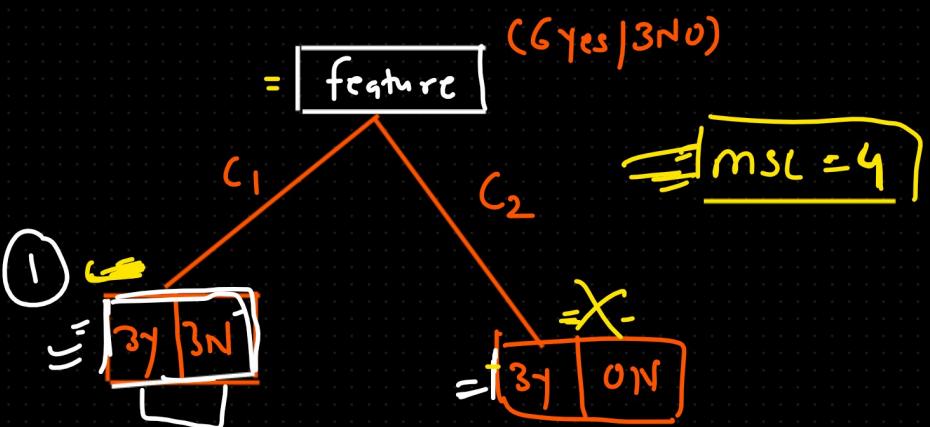
= ID3  $\Rightarrow$  Entropy  $\rightarrow$  Information gain —  
, CART  $\Rightarrow$  Gini impurity =



Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No



<u>feature</u>	0 P
$c_1 \longrightarrow$	Y
$c_2 \dots$	Y
$c_1 \longrightarrow$	Y
$c_2 \dots$	Y
$c_1 \longrightarrow$	Y
$c_1 \longrightarrow$	N
$c_2 \dots$	Y
$c_1 \longrightarrow$	N
$c_1 \longrightarrow$	N



ID<sub>3</sub> = entropy

CART = Gini coeff

$$\text{ENTROPY} = - \sum_{i=1}^N p_i \times \log_2(p_i) \quad (\text{2 class})$$

$$H(S) = - p_Y \times \log_2(p_Y) - p_N \times \log_2(p_N)$$

$$\log_2 2 = 1$$

$$H(S)(3y | 3N) = - \frac{3}{6} \times \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right)$$

$$= -\frac{1}{2} \times \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right)$$

$$= \frac{-1}{2} \left[ \log_2(1) - \log_2(2) \right] - \frac{1}{2} \left[ \log_2(1) - \log_2(2) \right]$$

$$= -\frac{1}{2} [0-1] - \frac{1}{2} [0-1]$$

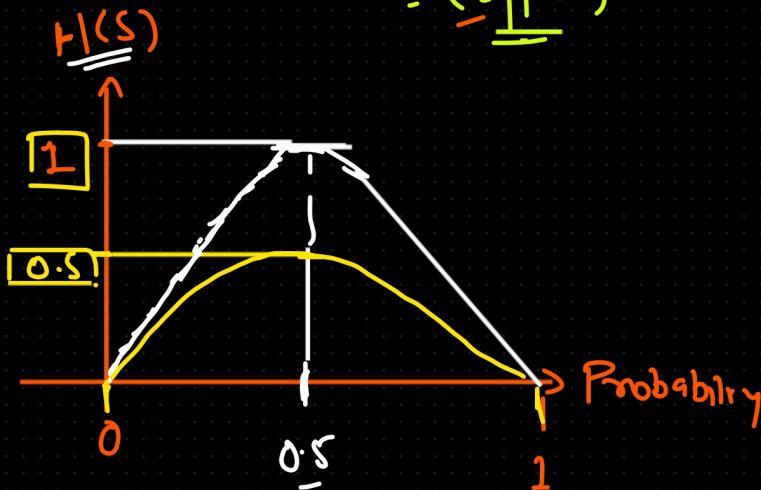
$$= -\frac{1}{2} [-1] - \frac{1}{2} [-1]$$

$$H(S) = \frac{1}{3} + \frac{2}{3} = \boxed{1}$$

$$H(S)(3Y|0N) = -\frac{3}{3} \log_2 \left( \frac{3}{3} \right) - \frac{0}{3} \log_2 \left( \frac{0}{3} \right)$$

$$H(S) = -\frac{3}{3} \log_2 (1) = \boxed{0}$$

feature with less impurity is considered to be the root feature?



$$=\underline{\underline{(2|3N)}}$$

$$H(S) = ?$$

$$-\sum_{i=1}^N p_i \times \log_2(p_i)$$

$$\underline{\underline{2+3=5}}$$

$$0.97 \rightarrow 1$$

$$\Rightarrow -p_1 \log(p_1) - p_N \log(p_N)$$

$$\Rightarrow -\frac{2}{5} \log_2(2/5) - \frac{3}{5} \log_2(3/5)$$

$$\Rightarrow \underline{\underline{0.97}}$$

Cross Entropy or Cross Impurity  $\Rightarrow$

$$1 - \sum_{i=1}^k (p_i)^2$$

$$\textcircled{1} \quad 3y|3N \Rightarrow$$

$$\textcircled{2} \quad 3y|0N \Rightarrow$$

$$\text{G.I.} \Rightarrow 1 - \left[ p_y^2 + p_N^2 \right]$$

$$\textcircled{1} \quad \underline{\underline{3y|3N}} \Rightarrow$$

$$1 - \left[ \left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2 \right]$$

$$1 - \left[ \left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2 \right]$$

$$\Rightarrow 1 - \left[ \left( \frac{1}{2} \right)^2 + \left( \frac{1}{2} \right)^2 \right]$$

$$1 - \left[ (1)^2 + 0 \right]$$

$$1 - 1 = 0$$

$$\Rightarrow 1 - \left[ \frac{1}{2} + \frac{1}{2} \right]$$

$$\Rightarrow 1 - \left[ \frac{2}{2} \right]$$

$$\Rightarrow 1 - 1 = 0$$

$\underline{\text{S0t}}$   $\underline{\text{3y}}$   $\underline{\text{13N}}$

impure entropy



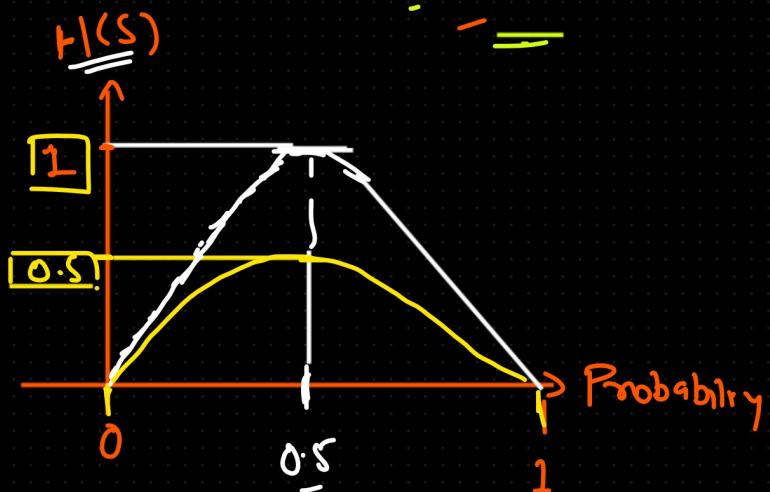
0.5

①

$$3y | \bar{0}N \Rightarrow 0 \Rightarrow \underline{\text{Impurity}} = 0$$

$G_{\text{mix}} = 0$

entropy = 0



$\boxed{\text{Entropy} = }$  slower  
 $\boxed{\text{Gini-coeff} = 0.5}$   
 $\boxed{\text{faster}}$

Which One Should Use  
ID3 = Entropy

CART = Gini-coeff

$\Leftarrow$  Large Dataset  $\rightarrow$  ~~1 2 3~~  $\Rightarrow$  Slow

= CART

$$\Rightarrow \boxed{1 - \sum_{i=1}^n (p_i)^2}$$

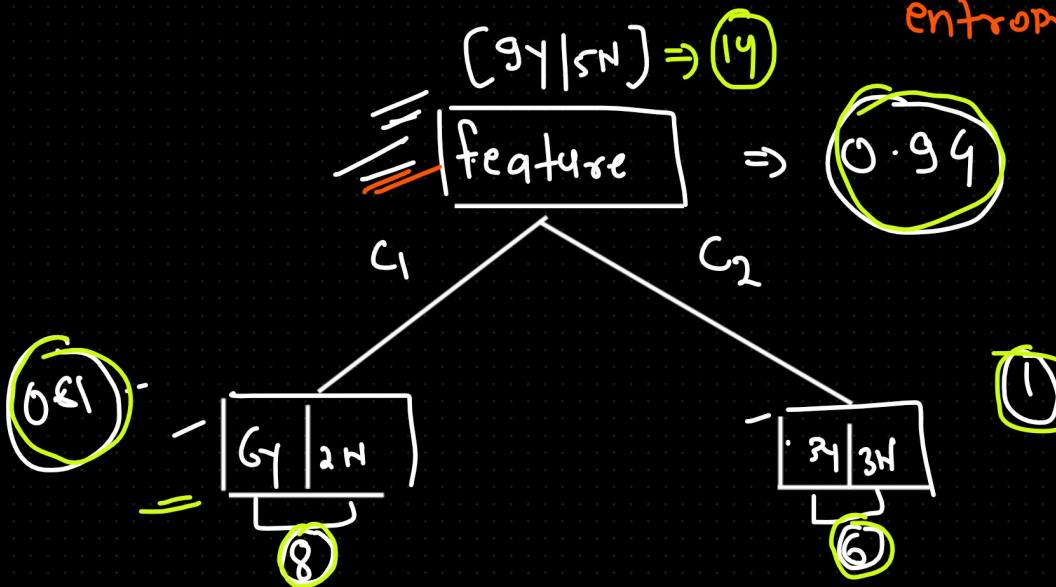
Information gain

$$\text{G.I.} \Rightarrow 1 - \left[ p_y^2 + p_N^2 \right]$$

$$Gain = H(S) - \sum_{V(G)-Val} \frac{|S_V|}{|S|} H(S_V)$$

root Node  
entropy

child node entropy  
or  
entropy after split



$$H(S) = \text{root feature entropy} = -P_Y \log_2(P_Y) - P_N \log_2(P_N)$$

$$= -\frac{9}{14} \log_2(\frac{9}{14}) - \frac{5}{14} \log_2(\frac{5}{14})$$

$$= -(0.64) \log_2(0.64) - 0.35 \log_2(0.35)$$

$$H(S)(\text{root}) \approx 0.9402$$

$$H(Y|Z=1) = -\frac{6}{8} \log_2(6/8) - \frac{2}{8} \log_2(2/8) = 0.81$$

$$H(Y|Z=0) = -\frac{3}{6} \log_2(3/6) - \frac{3}{6} \log_2(3/6) = 1$$

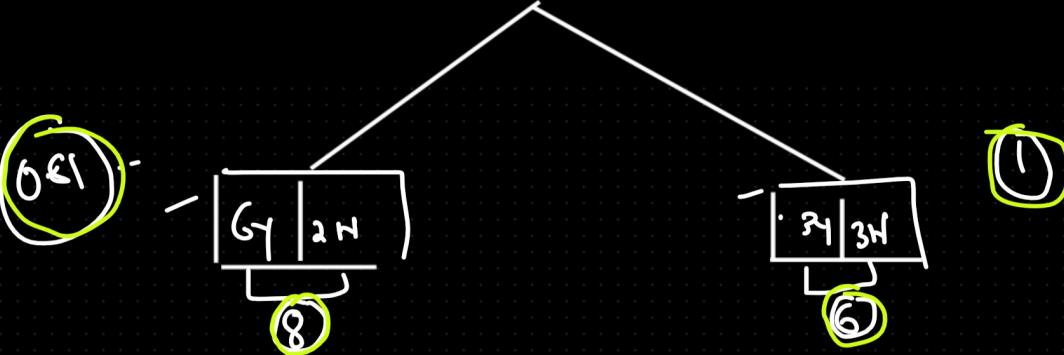
↓

$$\text{gain} \Rightarrow H(S) = \sum_{v \in V^{\text{int}}} \frac{|S_v|}{|S|} H(S_v)$$

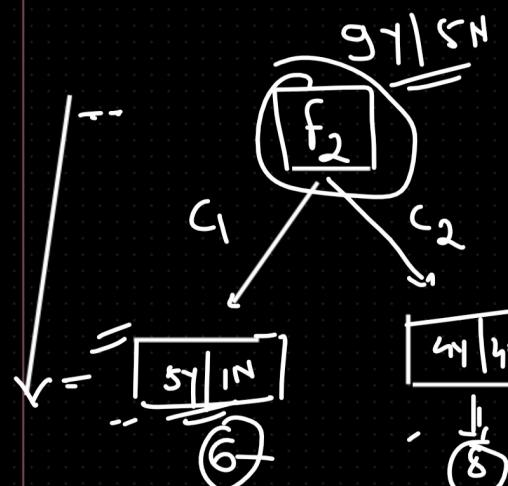
$$\Rightarrow 0.94 - \left[ \frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$\Rightarrow 0.94 - [0.462 + 0.42] = 0.0485 \approx 0.049$$

$$\begin{aligned} & (H(Y|S=1)) \Rightarrow 0.94 \\ \hline & \boxed{\text{feature}} \\ & \hline C_1 & C_2 & \Rightarrow 0.94 \end{aligned}$$



$$\text{gain} \Rightarrow H(S) - \sum_{v \in V_1} \frac{|S_v|}{|S|} H(S_v)$$



$$0.94 - \left[ \frac{6}{14} \times 0.650 + \frac{8}{14} \times 1 \right]$$

$$0.94 - \left( \frac{6 \times (0.650)}{14} + \frac{8}{14} \right)$$

$$\frac{0.94 - 0.85}{}$$

$$mgs = 5$$

$$f_2 \Rightarrow -0.09$$

$$f_1 = 0.049$$

$$f_2 > f_1$$

root?

$f_2$  will be my root feature = ?

Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No

①



Task

I/4

$$\left\{ \begin{array}{l} \text{Outlook} = 0.01 \\ \text{temp} = 0.5 \\ \text{Humidity} = 0.67 \\ \text{wind} = 0.012 \end{array} \right.$$



①

Build the complete DT on top of this Data.

②

C4.5 algo it's update version of ID3

①

→ | Humidity

18%

$$\left\{ \begin{array}{l} \text{outlook} \\ \text{temp} \\ \text{wind} \end{array} \right. \rightarrow \boxed{\text{Humidity}}$$

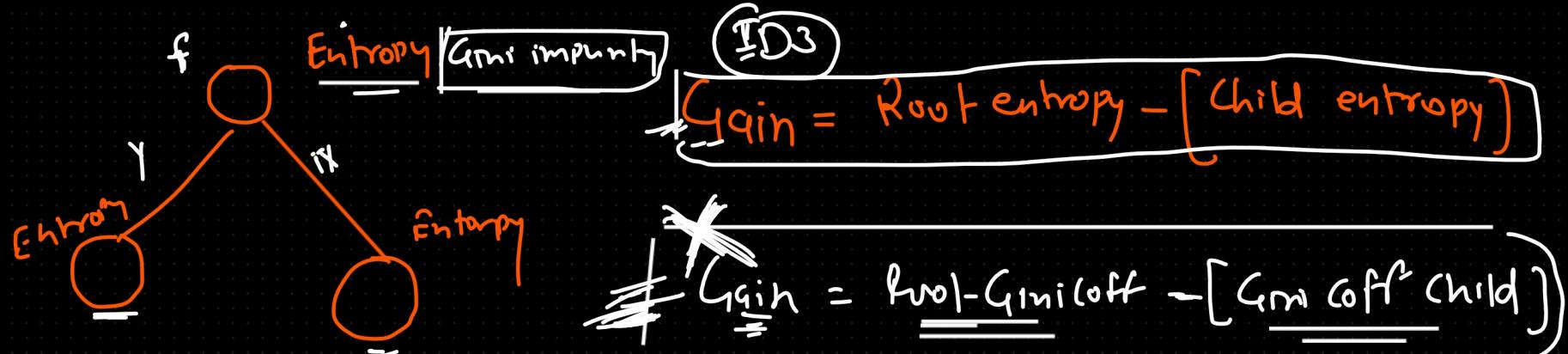
I/4

$$\left\{ \begin{array}{l} \text{outlook} \\ \text{temp} \\ \text{wind} \end{array} \right. \rightarrow \boxed{\text{Humidity}}$$

$$\left\{ \begin{array}{l} \text{outlook} \\ \text{temp} \\ \text{wind} \end{array} \right. \rightarrow \boxed{\text{Humidity}}$$

$$\left\{ \begin{array}{l} \text{outlook} \\ \text{temp} \\ \text{wind} \end{array} \right. \rightarrow \boxed{\text{Humidity}}$$

$$\boxed{\text{Humidity}}$$



~~Gain = Root-Gini<sup>off</sup> - [Gini<sup>off</sup> child]~~

CART

Calculation of all Gini<sup>off</sup>

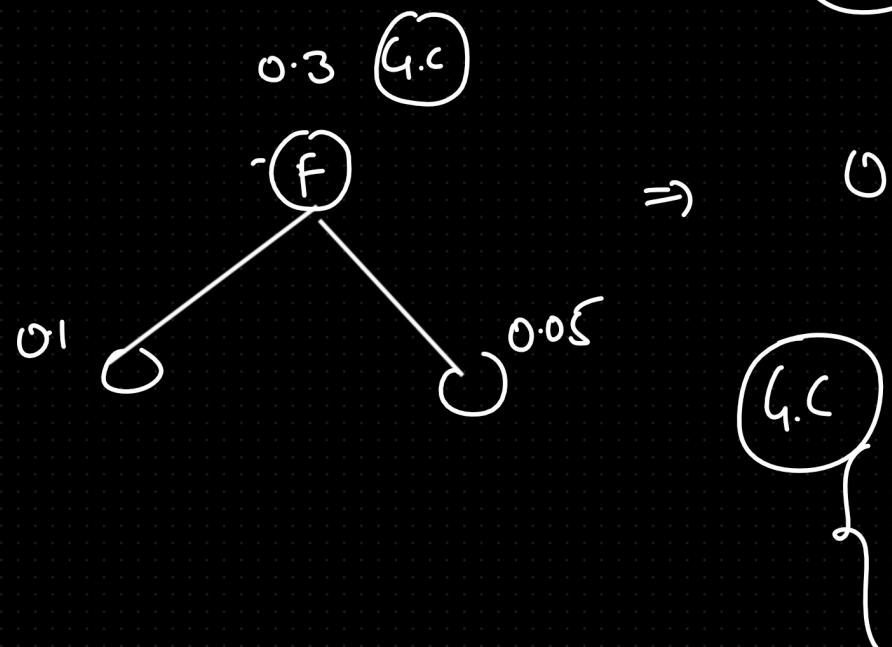
Summation

$0.3 + 0.1 + 0.05 = 0.45$

$0.4 + 0.05$

$f_1 = 0.45$   
 $f_2 = 0.35$   
 $f_3 = 0.01$

less impurity



① DTC  $\Rightarrow$  Implementation

② DTR -

③ Pruning | Post-Pruning

① DTC

① ID3  $\Rightarrow$  Entropy | IG  
② CART  $\Rightarrow$  Gini Coeff.

Entropy  $\Rightarrow$  formula

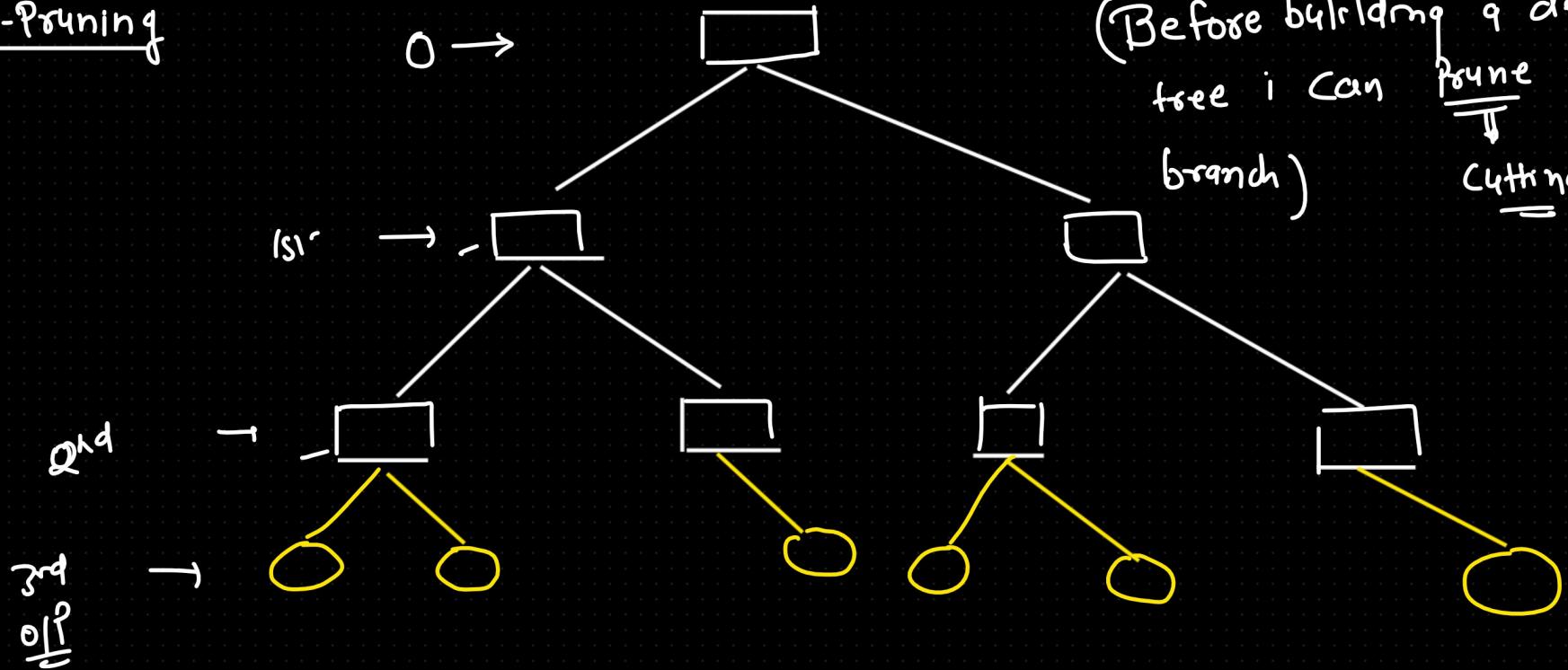
Gini formula

I.G. formula

Assignment.  $\Rightarrow$  Build a complete DT  
on the given Data

Explore C4.5 algo

Pre-Pruning

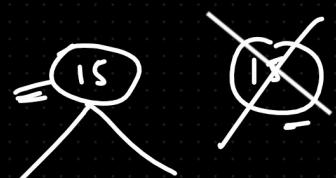


if we going to build DT till complete Depth = ?

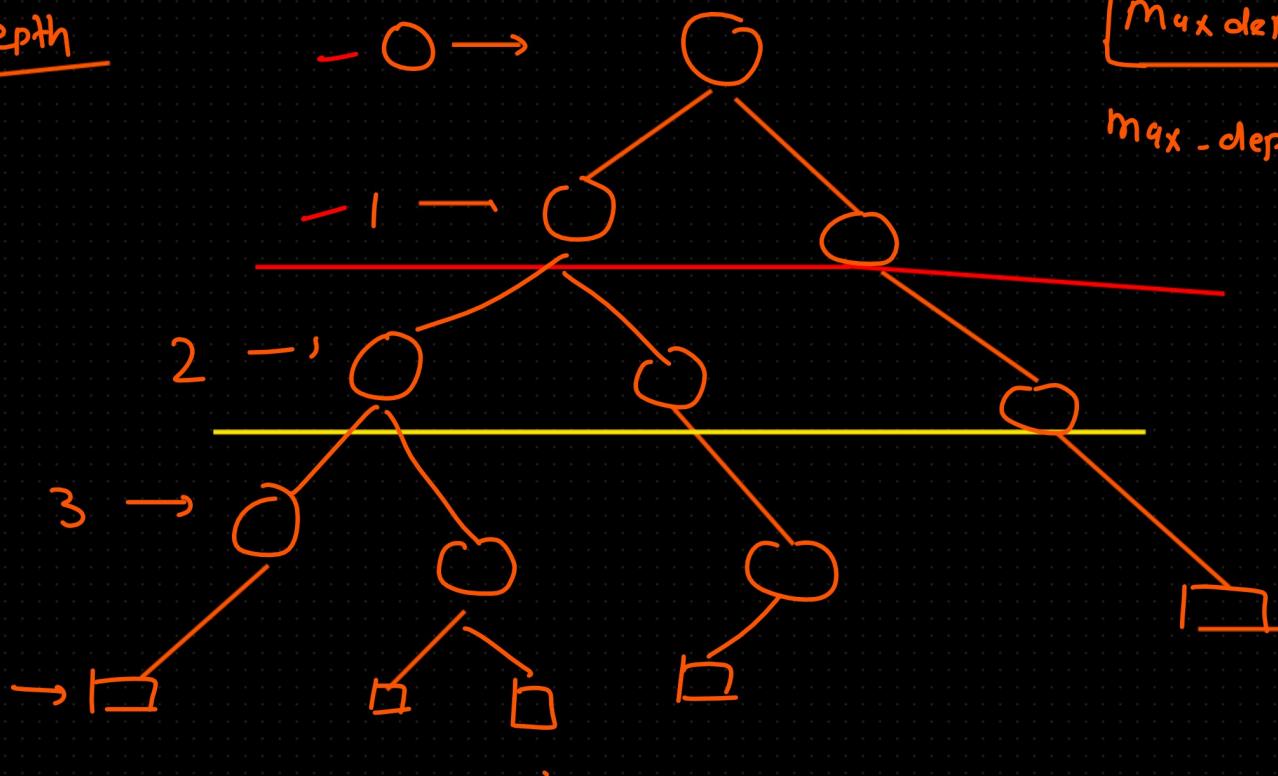
- ① Overfitting (train accuracy will be good test accuracy will be bad)
- ② Computational expensive

Preprune = {  
 ① Max\_Depth  
 ② Minimum Sample leaf  
 ③ Minimum Sample Split  
 ④ max\_feature  
 ⑤
 }

$$\text{MSS} = 10 \\ \underline{\text{mss}} = 18$$



① Max-Depth



②

minimum + sample - leaf

(constraint)

$$msL = 14$$

= 15

12

10

③ minimum sample split



④ max\_feature  $\Rightarrow$   $[f_1, f_2, f_3, f_4, f_5, f_6]$