

Day 4 - Machine Learning Algorithms

Agenda

- ① Decision Tree CLASSIFICATION
- ② Decision Tree REGRESSION
- ③ PRACTICAL IMPLEMENTATION
- ④ Ensemble Techniques

Agenda

{ DAY 1, DAY 2, DAY 3 }

⇓
Experience

Decision Tree {solving many usecase}

↳ Regression
↳ Classification

Regression
& Classification

⇓⇓⇓⇓
Decision Trees ⇓

if (age ≤ 18):

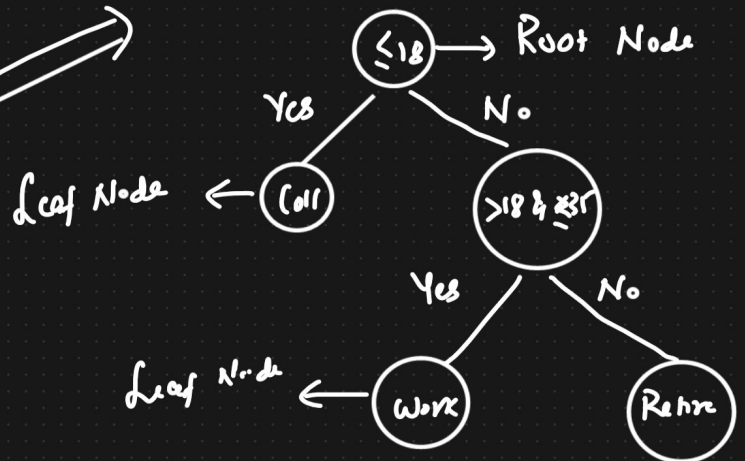
Print ("College")

elif (age > 18 and age ≤ 35):

Print (work)

else :

Print ("Retire")



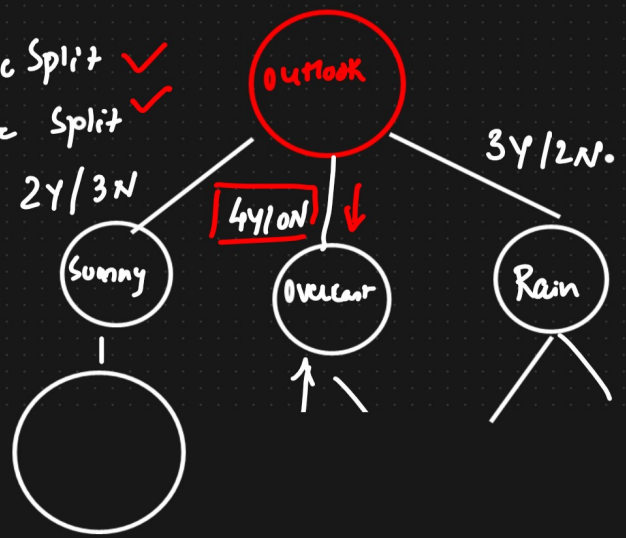
DECISION TREE

Nest if else \Rightarrow Decision Tree

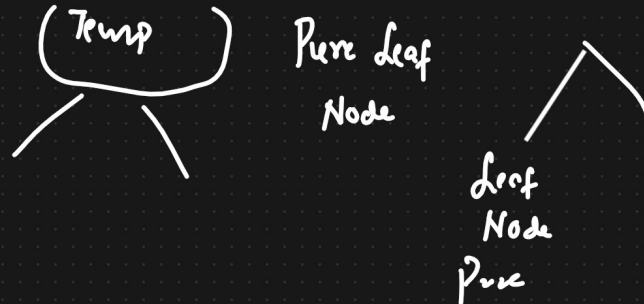
CLASSIFICATION 94/5 N

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	No
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Pure Split ✓
Impure Split ✓



- ① Purity \rightarrow Pure Split ??
- Entropy
 - Gini Impurity



- ② How the features are selected
- Information Gain?

① Entropy

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

① Gini Impurity

$$G.I = 1 - \sum_{i=1}^n (p_i)^2$$

+ } Binary
- } Binary

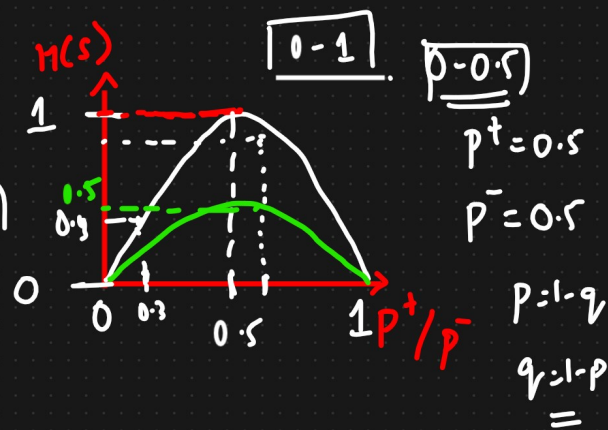
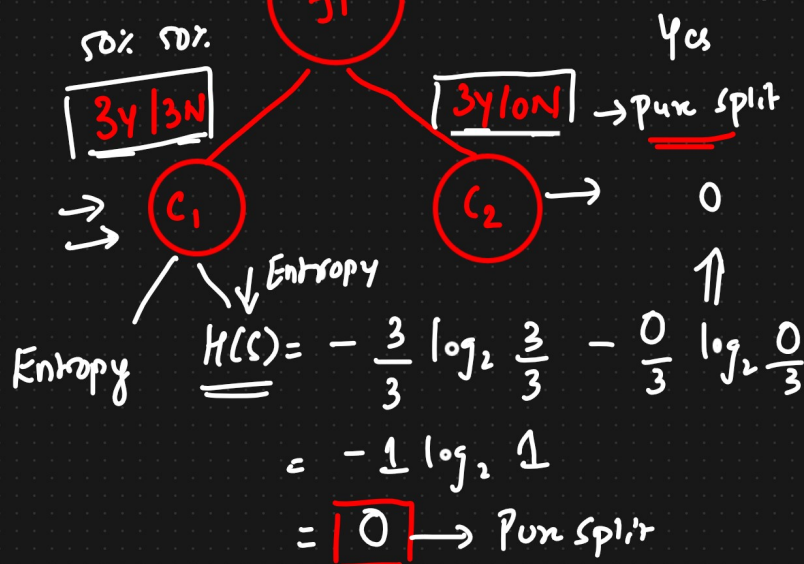
50% 50%

$64/3N$

f_1

$P_+ = \text{Probability of Yes}$

$+ \Rightarrow \text{Yes} \quad - \Rightarrow \text{No}$



$$H(S) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

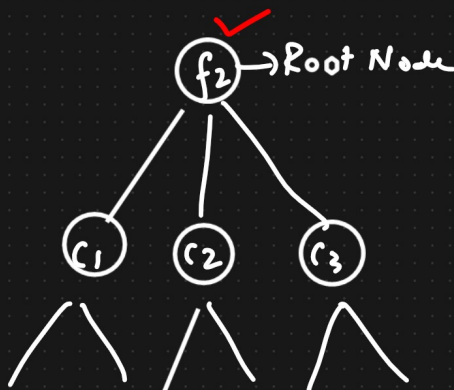
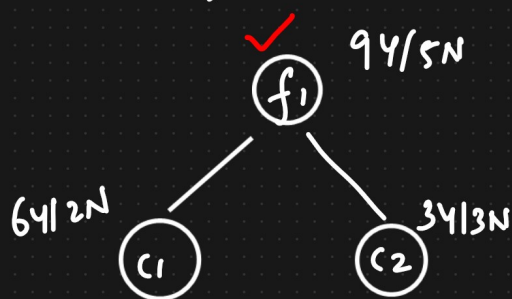
$$= 1 \checkmark$$

Impure Split

Purity Test \rightarrow Entropy

$f_1 \quad f_2 \quad f_3$

- ② Which feature to take to split??



Information Gain

$$\text{Gain}(S, f_1) = H(S) - \sum_{v \in \text{Val}} \frac{|S_v|}{|S|} H(S_v) \quad \checkmark$$

$$H(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

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

$$\approx \underline{\underline{0.94}}$$

$$H(c_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$\underline{H(c_1) = 0.81} \quad \underline{H(c_2) = 1}$$

$$\text{Gain}(S, f_1) = 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$\text{Gain}(S, f_1) = \underline{\underline{0.049}}$$

$$\text{Gain}(S, f_2) = \underline{\underline{0.051}}$$

$$\text{Gain}(S, f_2) >> \text{Gain}(S, f_1)$$

Using which feature
Should I start splitting
first

⑧ Gini Impurity \checkmark

$$G.I = 1 - \sum_{i=1}^n (p_i)^2 \rightarrow$$

$$= 1 - [(p_+)^2 + (p_-)^2]$$

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

$$= 1 - \left[\frac{1}{2} \right] = \underline{\underline{0.5}}$$

{Entropy}

n = 2 output {Yes
No

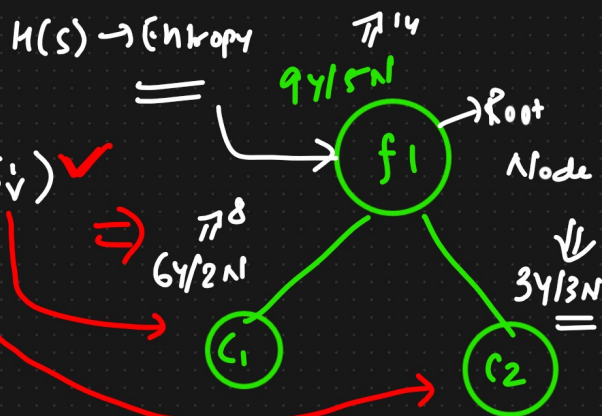
2Y/2N \Rightarrow Impure Split



Entropy = 1

Gini Impurity = 0.5

Entropy \rightarrow log

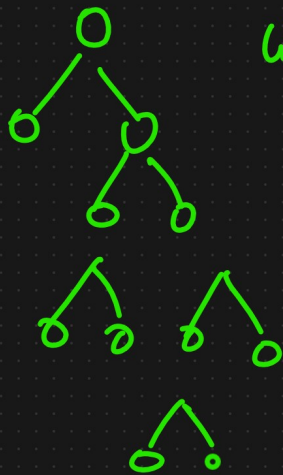


100

Gini Impurity → Simple Maths

Fast

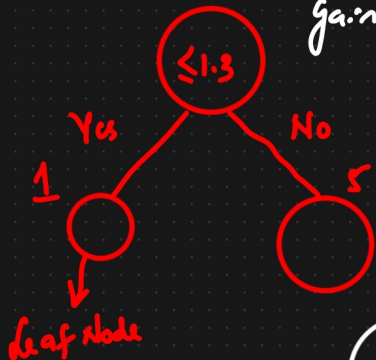
Gini >> Entropy



Continuous

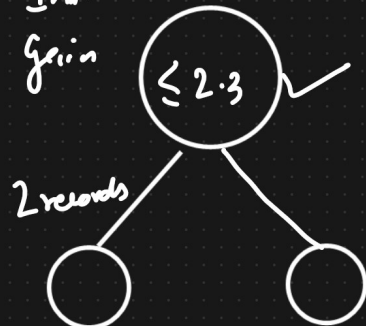
f_1	o/p	$\Rightarrow \boxed{f_1}$
2.3		$\frac{1.3}{2.3}$
1.3		$\frac{3}{4}$
4		
5		
7		
3		

%

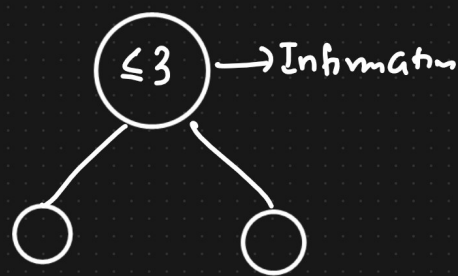


Information gain

Information gain

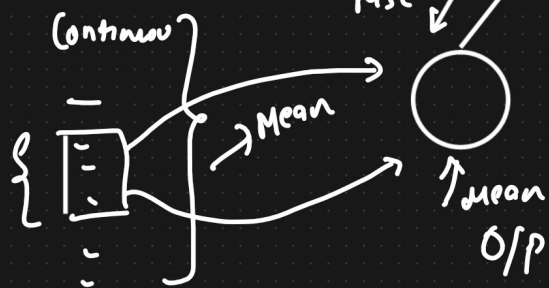


Best Info gain



Decision Tree Regressor

f_1 f_2 o/p



MSE

f_1 → Mean

$\boxed{\text{MSE or MAE}}$

$$\frac{1}{2m} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Hyperparameters

Decision → Overfitting

- { ① Post Pruning }
- { ② Pre Pruning }



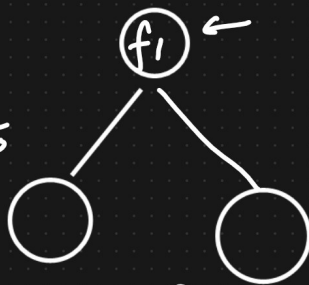
Overfitting

Decision Tree Regressor

f_1
 C_1
 C_1
 $-$
 $-$
 $-$

0/7
 20
 24
 26
 28
 30

MSE



MSE = 37

{

MSE

}

post pruning

pre pruning

Hyperparameter

GridSearchCV

max_depth, max_leaf

