# Decision Tree

① Decision Tree Classifier  [classification]

② Decision Tree Regressor  [Regression]

## ID3



Root Node ←

## CART



## Decision Tree Classifier

Two techniques

① ID3  [Iterative Dichotomisu 3]

② CART  [classification And Regression Tree]
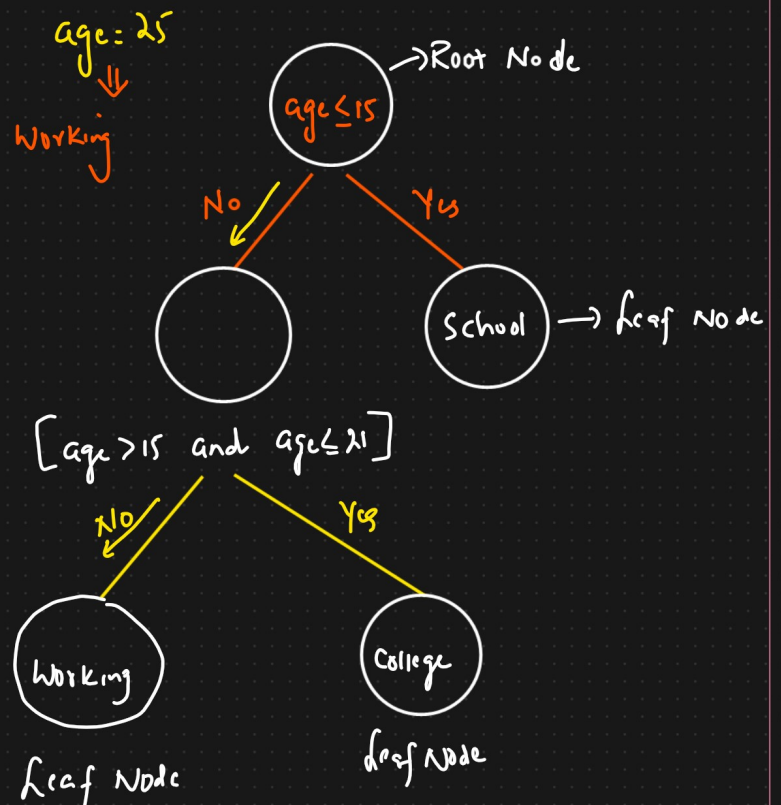
Multinuhd  if else clause

age = 14

age: 25
⇓
Working

```
if (age ≤ 15):
    Print ("School")

elif (age>15 and age ≤21):
    Print (" College")

else:
    Print (" Working").
```



age ≤ 15 → Root Node

No / Yes

School → Leaf Node

[age >15 and age ≤ 21]

No / Yes

Working
Leaf Node

College
leaf Node

Dataset → Problem statement

| Day | Outlook | Temperature | Humidity | Wind | Play Tennis |
|-----|---------|-------------|----------|------|-------------|
| 1 | Sunny . | Hot | High | Weak | No |
| 2 | Sunny | Hot | High | Strong | No |
| 3 | Overcast . | Hot | High | Weak | Yes |
| 4 | Rain . | Mild | High | Weak | Yes |
| 5 | Rain | Cool | Normal | Weak | Yes |
| 6 | Rain | Cool | Normal | Strong | No |
| 7 | Overcast | Cool | Normal | Strong | Yes |
| 8 | Sunny | Mild | High | Weak | No . |
| 9 | Sunny | Cool | Normal | Weak | Yes |
| 10 | Rain | Mild | Normal | Weak | Yes |
| 11 | Sunny | Mild | Normal | Strong | Yes |
| 12 | Overcast | Mild | High | Strong | Yes |
| 13 | Overcast | Hot | Normal | Weak | Yes |
| 14 | Rain | Mild | High | Strong | No |

Impure split

Outlook  9Yes/5No

2 Yes | 3No          3Yes/2No

4Yes| 0No

Sunny      Overcast      Rain

Leaf Node

① Purity Split check — Pure Split or Impure Split

→ Entropy
→ Gini Impurity.

} Measure of Purity

② What feature you need to select to start the split — Information Gain.

① Purity Check
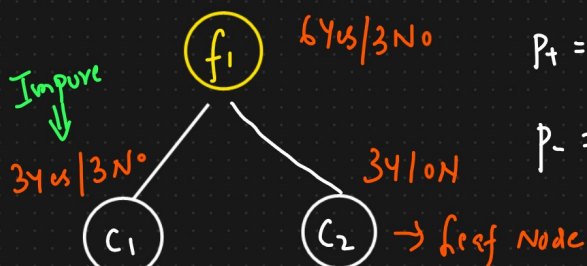
Binary classification

① Entropy

② Gini Impurity

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

$GI = 1 - \sum\limits_{i=1}^{n} (p)^2$

$P_+ = $ probability of positive category

$P_- = $      "          " negative category

Impure

f1      6Yes/3No

3Yes | 3No°           3Y | 0N

C1              C2  → Leaf Node

$P_+ = \dfrac{3}{6} = \dfrac{1}{2}$

$P_- = \dfrac{3}{6} = \dfrac{1}{2}$

$H(C_1) = -P_+ \log_2 P_+ - P_- \log_2 P_-$

$$= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log \frac{3}{6}.$$

$\boxed{H(c_1) = 1} \Rightarrow$ Impure Split

$$H(c_2) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$

$H(c_2) = 0 \Rightarrow$ Pure Split



② Gini Impurity

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

$$G.I(c_1) = 1 - \left[ (p_+)^2 + (p_-)^2 \right]$$

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

$$= 1 - \frac{1}{2} = 0.5 \Rightarrow \text{Impure Split}$$

3 Yes / 0 No

$$G.I(c_2) = 1 - \left[ \left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2 \right]$$

$$= 1 - 1 = 0 \Rightarrow \text{Pure Split}$$

Multiclass Classification Problem $\div$ 3 Categories In O/P

$$H(s) = -P_{c_1} \log_2 P_{c_1} - P_{c_2} \log_2 P_{c_2} - P_{c_3} \log_2 P_{c_3}$$

$$G.I = 1 - \left[ (P_{c_1})^2 + (P_{c_2})^2 + (P_{c_3})^2 \right]$$
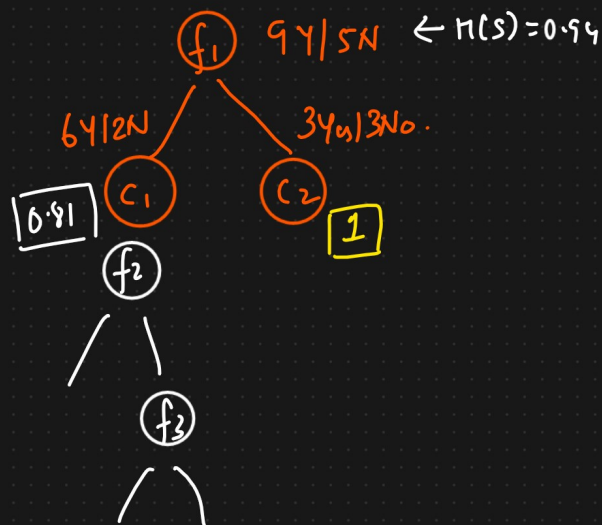
② <u>Information Gain</u> $\rightarrow$ Which feature to Select to Start the split?

$$\text{Gain}(S, f_i) = H(s) - \sum_{v \in val} \frac{|S_v|}{|S|} H(S_v) \rightarrow \text{Entropy of Categories}$$

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

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \left(\frac{5}{14}\right)$$
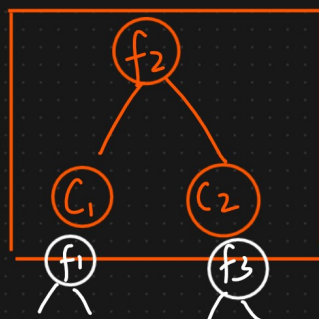
$$\approx 0.94$$

$$H(c_1) = -\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \frac{2}{8} \approx \boxed{0.81}$$

$$H(c_2) = \underline{\underline{1}}$$

f1  f2  f3  O/P



(f1) 9Y/5N  ← H(s) = 0.94

6Y|2N       3Yes|3No.

0.81 (c1)   (c2) 1

(f2)

$$\text{Gain}(S, f_1) = H(s) - \sum_{v \in val} \frac{|Sv|}{|S|} H(Sv) \rightarrow \text{Entropy of Categories}$$

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

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

⇒ Information Gain = 0.051

$$\boxed{\text{Gain}(S, f_2) = 0.051} > \boxed{\text{Gain}(S, f_1) = 0.049}$$

We need to Splitting by using f2 features

Entropy vs Gini Impurity

When dataset is small → Entropy   [log formula]

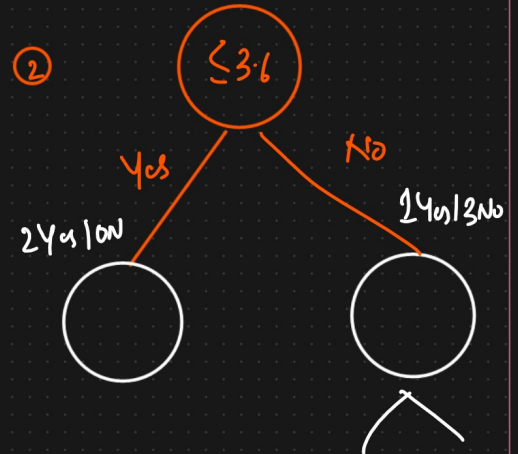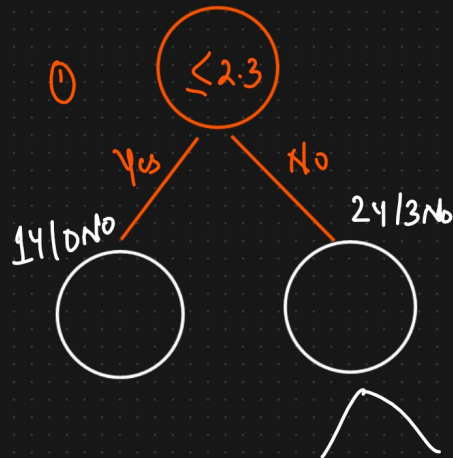When dataset is huge → Gini Impurity [Simple Maths]

⑧ What if my feature is continuous.

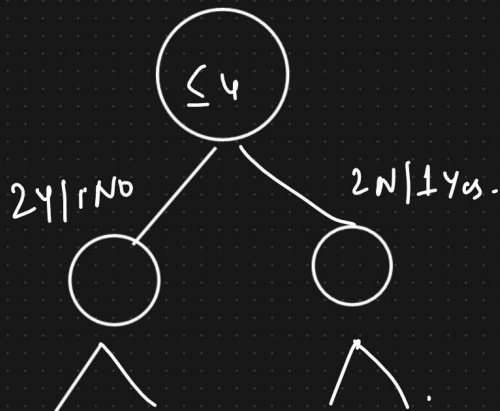| $f_1$ | %p |
|-------|-----|
| →2.3 | Yes |
| →3.6 | Yes |
| 4 | No |
| 5.2 | No |
| 67 | Yes |
| 7.8 | No |

① Sort the feature $f_1$

① Threshold = 2.3

⓪ $\leq 2.3$
 Yes → 1Y/0No
 No → 2Y/3No

② $\leq 3.6$
 Yes → 2Yes/0N
 No → 2Yes/3No

③ Threshold = 4

$\leq 4$
 2Y/1No
 2N/1Yes.

Time Complexity ↑↑

DATASET ↑↑.