# Decision Tree Classifier

Decision Tree Classifier
- → ID3
- → CART ✓

a) Entropy and Gini Index → Purity Split
b) Information Gain → features to select for

DT construction

```
age = 14

if (age ≤ 15):
    Print ("The person is in School)

elif (age >15 and age ≤21):
    Print ("The person may be College)

else:
    Print ("The person has passed)
```
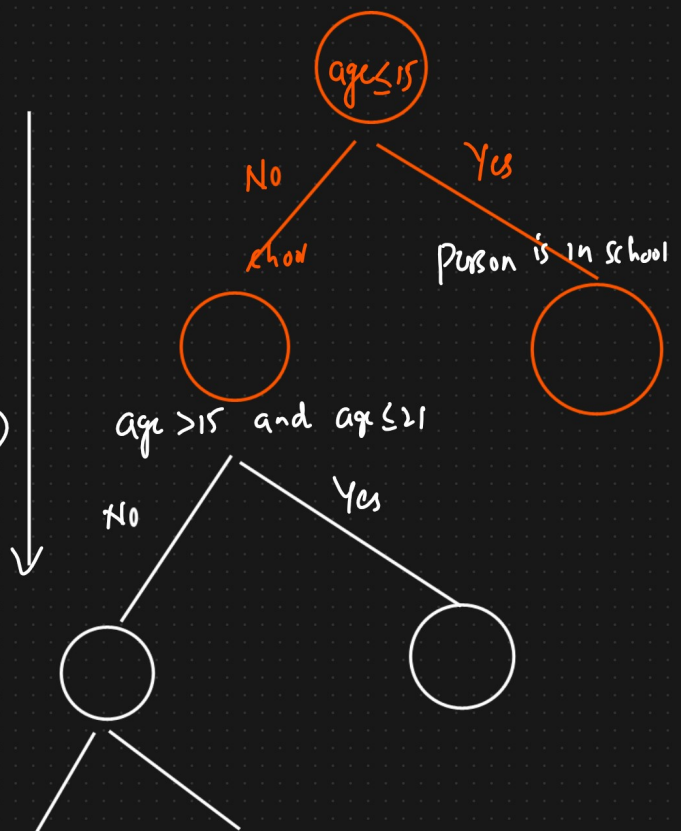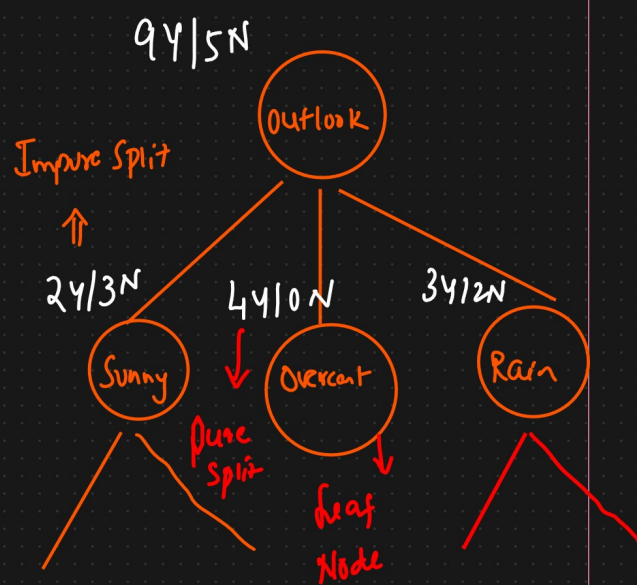
age ≤ 15

No        Yes

know        Person is in school

age >15 and age ≤ 21

No        Yes

Datant

Binary Clanification

| 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 |

Outlook

Impure Split

↑

2Y/3N    4Y|0N    3Y|2N

Sunny    Overcast    Rain

Pure Split

Leaf Node

① Punty → Pure or Impure Split

└→ Entropy
└→ Gini Impurity

② What feahive you necd Sclect for

Splitting → Information Gain }

1
0

{Binary clanification}
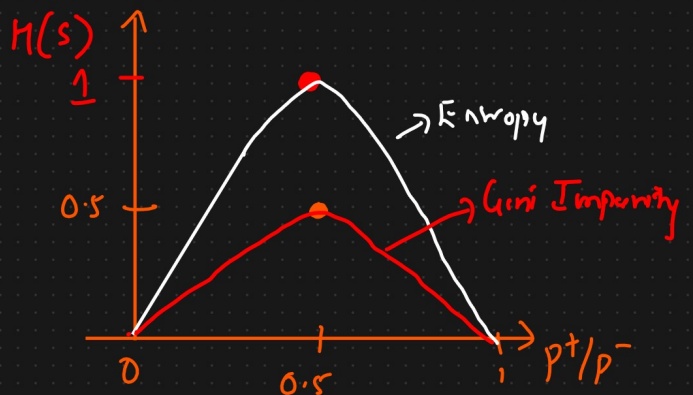
1) Entropy

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

6Yes/3No

f₁

Impure Split
↑
3Y/3N

Pure Split
↓
3Y|0N

C₁    C₂

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

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

② Gini Impurity

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

H(S)
1

→ Entropy

0.5

→ Gini Impurity

0    0.5    1    P⁺/P⁻

$$= 1 \Rightarrow \text{Impure Split}$$

$$H(C_2) = -\frac{3}{3} \log_2 \frac{3}{3} - 0 \log_2 0$$

$$= -1 \log_2 1 \Rightarrow 0 \Rightarrow \text{Pure Split}$$

② **Gini Impurity**

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

$$= 1 - \left((p_+)^2 + (p_-)^2\right)$$

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

$$= 0.5 \Rightarrow \text{Impure Split}$$
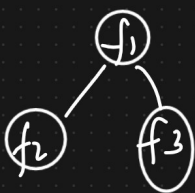
3Y|0N

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

$$= 1 - 1$$

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

f1    f2    f3

Decision Tree Split



$\Rightarrow$ Information Gain $\}$

**Information Gain**

$\nearrow$ Entropy of the root node

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

$\nearrow$ Root Node

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

f1  f2  f3  O/P

14 = 9Y/5N

8 = 6Y|2N

3Y/3N = 6



Impure split

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

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

$$\approx 0.94$$

$$\boxed{H(c_1) = 0.81}$$

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

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

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



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

Information is Basically calculated.
(Gain)

## Entropy Vs Gini Impurity

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

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

$$\boxed{O/P = 3 \text{ categories}}$$

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

Whenever dataset is small $\rightarrow$ Entropy
large $\rightarrow$ Gini Impurity $\Big\}$