# Recap of Previous Lecture

**Topic** Ridge

**Topic** lasso

**Topic** Regression

**Topic**

**Topic**

# Topics to be Covered

**Topic** Decision Tree

**Topic** Decision Tree Regressor

**Topic** Decision Tree Classifier

**Topic**

**Topic**

# Decision Tree Classifier And Regressor
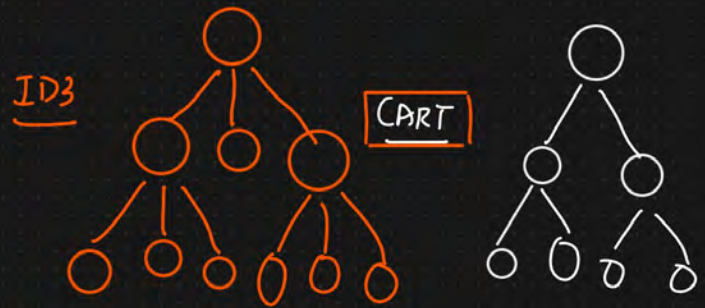
## Agenda

① Decision Tree Classifier [classification]

② Decision Tree Regressor [Regression]



$\underline{ID3}$      $\boxed{CART}$

## ① Decision Tree Classifier

Two Types

① ID3 [Iterative Dichotomiser 3]

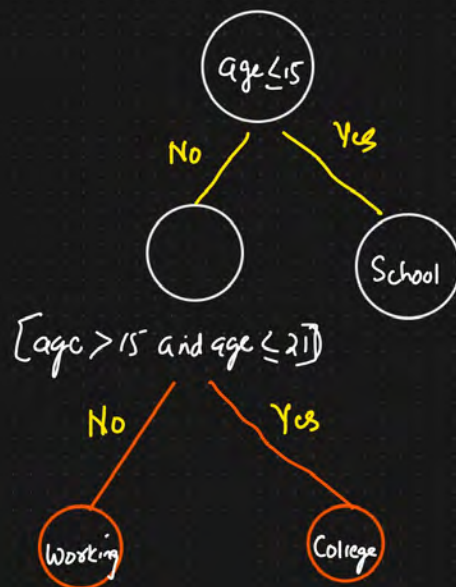② CART [Classification And Regression Trees]

age = 14

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

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

else (age > 21):
    Print ("Working")
```
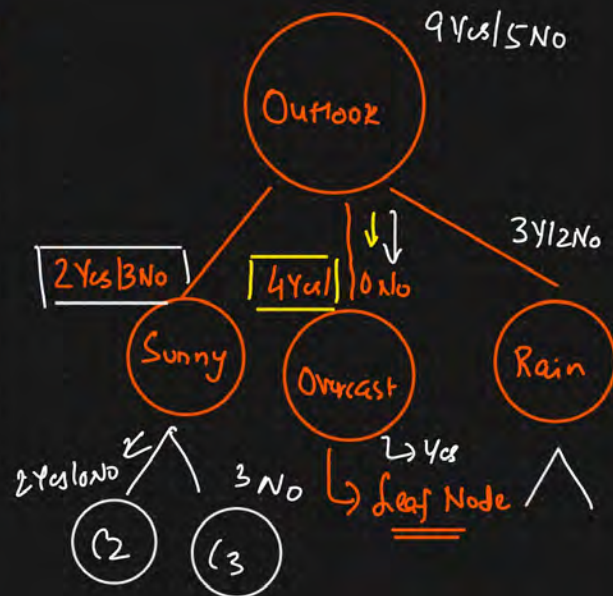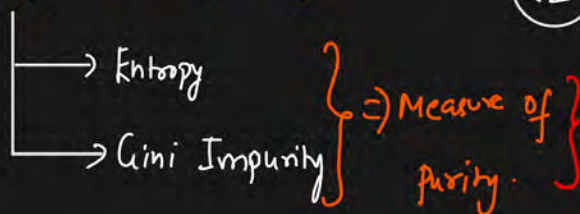
# Dataset : Predict Play Tennis or Not

| 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 |

↑ O/P

**9 Yes / 5 No**

Outlook

3 Y / 2 No

2 Yes | 3 No    4 Yes | 0 No

Sunny    Overcast    Rain

2 Yes | 0 No    3 No    ↳ Yes

C2    C3    ↳ Leaf Node

① **Purity 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}.
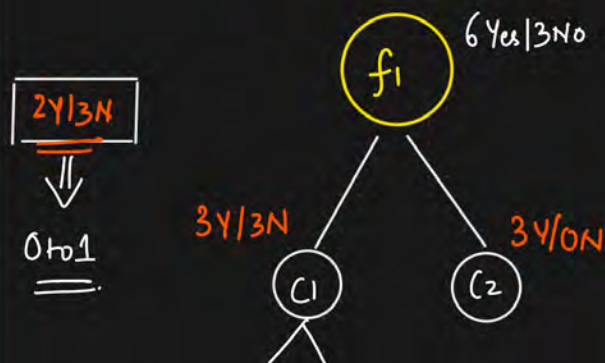
**2 Yes / 2 No**

$$P_t = \frac{2}{4} = 0.5$$

① **Entropy**

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

$P_+$ = probability of +ve Category

$P_-$ = probability of negative category

ii) **Gini Impurity**

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

6 Yes | 3 No

f1

2Y | 3N
⟱
0 to 1

3Y | 3N          3Y | 0N

C1          C2

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

H(s)

1

0.5

⟶ Entropy

⟶ G.I

0.5    1    $P_+, P_-$

**Entropy**
[0 - 1]

**Gini**
[0 - 0.5].

$$= -\frac{1}{2}\log_2(\tfrac{1}{2}) - \frac{1}{2}\log_2(\tfrac{1}{2})$$

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

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

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

$$= 0 \Rightarrow \text{Pure Split}$$
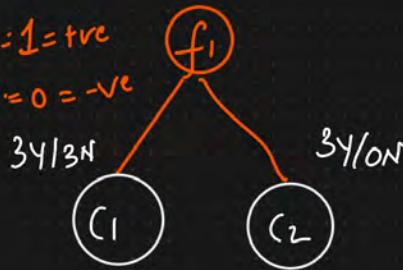
② **Gini Impurity**

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

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

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

$$= 1 - \left[\tfrac{1}{4} + \tfrac{1}{4}\right] = 0.5$$

$$\Downarrow$$

Impure Split.

$Yu = 1 = \text{tve}$
$No = 0 = -ve$
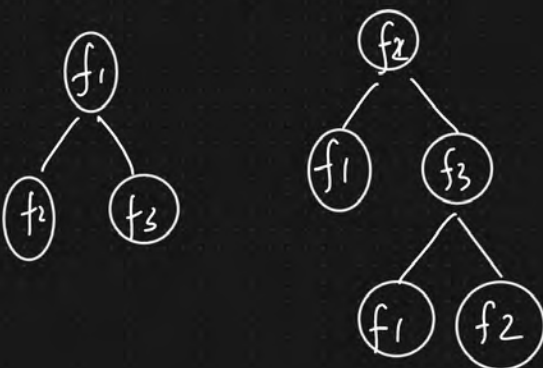
$f_1$

$3Y/3N$  →  $C_1$
$3Y/0N$  →  $C_2$

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

$$= 1 - 1$$

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

② What feature you need to select to start the split? → **Information Gain**

$f_1$   $f_2$   $f_3$       O/p

$f_1$
$f_2$   $f_3$

$f_2$
$f_1$   $f_3$
$f_1$   $f_2$

# Information Gain

Gain(S, f·) → Entropy of the root Node

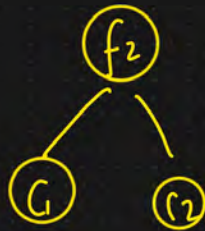$$Gain(S, f \cdot) = H(s) - \sum_{v \in val} \frac{|Sv|}{|S|} H(Sv)$$

f1    f2    f3    O/p

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

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

$$\boxed{H(s) \approx 0.94}$$

f1  9 Yes/5No

6 Y/2No        3 Yes/3No

C1              C2

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

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

$$H(C_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \approx \boxed{0.81}$$

$$H(C_2) = \underline{1}_{//}$$

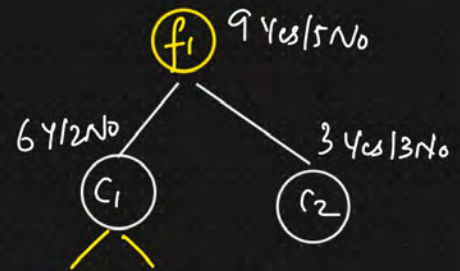f2

C1      C2

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

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

Higher the IG ↑

## Entropy Vs Gini Impurity

Whenever   dataset is Small → Entropy
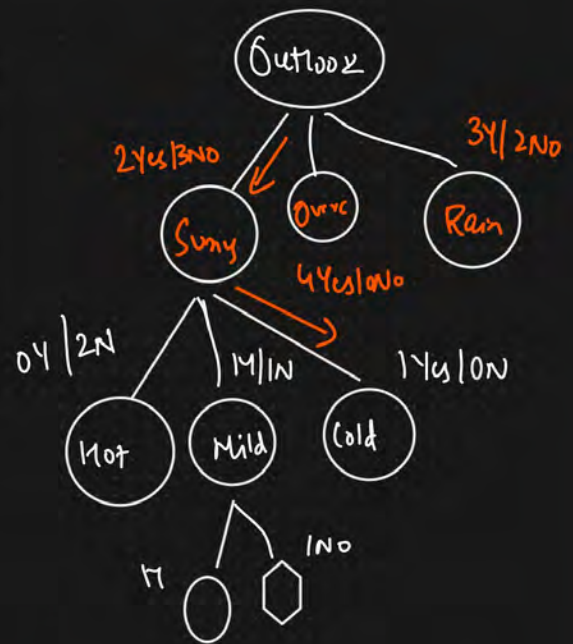
dataset is huge → Gini Impurity }

| 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
2Yes/3No
Suny
Ovrc
Rain
3Y/2No
4Yeslovo
0Y/2N
M/IN
1Yes/0N
Hot
Mild
Cold
M
1NO

→ Outlook    Temperature         [Play]

Sunny        Cold                [Yes]

---

Decision Tree Post Prunning And Pre prunning [Post Prunning And Pre prunning].

f1 →Level 0    Overfitting

Post prunning  {Reduce Overfitting}

⇓              ⇓

Small datant   Yes

→ Level 1
9Y/3NoY
→ Level 2
→ Level 3
→ Level 3.
9Yeslovo    0Yes/3No.

Select the best parameters
⇑

Pre Prunning → DATASET → Suitable parameter ⇒ Hyperparameter Tunnng.

⇓

Max depth, split,

Criterion

Gini       Entropy

max depth = 4          $\Rightarrow$ Accuracy $\Rightarrow$ Hyperparameter
                                    Tuning.

$\longrightarrow$ max depth = 3

④ Decision Tree Regressor

DATASET

$\downarrow$ OIP

| Exp | Career Gap | Salary |
|-----|-----------|--------|
| $\rightarrow$ 2 | Yes | 40K |
| 2.5 | Yes | 42K |
| $\rightarrow$ 3 | No | 52K |
| 4 | No | 60K |
| 4.5 | Yes | 56K |

$y = 50K$

[40K, 42K, 52K, 60K, 56K]

$\leq 2$  $\leftarrow 0$

$\downarrow$ Split {Variance Reduction}.

$\leq 2.5$ $\Rightarrow$ Var (0.04).

yes

| 40K |   | 42K, 52K, 60K, 56K |

$<$

Variance Reduction

Mean of y

$\underset{\text{(Root)}}{\text{Variance}} = \frac{1}{n} \sum_{i=1}^{n} (y - \bar{y})^2 \; [MSE].$

$= \frac{1}{5} \left[ (10)^2 + (8)^2 + (2)^2 + 10^2 + 6^2 \right]$

Variance Reduction :

$Var(Root) - \sum W_i \, Var(child)$

$60.8 - \left[ \frac{1}{5} * 100 + \frac{4}{5} * 51 \right]$

Variance$\left(\substack{\text{Root}}\right) = 0 \ 0 \cdot 8$

Var$(\text{Right}) = \underline{51} \ = \ 0$.

Var$(\text{c c t}) = \frac{1}{1} \left[ (40 - 50)^2 \right]$

Var$(\text{Left}) = 100$