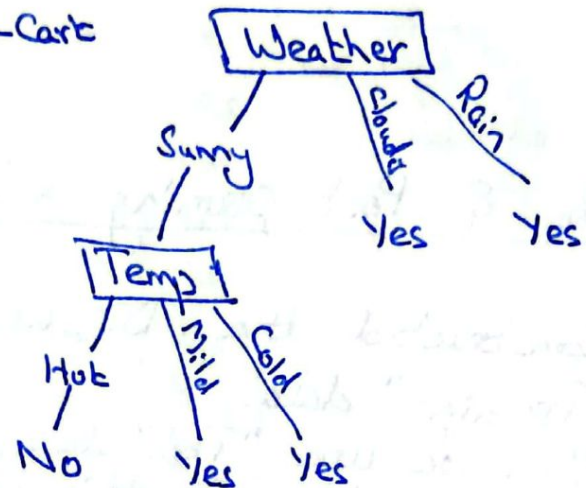→ Decision Tree

1) Can be used for both classification & regression, however mainly used for classification.

2) Some important points to consider:
   a) Tree Structure ⎤⎡— ID3 —
   b) Decision Nodes ⎦⎣— Cart
   c) Leaf Nodes
   d) Split (pure, impure)
   e) Entropy
   f) Information gain
   g) Pruning (Pre, post)

Weather
- Sunny → Temp
- cloudy → Yes
- Rain → Yes

Temp
- Hot → No
- Mild → Yes
- Cold → Yes

Questions: To check the purity of split we use
↳ Entropy
↳ Gini Index (for large datasets. Have simple math so computationaly efficient)
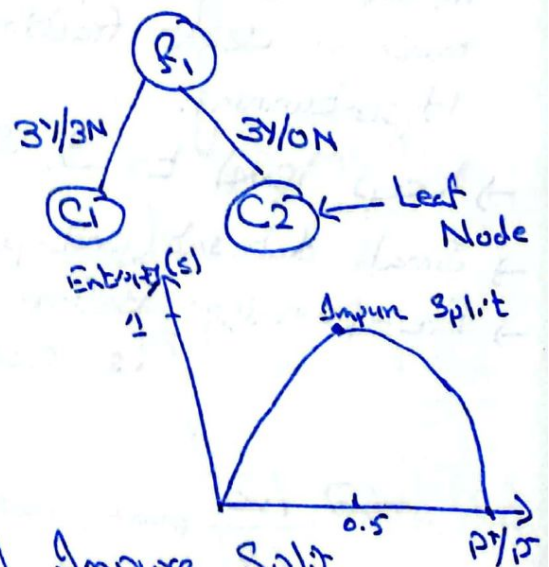: Which feature to chose for splitting.
↳ Information Gain.

Entropy

Assuming Binary (Yes/No)

$$Entropy(s) = P^+ \log_2 P^+ - P^- \log_2 P^-$$

$$S_{c_2} = \frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$
$$= -1 \log_2 1 = \boxed{0} \quad Pure\ Split$$

$$S_{c_1} = \frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$
$$= \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = \boxed{1} \quad Impure\ Split$$
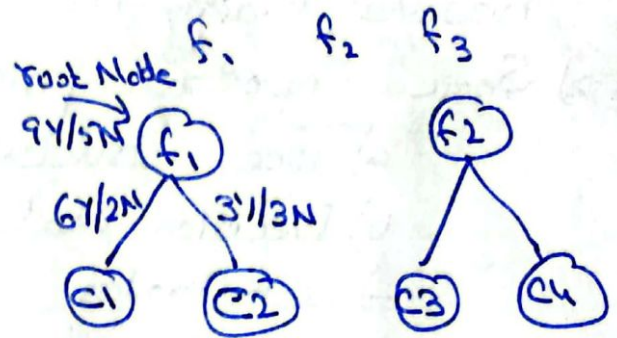
$P_1$
- 3Y/3N → C1
- 3Y/0N → C2 ← Leaf Node

Entropy(s)

Entropy will always be between 0~1.

→ How to split, which feature to be the node.

__Information Gain__

$$\text{Gain}(S, f_i) = \text{Entropy}(S) - \sum_{v \in} \frac{|S_v|}{|S|} \text{Entropy}(S)$$

$f_1 \quad f_2 \quad f_3$

root Node

$9Y/5N \xrightarrow{} (f_1)$

$6Y/2N \quad 3I/3N$

$(C1) \quad (C2)$

$(f_2)$

$(C3) \quad (C4)$

→ __Pre & Post pruning DT__

→ Constructed the DT using "Training" data.

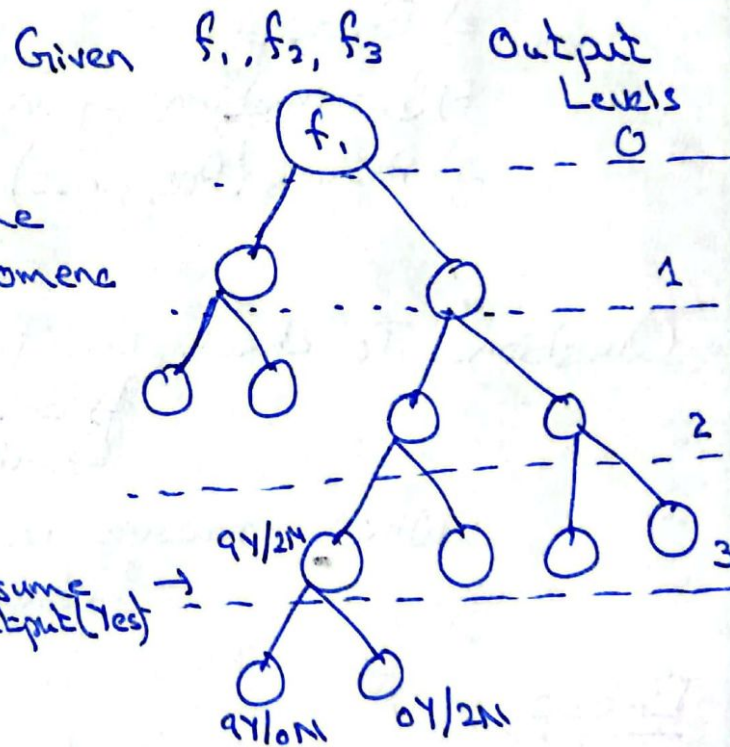→ When we use "Test" data, the accuracy is lower & this phenomena is known as overfitting.

→ To resolve overfitting
  └→ Post Pruning
  └→ Pre Pruning

→ Hyperparameters (Max depth, minimum samples split/leaf, minimum weight fraction).& A Hypertuning.

→ Keep level to 3. So reduced the overfitting issue

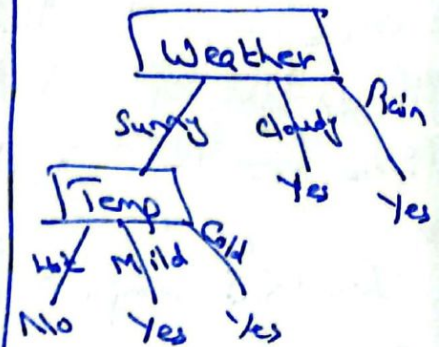→ Small dataset (post-pruning), Big datasets (pre-pruning)

→ Pre-pruning: Before constructing DT, at what depth is accuracy highest.

Given $f_1, f_2, f_3$    Output Levels

$(f_1)$    0

   1

   2

$9Y/2N$    3

Assume
Output (Test) →

$9Y/0N \quad 0Y/2N$

# Numerical Example:

| Day | Outlook | Temp | 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 | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

Values(Outlook) = Sunny, Overcast, Rain

$S[9+, 5-]$  Entropy$(s) = \frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$

$= \boxed{0.94}$

$S_{sunny}[2+, 3-]$  Entropy$(S_{sunny}) = \frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = \boxed{0.971}$

$S_{overcast}[4+, 0-]$  Entropy$(S_{overcast}) = \frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} = \boxed{0}$

$S_{Rain}[3+, 2-]$  Entropy$(S_{Rain}) = \frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = \boxed{0.97}$

$Gain(S, Outlook) = Entropy(s) - \sum\limits_{S \in (overcast, sunny, Rain)} \frac{|S_v|}{|S|} Entropy(S_v)$

$= 0.94 - \frac{5}{14}(0.971) - \frac{4}{14}(0) - \frac{5}{14}(0.971)$

$= 0.2464$

## Attribute: Temp

Values (Temp) = Hot, Mild, Cool

$S = [9+, 5-]$     Entropy $(S) = 0.94$

$S_{Hot} [2+, 2-]$     Entropy $(S_{Hot}) = \frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1$

$S_{Mild} [4+, 2-]$     Entropy $(S_{Mild}) = \frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.9183$

$S_{Cool} [3+, 1-]$     Entropy $(S_{Cool}) = \frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.8113$

$$Gain (S, Temp) = Entropy (S) - \sum_{v \in (Hot, Mild, Cool)} \frac{|S_v|}{|S|} Entropy (S_v)$$

$$= 0.94 - \frac{4}{14}(1) - \frac{6}{14}(0.9183) - \frac{4}{14}(0.8113)$$

$$= 0.0289$$

## Attribute: Humidity

Value (Humidity) = High, Normal

$S = [9+, 5-]$     Entropy $(S) = 0.94$

$S_{High} [3+, 4-]$     Entropy $(S_{High}) = 0.9852$

$S_{Normal} [6+, 1-]$     Entropy $(S_{Normal}) = 0.5916$

$$Gain (S, Humidity) = Entropy (S) - \sum_{v \in (High, Normal)} \frac{|S_v|}{|S|} Entropy (S_v)$$

$$= 0.94 - \frac{7}{14}(0.9852) - \frac{7}{14}(0.5916)$$

$$= 0.1516$$

## Attribute: Wind

Values (Wind) = Strong, Weak

$S = [9+, 5-]$     Entropy $(S) = 0.94$

$S_{strong} [3+, 3-]$     Entropy $(S_{strong}) = 1$

$S_{Weak} [6+, 2-]$     Entropy $(S_{Weak}) = 0.8113$

$$Gain (S, Wind) = 0.94 - \frac{6}{14}(1) - \frac{8}{14}(0.8113)$$

$$= 0.0478$$

So

$Gain(S, Outlook)$ = 0.2464

$Gain(S, Temp)$ = 0.0289

$Gain(S, Humidity)$ = 0.1516

$Gain(S, Wind)$ = 0.0478

$\{D1, D2, ..., D14\}$
$[9+, 5-]$

Outlook

Sunny

Overcast

Rain

$\{D1, D2, D8, D9, D11\}$
$[2+, 3-]$

Humidity

High

Normal

$\{D1, D2, D8\}$
No

$\{D9, D11\}$
Yes

$\{D3, D7, D12, D13\}$
$[4+, 0-]$

Yes

$\{D4, D5, D6, D10, D14\}$
$[3+, 2-]$

Wind

Strong

Weak

$\{D6, D14\}$
No

$\{D4, D5, D10\}$
Yes

| Day | Temp | Humidity | Wind | Play Tennis |
|-----|------|----------|------|-------------|
| D1 | Hot | High | Weak | No |
| D2 | Hot | High | Strong | No |
| D8 | Mild | High | Weak | No |
| D9 | Cool | Normal | Weak | Yes |
| D11 | Mild | Normal | Strong | Yes |

## Attribute: Temp

$Values(Temp)$ = Hot, Mild, Cool

$S_{sunny}$ = $[2+, 3-]$

$Entropy(S_{sunny})$ = 0.97

$S_{Hot} [0+, 2-]$     $Entropy (S_{Hot}) = 0$

$S_{mild} [1+, 1-]$     $Entropy (S_{mild}) = 1$

$S_{cool} [1+, 0-]$     $Entropy (S_{cool}) = 0$

$Gain (S_{sunny}, Temp) = Entropy (S) - \sum_{v \in (H,M,C)} \frac{|S_v|}{|S|} Entropy (S)$

$= 0.97 - \frac{2}{5}(0) - \frac{2}{5}(1) - \frac{1}{5}(0)$

$= \boxed{0.570}$

## Attribute: Humidity

$Values (Humidity) = High, Normal$

$S_{sunny} = [2+, 3-]$     $Entropy (S) = 0.97$

$S_{High} [0+, 3-]$     $Entropy (S_{High}) = 0$

$S_{Normal} [2+, 0-]$     $Entropy (S_{Normal}) = 0$

$Gain (S_{sunny}, S_{Humidity}) = Entropy (S) - \frac{3}{5} Entropy (S_{sunny}) -$

$\frac{2}{5} Entropy (S_{Normal})$

$= 0.97 - \frac{3}{5}(0) - \frac{2}{5}(0)$

$= 0.97$

## Attribute: Wind

$Value (Wind) = Strong, Wind$

$S_{sunny} = 0.97$

$S_{strong} [1+, 1-]$     $Entropy (S_{strong}) = 1$

$S_{weak} [1+, 2-]$     $Entropy (S_{weak}) = \frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$

$= 0.9183$

$Gain (S_{sunny}, Wind) = 0.0192$

For Rain

| Day | Temp | Humidity | Wind | Play Tennis |
|-----|------|----------|------|-------------|
| D4 | Mild | High | Weak | Yes |
| D5 | Cool | Normal | Weak | Yes |
| D6 | Cool | Normal | Strong | No |
| D10 | Mild | Normal | Weak | Yes |
| D14 | Mild | High | Strong | No |

Attribute: Temp

Value(Temp) = Hot, Mild, Cool

$S_{Rain} = [3+, 2-]$

$S_{Hot} \leftarrow [0+, 0-]$

$S_{Mild} \leftarrow [2+, 1-]$

$S_{Cool} [1+, 1-]$

Entropy $(S_{sunny}) = 0.97$

Entropy $(S_{Hot}) = 0$

Entropy $(S_{Mild}) = 0.9183$

Entropy $(S_{Cool}) = 1$

Gain $(S_{Rain}, Temp) = 0.0192$

Attribute: Humidity

$S_{Rain} [3+, 2-]$  Entropy $(S_{sunny}) = 0.97$

$S_{High} [1+, 1-]$  Entropy $(S_{High}) = 1$

$S_{Normal} [2+, 1-]$  Entropy $(S_{Normal}) = 0.9183$

Gain $(S_{Rain}, Humidity) = 0.0192$

Attribute: Wind

$S_{Rain} [3+, 2-]$

$S_{Strong} [0+, 2-]$

$S_{Weak} [3+, 0-]$

Entropy $(S_{sunny}) = 0.97$

Entropy $(S_{Strong}) = 0$

Entropy $(S_{Weak}) = 0$

Gain $(S_{Rain}, Wind) = 0.97$