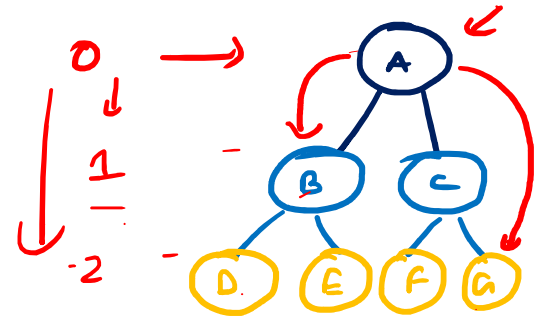→ Decision Trees → Hierarchical Tree

→ Supervised Learning , F + T

→ Classification , Regression

Root node (A) → The starting of the tree.

Splitting → Process of dividing a node into subnodes

Decision Node → Any node that is participating in a condition or has child branches (A, B, C)

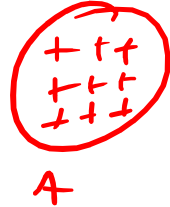Parent Node → Any node that has a child. (A, B, C)

Child Node → Any Subnode that has a parent (B, C, D, E, F, G)

leaf Node → Any node that is NOT a decision node or not has a child or the tree ends there (D, E, F, G)

Depth → # of generations . Root node = 0 depth

Siblings → Nodes at Same depth

0 →

1

-2

A
B   C
D  E  F  G

→ Pure Node

Any Node that has Rows of the
    Same class

```
  + + +
  + + +
  + + +
```
A

→ Impure Node

Any Node that has rows of
    multiple class

```
  + - +
  + - -
```
A

→ | Algorithm | Splitting criteria |

CART →Regression Algo
C4.5 ⎤
I D3  ⎥ → classification
CHAID ⎦

⇓

Decision Trees

Gini
Entropy
Entropy
Chi-Square

# Splitting Criterias

Features | Target

| S. No. | Name | Performance | Height | Class | Plays Cricket |
|--------|------|-------------|--------|-------|---------------|
| 1 | A | Above Avg. | > 5.5 | 1X | 0 |
| 2 | B | Below Avg | ≤ 5.5 | X | 0 |
| 3 | C | Below Avg | ≤ 5.5 | X | 1 |
| 4 | D | Above Avg | > 5.5 | X | 0 |
| 5 | E | Below Avg | > 5.5 | 1X | 1 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

20

→ 10/20
→ +ve Class

→ -ve Class
↳ 10/20

Treasure DM
~M

Treasure Ms

Reass

Gook

look

Loon

Toy h

Look

$\rightarrow$



$\leq r.r.gt$     $> r.r.gt$

(8)     (12)

Hieght

BA     AA

(14)     (6)

Performen

$!X$     $X$

(10)     (10)

Class

Now to dicide which one feature to choose as the Best split

$\hookrightarrow$ based on splitting criteria you choose

① Gini Index $\delta$ Gini Impurity (G.I.)

G.I.          Hight          Pred.
                  5             ③

$\downarrow$ J

how much purity
the split will
provide

$\downarrow$
how much Impurit
is being left after
split.

Gini Impurity $= 1 - $ Gini

Prob of the Class in that node
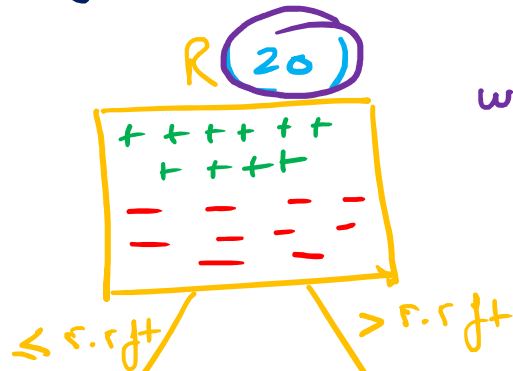
$$\text{Gini} = P_1^2 + P_2^2 + P_3^2 + P_4^2 + \cdots + P_n^2 \quad, \quad n = \# \text{ of Categories in the target}$$

of a
Node

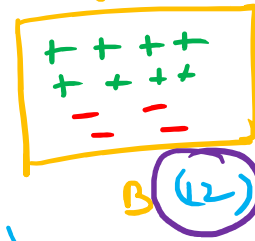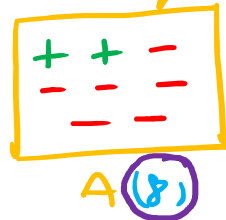Gini of a split $=$ Weighted Avg $\left( \text{Node A}, \text{Node B} \right) = \dfrac{W_A G_A + W_B G_B}{2}$

→ Split on Hieght based Gini Impurity

$$G.I._{Hieght} = W_A \, GI_A + W_B \, GI_B$$

&
weighted GI. $\Rightarrow \dfrac{8}{20} * 0.375 + \dfrac{12}{20} * 0.442$

$$\boxed{G.I._{hieght} \Rightarrow 0.415}$$



R (20)

≤ 5.5 ft    > 5.5 ft

+ve = 2
$P(+ve) = \dfrac{2}{8} = 0.25$
$P(-ve) = 0.75$

A (8)

$P(+ve) = 8/12 = 0.67$
$P(-ve) = 0.33$

B (12)

$$G.I. = 1 - Gini = 1 - (P_1^2 + P_2^2 + P_3^2 + \ldots P_n^2)$$

$$G.I._A = 1 - \left[ P_{(+ve)}^2 + P_{(-ve)}^2 \right]$$

$$= 1 - \left[ (0.25)^2 + (0.75)^2 \right]$$

$$\Rightarrow 0.375$$

$$G.I._B = 1 - \left[ P_{(+ve)}^2 + P_{(-ve)}^2 \right]$$

$$= 1 - \left[ 0.67^2 + 0.33^2 \right]$$

$$\Rightarrow 0.442$$

→ Split on Class based on Gini Impurity



$$G.I._A = 1 - [0.8^2 + 0.2^2]$$
$$\Rightarrow 0.32$$

R (20)

+ + + + + + +
+ + +
− − − − −
− − − − −

1X

P(+ve) = 0.8
P(−ve) = 0.2

A(10)

+ + + + −
+ + + + −

X

+ + − − −
− − − − −

B (10)

P(+ve) = 0.2
P(−ve) = 0.8

$$G.I._B = 1 - [0.2^2 + 0.8^2]$$
$$\Rightarrow 0.32$$

$$G.I._{Class} \Rightarrow W_A \, GI_A + W_B \, GI_B$$

$$\Rightarrow \frac{10}{20} * 0.32 + \frac{10}{20} * 0.32 = 0.32 \left(\frac{1}{2} + \frac{1}{2}\right)$$

$$G.I._{Class} \Rightarrow 0.32$$

| Feature | G.I |
|---|---|
| Height | 0.415 |
| Class | 0.32 |
| P.tf | 0.16 |