# Decision Trees
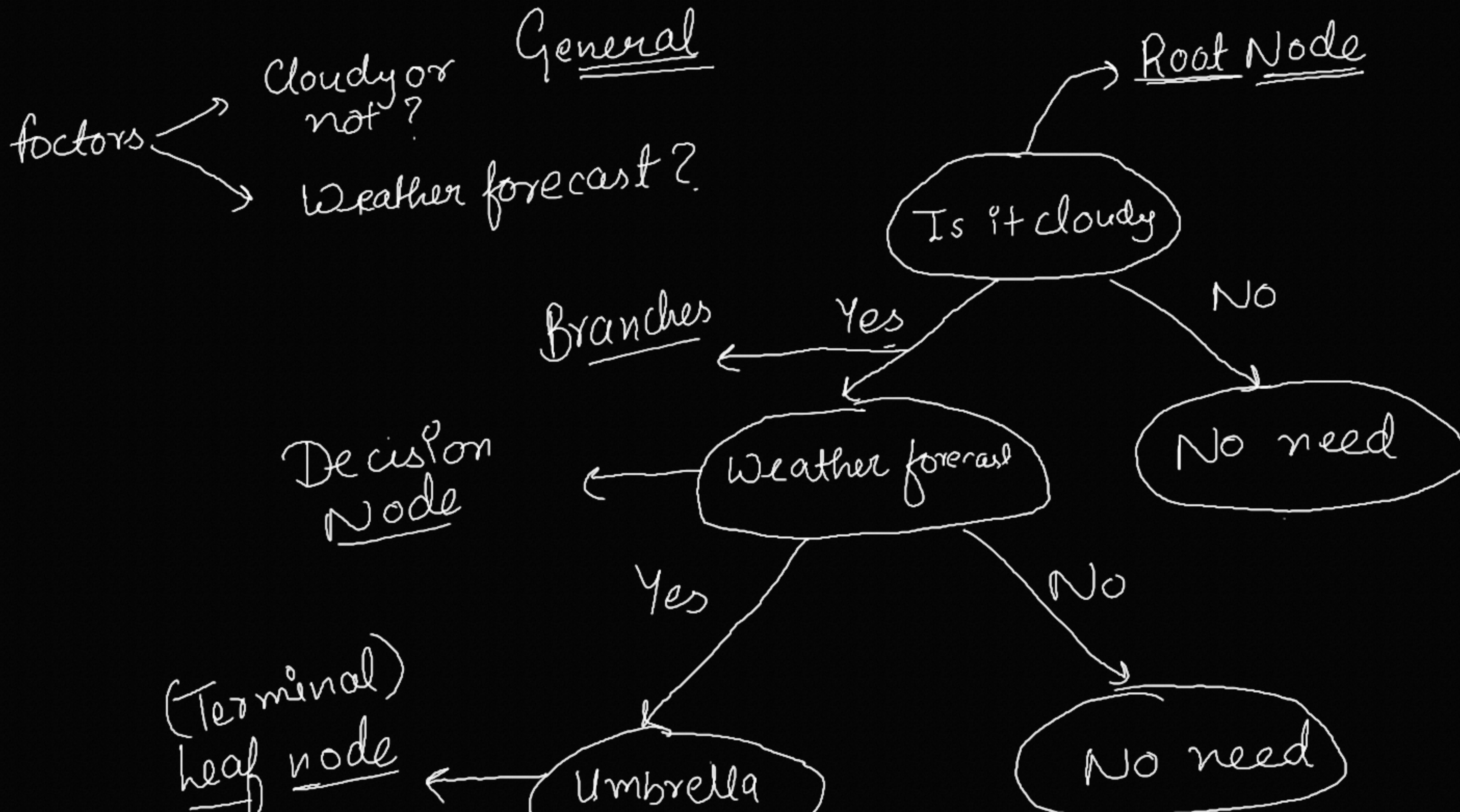
* Classification & Regression
* Core of Ensemble Learning
* Looks like <u>Nested - If else</u>
  Structure is like a <u>tree</u>
* <u>History</u>

1986
Ross Quinlan $\longrightarrow$
ID3

1986
CART
Leo, Jerome, Richard, Charles
Classification $\longrightarrow$ Gini, Regression $\longrightarrow$ MSE

1990's
(<u>Prunning</u>)

factors → Cloudy or not?

**General**

Weather forecast?

**Root Node**

Is it cloudy

**Branches** ← Yes / No

**Decision Node** ← Weather forecast

No need

Yes / No

**(Terminal) Leaf node** ← Umbrella

No need

+ Criteria to choose Nodes

Classification

Gini Impurity

Information Gain

Lower value is best value

Greater value is best value

*Entropy (concept)

Regression

MSE          MAE

# Mathematical (classification)

| Person | Income ($X_1$) | Age ($X_2$) | Buy House? ($Y$) |
|--------|--------|-----|------------|
| 1 | High | Young | No |
| 2 | High | Old | Yes |
| 3 | Low | Young | No |
| 4 | Low | Old | Yes |

① Calculate Total Gini Impurity

$$\boxed{Gini = 1 - \sum p_i^{\,2}}$$

$P(\text{Buy House} = Yes) = \dfrac{2}{4} = 0.5$

$P(\text{Buy House} = No) = \dfrac{2}{4} = 0.5$

$Gini = 1 - \left( (0.5)^2 + (0.5)^2 \right)$

$= \underline{\underline{0.5}}$  (overall $\underline{gini}$)

# ② Calculate Gini for splits

## (a) Income

(i) Income = "High" $\Rightarrow$ P1 $\longrightarrow$ No

P2 $\longrightarrow$ Yes

$$P(Yes) = \frac{1}{2} = 0.5 \qquad P(No) = \frac{1}{2} = 0.5$$

$$Gini(Income = High) = 1 - (0.5^2 + 0.5^2)$$

$$= \underline{\underline{0.5}}$$

(ii) Income = Low $\Rightarrow$ P3 $\longrightarrow$ No $\Rightarrow$ $P(Yes) = 0.5$

P4 $\longrightarrow$ Yes $\qquad P(No) = 0.5$

$$Gini(Income = Low) = \underline{\underline{0.5}}$$

\* Weighted Gini (Income) $= 0.5 \times \underbrace{\dfrac{2}{4}}_{\text{High}} + 0.5 \times \underbrace{\dfrac{2}{4}}_{\text{Low}}$

$$= \underline{0.5}$$

(b) <u>Age</u>

(i) Age = Young $\Rightarrow$ P1 $\longrightarrow$ No

P3 $\longrightarrow$ No

$P(Yes) = 0$

$P(No) = \dfrac{2}{2} = 1$

$$\text{Gini (Age = Young)} = 1 - [1^2 + 0^2]$$

$$= \underline{\underline{0}}$$

iii) Age = Old $\Rightarrow$ P2 $\longrightarrow$ Yes      P(Yes) = 1

               P4 $\longrightarrow$ Yes      P(No) = 0

$$Gini \, (Age = Old) = 1 - (1^2 + 0^2)$$

$$= \underline{\underline{0}}$$

$$\#\ Weighted \ Gini \ (Age) = 0 \times \frac{2}{4} + 0 \times \frac{2}{4}$$

$$= \underline{\underline{0}}.$$

③ Choose the best split

$Gini\ (Income) = 0.5$
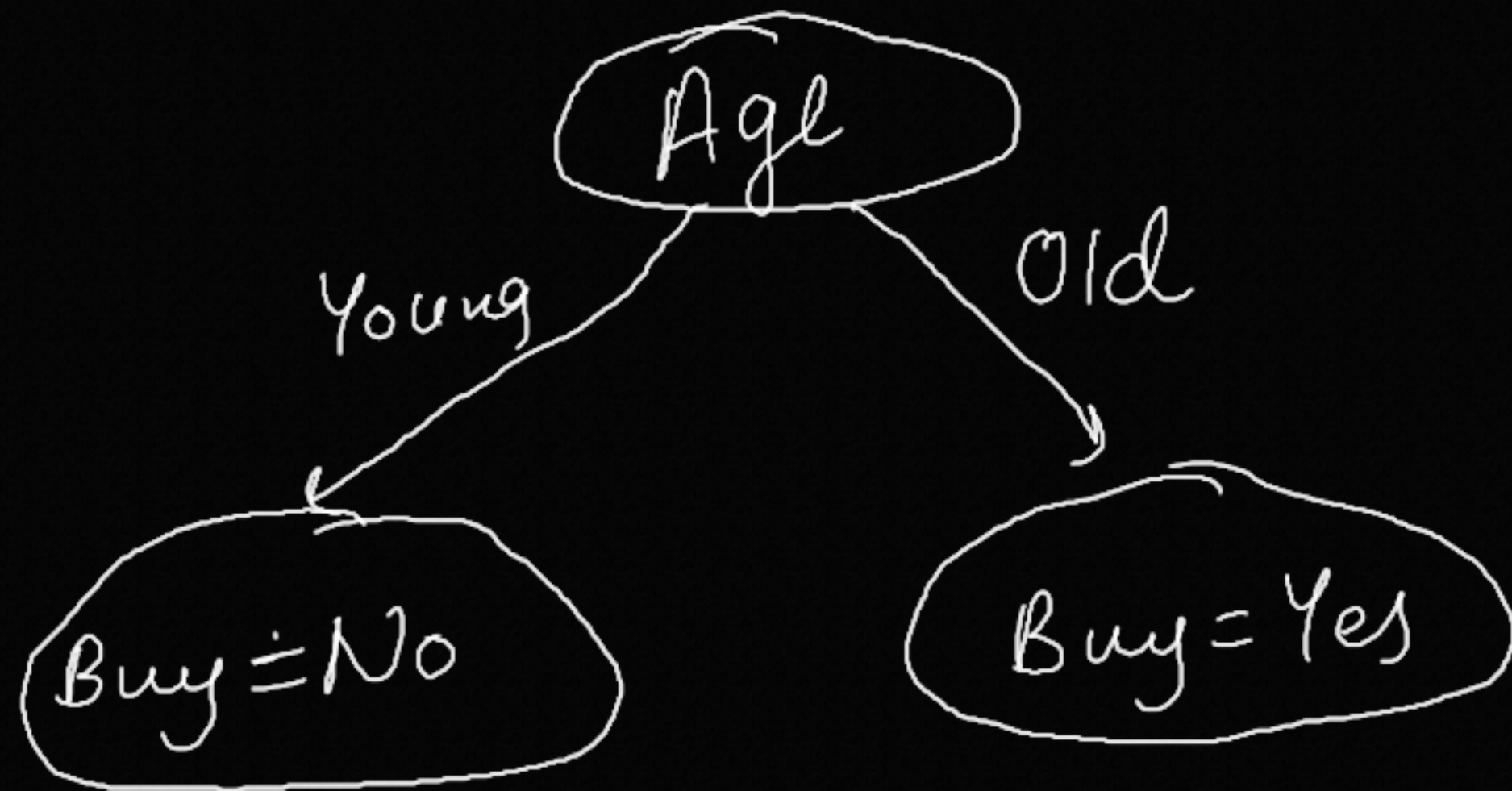
$Gini\ (Age) = 0$ ✓

Lower is best.

When $gini = 0 \longrightarrow$ perfect / pure split

final.

Age

Young — Buy = No

Old → Buy = Yes

(DT will end when you get pure split)

④ Repeat all the steps again
till you get
pure split

(Information Gain)

Age
Y / \ Old
No   Yes

| $X_1$ | $X_2$ | Y |
|------|-------|---|
| High | Young | No |
| Low | Young | No |

Same approach

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| 1 | 2 | 10 |
| 2 | 1 | 20 |
| 3 | 2 | 15 |
| 4 | 3 | 25 |
| 5 | 4 | 30 |

Regression (Maths)

① Calculate Total MSE

Mean $= 20 (\bar{x})$

$$MSE = \frac{1}{n} \Sigma (\bar{x} - x_i)^2$$

$$MSE = \frac{1}{5}\left[ (20-10)^2 + (20-20)^2 + (20-15)^2 + (20-25)^2 + (20-30)^2 \right]$$

$$= \underline{\underline{50}}$$

② Calculate **MSE** _at_ _each feature_

(a) $\underline{X_1}$ $\Rightarrow$ Threshold $= 1.5$

$\underline{X_1 \leq 1.5}$
(Left subset)

$\underline{X_1 > 1.5}$
(Right subset)

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| 1 | 2 | 10 |

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| 2 | 1 | 20 |
| 3 | 2 | 15 |
| 4 | 3 | 25 |
| 5 | 4 | 30 |

Left subset

Mean = 10

MSE = 0

Right subset

Mean = 22.5

MSE = 31.25

$$* \text{ Weighted MSE}(X_1) = \frac{1}{n}\left[\text{no of rows (left)} \times \text{MSE} + \text{no of rows (right)} \times \text{MSE}\right]$$

$$= \frac{1}{5}\left[1 \times 0 + 4 \times 31.25\right]$$

$$= 25$$

(b) $X_2$ $\Rightarrow$ Threshold = 2.5

$X_2 \leq 2.5$
(Left)

$X_2 > 2.5$
(Right)

| $X_1$ | $X_2$ | $Y$ |
|---|---|---|
| 1 | 2 | 10 |
| 2 | 1 | 20 |
| 3 | 2 | 15 |

Mean = 15

MSE = 16.67

| $X_1$ | $X_2$ | $Y$ |
|---|---|---|
| 4 | 3 | 25 |
| 5 | 4 | 30 |

Mean = 27.5

MSE = 6.25

$*$ Weighted MSE $(X_2) = \frac{1}{5} \left[ 3 \times 16.67 + 2 \times 6.25 \right]$

$= 12.25$

③ Choose best split

$X_1 = 25$

$X_2 = 12.25$ ✓

Lower is better

Tree

$X_2$

$X_2 \leq 2.5$     $X_2 > 2.5$

Mean of Y in Subset

Mean of Y in s

Answer

④ Repeat steps.

## Pros

1. Easy to understand like if-else.

2. No feature scaling reqd.

3. Can handle Non linear data.

4. feature selection can do.

## Cons

1. Overfitting ✓

2. Computationally extensive