

02-11-2025

Agenda:

→ code implementation :-

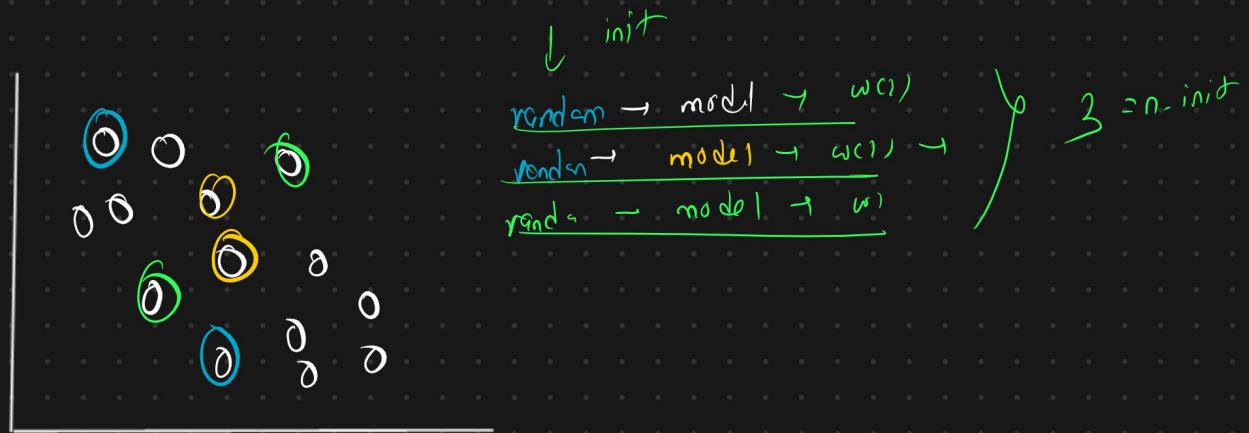
→ naive Bayes

↳ (categorical) NB

↳ Gaussian NB

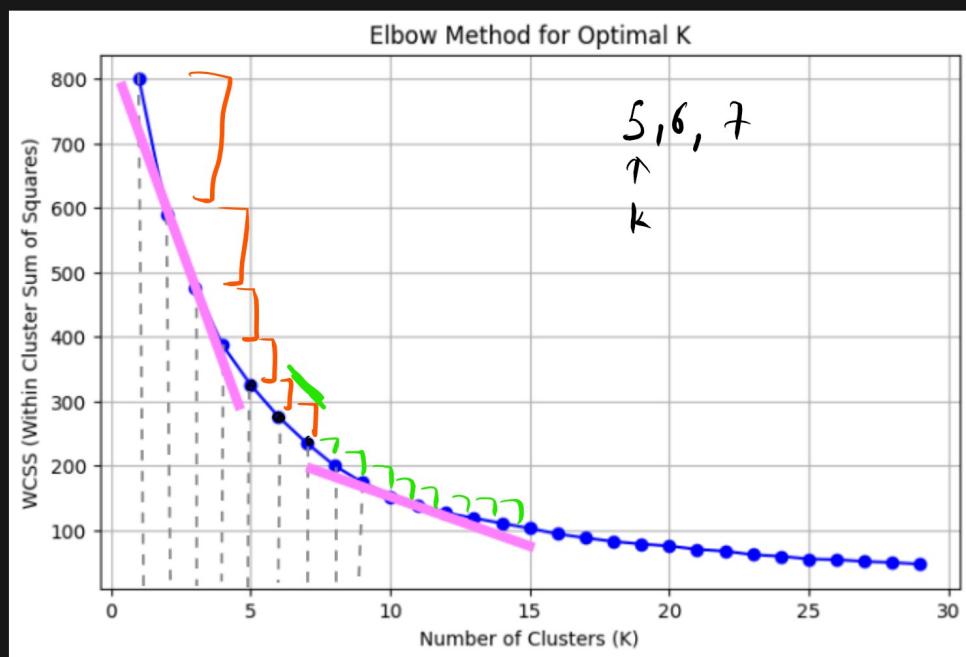
→ k-mean clustering

→ Decision Trees (ML Algo)

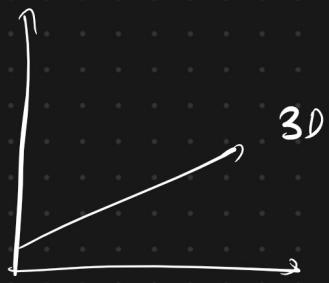
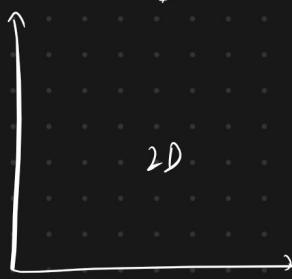


→ divide the dataset optimally

→ i want to divide the dataset into 2 categories



x_1 x_2 x_3 x_4 \rightarrow Red
Blue

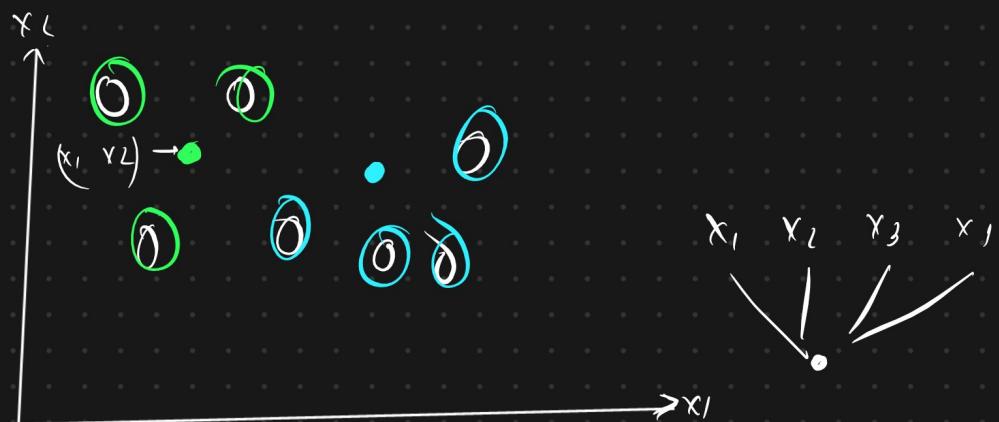


PCA \rightarrow algo \rightarrow High dimension data \rightarrow low dimensional data

Principal component analysis

$x_1, x_2, x_3, x_4 \rightarrow x_1, x_2$ (2D)

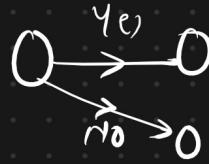
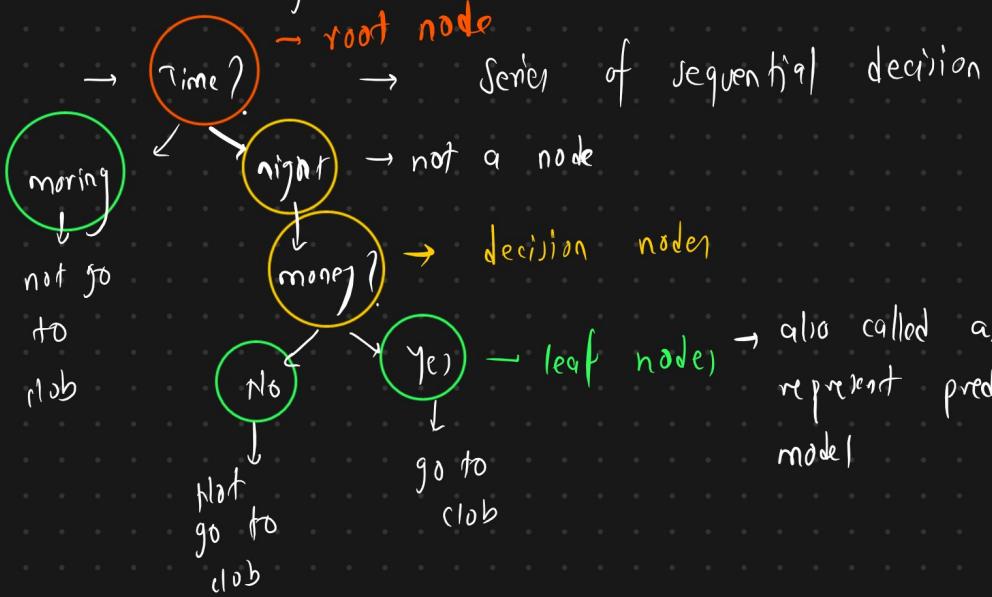
\rightarrow
lose
information



→ Decision Trees:

- supervised
- ease of interpretation
- both regression & classification problems.

Task: Go to night club

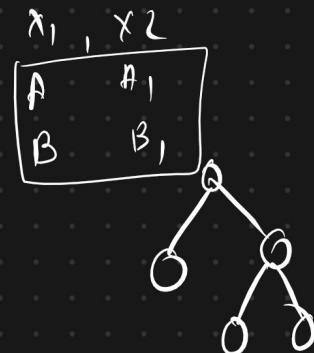


For decision tree, the resulting flow-like structure is navigated via conditional control statements, or if-then rules.

if-else statement → why it is a ML algo?

x_1	x_2	x_3	x_4	x_5	x_6

→ million of rows



Decision Tree classifier

ID	Weather	Temp	Play
1	Sunny	Hot	No
2	Sunny	Mild	No
3	Cloudy	Hot	Yes
4	Rain	Mild	Yes
5	Rain	Cool	Yes
6	Cloudy	Cool	Yes
7	Rain	Hot	No
8	Sunny	Cool	Yes

impurity = how good a split is

→ Entropy: measure of uncertainty: $-\sum p_i \log_2(p_i)$

→ Gini impurity: $1 - \sum p_i^2$

information gain:

: How much uncertainty (entropy) is reduced after splitting on a feature.

$$\text{Information Gain} = \text{Entropy}_{\text{parent}} - \sum \frac{n_{\text{child}}}{n_{\text{total}}} \times \text{Entropy}_{\text{child}}$$

Gender: $\left[\begin{matrix} M & F & F & F \\ 0 & 0 & 0 & 3 \end{matrix}, \begin{matrix} M & F & F & M \\ 1 & 1 & 1 & 0 \end{matrix}, \begin{matrix} S \\ 5 \end{matrix}, \begin{matrix} F \\ 6 \end{matrix} \end{matrix} \right] \rightarrow \underline{\text{impurity}} \xrightarrow{\text{high}} \text{low}$

$$\boxed{\begin{matrix} M-3 \\ F-6 \end{matrix}} \rightarrow \text{Entropy} \quad \leftarrow \text{Impurity}$$

Gini - impurity

Gender → $G_1 [M, F, F, F, F] \rightarrow \text{Entropy}$ ↗ High Information
 $\rightarrow G_2 [M, M, F] \rightarrow \text{Entropy}$ ↘ low IG

- Better impurity a) compared to previous split
- Not the best IG, still higher side.

~~X~~ weather Temp play

Step 1: calculate parent entropy (x_1)

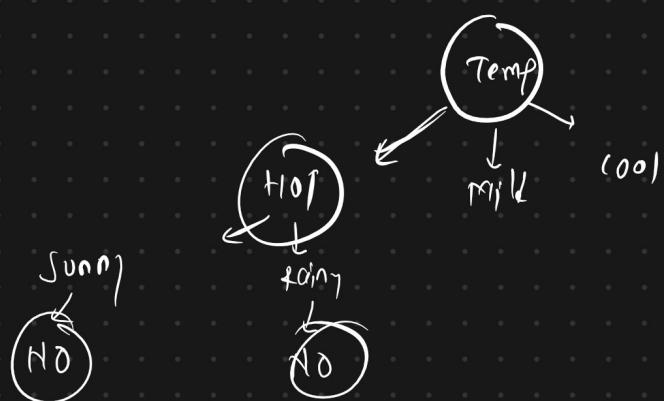
Step 2: Try each feature

$$\begin{array}{c} \text{weather} \\ \downarrow \\ \text{entropy} (x_2) \\ \downarrow \\ \text{IG (parent - child)} \\ \downarrow \\ (x_1 - x_2) \\ \rightarrow A_1 (1, 2, 3) \end{array}$$

$$\begin{array}{c} \text{Temp} \\ \downarrow \\ \text{entropy} (x_3) \\ \downarrow \\ \text{IG (parent - child)} \\ \rightarrow (x_1 - x_3) \\ \rightarrow A_2 (5, 6) \end{array}$$

→ whichever is having the higher IG becomes the root node

→ $A_2 \rightarrow \text{Temperature}$



ID	weather	Temp	play
1	sunny	hot	No
2	sunny	mild	No
3	cloudy	hot	Yes
4	rain	mild	Yes
5	Rain	cool	Yes
6	cloudy	cool	Yes
7	Rain	Hot	No
8	sunny	cool	Yes

Step 1: calculate parent entropy:

$$\begin{array}{l|l} \text{prob} & \\ \hline \text{Yes} \rightarrow 5 & S/P = 0.625 \\ \text{No} \rightarrow 3 & S/P = 0.375 \end{array}$$

$$\text{Entropy: } - (p_i \log_2 p_i) \rightarrow - (0.625 \log_2 0.625 + 0.375 \log_2 0.375) \rightarrow \underline{\underline{0.954}}$$

Step 2: Try each feature

Feature 1: weather		Entropy		
weather	Count	Yes	No	
sunny	3	1	2	$-(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}) = 0.918$
cloudy	2	2	0	0
Rain	3	2	1	0.918

weighted average entropy after split:

$$\frac{3}{8} \times 0.918 + \frac{2}{8} \times 0 + \frac{3}{8} \times 0.918 \rightarrow 0.688$$

$$\begin{aligned}\text{Information gain (weather)} &= 0.954 - 0.688 \\ &= \underline{\underline{0.266}}\end{aligned}$$

Feature 2: Temp:

Temp	Count	Yes	No	Entropy
Hot	3	1	2	0.918
Mild	2	1	1	1
Cool	3	3	0	0

Weighted Entropy:

$$\frac{3}{8} \times 0.918 + \frac{2}{8} \times 1 + \frac{3}{8} \times 0 = 0.594$$

$$\begin{aligned}\text{Information gain (Temp)} &= 0.954 - 0.594 \\ &= 0.360\end{aligned}$$

$I^G(\text{Weather})$

$I^G(\text{Temp})$

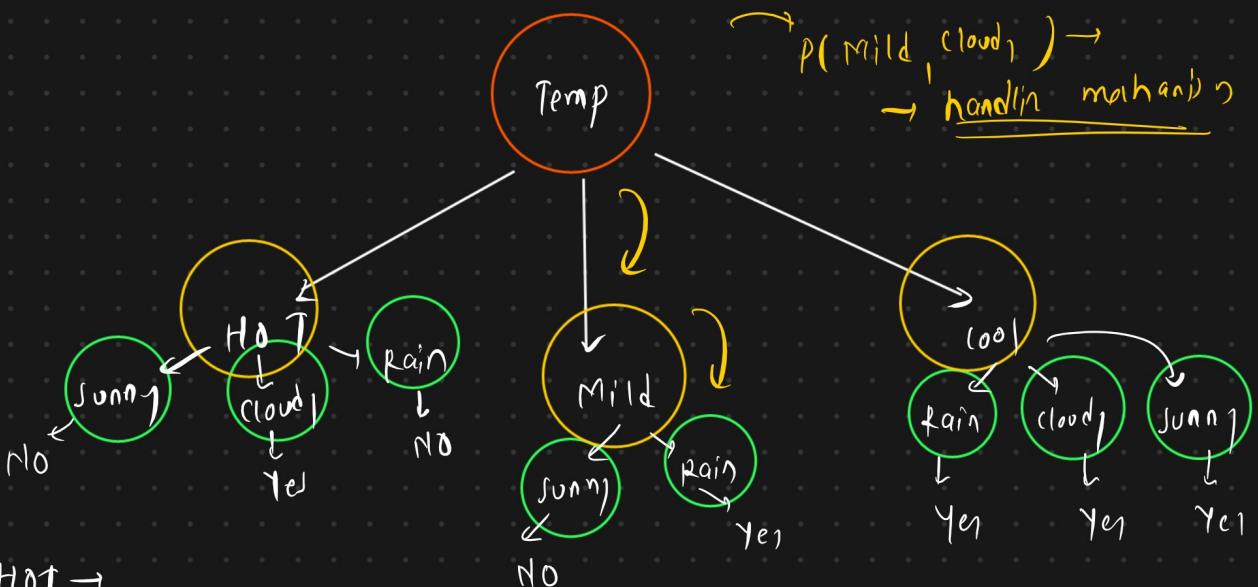
0.266

<

0.360



Temp becomes our root node.



Hot →

weather

sunny

cloudy

Rain

play

No

Yes

No

Yes → 1

No → 2

Hot	Count	Yes	No	Entropy
	3	1	2	0.918 (same as before)

↓
 we must split again to make
 it pure

Mild →

=

sunny

Rain

play

No

Yes

Entropy = 1.0 (mixed) (same as before) → split

