

Class

Entropy

Data Mining

GIA

Decision tree Algo

age' select w.r.t 20. Info. gain ~~max~~ min, Gain ratio etc,

Info. Gain

Class same w.r.t bias w.r.t output of Df entropy

w.r.t age 20,

$$H(S) = - \sum_{i=1}^n p_i \log_2 p_i$$

$$= - \frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$= 0.902$$

∴ Perfectly imbalance

(some class 1 0)

$$IG(D, A) =$$



Attribute

$$IG(S, \text{Outlook}) = H(S) - \sum_{v \in S, \text{Out}} \frac{|S_v|}{|S|} \cdot H(S_v)$$

$$= \left(-\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \right)$$

Condindity
Gained Data
Value

$$- \left(\frac{5}{14} \times 0.97 \right) + \left(\frac{9}{14} \times 0 + \frac{5}{14} \times 0.97 \right)$$

(Sunny) H (Rain) P

$$(0 \times \frac{1}{2}) = 0.247$$

$$IG(S, \text{Humidity}) = 0.0902 \left\{ \left(\frac{7}{14} \times 0.985 \right) + \left(\frac{7}{14} \times 0.0916 \right) \right\}$$

$$= 0.1519$$

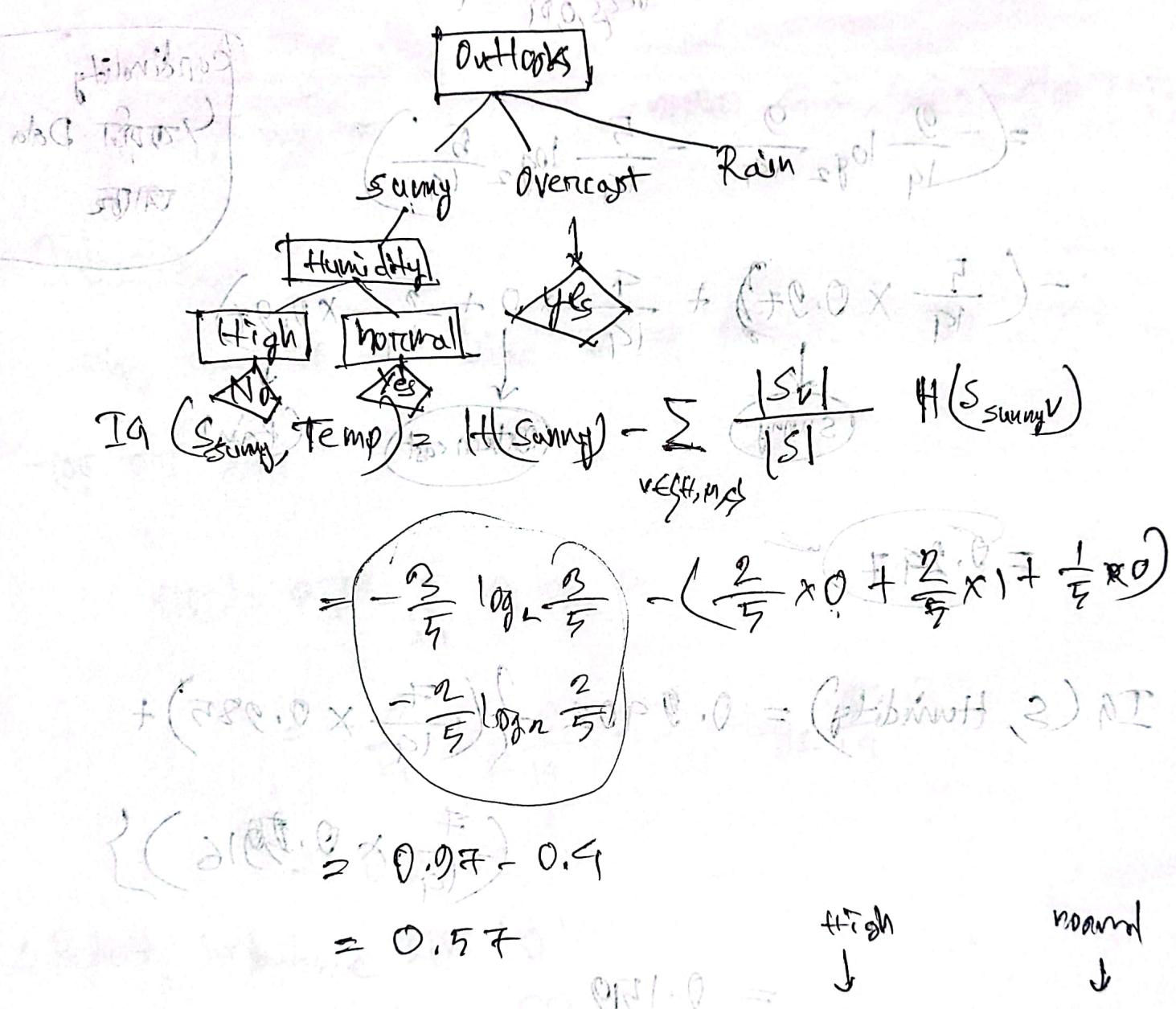
$$IG(S, \text{Temp}) = 0.0902 - \left\{ \frac{9}{14} \times \left(-\frac{2}{9} \log_2 \frac{2}{9} - \frac{2}{9} \log_2 \frac{2}{9} \right) \right.$$

$$+ \frac{5}{14} \left(-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right) + \frac{9}{14} \left(-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right) \left. \right\}$$

$$\approx 0.1519$$

With 5 bits second letter is (fifth digit, first) R

④ Outlook: 2nd. highest DTF off root node.



$$IG(Sunny, Humidity) = 0.97 - \left(\frac{3}{5} \times 0 + \frac{2}{5} \times 0.92 \right)$$

$$IG(Sunny, Wind) = 0.97 - \left(\frac{2}{5} \times 1 + \frac{3}{5} \times 0.92 \right)$$

$$= 0.97 - 0.952$$

$$= 0.018$$

The IG_{Wind}(Sunny, Humidity) is accepted, because it is high.

Class 3

Data mining

Decision tree

Gini Index:

$$G(S) = 1 - \sum_{i=1}^c (P_i)^2$$

$$G(\text{Outlook}) = 0.98 \times \frac{5}{14} + 0 + 0.98 \times \frac{5}{14}$$

$$\begin{aligned} &= 0.342 \\ &\approx 0.342 \\ &\approx 0.96 \end{aligned}$$

$$G(\text{Outlook} \rightarrow \text{Sunny}) = 1 - \left\{ \left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right\} = 0.98$$

$$G(\text{Outlook} \rightarrow \text{Overcast}) = 1 - \left\{ \left(\frac{9}{14} \right)^2 \right\} = 0.$$

$$G(\text{Outlook} \rightarrow \text{Rain}) = 1 - \left\{ \left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right\} = 0.98$$

Wind

$$G(\text{Wind}) = 0.375 \times \frac{8}{14} + 0.5 \times \frac{6}{14} = 0.43$$

$$G(\text{Weather} \rightarrow \text{Wind} \rightarrow \text{Weak}) = 1 - \left\{ \left(\frac{6}{8} \right)^2 + \left(\frac{2}{8} \right)^2 \right\} = 0.375$$

$$G(\text{Wind} \rightarrow \text{Strong}) = 1 - \left\{ \left(\frac{3}{8} \right)^2 + \left(\frac{3}{6} \right)^2 \right\} = 0.5$$

$$\therefore G(\text{Humidity}) = 0.36$$

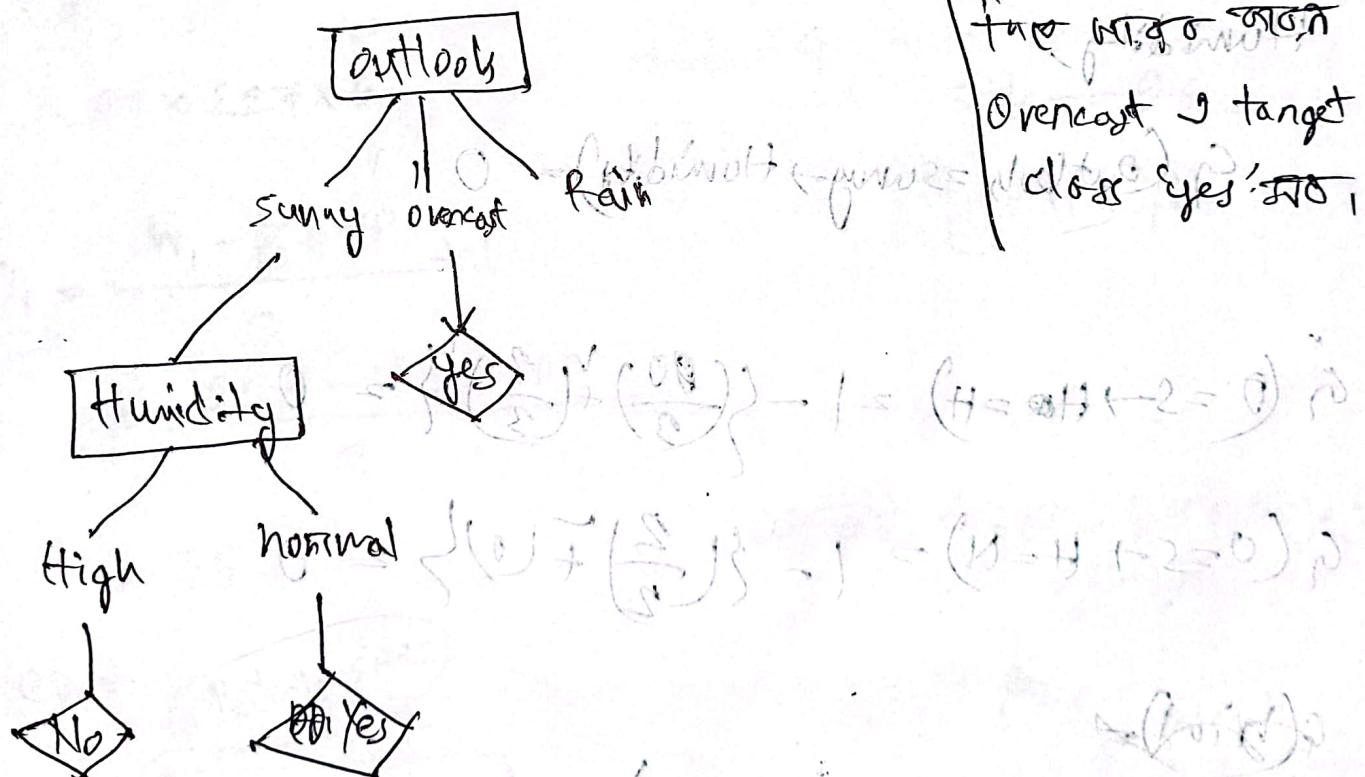
$$G(t \rightarrow \text{Hot}) = 1 - \left\{ \left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right\} = 0.5$$

$$G(t \rightarrow \text{Mild}) = 1 - \left\{ \left(\frac{4}{6} \right)^2 + \left(\frac{2}{6} \right)^2 \right\} = 0.54$$

$$G(t \rightarrow \text{Cold}) = 1 - \left\{ \left(\frac{3}{9} \right)^2 + \left(\frac{1}{9} \right)^2 \right\} = 0.375$$

$$G(\text{Temp}) = 0.5 \times \frac{9}{14} + 0.44 \times \frac{8}{14} + \frac{9}{14} \times 0.375 \\ = 0.94$$

We take outlook.



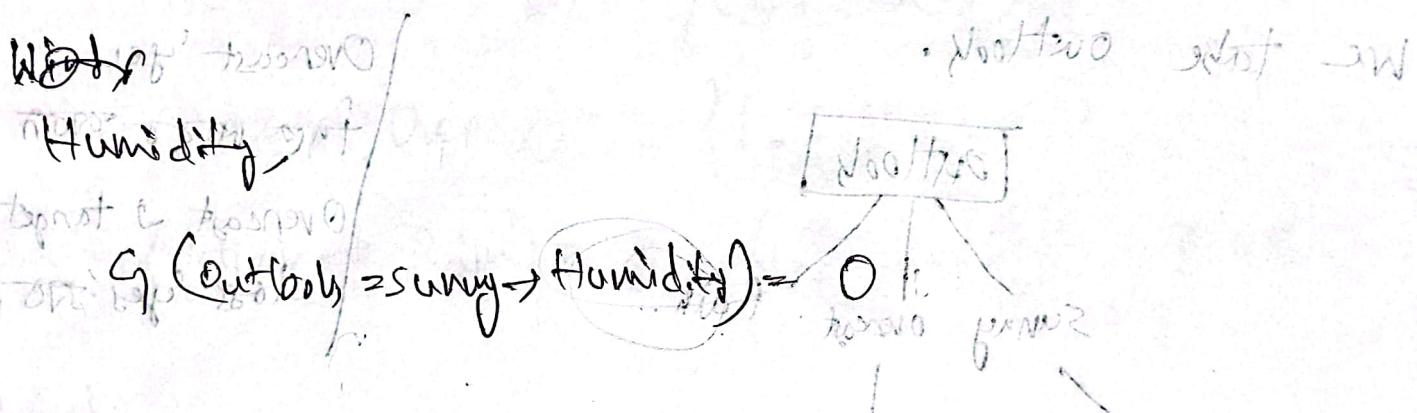
$$G_1(\text{Outlook} = \text{Sunny} \rightarrow \text{Temp}) = 0.02 \times 0 \times \frac{2}{5} + 0.4 \times \frac{2}{5} + 0.4 \times \frac{1}{5}$$

PP.0 = $\{ \rightarrow 0.2(5) \}, 1 - (0.1 \times 5) = 0.8$

$$G_1(0=S \rightarrow T=H) = 0.1 - \left\{ \left(\frac{0}{2} \right)^v + \left(\frac{2}{2} \right)^v \right\} = 0.00$$

$$G_1(0=S \rightarrow T=M) = 0.1 - \left\{ \left(\frac{1}{2} \right)^v + \left(\frac{1}{2} \right)^v \right\} = 0.5$$

$$G_1(0=S \rightarrow T=C) = 1 - \left(\frac{1}{1} \right)^v = 0$$



$$G_1(0=S \rightarrow H=H) = 1 - \left\{ \left(\frac{0}{3} \right)^v + \left(\frac{3}{3} \right)^v \right\} = 0$$

$$G_1(0=S \rightarrow H=N) = 1 - \left\{ \left(\frac{2}{2} \right)^v + 0 \right\} = 0$$

Wind

$$G_1(\text{Outlook} = \text{sunny} \rightarrow \text{Wind}) = 0.49$$

} Segmentation
Image Compression
Morphological

3010

1130

bait → nitroloids → nitrates + gases

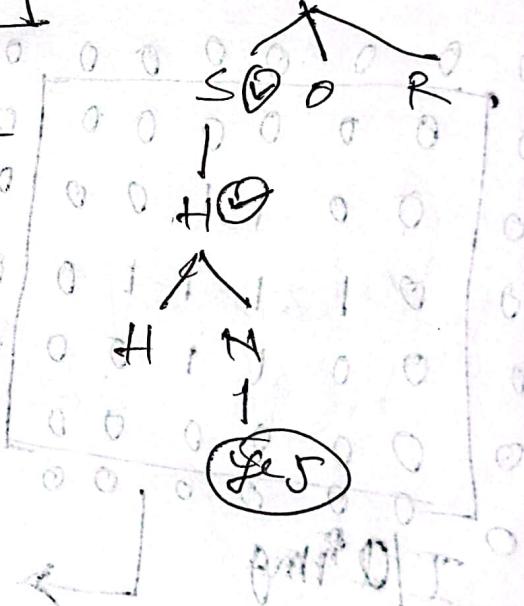
Class 4

Data Mining

Decision tree

Sunny, ~~Mild~~, Strong, Normal, Yes

Unknown Data Match



↳ Probability
↳ Conditional Probability

$$P(A|B)$$

Event Ω is a set of possibility

$$P(A|B)$$

feature / attribute

class

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \rightarrow P(A|B) * P(B) = P(A \cap B)$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)} \rightarrow P(B|A) * P(A) = P(B \cap A)$$

$$P(A|B) = \frac{P(A) * P(B|A)}{P(B)}$$

\rightarrow Bayesian formula.

Naive Bayesian

$$P(A|B)$$

$$C_1$$

$$C_2$$

Class

$$A = \{\text{yes}, \text{No}\}$$

$$B = \{B_1, B_2, \dots, B_n\}$$

$$P(A|B) = \frac{P(A) * P(B|A)}{P(B)}$$

$$P(\text{yes}|B) = \frac{P(\text{yes}) * P(B|\text{yes})}{P(B)}$$

$$P(\text{No}|B) = \frac{P(\text{No}) * P(B|\text{No})}{P(B)}$$

$$P(\text{Yes}) \cdot P(B_1 | \text{Yes}) = P(B_n | \text{Yes})$$

$$\therefore P(\text{Yes} | B) = \frac{P(\text{Yes}) \cdot P(B_1) \cdot P(B_2) \cdots P(B_n)}{P(B_1) \cdot P(B_2) \cdots P(B_n)}$$

$$\therefore P(\text{No} | B) = \frac{P(\text{No}) \cdot P(B_1 | \text{No}) \cdot P(B_2 | \text{No}) \cdots P(B_n | \text{No})}{P(B_1) \cdot P(B_2) \cdots P(B_n)}$$

\therefore 2 st class than 2nd 10 times closer to same

20,

$$B = \{32, \text{High}, \text{No, excellent}\}$$

$$P(\text{Yes} | 32, \text{High, No, excellent}) =$$

without for solve

$$P(\text{Yes}) | P(32 | \text{Yes}), P(\text{High} | \text{Yes})$$

$$P(\text{High} | \text{Yes}), P(\text{excellent} | \text{Yes})$$

$$= \left(\frac{12}{19}\right) * \left(\frac{9}{10}\right) * \left(\frac{2}{5}\right) * \left(\frac{3}{9}\right) * \left(\frac{2}{3}\right)$$

$$= 0.00705$$

$$P(\text{No} | 32, \text{High, No, excellent}) = P(\text{No}) | P(32 | \text{No}), P(\text{High} | \text{No})$$

This cont
not accepted

$$= \left(\frac{12}{19}\right) * \left(\frac{1}{5}\right) * \left(\frac{2}{5}\right) * \left(\frac{9}{10}\right) * \left(\frac{3}{5}\right)$$

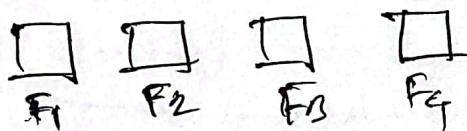
Name by \rightarrow Adv / Disadv:

④ k-fold, cross validation, hold out method, Bootstrap

\downarrow
Define

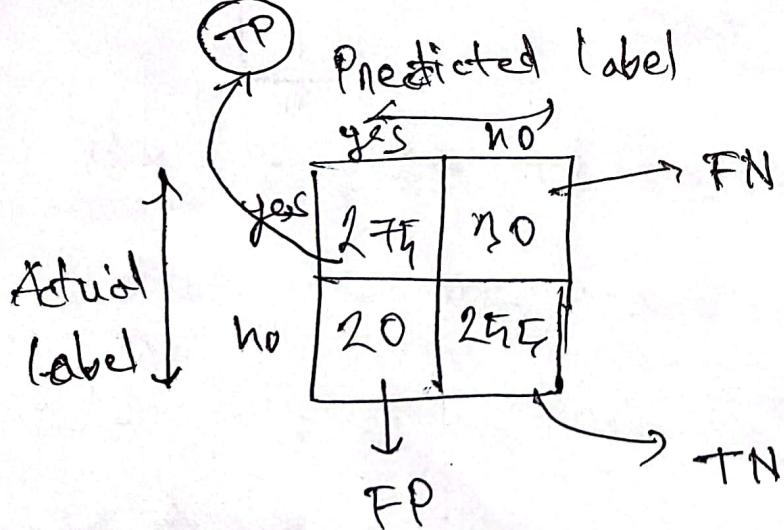
k-fold

\circlearrowleft k-fold



<u>Train</u>	<u>Test</u>
F_1, F_2, F_3	F_4
F_1, F_2, F_4	F_3
F_1, F_3, F_4	F_2
F_1, F_2, F_3	F_4

Confusion



$$\text{Acc} = \frac{TP + TN}{TP + TN + FN + FP}$$

$$= 0.91$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{recall / sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Predicted
ஏது போன்றை

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

$$F_1 \text{ score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

True Positive (TP)
False Positive (FP)
True Negative (TN)
False Negative (FN)

9 class

		Predicted			
		0	1	2	3
Actual	0	TP ₀	FN ₀	FN ₀	FN ₀
	1	FP ₀	TP ₁	TN ₀	TN ₀
2	0	FP ₀	TN ₀	TP ₂	TN ₀
	1	FP ₀	TN ₀	TN ₀	TP ₃

Prediction = 0 \rightarrow TP
Prediction = 1 \rightarrow FN

Class no 1

		Predicted			
		0	1	2	3
Actual	0	FP ₁			
	1	FN ₁	TP ₁	FN ₁	FN ₁
2	0	FP ₁			
	1	FP ₁			

Class no 2

		Predicted			
		0	1	2	3
Actual	0	FP ₂			
	1	FP ₂			
2	0	FN ₂	FN ₂	TP ₂	FN ₂
	1	FP ₂			

Class 6

Data Mining

9T

9T + 9T

9T + 9T

clustering → (Unsupervised)

soft cluster → 1 ft multiple clusters

hard → 1 ft 1 ft cluster

Partition based

k-means

Problem → k value को select करें, या Problem

elbow

यह Data include करें कि, किसी जगह डिवाइस

9T ← 0, calculate SSE,

HIT ← 1 = 9T + 9T + 9T + 9T

Solu
Elbow
method.

dense

9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T

Partition

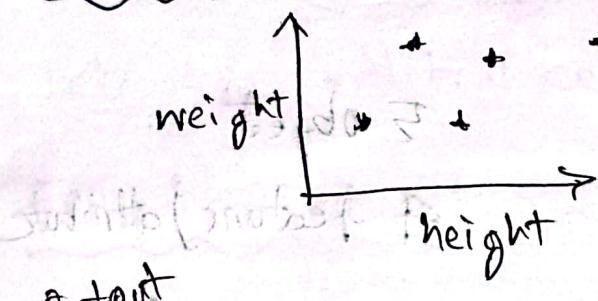
Hierarchical

Density

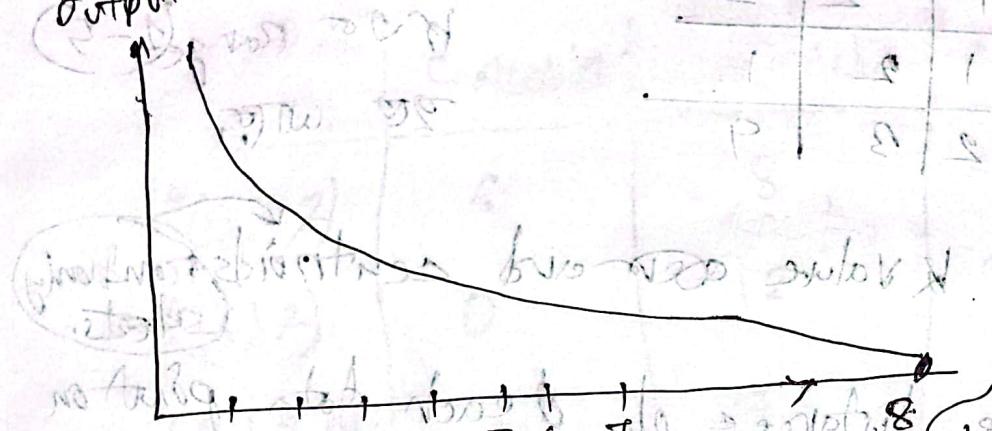
9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T

9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T
9T	9T	9T	9T	9T	9T

Elbow



Output



$$W_{\text{tot}} = d_1^r + d_2^r + d_3^r + d_8^r$$

(*Chionanthus* LK) \rightarrow $b=1$, \nearrow

For $y = 2^{-\frac{1}{2}x}$, the graph is decreasing and passes through $(0, 1)$.

$$W_{ex,1} + W_{ex,2} = W_{ex}$$

2 group

मुद्दात आर्टिंग slowly decrease with elbow point.

K-means

ID	A ₁	A ₂	A ₃	A ₄
1	5	9	3	2
2	7	6	1	2
3	0	1	2	3
4	2	1	2	1
5	1	2	3	9

5 objects

4 feature / attribute

k of range (1-5)

20 wts,

Step-1 Initially find the k value and centroids randomly selects.

Step-2 calculate the distances of each data point on objects from no. of centroids (all centroids)

Step-3 Assign the points come in to their corresponding clusters.

Step-4 repeat step (2,3) until the cluster are stable

K value = centroid

Centroids (2, 9), $k=2$

Let's, the initial centroids be $(7, 6, 1, 2)$ & $(2, 1, 2, 1)$

For iteration 1

Data obj	Centroid 1	Centroid 2	cluster
	cluster 1	cluster 2	cluster 3
(5, 9, 3, 2)	6	8 0	cluster 1.
(7, 6, 1, 2)	0	12	
(0, 1, 2, 3)	19	9	cluster 2
(2, 1, 3, 1)	12	0	cluster 2
(1, 2, 3, 4)	19	6	cluster 2

New centroid,

$$\text{new-centroid 1} = \left(\frac{7+5}{2}, \frac{9+6}{2}, \frac{1+1+2}{3}, \frac{2+1}{2} \right) \\ = (6, 7.5, 2, 1.5)$$

$$\text{new-centroid 2} = \left(\frac{0+2+1}{3}, \frac{1+1+2}{3}, \frac{2+1+3}{3} \right) \\ = (1, 1.33, 2.33, 2.67)$$

complete the iteration, if iteration match all data points then change in centroid becomes zero

Training \rightarrow (Draw and loss function)

class 7

Data Mining

Density based clustering (DBSCAN)

K-mean
Limitation
- Advantage
Application

Density based clustering

Location wise clustering

∴ two parameters \rightarrow epsilon and minPoints.

① core point, border point, and noise.

मात्र एक नहीं point जैसे But
minA = 3 का लिया गया

if 3 नहीं तो 2
total center point
एक core point

$$\begin{cases} \epsilon = 1 \\ \text{minA} = 3 \end{cases}$$

3 for point
जैसा

मात्र
एक

② trial and error method एवं ε (E) द्वारा

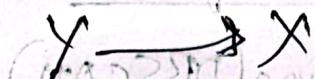
जैसा कि $E = 1$

minPoint \rightarrow Dimension of 1

feature/attribute

Reachability and Connectivity

X is direct density-reachable from point Y.

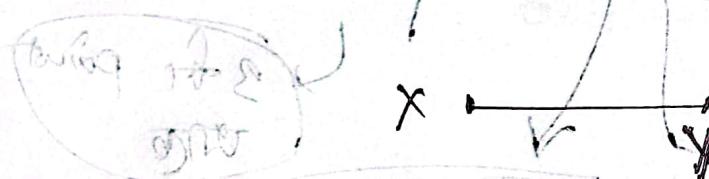


④ Y is P₂ indirectly density-reachable.

⑤ X & Y are Density Reachable.

Density connected

X is Y & Y is Z same density Reachable.



Algorithm steps

①

Visit (1) to (2) to (3) to (4) to (5) to (6) to (7) to (8)

start from 1 to 2 to 3 to 4 to 5 to 6 to 7 to 8

Example:

	u_1	u_2	u_3	u_4	u_5	u_6	u_7	u_8	u_9
u_1	0								
u_2	1	0							
u_3	1.41		0	.					
u_4				0					
u_5					0				
u_6						0			
u_7							0		
u_8								0	
u_9									0

euclidean $\rightarrow d(u_1, u_2) = \sqrt{(1-0)^2 + (0-0)^2} = 1$

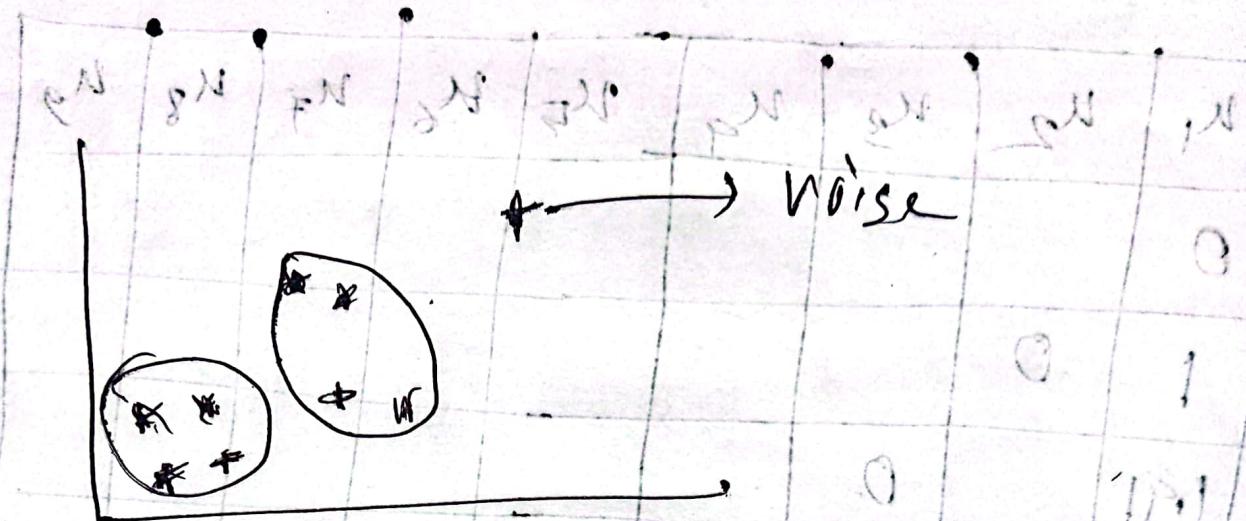
- ① $u_1, u_2 \in$ १st मात्र \Rightarrow उनकी distance ऐसा होगा कि वे एक नियंत्रित निकट स्वरूप हों।

min point = 3

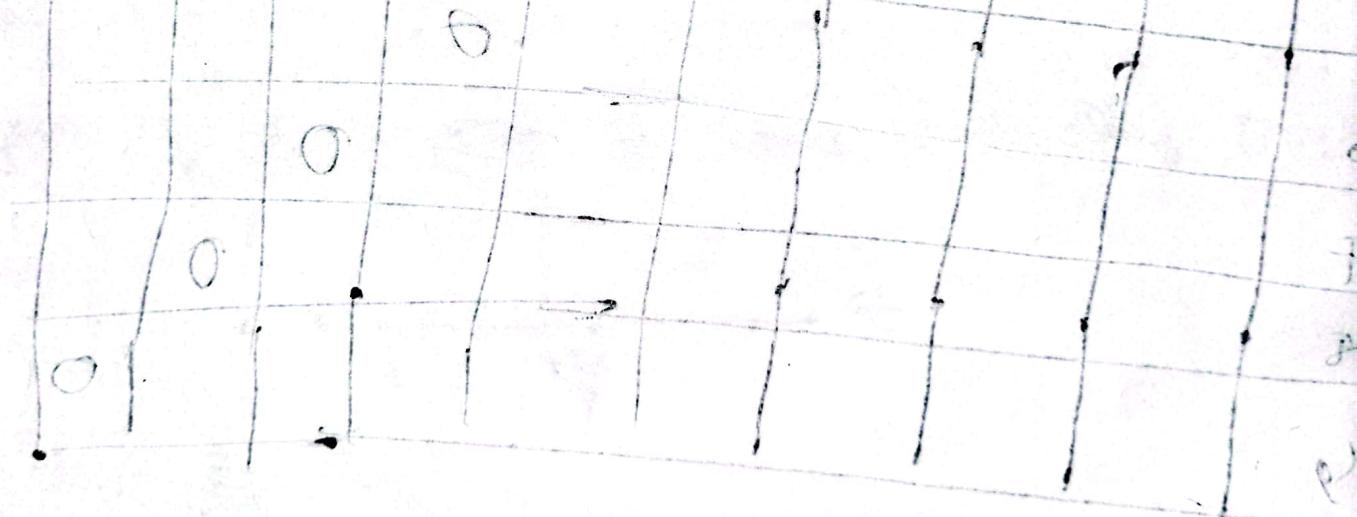
- ② $u_5, u_6 \rightarrow$ border point $\Rightarrow u_9 = \text{noise}$

- ③ u_3, u_4 Directly and density not reachable \Rightarrow outlier

क्लूसिंग के तरीके



⑧ क्लूसिंग.



$$1 = \sim(6-0) + \sim(1-0) = 5, \text{ जहाँ } \leftarrow \text{ ना}$$

प्रति एक विशिष्ट रेखा पर 0 अंक हो

1000 लोगों का बच्चादासीय रेखा

2 = दो

OLAP operation

- Roll up
- Drill Down → multidimensional
- Slice → 3D ,
- Dice → sub-cube तथा उत्पाद
- Pivot → Transpose matrix or convert उत्पाद, (Rotated उत्पाद तथा उत्पाद)

④ Dimensions

④ Facts and Measure

↳ Dimension वे वर्ग हैं जो measure.

Schemas

Star, snowflake ,

Star schema → 1 fact table.

Snowflake →

Galaxy schema → more than 1 fact table.

Figure fact schema explain तथा

