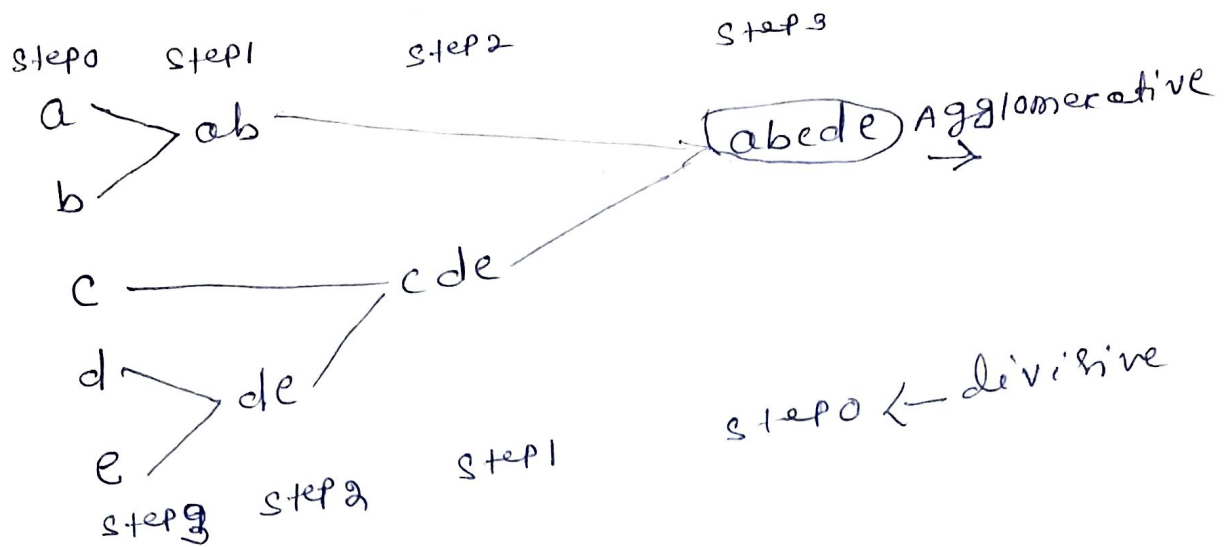


# Hierarchical clustering (unsupervised) ①

- 1) Agglomerative (bottom-up)
- 2) Divisive (top-down)



cluster distance measure.

Ex

	a	b	c	d	e
	1	2	4	5	6

$$C_1 = \{a, b\}$$

$$C_2 = \{c, d, e\}$$

Distance Matrix.

	a	b	c	d	e
a	0	1	3	4	5
b	1	0	2	3	4
c	3	2	0	1	2
d	4	3	1	0	1
e	5	4	2	1	0

single link

$$\begin{aligned} \text{dist}(C_1, C_2) &= \min \{d(a, c), d(a, d), d(a, e), \\ &\quad d(b, c), d(b, d), d(b, e)\} \\ &= \min \{3, 4, 5, 2, 3, 4\} = 2 \end{aligned}$$

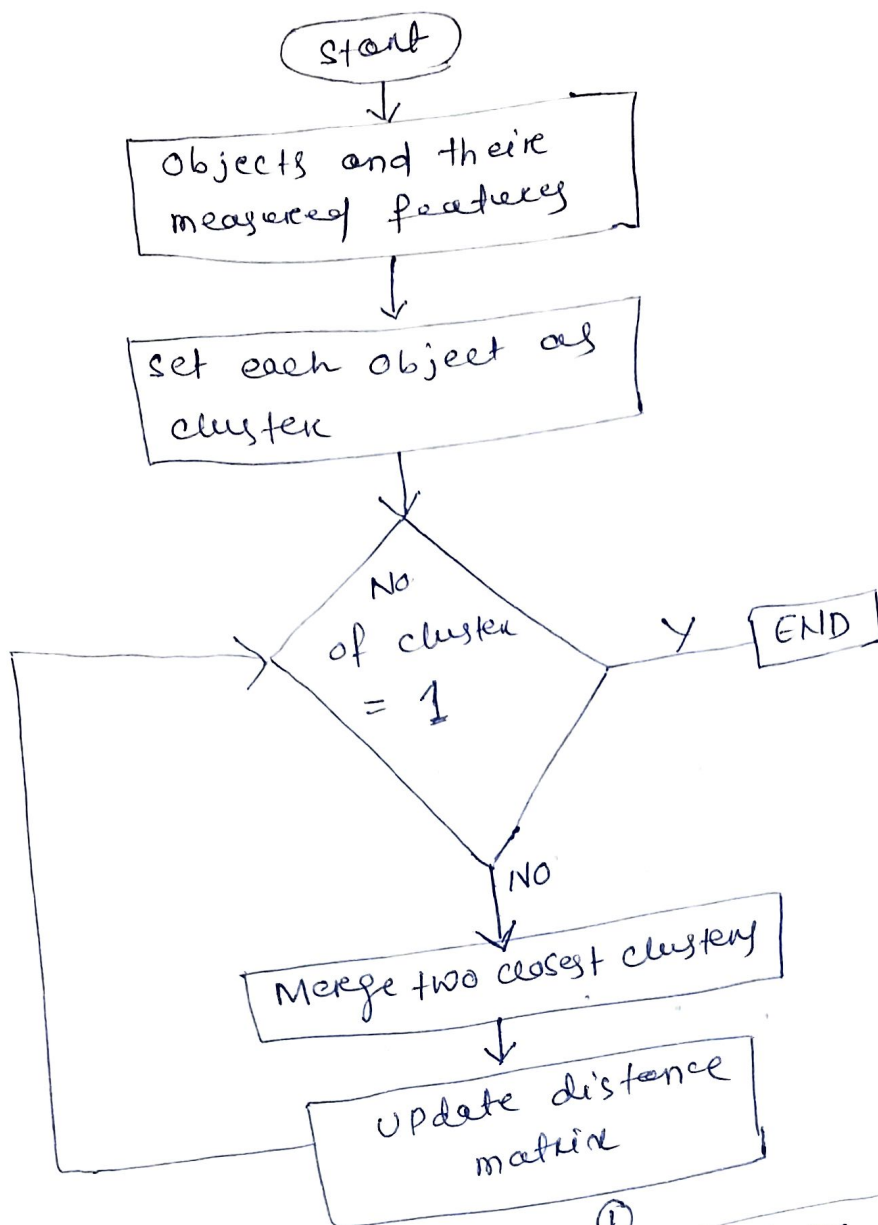
complete link

$$\begin{aligned} \text{dist}(C_1, C_2) &= \max \{d(a, c), d(a, d), d(a, e), \\ &\quad d(b, c), d(b, d), d(b, e)\} \\ &= \max \{3, 4, 5, 2, 3, 4\} = 5 \end{aligned}$$

Average link

$$\begin{aligned} \text{dist}(C_1, C_2) &= \frac{3+4+5+2+3+4}{6} \\ &= \frac{21}{6} = 3.5 \end{aligned}$$

# Algorithm (Agglomerative)



Example

	$x_1$	$x_2$
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5

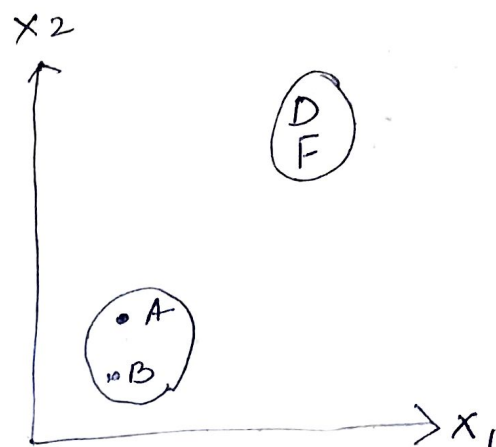
①

	A	B	C	D	E	F
A	0	0.71	5.66	3.61	4.24	3.2
B	0.71	0	4.95	2.92	3.54	2.5
C	5.66	4.95	0	2.24	1.41	2.5
D	3.61	2.92	2.24	0	1	0.5
E	4.24	3.54	1.41	1	0	1.12
F	3.2	2.5	2.5	0.5	1.12	0

$$d_{AB} = \sqrt{(1-1.5)^2 + (1-1.5)^2} = 0.7071 \quad \left| \quad d_{DF} = \sqrt{(3-3)^2 + (4-3.5)^2} = 0.5 \right.$$

②

	A	B	C	DF	E
A	0	0.71	5.66	3.20	4.24
B	0.71	0	4.95	2.5	3.54
C	5.66	4.95	0	2.24	1.41
DF	3.20	2.50	2.24	0	1
E	4.24	3.54	1.41	1.00	0



single linkage

$$d_{DF \rightarrow A} = \min(d_{DA}, d_{FA}) = \min(3.61, 3.20) = 3.20$$

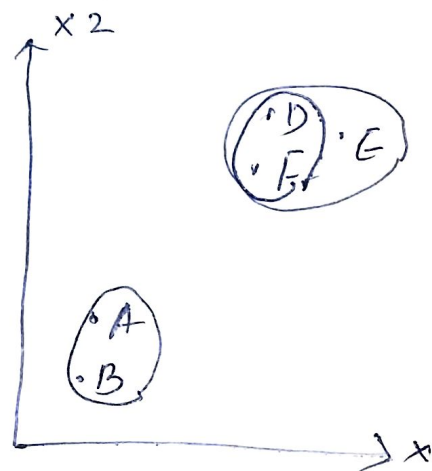
$$d_{DF \rightarrow B} = \min(d_{DB}, d_{FB}) = \min(2.92, 2.50) = 2.50$$

$$d_{DF \rightarrow C} = \min(d_{DC}, d_{FC}) = \min(2.24, 2.50) = 2.24$$

$$d_{DF \rightarrow E} = \min(d_{DE}, d_{FE}) = \min(1, 1.12) = 1$$

③

	AB	C	DF	E
AB	0	4.95	2.5	3.5
C	4.95	0	2.24	1.41
DF	2.5	2.24	0	①
E	3.54	1.41	①	0



$$d(C \rightarrow AB) = \min(d_{CA}, d_{CB}) = \min(5.66, 4.95) = 4.95$$

$$d_{DF \rightarrow AB} = \min(d_{DF \rightarrow A}, d_{DF \rightarrow B}) = \min(3.2, 2.5) = 2.5$$

$$d_{E \rightarrow AB} = \min(d_{EA}, d_{EB}) = \min(4.24, 3.54) = 3.54$$

④

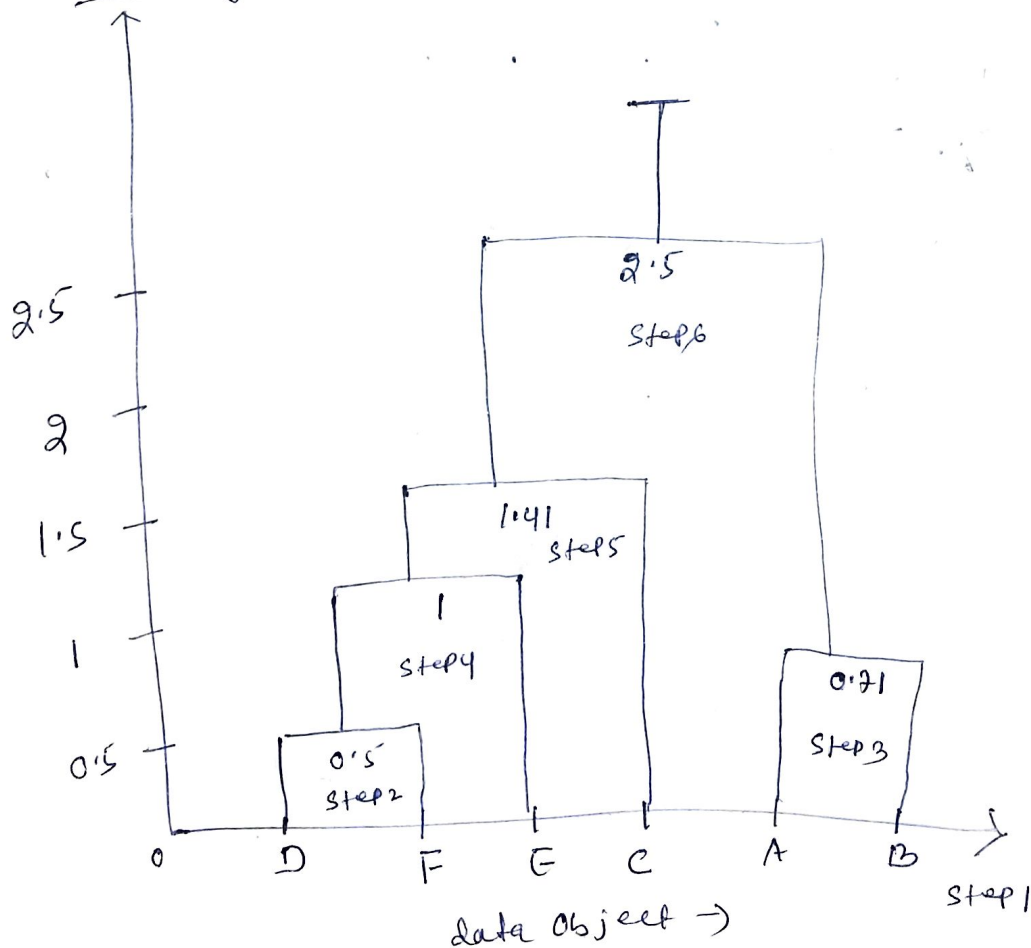
	AB	C	DFE
AB	0	4.95	2.50
C	4.95	0	(1.41)
DFE	2.50	(1.41)	0

$$\begin{aligned}
 d_{AB \rightarrow DFE} &= \min(d_{AB \rightarrow DF}, d_{AB \rightarrow E}) \\
 &= \min(2.5, 3.5) = 2.5 \\
 d_{C \rightarrow DFE} &= \min(d_{C \rightarrow DF}, d_{C \rightarrow E}) \\
 &= \min(2.24, 1.41) = 1.41
 \end{aligned}$$

	AB	DFEC
AB	0	(2.50)
DFEC	(2.50)	0

$$\begin{aligned}
 d_{AB \rightarrow DFEC} &= \min(d_{AB \rightarrow DFE}, d_{AB \rightarrow E}) \\
 &= \min(2.5, 4.95) = 2.5
 \end{aligned}$$

Dendrogram tree representation





DBSCAN (Density based spatial clustering of applications with Noise)

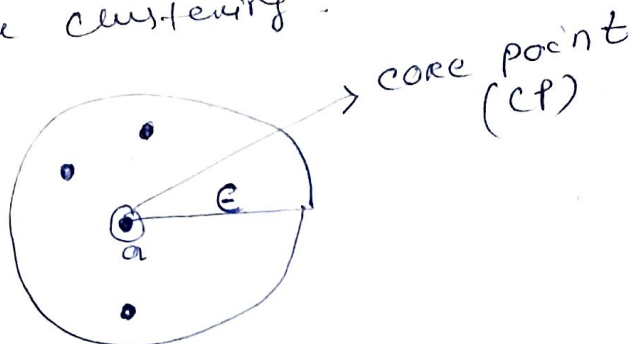
(3)

Density  $\rightarrow$  No. of points in a given area

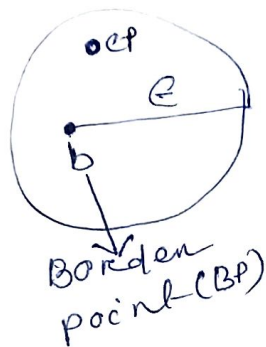
Inputs - Epsilon ( $\epsilon$ )

Minimum points = let it be 3

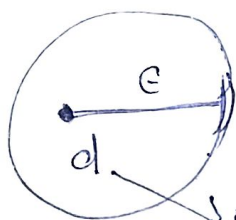
$\epsilon$  represents the radius to draw a circle for clustering.



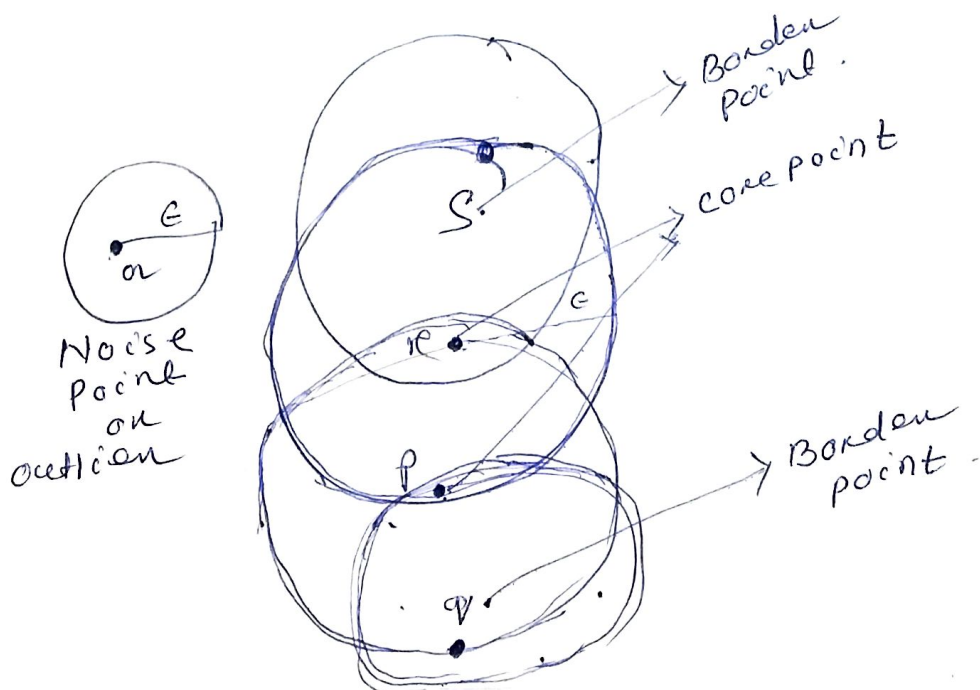
[It can detect Noise or Outliers easily.]



part of the cluster because it is within  $\epsilon$  of a core point but does not meet the min-point criteria.



No core point  
Not satisfying min-point criteria



Ex	B/N	B/N	B/N	CP	CP	CP
	A	B	C	D	E	F
A	0	0.7	5.7	3.6	4.2	3.2
B	0.7	0	4.9	2.9	3.5	2.5
C	5.9	4.9	0	2.2	1.4 ✓	2.5
D	3.6	2.9	2.2	0 ✓	1	0.5
E	1.2	3.5	1.4	1	0	1.5
F	3.2	2.5	2.5	0.5	1.1	0

Let  $E = 1.5$   
 min-point  $\geq 3$ .

C - BP (Border point).  
 as close to E a cp  
 B - Noise point  
 A - Noise point.