

Q: What is Clustering?

Regression & Classification

Dep. variable
 $D = \{x_i, y_i\}, y = f(x)$

$y \in \{0, 1\} \text{ or } 0, 1, 2, 3, \dots$

$y \in \mathbb{R}^q$

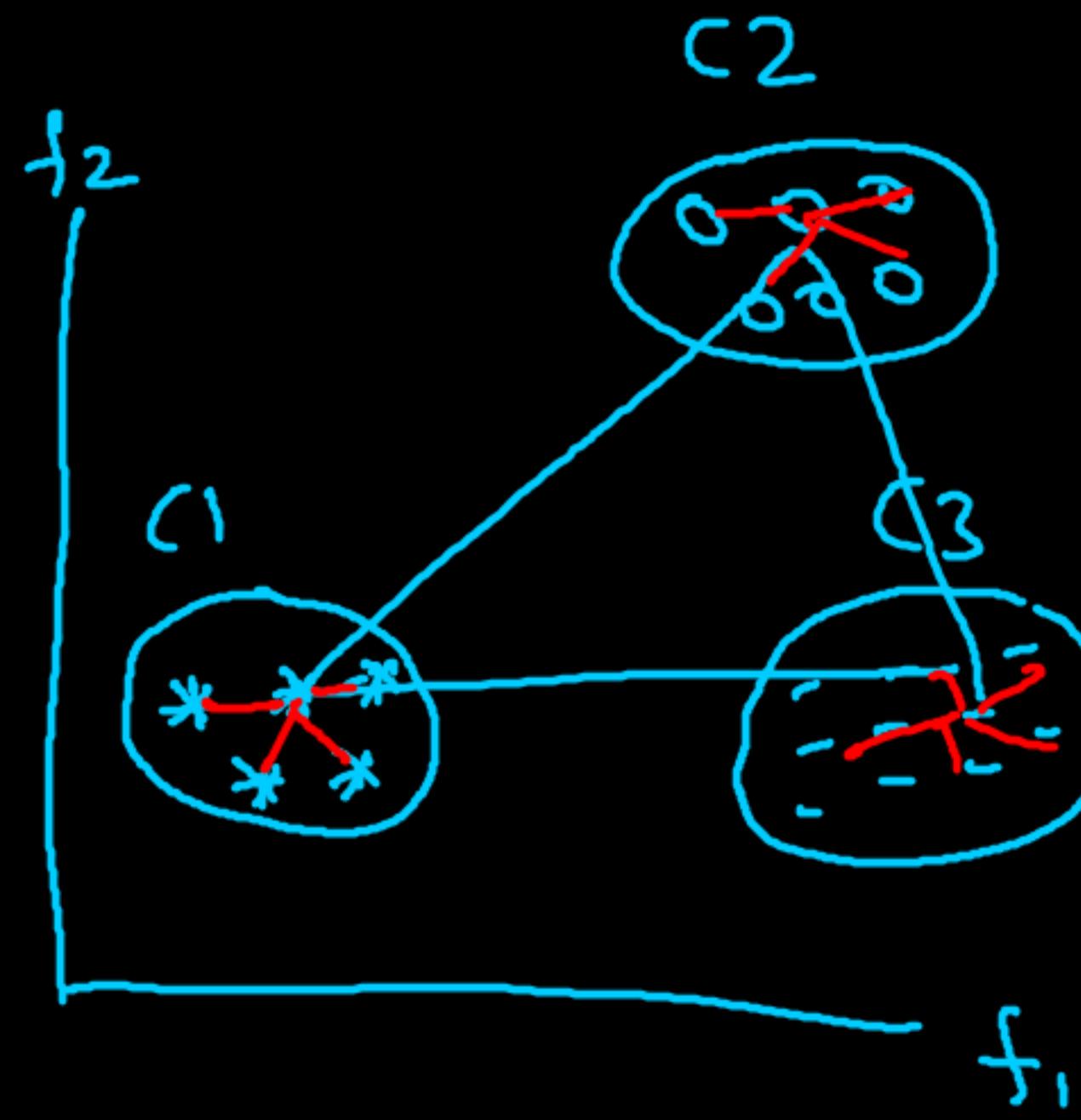
$x_i = \text{Ind. variable}$

"Clustering"

$D = \{x_i\}$, No y_i | No dep. variable
 Given Not given
 Output Result var.

TASK := Group/cluster

"Similar" data point



Task of clustering

→ Group | Cluster | Segment "similarity"

Point

→ points are close together
by the help of Euclidean dist

Point 1 :- points in a cluster are closer together.

Point 2 :- points in different cluster are far away.

Supervised ML — classification-report, accuracy-score, Adj-R-Sqr,
MAE, MAPE, MSE, RMSE

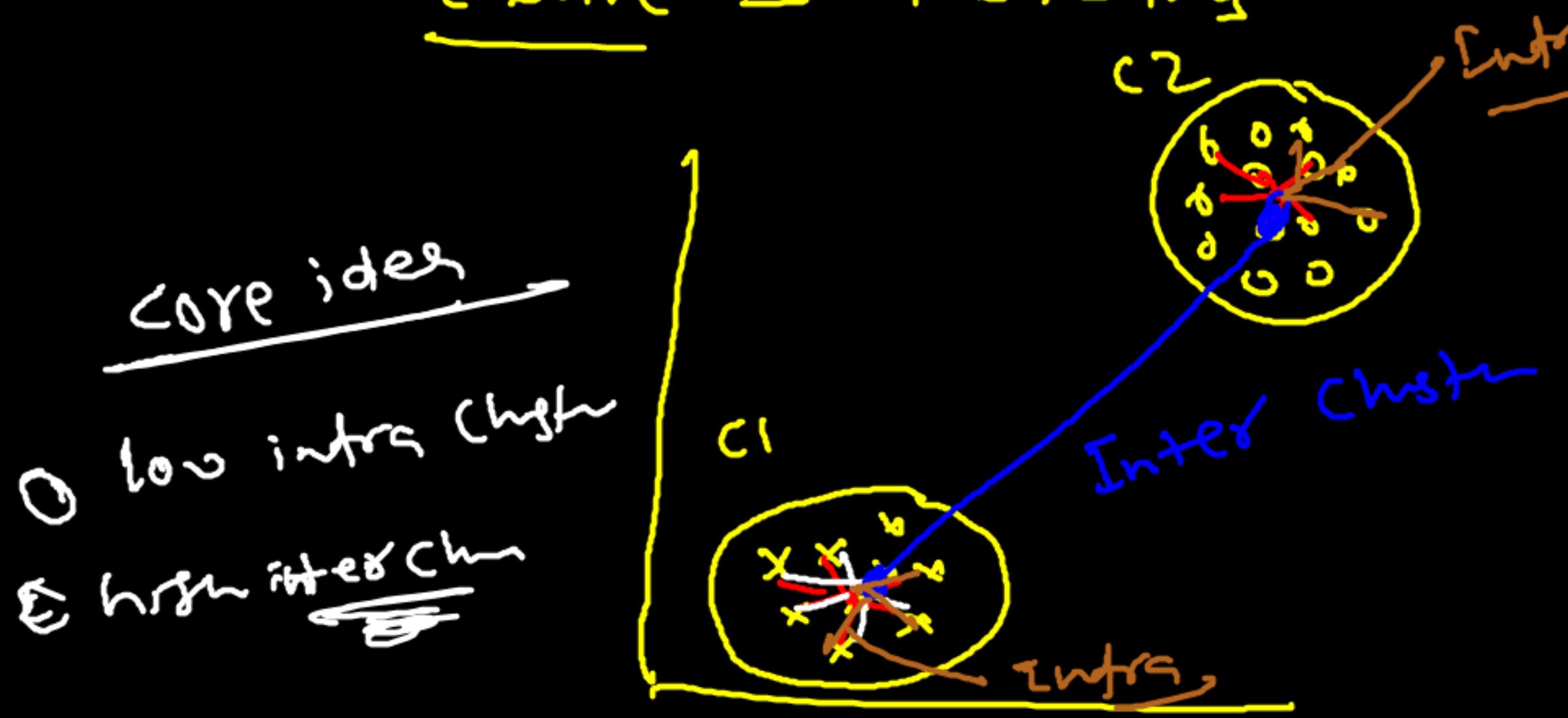
how to measure

Unsupervised ML — K-means, Hierarchical, DBSCAN
X X
General purpose

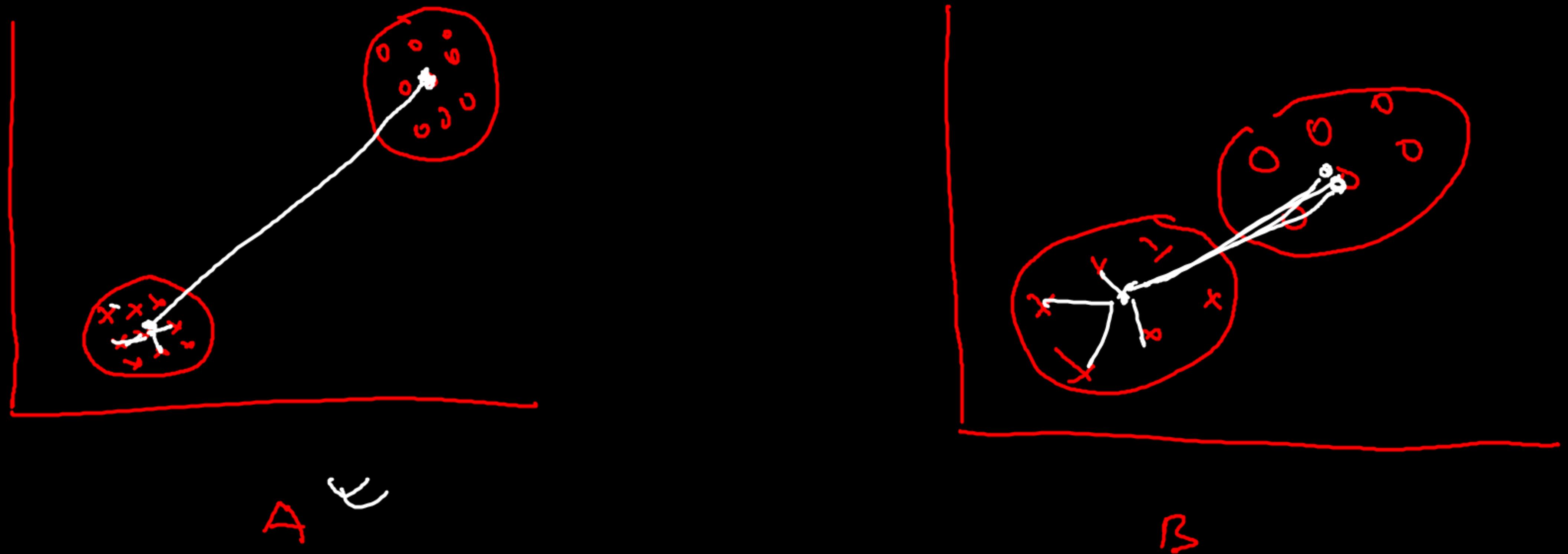
Clustering \rightarrow Data minus topic \rightarrow Def / output / loss / result

$D = \{x_i\}$: There is no y_i variable
Given

Measure - k-means



- ① Intra-cluster :- within a cluster
dist - very low
Small
- ② Inter-cluster :- across or
 between cluster
dist - high - Great



Q:- Which cluster is better?

Ans :- Inter-clust = very high = A \checkmark
Intraclust = dist = 11 Small = A \checkmark

K-Means Cluster

ViratJain raised hand

View

x

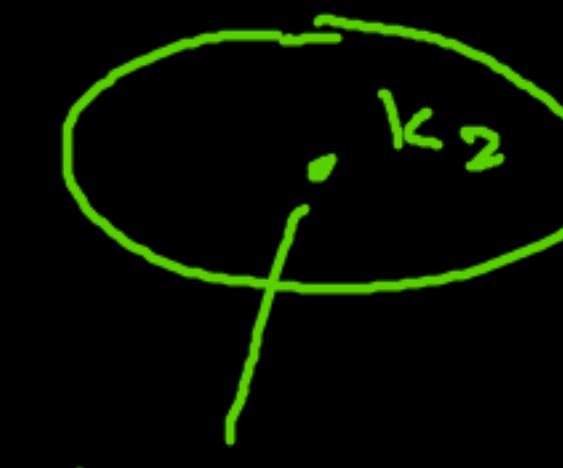
Height Weight

	Height	Weight	
1.	168	60	$\rightarrow K_1$
2.	185	72	$\rightarrow K_2$
3.	170	56	$\rightarrow K_1$
4.	190	78	$\rightarrow K_2$
5.	185	77	$\rightarrow ? \quad K_2$
6.	183	84	$\rightarrow ?$
7.	165	66	$\rightarrow ? \quad K_2$

$$K = 2$$

Elbow method → how many "K" is required

↳ Randomly Selected



$$\text{new cluster (K1)} \\ (168, 60), (170, 56)$$

What is $(170, 56) = \underline{K_1 \text{ or } K_2}$?

$$\frac{168+170}{2}, \frac{60+56}{2}$$

$$\left\{ \begin{array}{l} ED_1 = \sqrt{2^2 + 4^2} = \\ ED_2 = \sqrt{15^2 + 16^2} = \end{array} \right.$$

$$K_1 = (169, 58) \\ (169, 58) \cancel{\in}$$

$$\rightarrow \boxed{190, 78} = K_1 \text{ or } K_2 = ?$$

$K_1 = 169,58$

$K_2 = 185,72$

$$ED_1 = \sqrt{21^2 + 20^2}$$

$$ED_2 = \sqrt{5^2 + 6^2} \quad ED_2 < ED_1$$

$$\frac{190+185}{2}, \frac{78+72}{2}$$

$$\therefore \underline{K_2} = \boxed{187,5, 75}$$

$$\underline{187,5, 75} = ? \quad \begin{matrix} K_1 \\ K_2 \end{matrix}$$

$$K_1 = \sqrt{16^2 + 19^2}$$

$$K_2 = \sqrt{(2,5)^2 + 2^2} \quad \left(\frac{187,5+185}{2}, \frac{75+72}{2} \right)$$

K-Means :- Very very hard to solve mathematically

in computer science \rightarrow "Complexity Theory"



Exponentially Time complexity

"Approximation algorithm" - hard prob

hacks & heuristics

Lloyd's algorithm \rightarrow N.V. Simple

k-means : Lloyd's Algorithm

Rule 1 :- random initialize

→ randomly pick K points from D

$c_1 \ c_2 \ c_3 \dots$

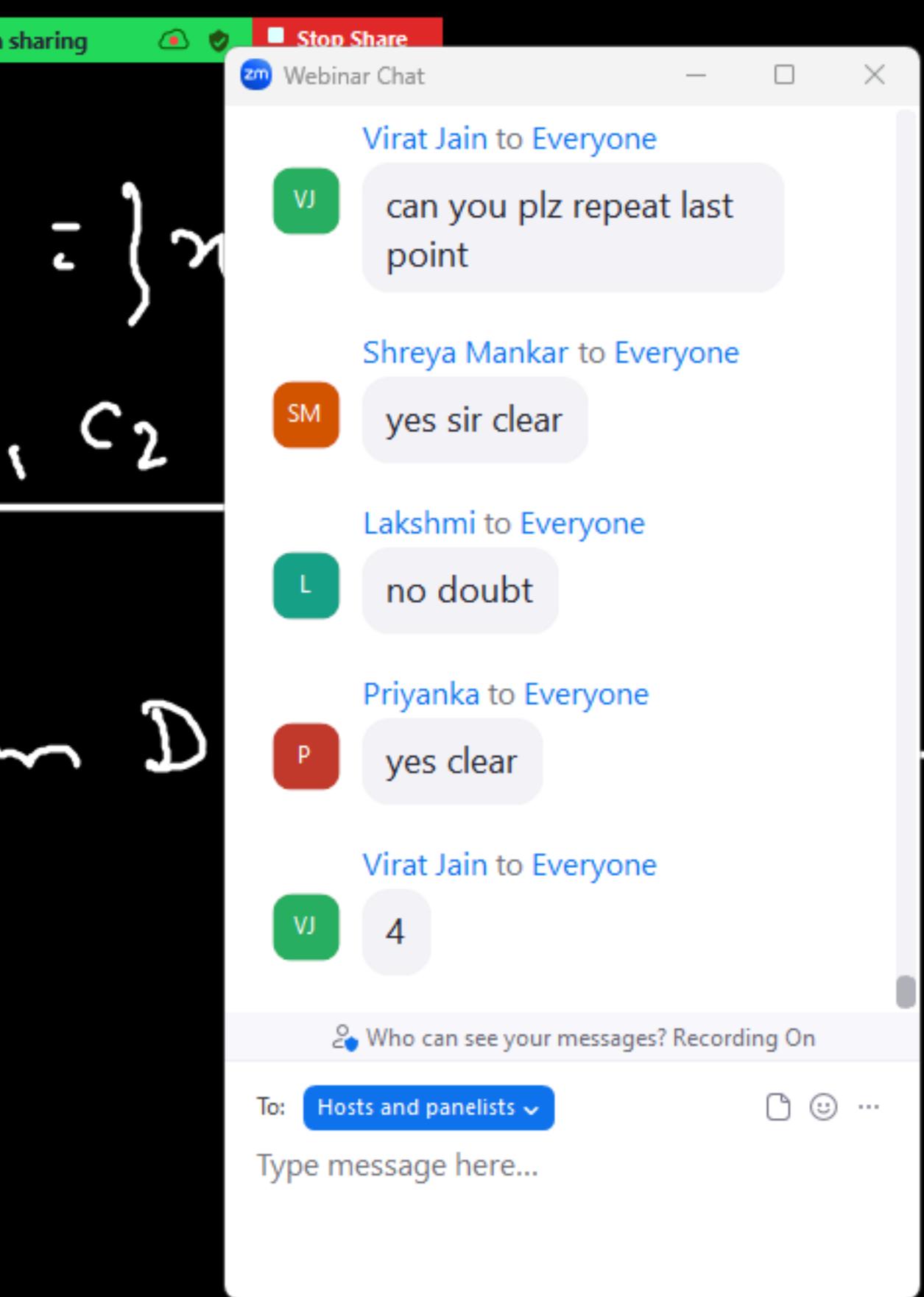
Rule 2 :- Assignment

iteration loop { for each points x_i in D
→ select the nearest one
Eucl-dist

$D = \{x\}$

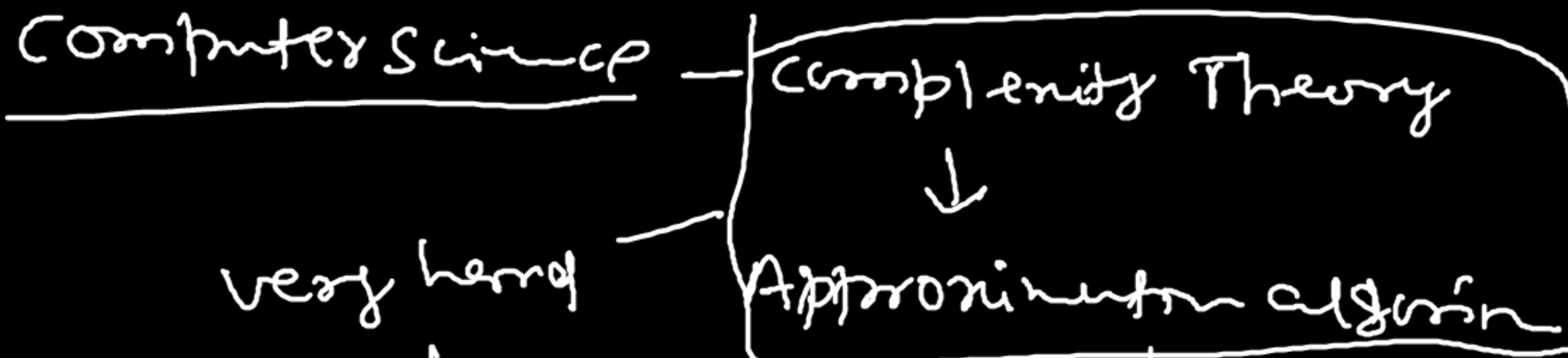
$c_1 \ c_2$

$c_3 \dots$



Rule 3 :- Recompute centers & update tree

Rule 4 :- Repeat step 2 & 3 until



Lloyd's method

dist between new
centers & old centers
is very small.

Virat Jain to Everyone

VJ

4

yatin mistry to Everyone

YM

how to know how
many clusters are there
in random points as per
step 1?

what if random points if
we get from only one
cluster only and not
aware other cluster
points?

Yushma Premchand to Everyone

YP

yes

Who can see your messages? Recording On

To:

Hosts and panelists

Type message here...

You are screen sharing Stop Share

Silhouette Score