

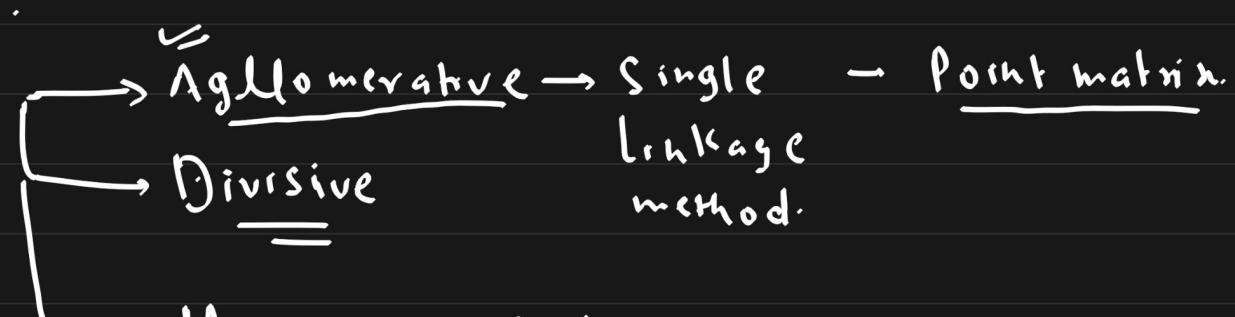
- 1 what is Unsupervised ML
- 2 what is clustering
- 3 Application of the clustering
- 4 Clustering algo

↳ k-Means → Elbow method

diff b/w inter cluster Dis
 { k-mean and } intra cluster dis
 k-mean++ WCSS.

→ how we can evaluate this
clustering.

↳ Hierarchical clustering



the concept of dendrogram
(for representation of cluster)

↳ DBScan clustering

→ eps d.s
- No. of Point
- Core Point

→ Bound, Point
→ Noise Point

Machine learning :-

Supervised

Unsupervised

Semi Supervised

Supervised :-

Data. →

X, y

Independent

Someone for supervising

dependent (target)

(Supervisor)

Machine Learning

Lm-reg ✓

Log reg ✓

Svm ✓

kNN ✓

DT ✓

RF ✓

AB ✓

GB ✓

XGBoost ✓

Nom.

Reg.

Cat.

Num.

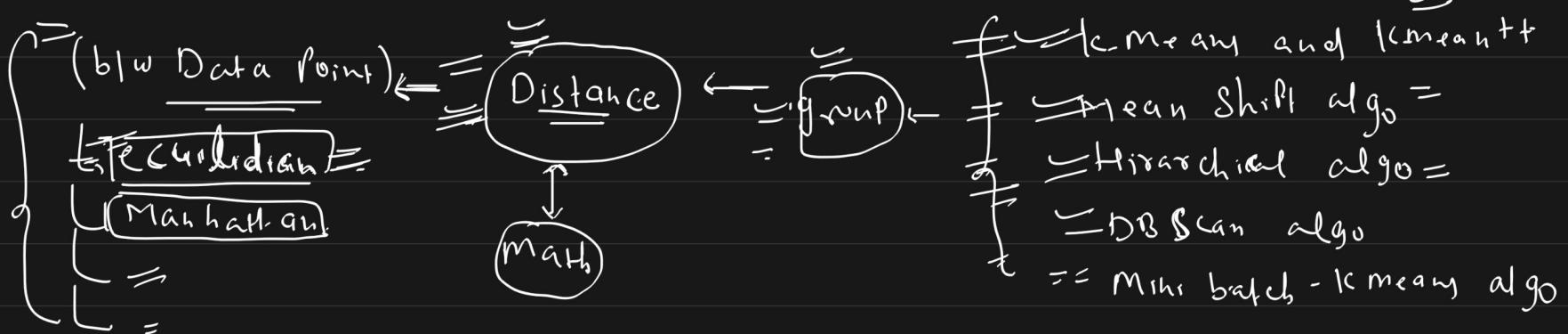
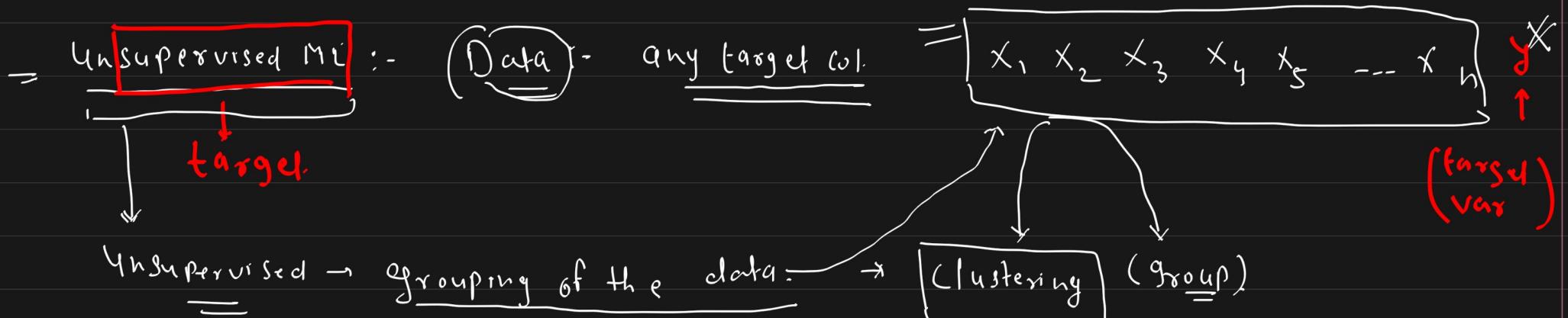
(Data.)
Supervised

&

Many - end to end
unsupervised Project

{
① Data. ✓
② EDA ✓
③ Preprocess ✓
④ Model ✓
⑤ Evaluation ✓

Supervised :- Regression
Classification.



\Rightarrow Clustering \rightarrow grouping of data \rightarrow Distance \rightarrow Ecu.
 \rightarrow Manf.

Supervisor \rightarrow Data - group - y =
 \Rightarrow x

Application:- Real time app:-

(1) Customer Seg. \rightarrow E-commerce \rightarrow Similarity \rightarrow Amazon \rightarrow Customer

(2) Image Seg.

(3) Data seg. (Sentiment Data)

↳ Review [Pos]
[Neg]
[Nut]

{ Female
Male
Payment.
Stuff
Review }

Segregate as similar pixel



Pixel Segmentation
grouping of the pixel.

(3)

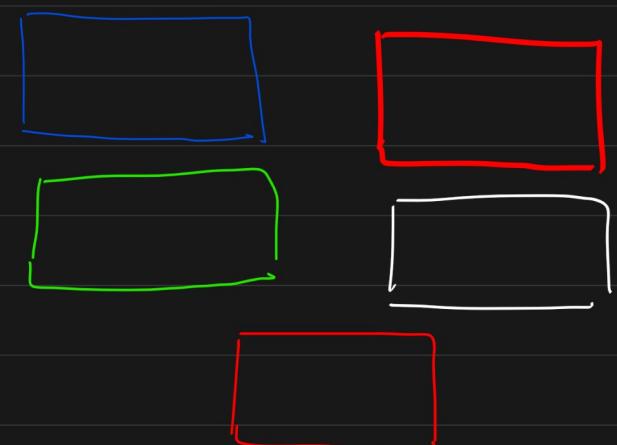
Structure for set :- Medical.

Agg.
education.

$(m \times n)$

$m = 1000$
 $n = 10$

$x_1 \ x_2 \ x_3 \ x_4 \ \dots \ x_n$

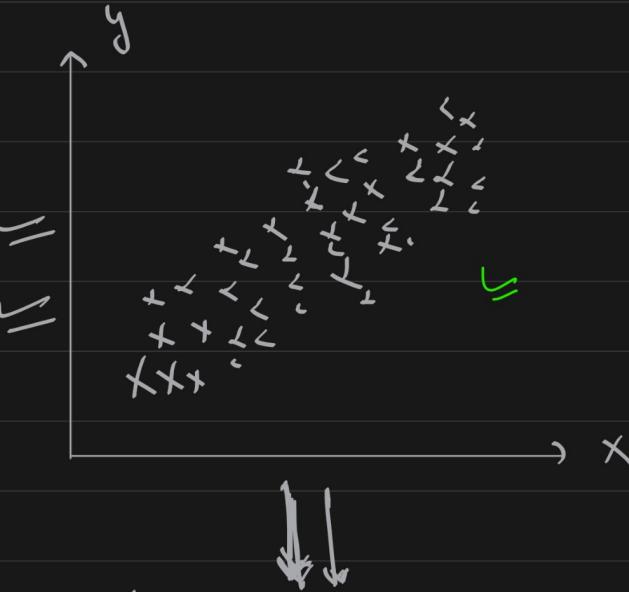


Clustering algo:- K-Means

(i) Concept of the centroid.

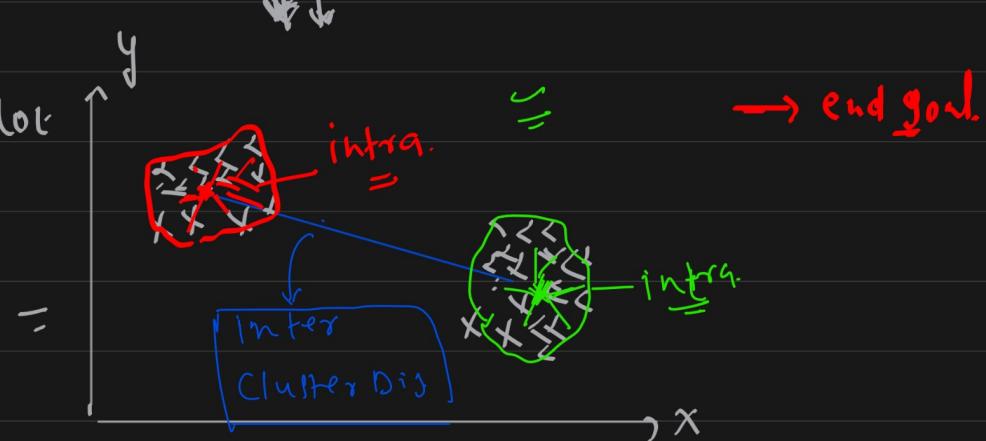
Distance

$\left\{ \begin{array}{l} \text{Wtch in cluster} \\ \text{blw cluster} \end{array} \right\}$ ~~(i)~~ intra. }
 $\left\{ \begin{array}{l} \\ \end{array} \right\}$ ~~(ii)~~ inter }

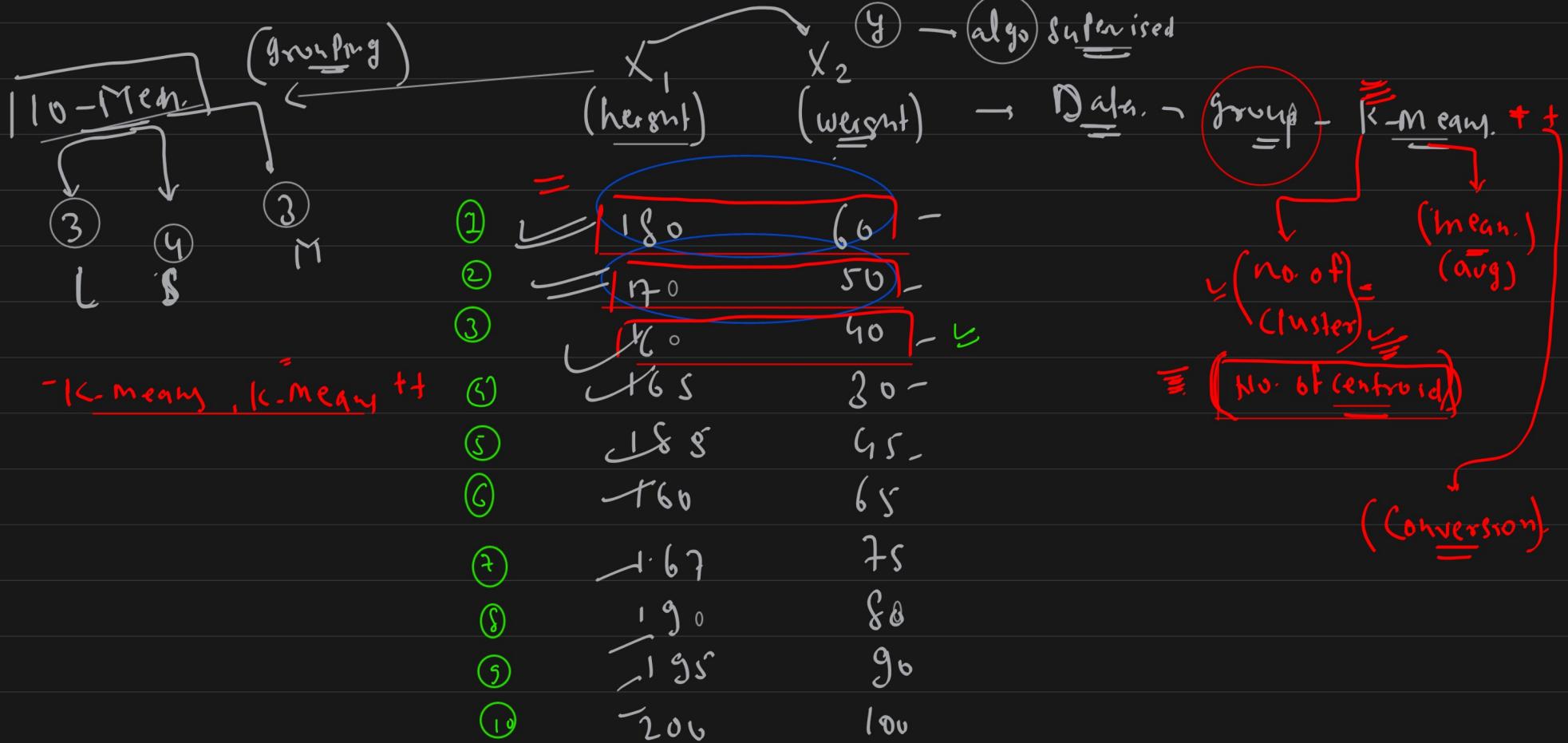


Evaluation of cluster \leftarrow $\begin{cases} \text{WCS} \rightarrow \text{elbow plot} \\ \text{Shallowee} \end{cases}$

Intra and
inter ! many more



K-Means algo :- Data \rightarrow group \rightarrow cluster - Supervisor not having any



= ① Calculate a centroid (data point) = (Random).

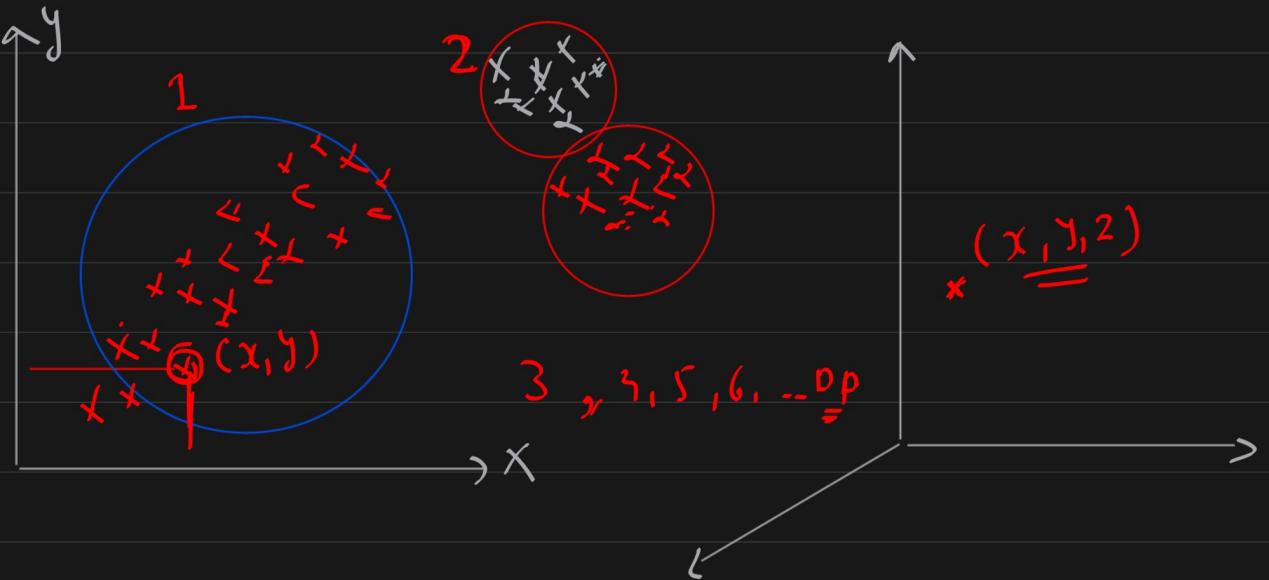
Minimum :- One cluster

Hierarchical

Maximum :- No. of data points (No. of Row) =

$$\text{L-1} \quad 20 \quad \text{3 cl}$$

$$(x_1, x_2, x_3, x_4) \quad \underline{\underline{(x_1, x_2)}} \quad \underline{\underline{(x_1, y_1, z)}}$$



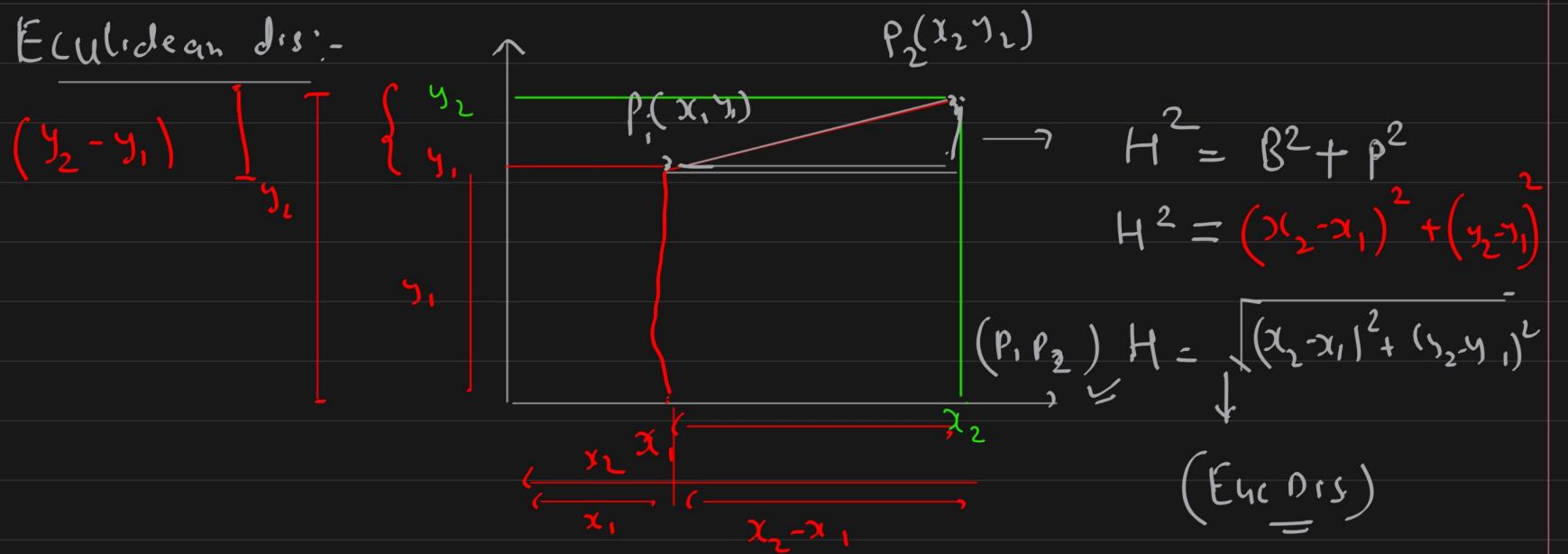
Minimum two :- $K = \text{No of centroid.}$ $\underline{\underline{K=2}}$

(Randomly)

$$C_1(180, 60)$$

Distance

$$C_2(170, 50)$$



1.) initialize a centroid
(Randomly)

- 2) Distance
3) fill the cluster
4) update the centroid.

$$c_1(180, 60)$$

$$3 \text{ Point} \rightarrow (160, 90)$$

$$= c_1 + 3 \sqrt{(160 - 180)^2 + (90 - 60)^2} = \approx (28 \cdot 28)$$

$$= c_2 + 3 \sqrt{(170 - 160)^2 + (50 - 90)^2} = \approx 14 \cdot 14$$



$$\bar{c}_1 = \left(\frac{180 + 160}{2}, \frac{60 + 40}{2} \right) = \left(170, 50 \right)$$

K-means

K-mean init =

mini batch - K-means.

$$c_1 = \bar{c}_1 = \left(\frac{180 + 160}{2}, \frac{60 + 40}{2} \right) = \left(170, 50 \right)$$

$$c_2 = \bar{c}_2 = \left(\frac{165 + 115}{2}, \frac{45 + 30}{2} \right) = \left(145, 37.5 \right)$$

update

$$c_1 \text{ to } q = \sqrt{(180 - 165)^2 + (60 - 45)^2} = 33.54$$

$$c_2 \text{ to } q = \sqrt{(165 - 115)^2 + (45 - 30)^2} = 52.90$$



update?

$$\bar{c}_1 = \left(\frac{180 + 160}{2}, \frac{60 + 40}{2} \right) = \left(170, 50 \right)$$

$$\bar{c}_2 = \left(\frac{165 + 115}{2}, \frac{45 + 30}{2} \right) = \left(145, 37.5 \right)$$

$$\left(\frac{165+175}{2}, \frac{45+30}{2} \right)$$

$$C_1(180, 60) \quad \left(165, 37.5 \right)$$



Elbo Plot :- \Rightarrow K-Means or K-Means++

No. of cluster how we can decide K-value?

$k=1$

$k=2$

$k=3$

$k=4$

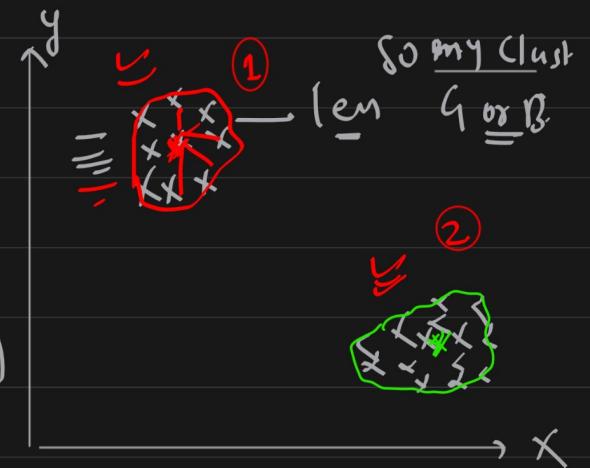
⋮

$k = \text{No. of Data Points}$

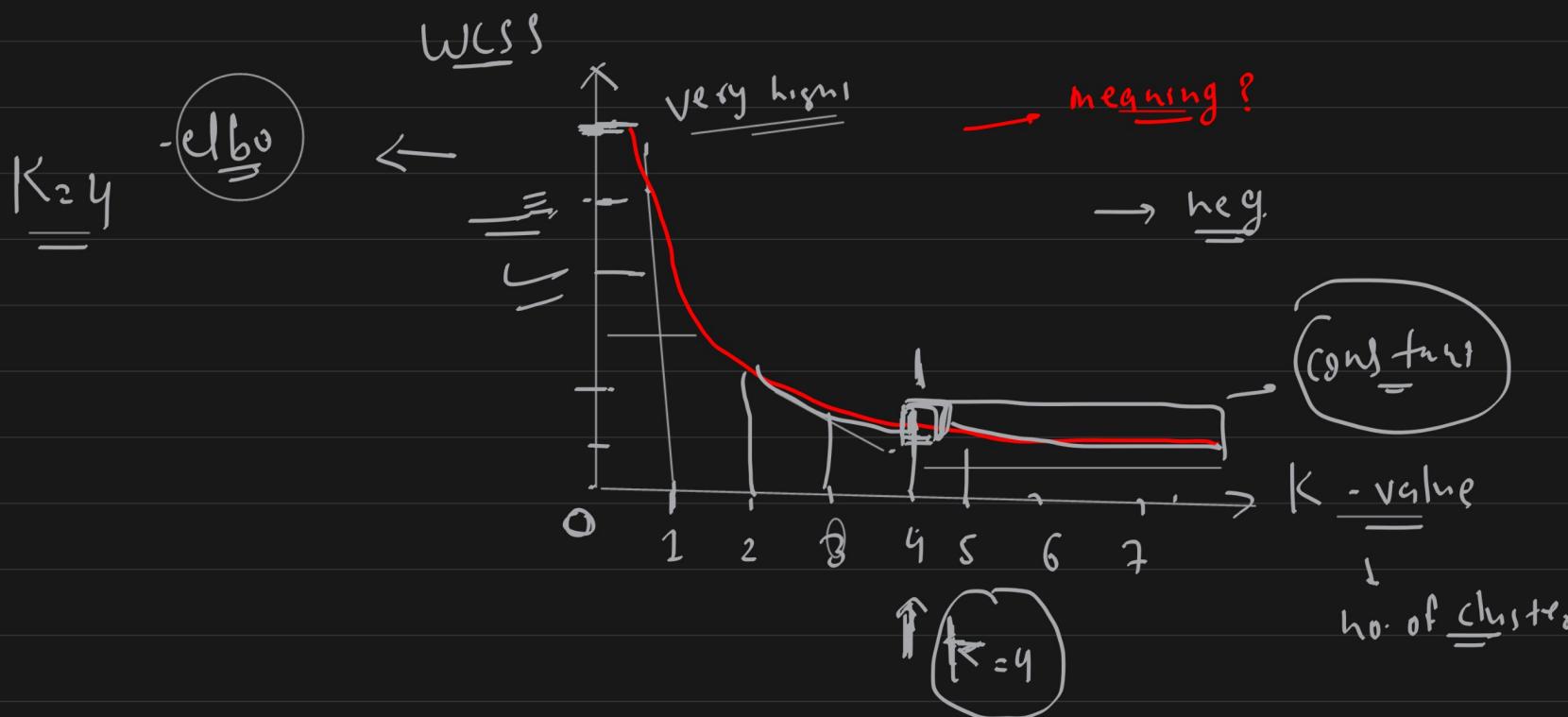
Distance \rightarrow intra and inter

Within cluster (WC)
 \equiv from centroid

WCSS → Within Cluster Sum of
(dense) \sum square
len intra cl.s (within clusters)
↓
good.



More dis → Bad (Sparse)



Data Point :-

①



WCSS \rightarrow Σ



WCSS \rightarrow [len]

= Evolution $k=4 \rightarrow$ correct?
Some other method

Over there (ratio) $\frac{\text{Within}}{\text{Between}}$ and inter
WC below cluster

- 1 Dunn index.
- 2 Silhouette Score
- 3 Jaccard index

WCSS (Matrix) $k=4 \rightarrow$ Validation - Schonette, Dunn, ...
Intra, inter = Mathematical - How much strong. (it is good or bad)

(Semi Supervised): all dots

Hierarchical

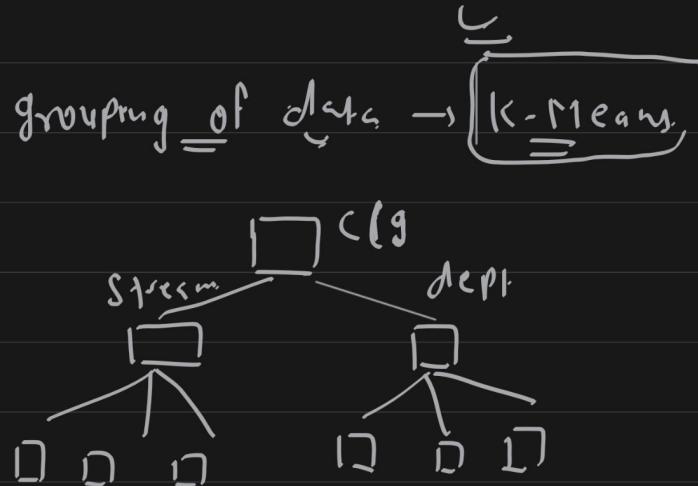
DBScan

Hierarchical clustering \rightarrow grouping of data \rightarrow k-Means

Hierarchy of cluster

↓
Cluster
↓

Distance
(Euclidean)



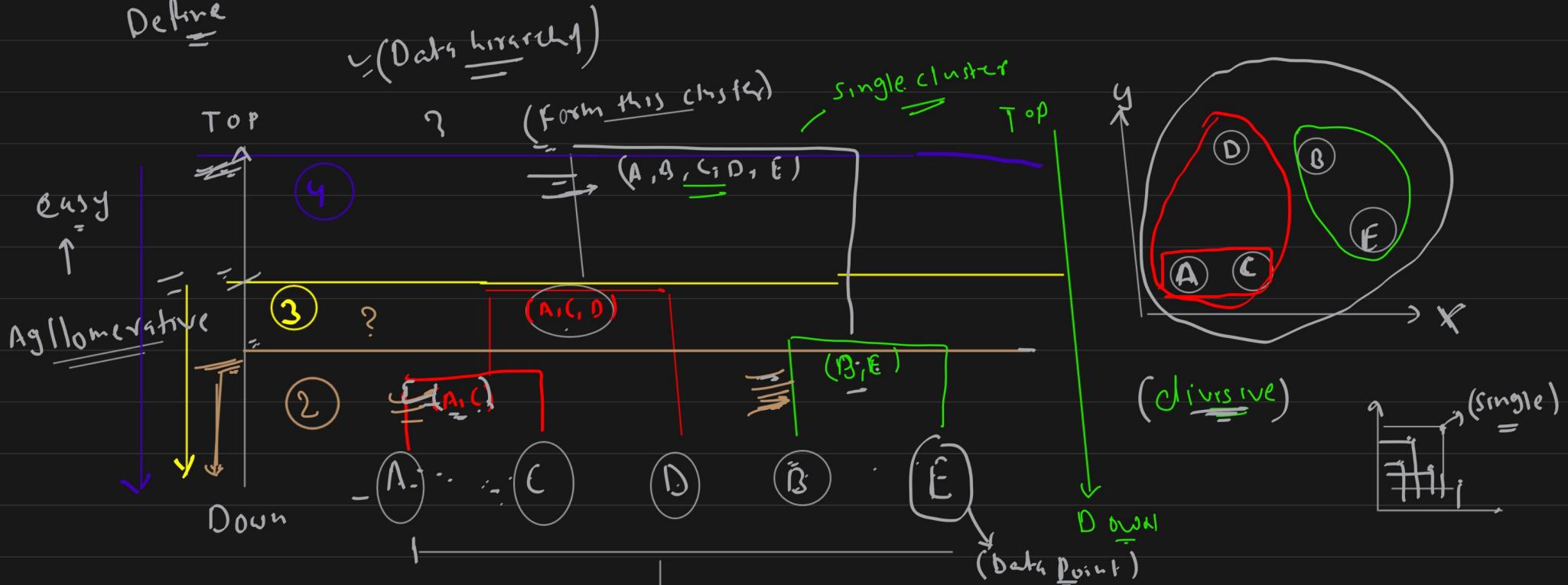
- Hierarchical :-
 = Agglomerative clustering → bottom to top → No. of cluster =
 = Divisive clustering → top to down → single cluster
 No. of data point

= Data :- N of row :- minimum cluster:- 1
 maximum cluster! No of data Point

Dendrogram :- Diagram which is representing cluster

K-Means and hierarchical

Define



↓

(single cluster) (individual)

practical

Hierarchical :- only once define hierar?

$$\begin{matrix} x_1 & x_2 \\ (\text{H}) & (\omega) \\ \swarrow & \searrow \\ \boxed{\begin{matrix} 180 & 10 \\ 1120 & 50 \\ 160 & 40 \end{matrix}} \end{matrix}$$

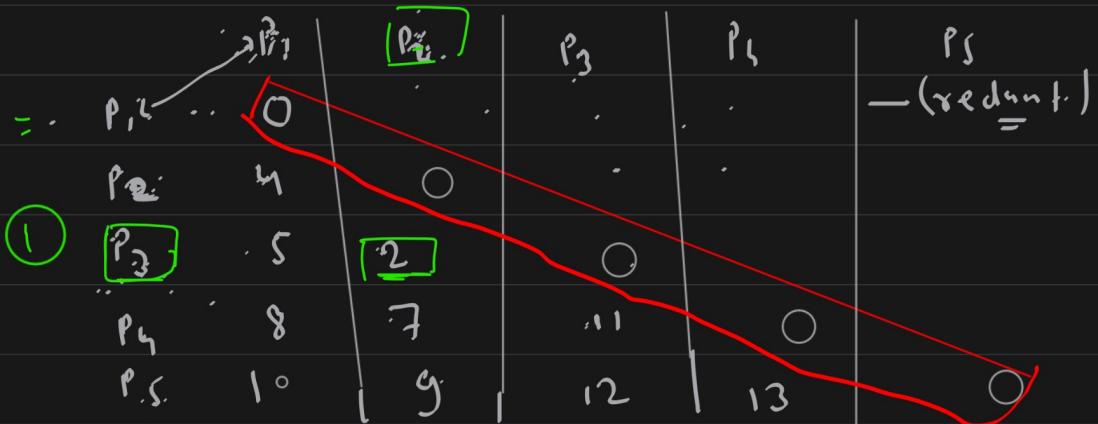
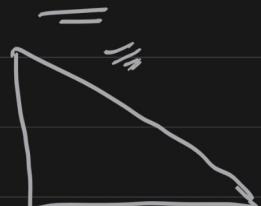
(24, 5, 6, 7, 8, ...)

	x_1	x_2
P_1		
P_2		
P_3		
P_4		
P_5		

Point matrix :-

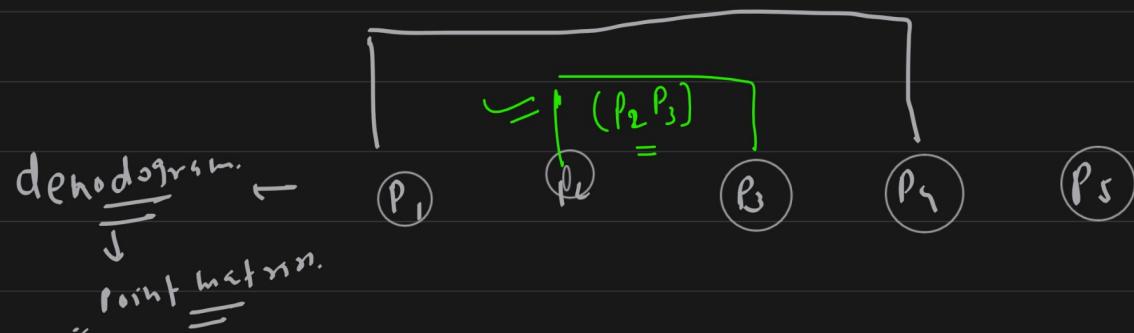
Assuming →

Heat map



matrix
(distance)

↓
Dendrogram



	ρ_1	(ρ_2, ρ_3)	ρ_4	ρ_5
ρ_1	○			
(ρ_2, ρ_3)	g	○		
ρ_4	8	18	○	
ρ_5	10	25	13	○

$d(\underbrace{\rho_1, (\rho_2, \rho_3)}_{d((\rho_1, \rho_2), (\rho_1, \rho_3))}) - \text{single linkage}$

$d(\underbrace{\rho_4, \rho_5}_{d(g, \rho_5)})$

$d(g)$

$d(\underbrace{\rho_4, (\rho_2, \rho_3)}_{d((\rho_4, \rho_2), (\rho_4, \rho_3))})$

$d(7, 11)$

18

$d(\underbrace{\rho_5, (\rho_4, \rho_3)}_{d((\rho_5, \rho_4), (\rho_5, \rho_3))})$

$d(\rho_5, \rho_4)$

$d(13, 12) = 25$