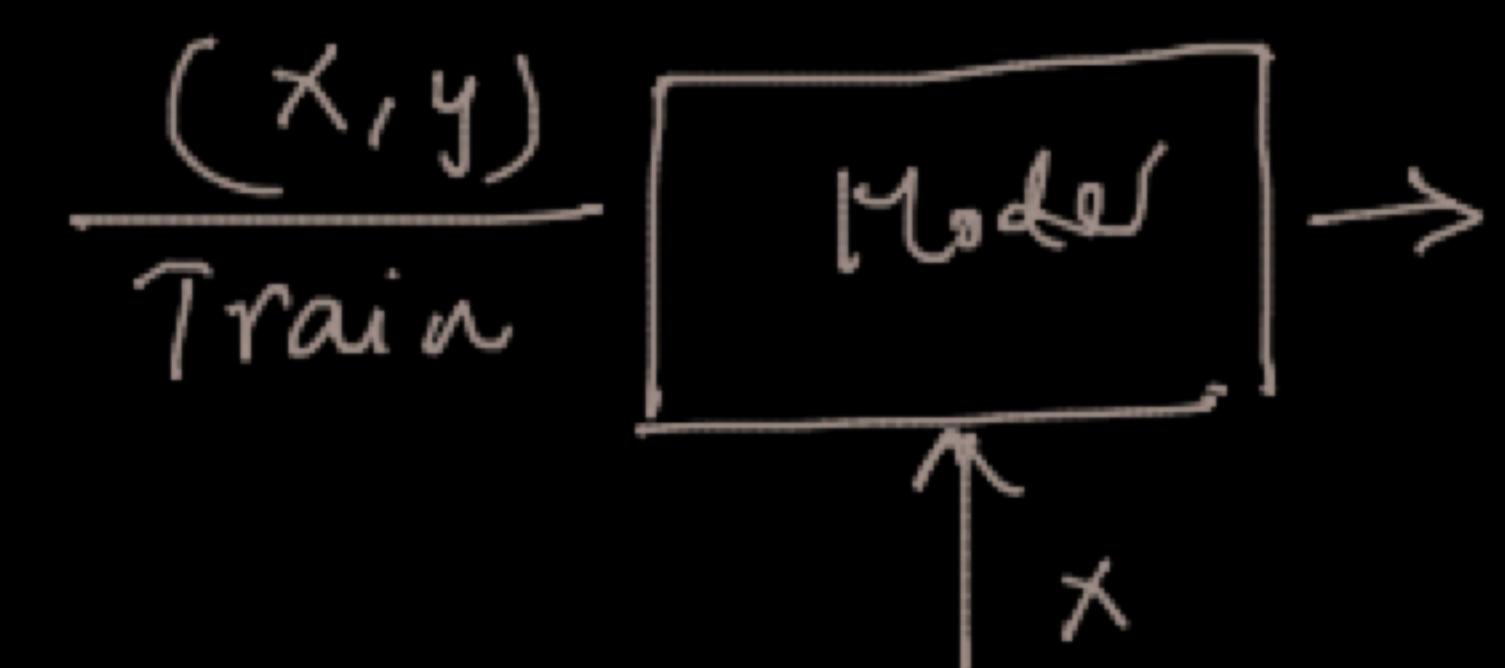


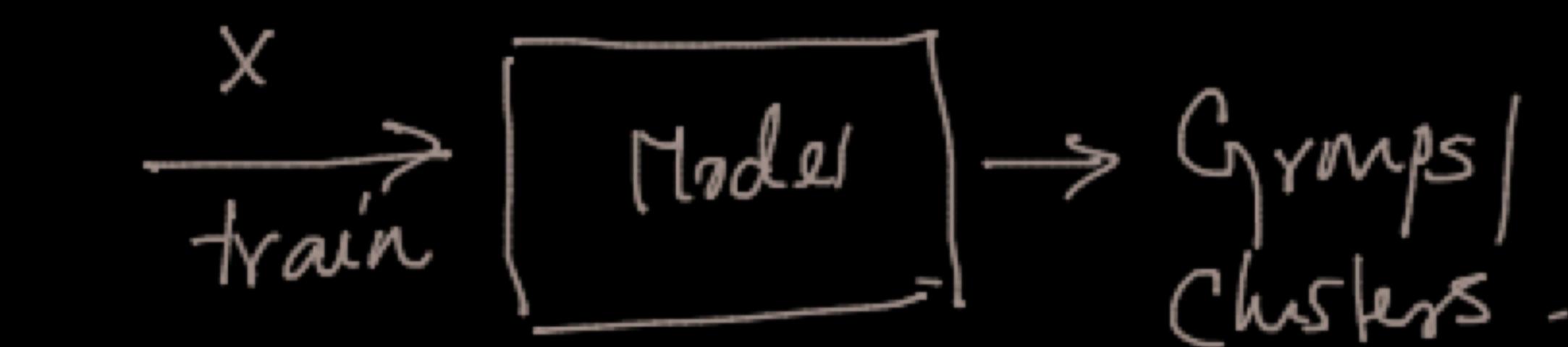
— Model is trained using both x & y .



Regression:

— 'y' is Continuous.
(Linear Regression)

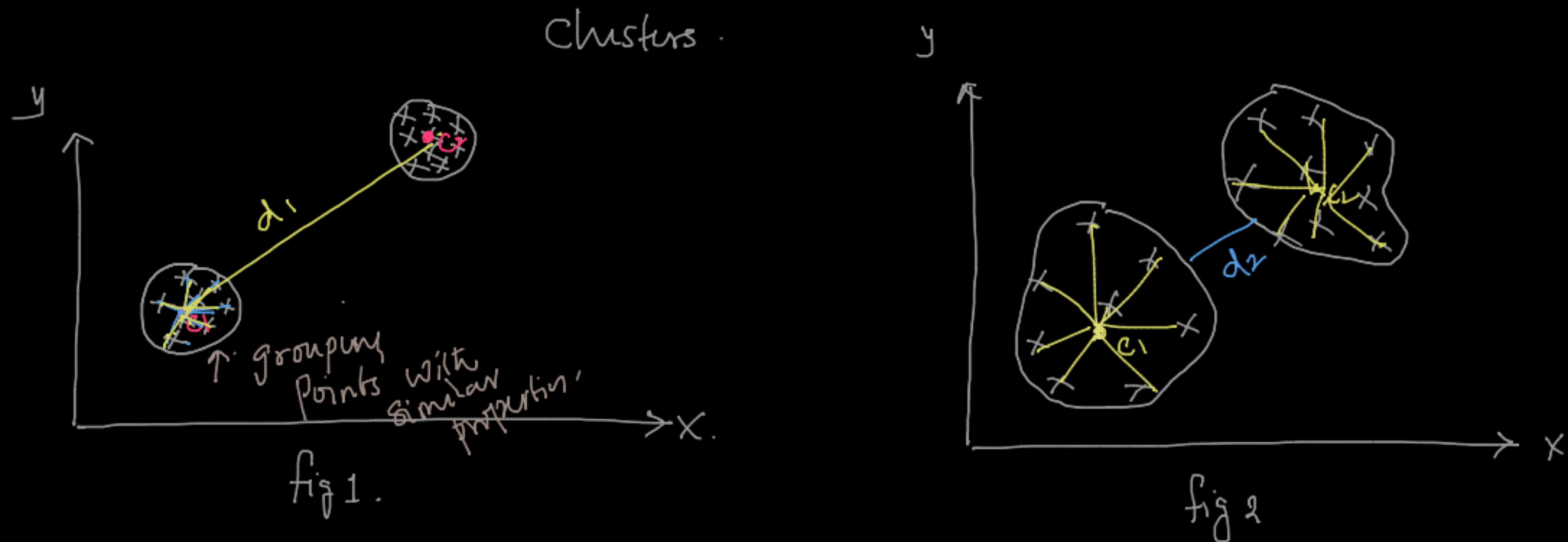
— Model is trained using only x -values.



Classification

— 'y' is categorical
(Logistic Regression)

- k-Means
- Hierarchical
- DBSCAN



Good clusters :

1. Intra cluster distance - Minimum.
2. Inter cluster distance \rightarrow Maximum.

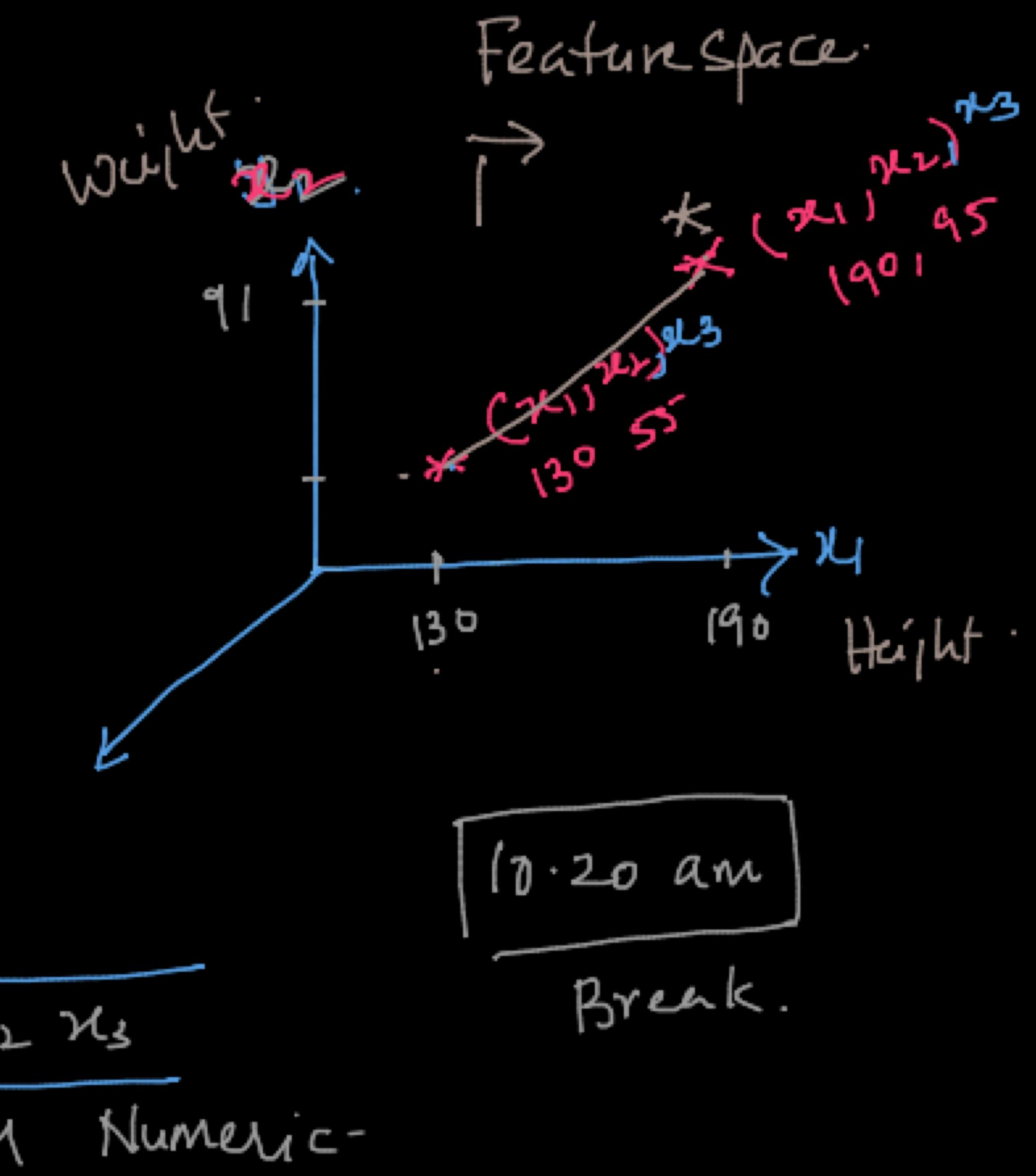
x_1	x_2	x_3
Height	Weight -	we

Similarity $\propto \frac{1}{\text{distance}}$

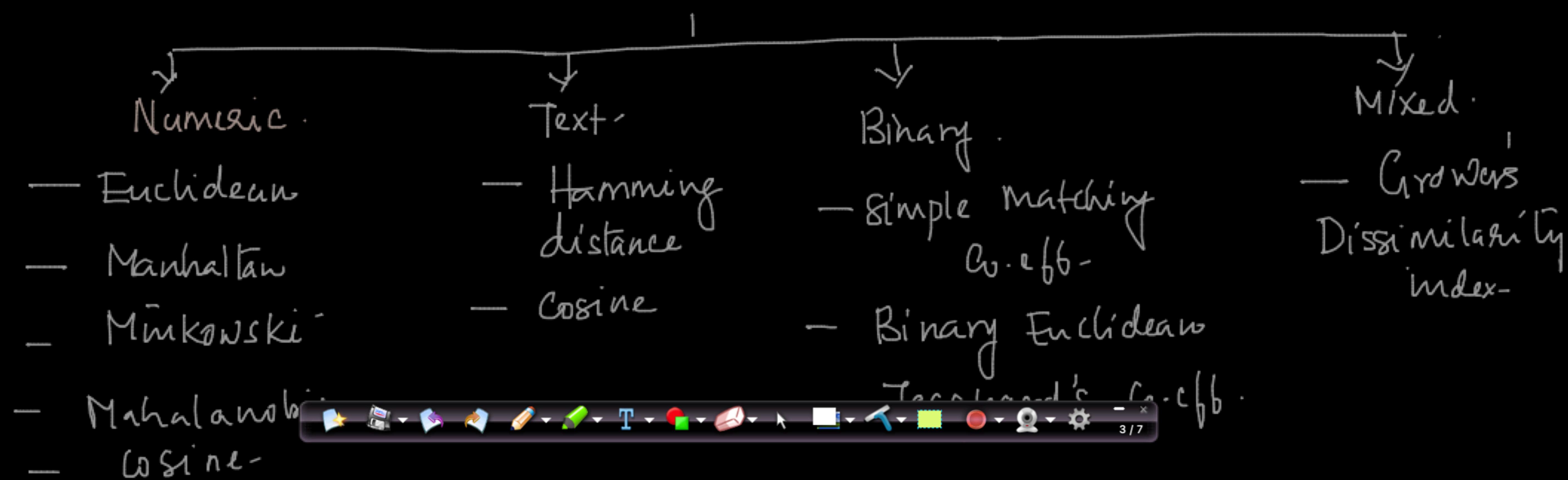
→ P ₁	190 ✓	·95	·63 ←	·
→ P ₂	130	55	95	· · ·
P ₃	189	100		

↳ High dimension

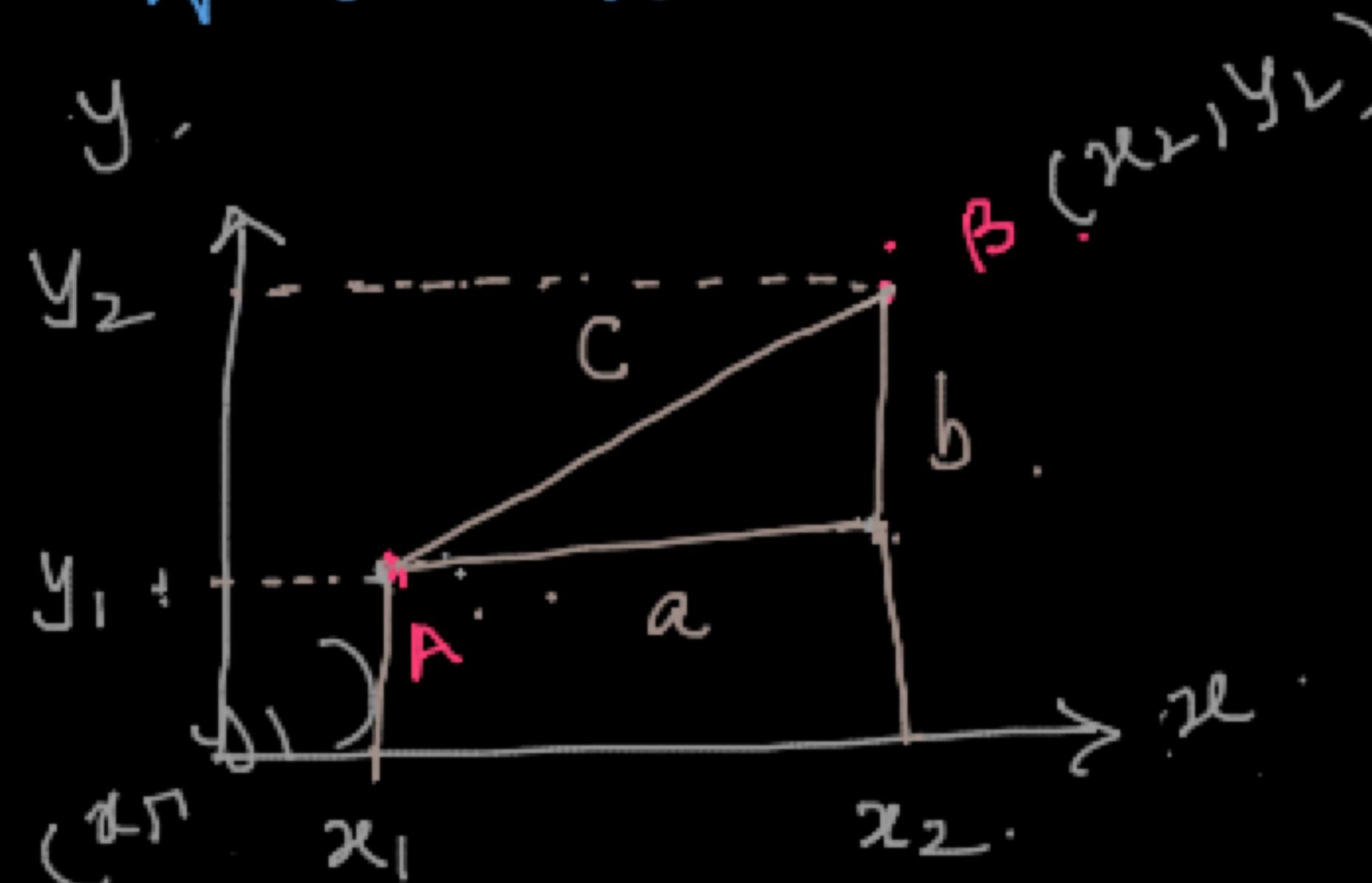
Binary.	
Voted.	Cured
Yes.	0
No.	1.



Distance Measures



Euclidean distance ✓

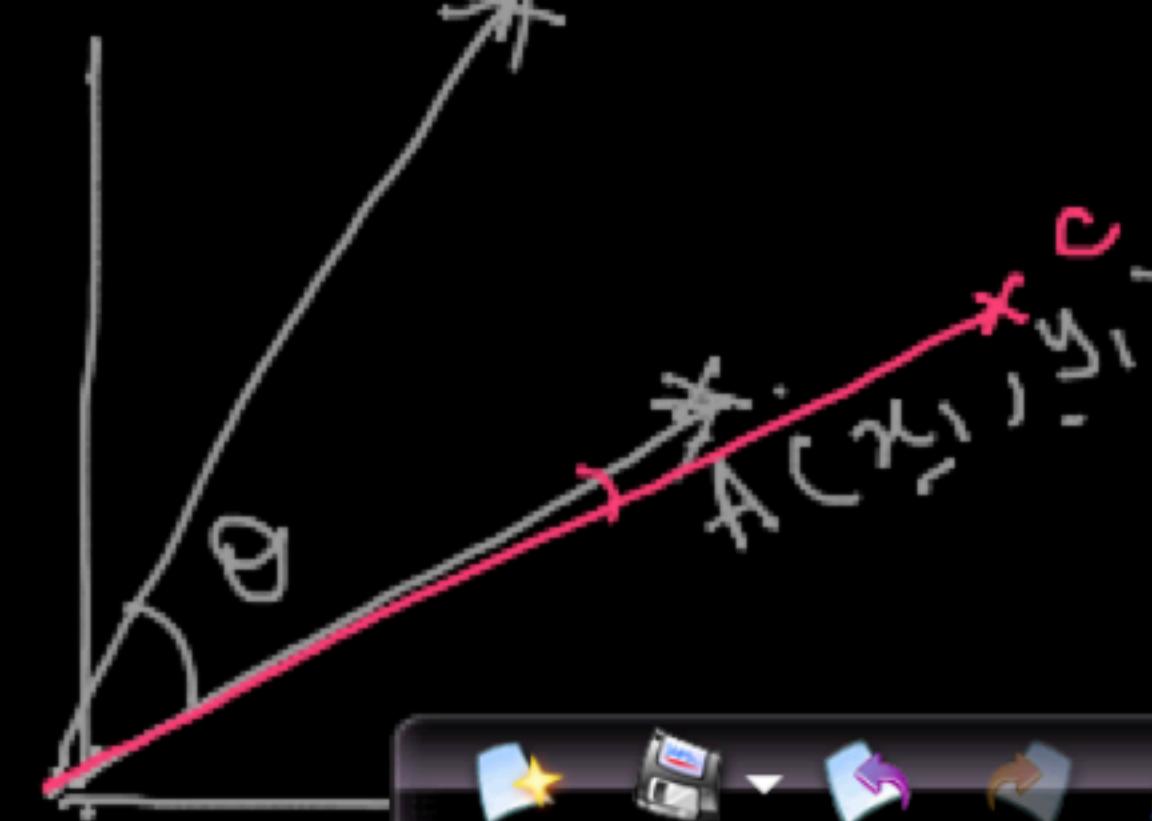


$$c = \sqrt{a^2 + b^2}$$

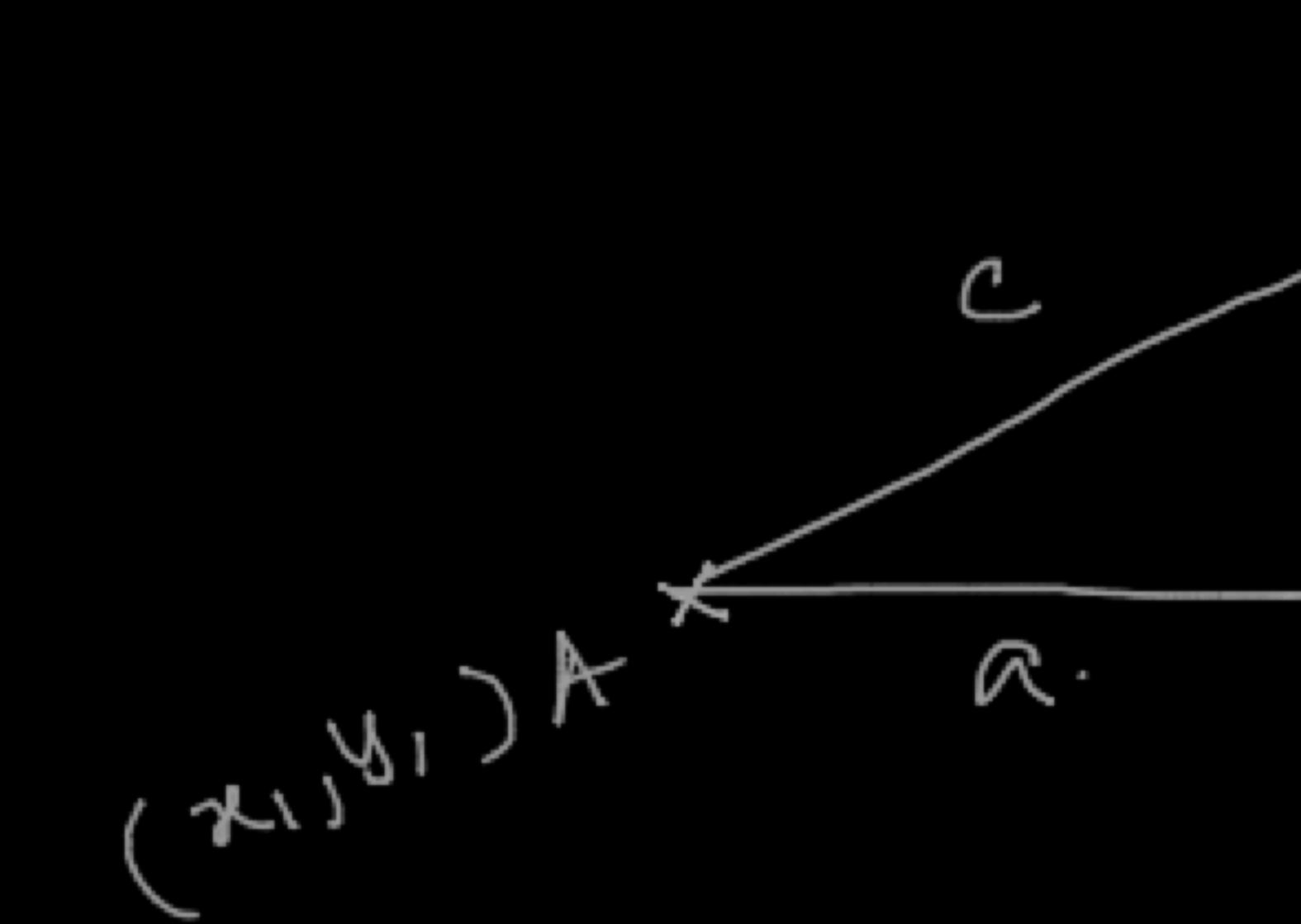
$$c = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$c = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2 + \dots}$$

Cosine Distance



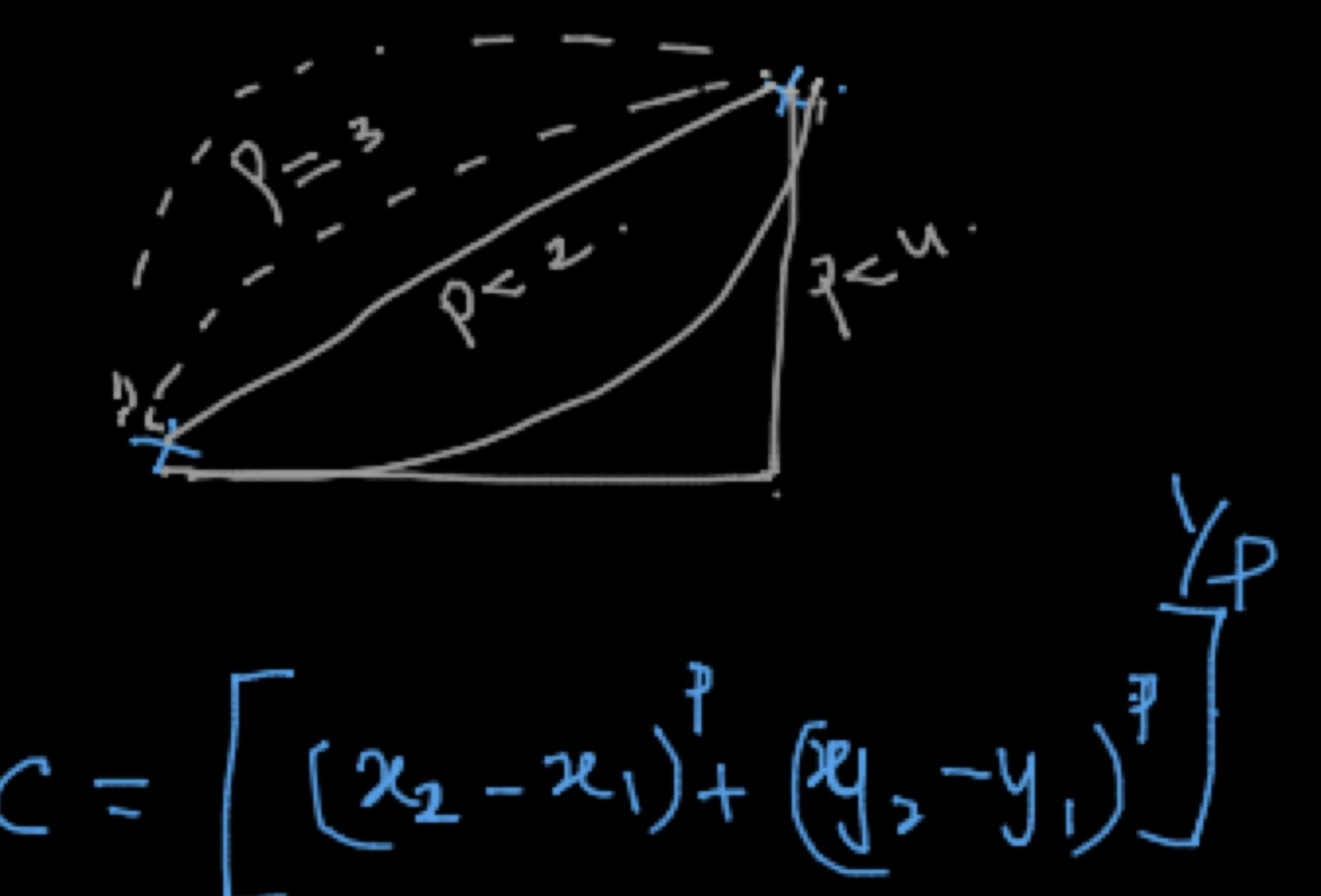
Manhattan Distance (y₂)



$$c = a + b$$

$$c = (x_2 - x_1) + (y_2 - y_1)$$

Minkowski



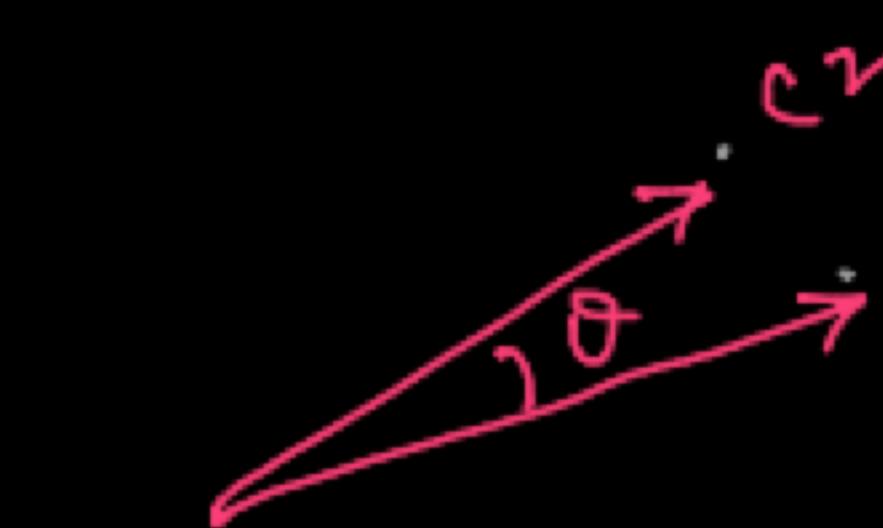
$$c = \left[(x_2 - x_1)^p + (y_2 - y_1)^p \right]^{\frac{1}{p}}$$

$$p = 1, 2, 3, \dots$$

$$c = (x_2 - x_1) + (y_2 - y_1)$$

p = 2 - Euclidean

$$c = \left[(x_2 - x_1)^2 + (y_2 - y_1)^2 \right]^{\frac{1}{2}}$$



$$\theta = 0$$

90°

$$\cos(\theta) \Rightarrow 1$$

0

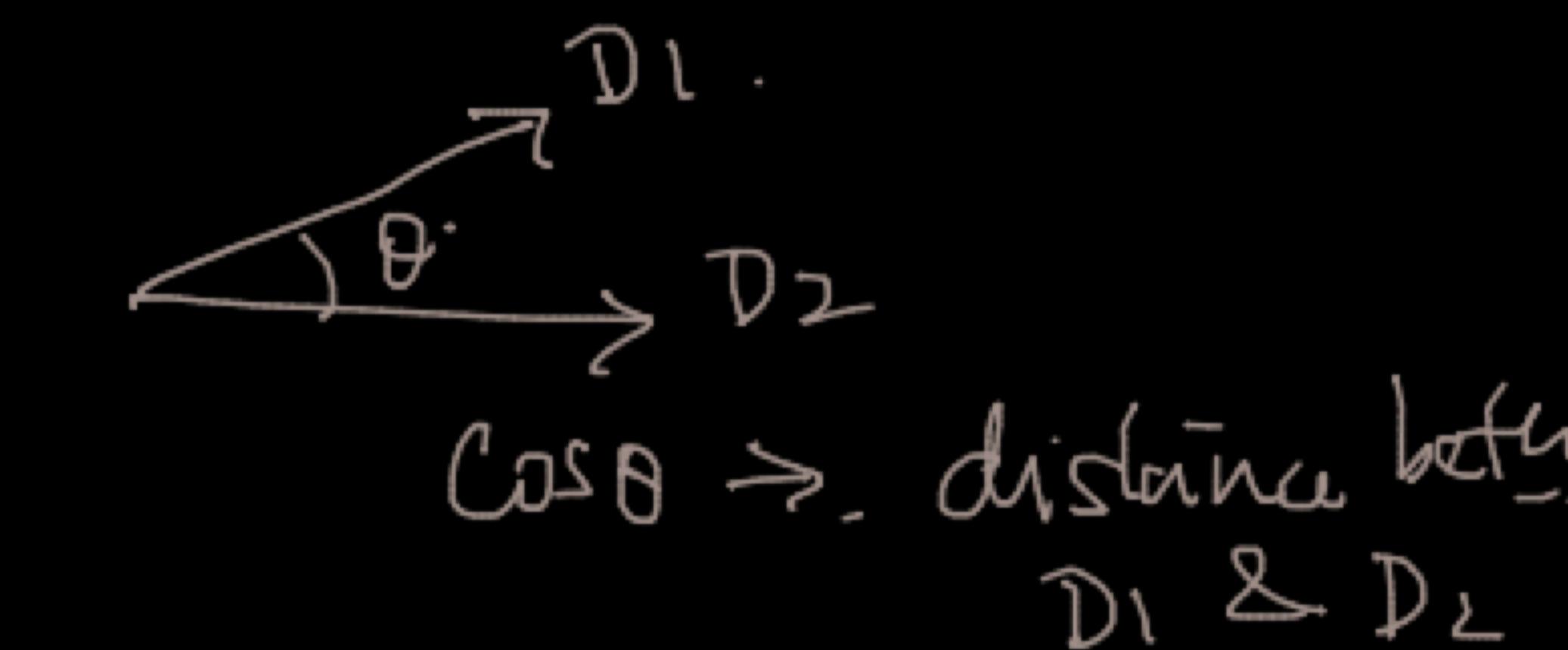
Text Values.

Corpus. {
 $D_1 : \text{The car is driven on the road} \rightarrow \text{document}$
 $D_2 : \text{The truck is driven on highway} \rightarrow \text{document}$

Vocab: The car truck is driven on road highway

$D_1 :$	1	1	0	1	1	1	0
$D_2 :$	1	0	1	1	1	0	1
			1			1	

Hamming distance = 2



Scaling \rightarrow Bring all columns
to the same range.

	SAT	Top 10	AR.	SFR	Expenses	GradRate	
\hookrightarrow	1310	29	22	6	450000	94	—
	1200	83	34	25	10000	73	—

- 1. Standardize.
- 2. Normalize.

$$\begin{array}{cc}
 x & y \downarrow \\
 \text{SFR} & \text{Expenses} \\
 \xrightarrow{x_1} & \xrightarrow{y_1} 450000 \\
 \xrightarrow{x_2} & \xrightarrow{y_2} 10000 \\
 \xrightarrow{13} & \xrightarrow{20000}
 \end{array}$$

Euclidean

$$\begin{aligned}
 d_1 &= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \\
 &= (25 - 6)^2 + (35000 - 45000)^2 \\
 &\Rightarrow [(19)^2 + (35000)^2]^{1/2}
 \end{aligned}$$

$$(13-6)^2 + (25000)^2$$

$$\left[\underset{\uparrow}{361} + \underset{\rightarrow}{1225000000} \right]$$

$$49 + \rightarrow$$

✓ ✓
 SFR. Expenses.

$$6 \rightarrow 45K.$$

$$13 \rightarrow 20K$$

$$25 \rightarrow 10K$$

$$10 \rightarrow 33K$$

$$\frac{\bar{x}_{SFR}}{\sigma_{SFR}}$$

SFR.	Exp.
-2.3	2.57
1.7	0.15
0.85	-1.32

Standardize → StandardScaler ✓

Z_SFR.

Z_Expenses.

$$\frac{6 - \bar{x}_{SFR}}{\sigma_{SFR}} = -3$$

$$\frac{45K - \bar{x}_{Exp}}{\sigma_{Exp}} = -3$$

$$\frac{13 - \bar{x}_{SFR}}{\sigma_{SFR}} = +3$$

$$= +3$$

Normalize

$$\frac{x - x_{min}}{x_{max} - x_{min}} \quad \text{MinMax Scaler} \checkmark$$

$$\frac{6 - 6}{25 - 6} = 0 \rightarrow 0$$

$$0:25$$

$$0:37$$

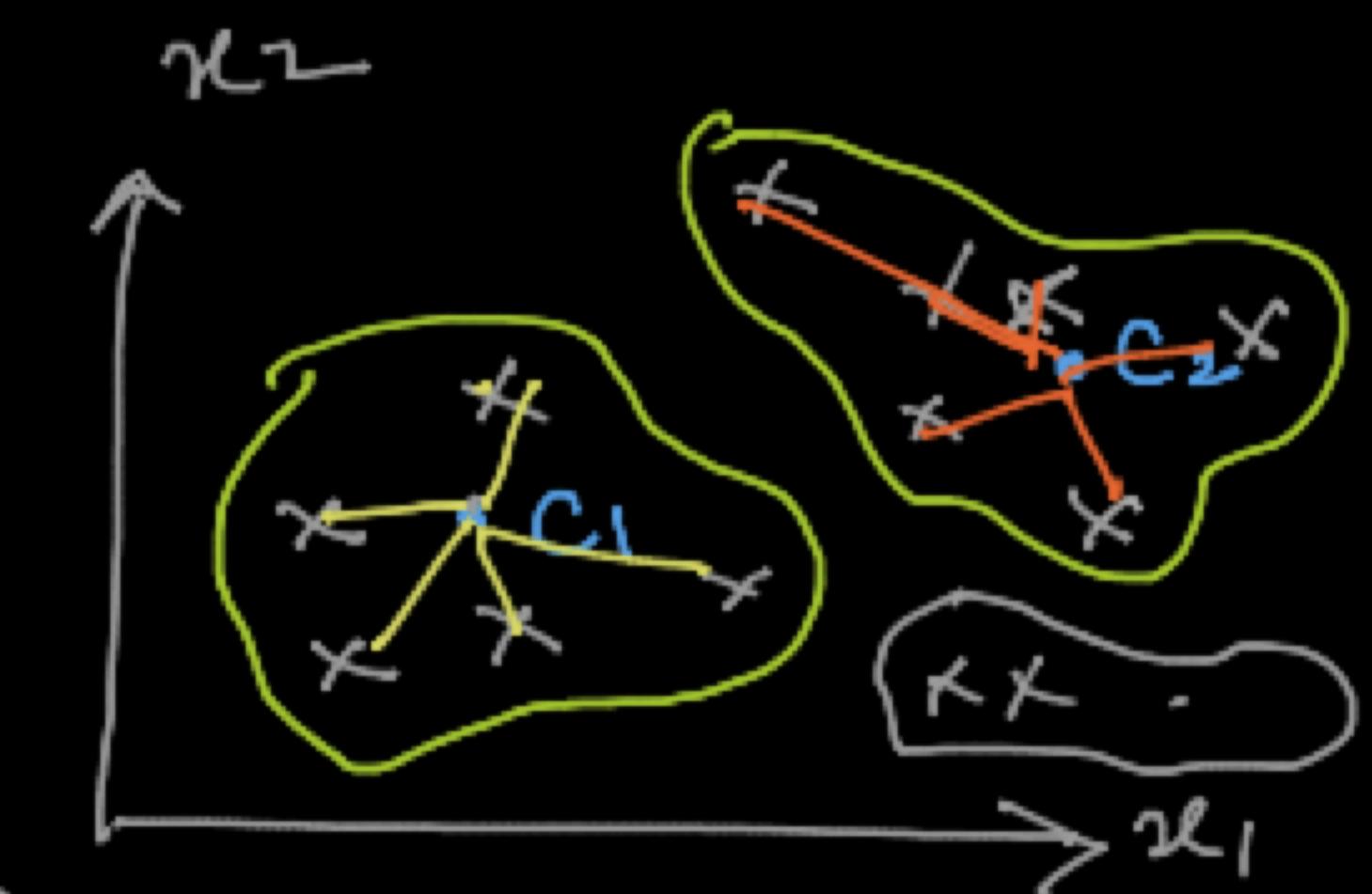
$$\frac{25 - 6}{25 - 6} = 1 \rightarrow 1$$

	✓	✓
	x_1	x_2
1	.	
2	.	
n	.	

K-Means Clustering

No. of groups $\leftarrow K = 2$
 (Hyperparam)

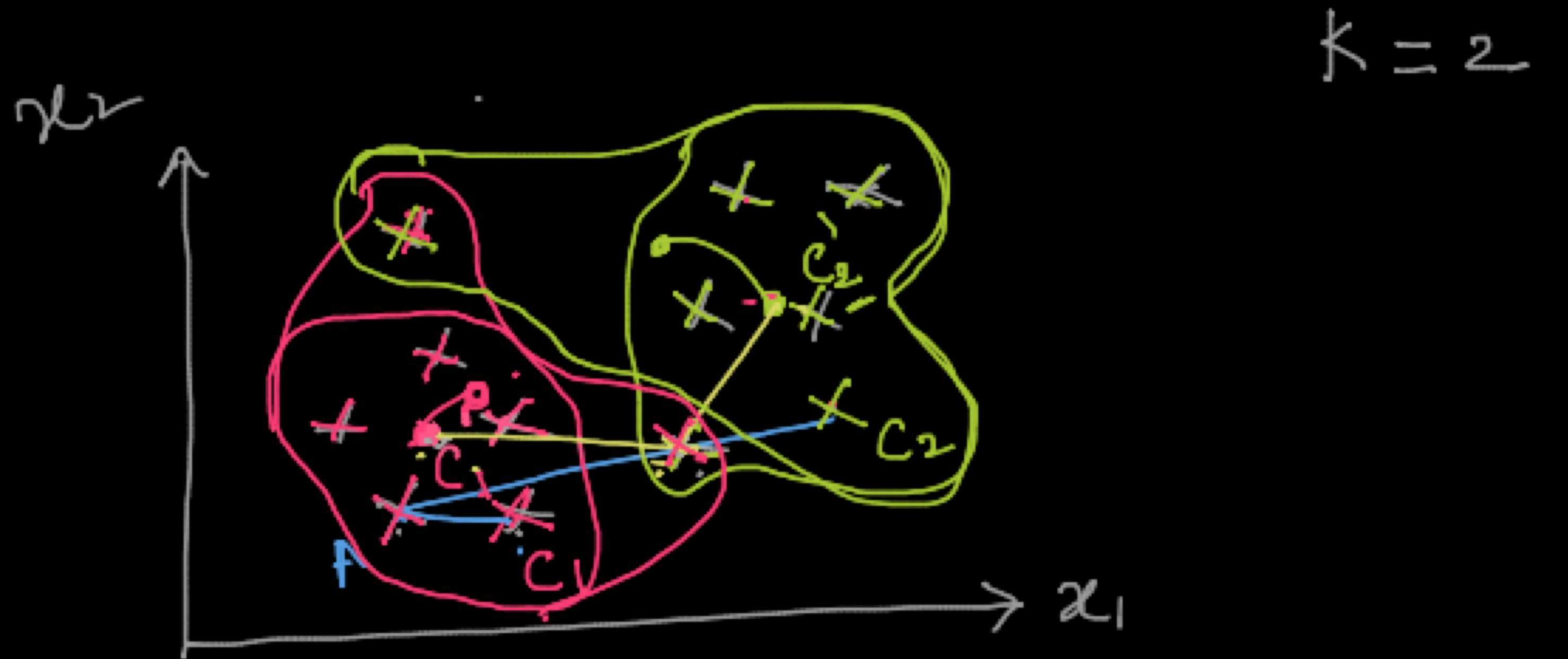
$$WCSS \text{ (Within Cluster Sum of Squares)} = \sum_{i=1}^{n_1} (x_i - c_1)^2 + \sum_{i=1}^{n_2} (x_i - c_2)^2 + \dots + \sum_{i=1}^{n_k} (x_i - c_k)^2$$



$$WCSS = \sum_{j=1}^k \sum_{i=1}^{n_k} (x_i - c_j)^2 \rightarrow \text{Lloyd's Algorithm}$$

Task: Find those centroids which will minimize WCSS

- NP Hard problem
- Approximations (Lloyd's algorithm)



$$C = \left(\frac{x_1 + x_2 + x_3}{3}, \frac{y_1 + y_2 + y_3}{3} \right)$$

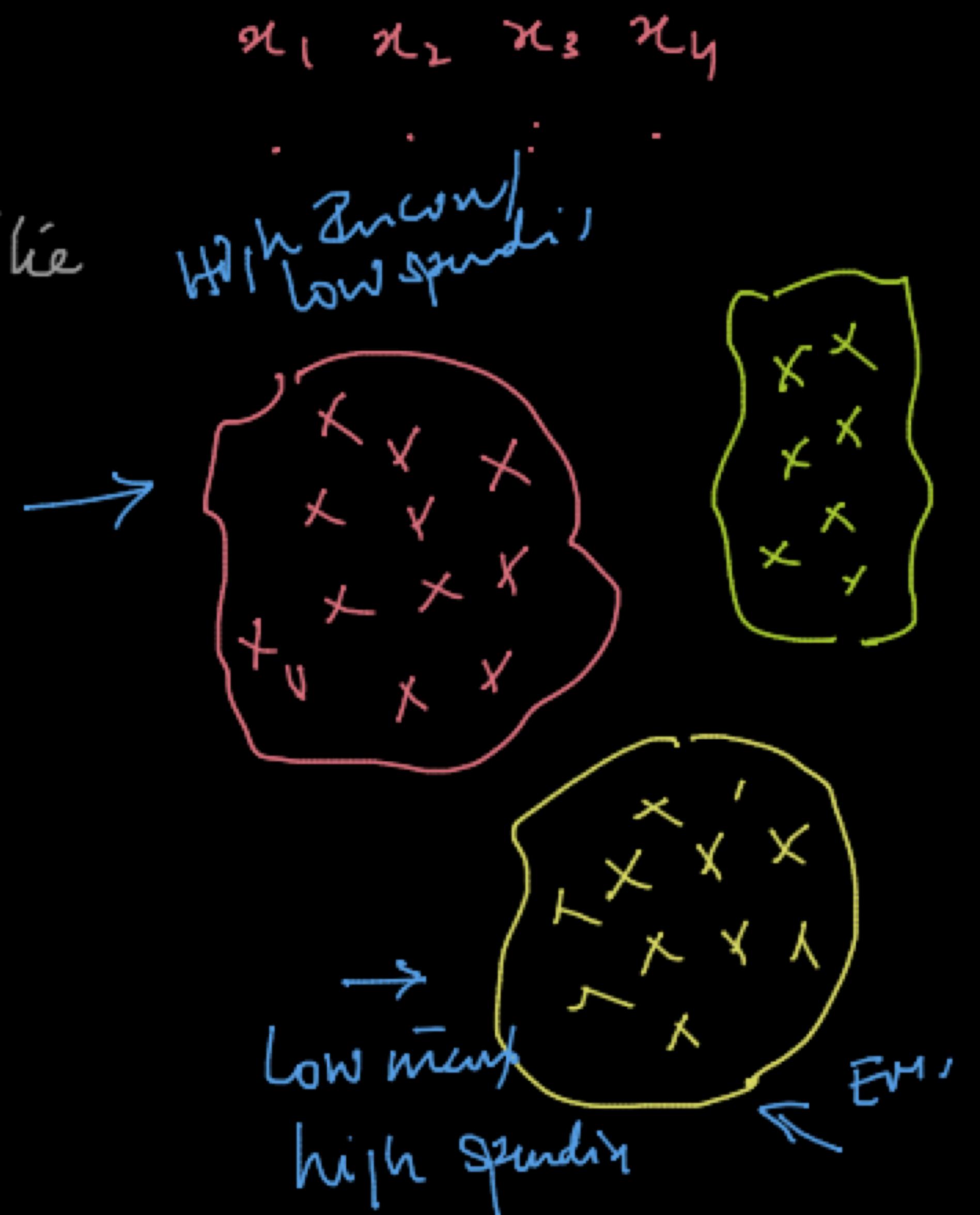
$$C = (x, y)$$

Initialization Sensitivity

→ Based on the choice of initial centroids (ie <sup>With known
low spread</sup>) final clusters would be different

K-Means++

→ Guidelines to choose the initial points



Lloyd's Algorithm

Step 1: Randomly choose k-points as centroids.

- repeat until convergence
- Step 2: Find the distance between every d_p and the k-centroids.
- Step 3: Assign the d_p to the closest centroid.
- Step 4: Recalculate the centroids based on the cluster created.

Convergence: The centroids don't change
— Stable clusters.

Step 1 : Define the objective function

— usually the loss function ✓

— Some form of residual / Error - ($y - \hat{y}$)

Step 2 : Find model parameters which minimize / maximize the
Value of objective fn.

|



Step 1 : Define the objective function

— usually the loss function ✓

— Some form of residual / Error - ($y - \hat{y}$)

Step 2 : Find model parameters which minimize / maximize the value of objective fn.

1. Which obj. fn-

— Squared Loss

2. What are we finding ?

→ β 's which minimize the squared loss

3. What algo ?

→ OLS.

4. Metrics

→ R^2 , MSE, RMSE, MAE.

