

Topics:

- { - clustering
- Dive deep into K-means + Variations

~~ques~~

1

 $\rightarrow O(t)$
 $x_{qj} = w_1 w_2 \dots w_k$ \curvearrowright #words in x_{qj} NOT train data

Run time compl = $O(k)$

precompute
from D_{train}

$\left\{ \begin{array}{l} p(y=1); p(y=0) \\ p(w_i|y=1); p(w_i|y=0) \end{array} \right.$

\downarrow

key = $i, 1$ val = 2

②

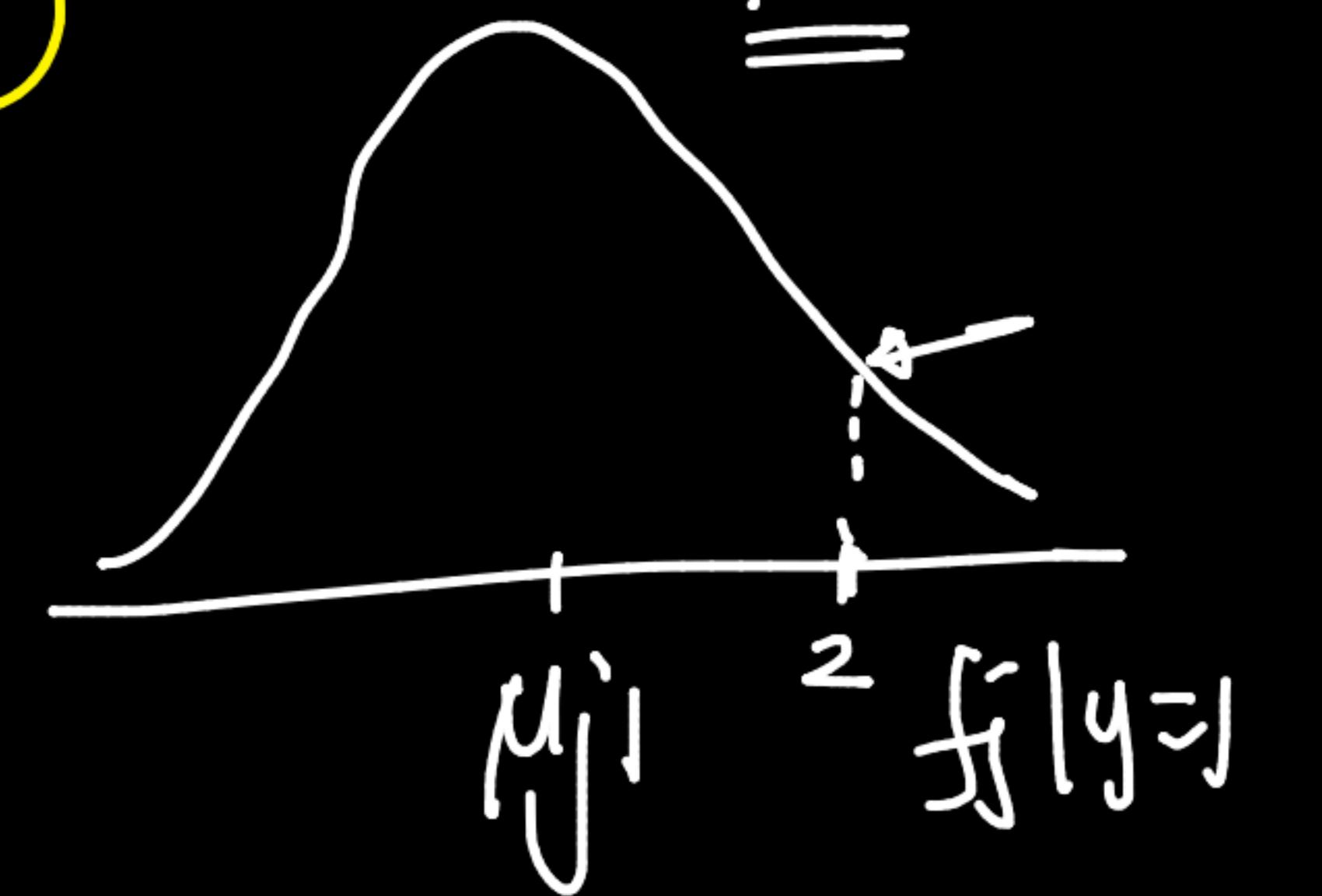
Gaussian NB:-

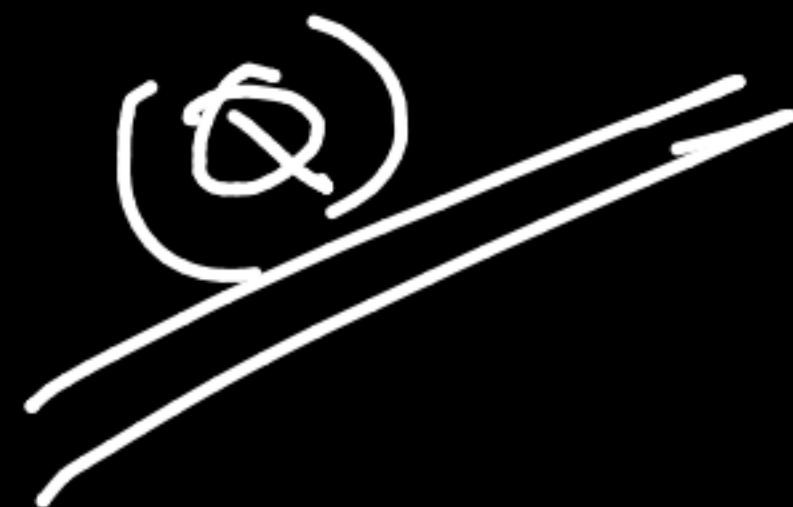
PDF

$$p(f_j | y=1)$$



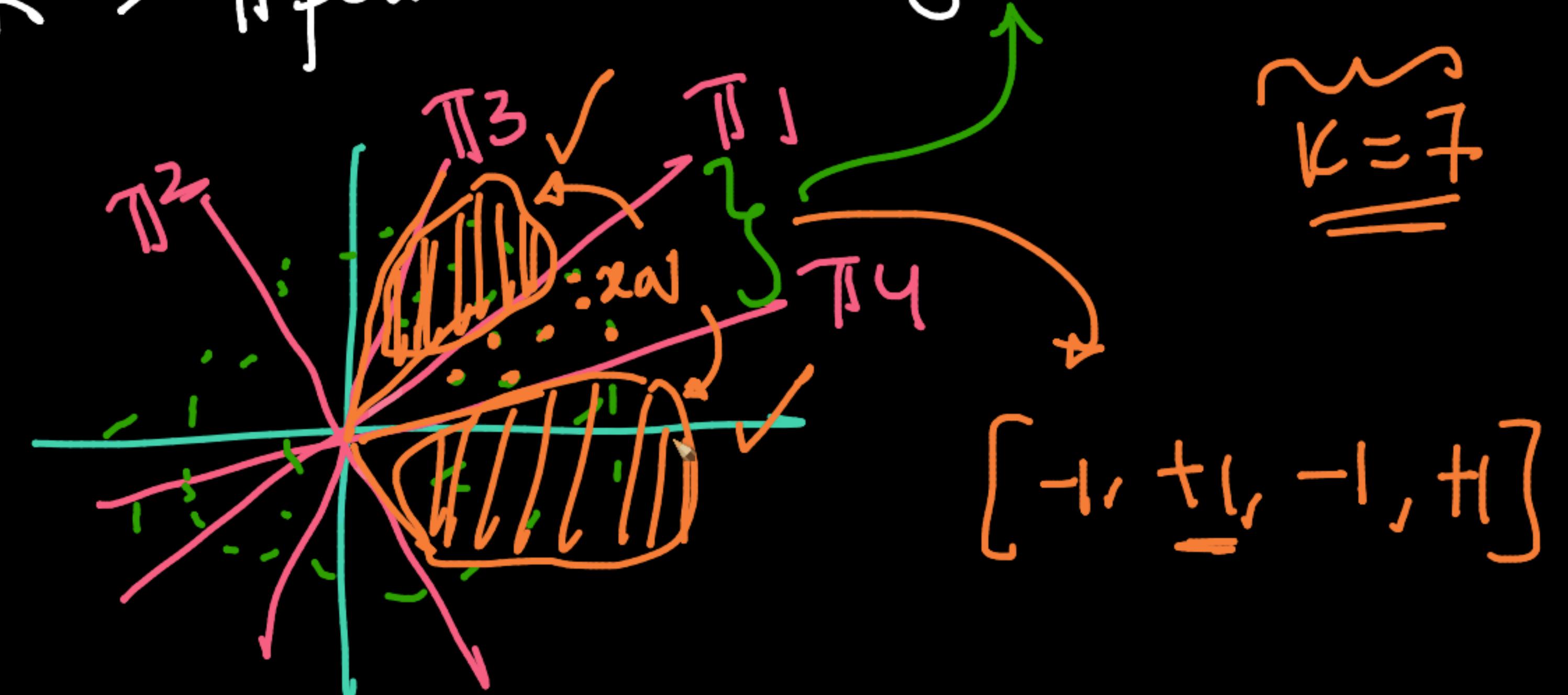
$$\hookrightarrow N(\mu_j, \sigma_j^2)$$

PDF



LSH: Cosine-Similarity (intuition)

$K > \# \text{points} \rightarrow \text{in a Segment}$



Clustering

Classification

$$\mathcal{D} = \left\{ (\underline{x}_i, \underline{y}_i) \right\}_{i=1}^n ; \underline{x}_i \in \mathbb{R}^d, \underline{y}_i \in \{0, 1\}$$

Regression

$$\underline{y}_i \in \mathbb{R}$$

Supervised

y_i exists

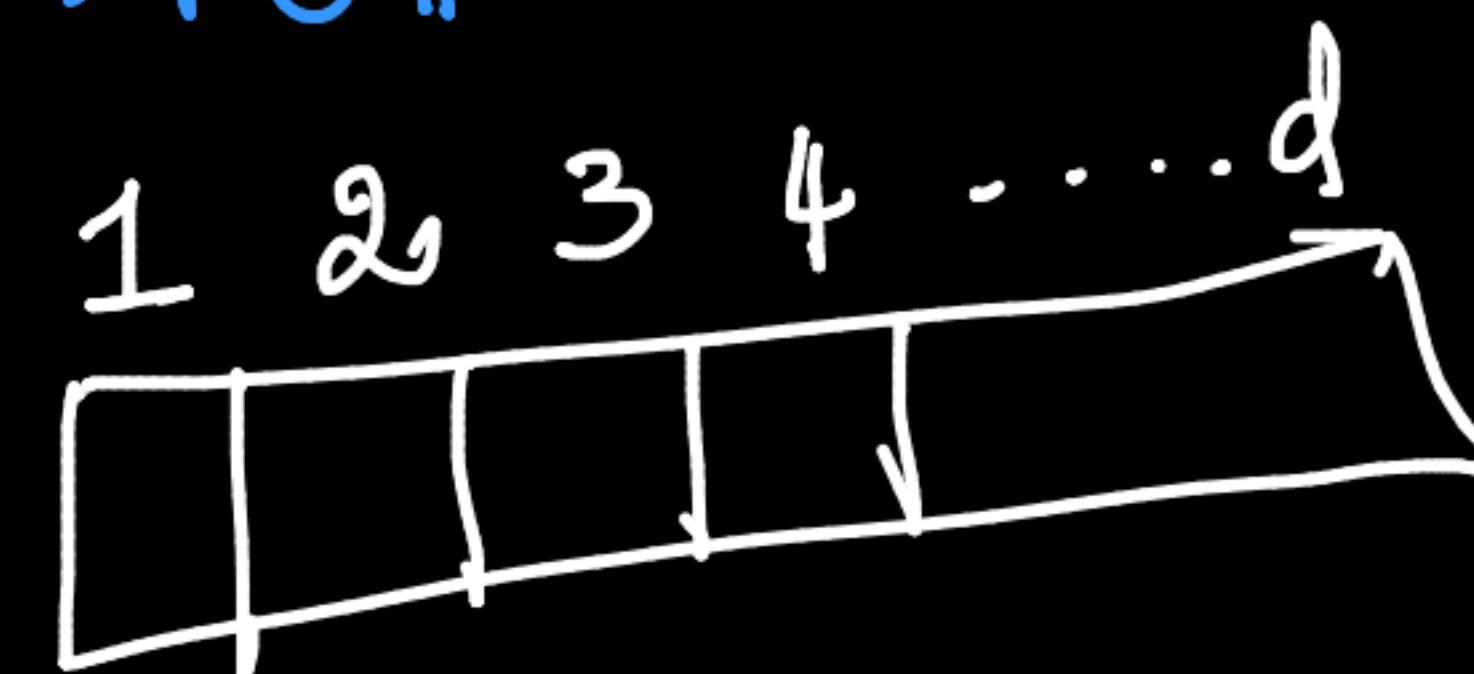
c.g.
~~if~~

100 MM Amazon Customers

10M - C₁
8M - C₂
..

$$\mathcal{D} = \{ x_i \}_{i=1}^n$$

$$x_i \in \mathbb{R}^d$$

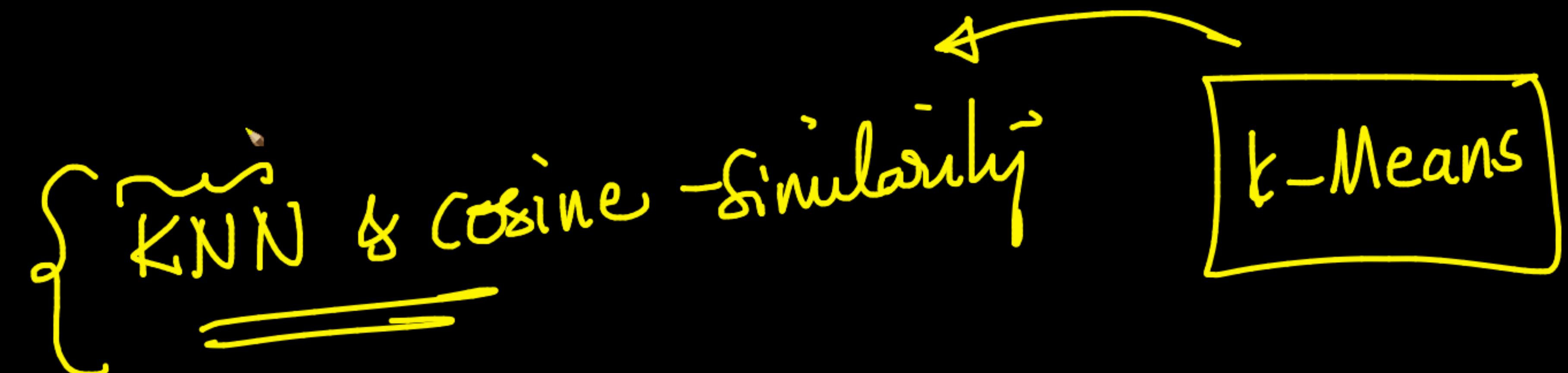


group similar
Customers

no y_i

e.g:
look english words → cluster similar words

LM newspaper articles → group similar articles
=



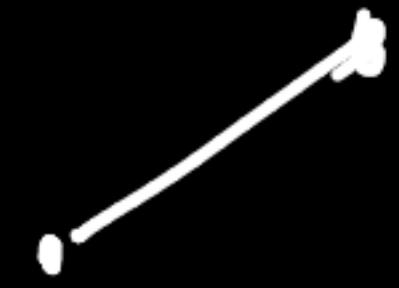
✓ [Clustering
Unsupervised Learning :-]

$x_i \rightarrow y_i$ ✓ { Semi-supervised learning : }
Self-supervised learning : } ↳ Images (CV & DL)

RCA: dim-reduction ✓
no y_i^T

Train ↗
Some pls have y_i^T
many DONOT have y_i^T

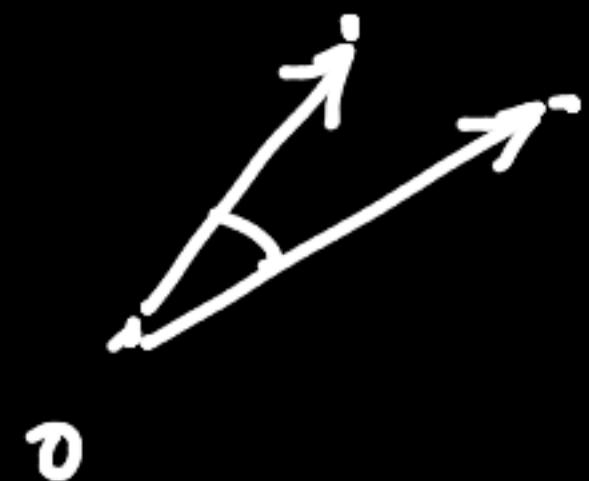
Practical



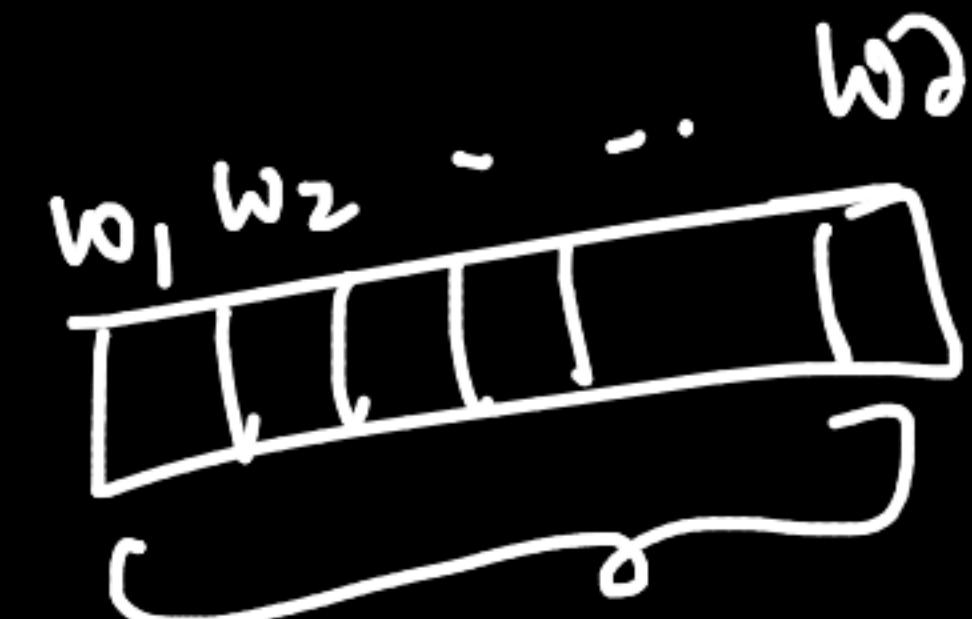
Euclidean - dist → low-dim

⋮

Manhattan dist → low-medium dim



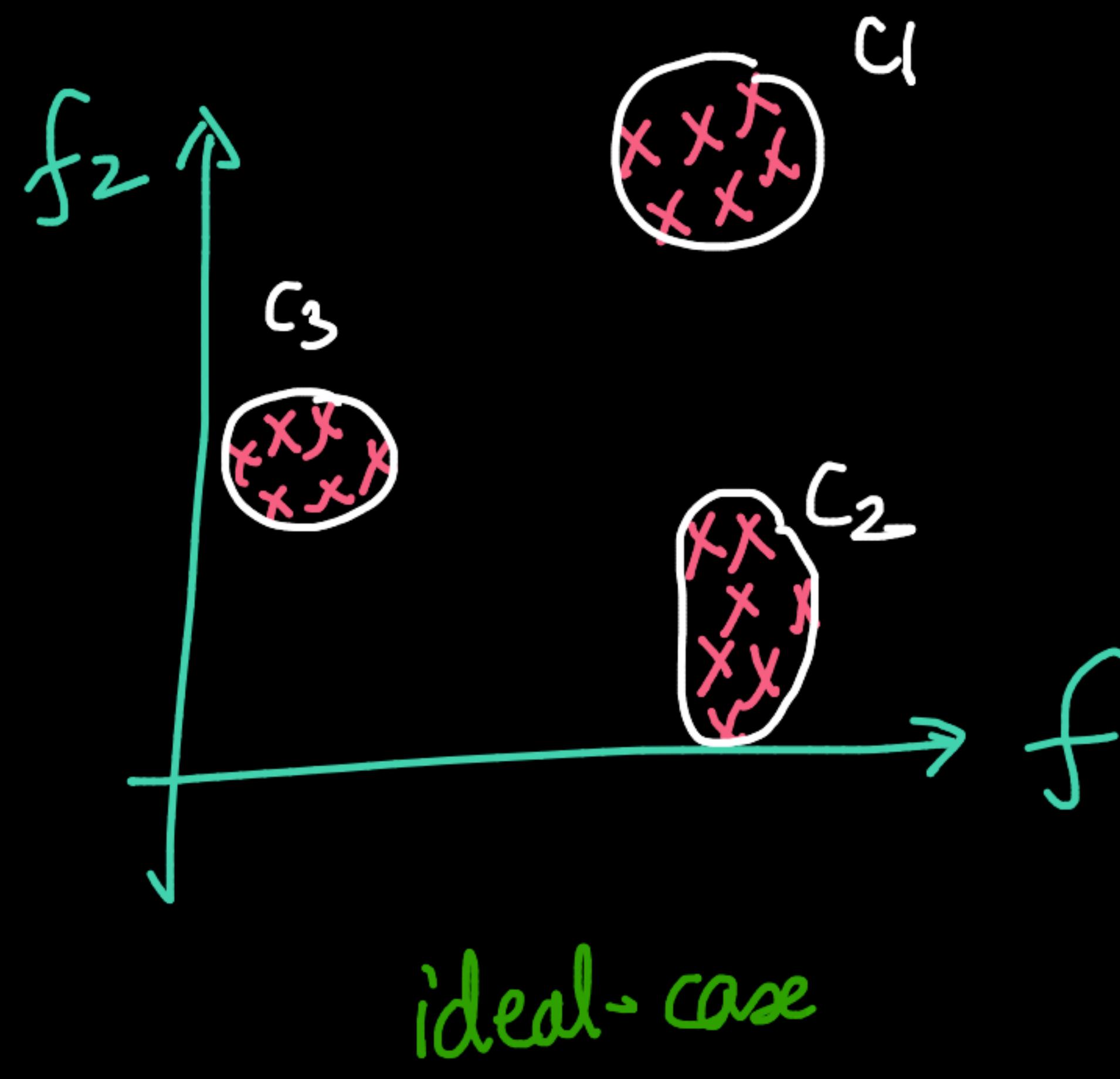
Cosine sim → high-dim
Text-data BOW
sparse



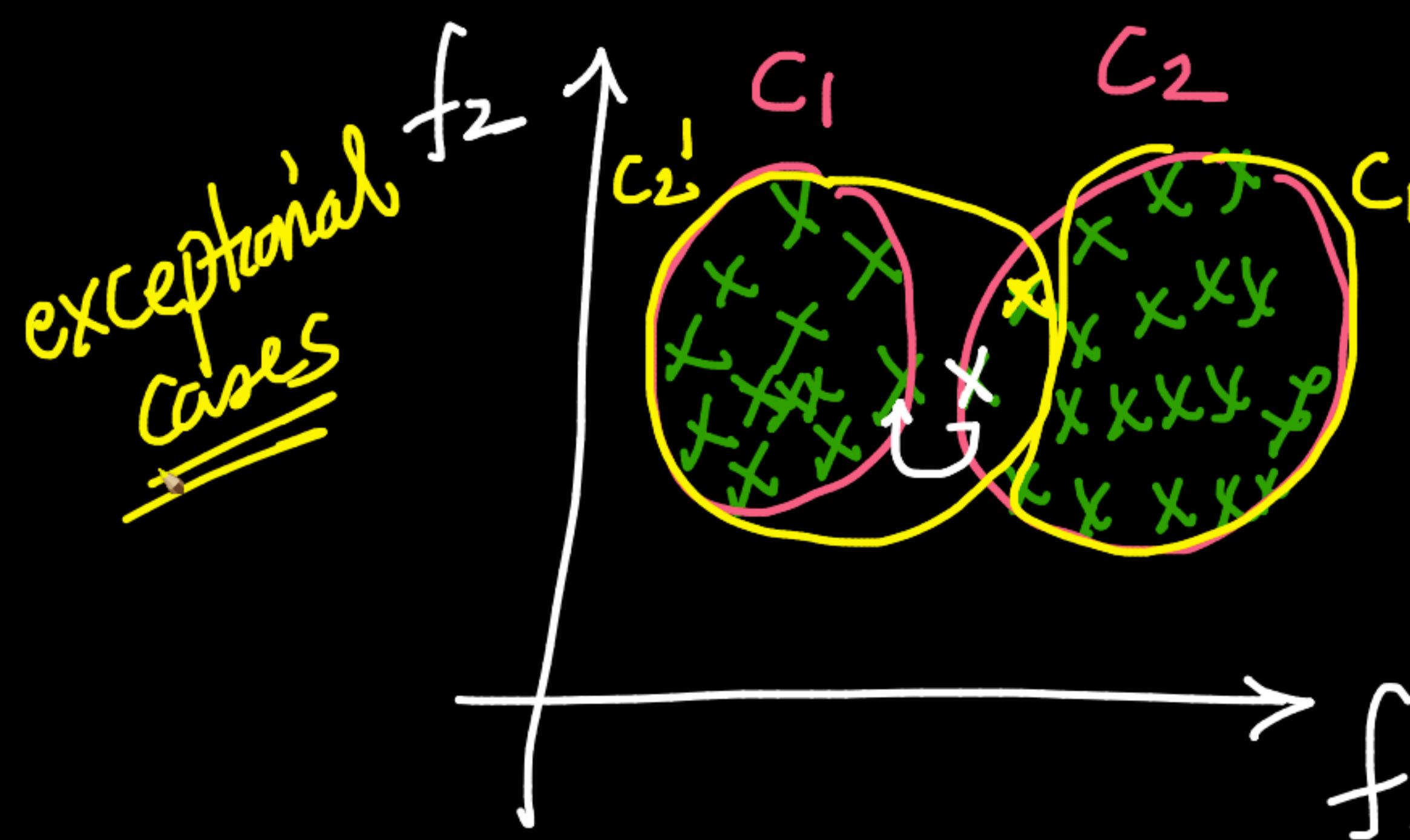
Rule of Thumb

Practically:

Distance-Metric is a hyper-passim



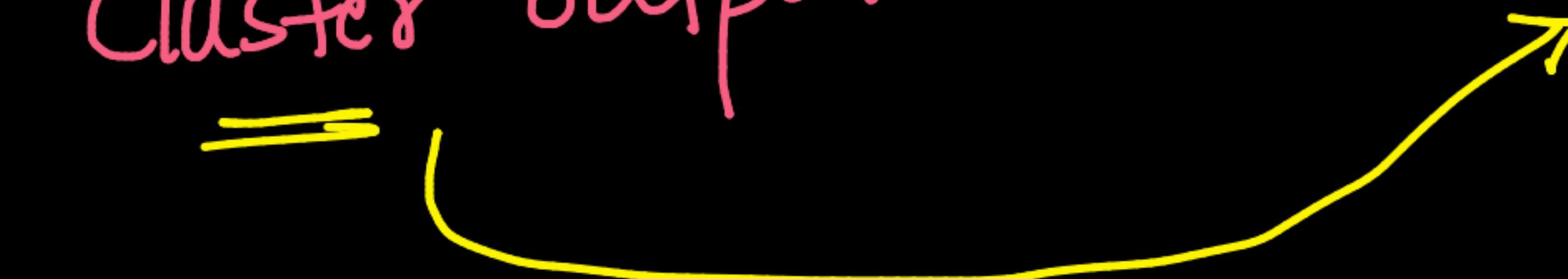
- INTUTION
- points in a cluster are close to each other
 - points in different clusters are far from each other



- different algorithms

- no ground truth

cluster output must make biz sense



DOMM Amazon

$x_i \in \mathbb{R}^d$

Special Deals

$c_1 c_2 \dots$
 m
high Spenders

c_{10}
 m

Input: $\{(x_i)_{i=1}^n; x_i \in \mathbb{R}^d\}$
Output: $c_1, c_2 \dots c_K$
 $\{ \} \quad \{ \} \quad \{ \}$

x_{new} → which cluster does x_{new} belong to?

personas
ecommerce

n $\in \mathbb{N}$



bif -sense
 =

→ Supervised
Euc. dist,
→ dist
→ distance

Manhattan dist; Cosine-Sim

COMMON-
sense
metrics
↳ Dunn-Index

Classfn \rightarrow AUC; Pr, Fe, F1 ...

Metric

Reg $\rightarrow R^2$

Metric

$$\mathcal{D} = \left\{ (x_i)_{i=1}^n ; x_i \in \mathbb{R}^d \right\}$$

dist b/w cluster centres

min dist

max dist

inter-cluster dist (\uparrow high)

intra-cluster dist

(\downarrow low)

avg dist from center

max dist

f_2

f_1

C_2

$x_j \in C_2$

$x_i \in C_1$

$\min d(x_i, x_j)$

s.t $x_i \in C_1$
 $x_j \in C_2$

k-clusters

Metric

$$\text{DUNN - Index} = \frac{\min_{i,j} d(i,j)}{\max_k d^*(k)}$$

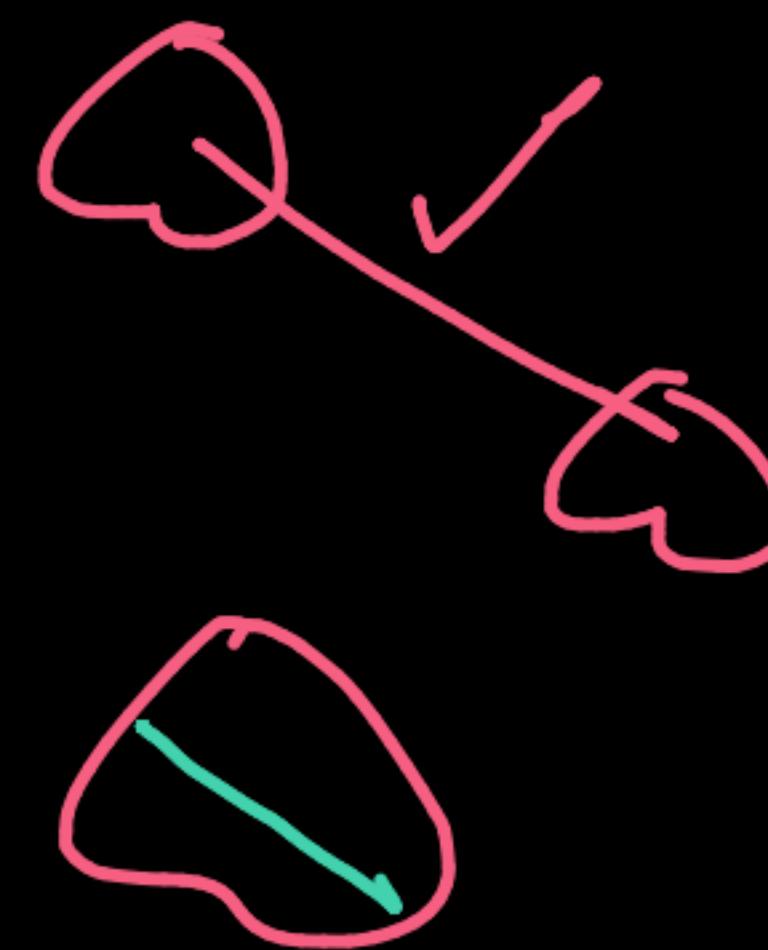
Inter-cluster dist

Intra-cluster dist

Diagram illustrating the Dunn Index formula:

- The numerator $\min_{i,j} d(i,j)$ is shown with a yellow bracket under the term $d(i,j)$. A yellow arrow points from the term $d(i,j)$ to the bracket.
- The denominator $\max_k d^*(k)$ is shown with a blue bracket under the term $d^*(k)$. A blue arrow points from the term $d^*(k)$ to the bracket.
- A pink box encloses the term $d(i,j)$.
- A blue box encloses the term $d^*(k)$.

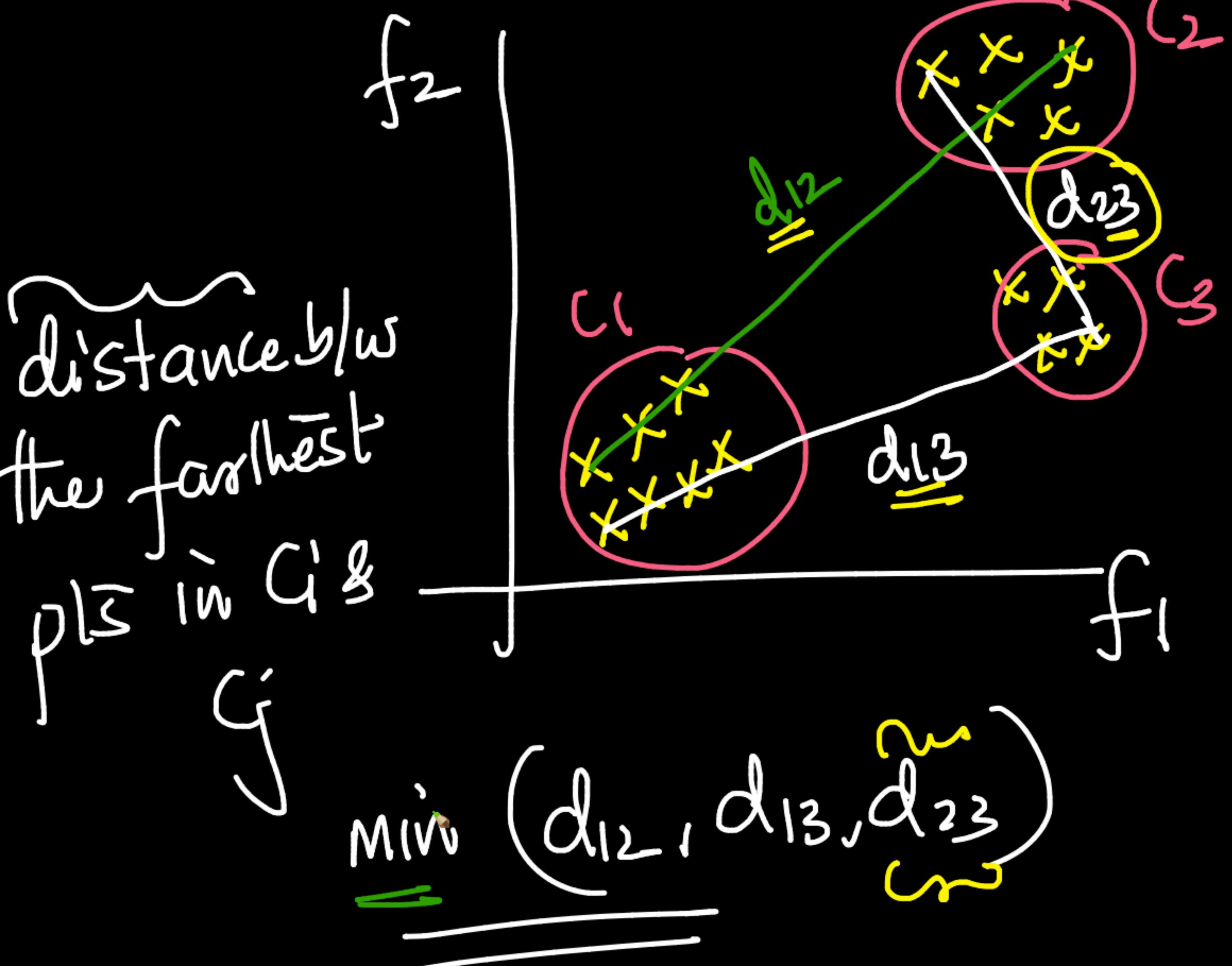
$$\uparrow \quad \textcircled{D} = \min_{i,j} d(i,j)$$
$$\downarrow \max_k d^l(k)$$



Inter-cluster
distance

$d(i,j)$ = distance b/w
the farthest
pts in $C_i \& C_j$

$\{k\text{-cluster}\}$
 kC_2 distance

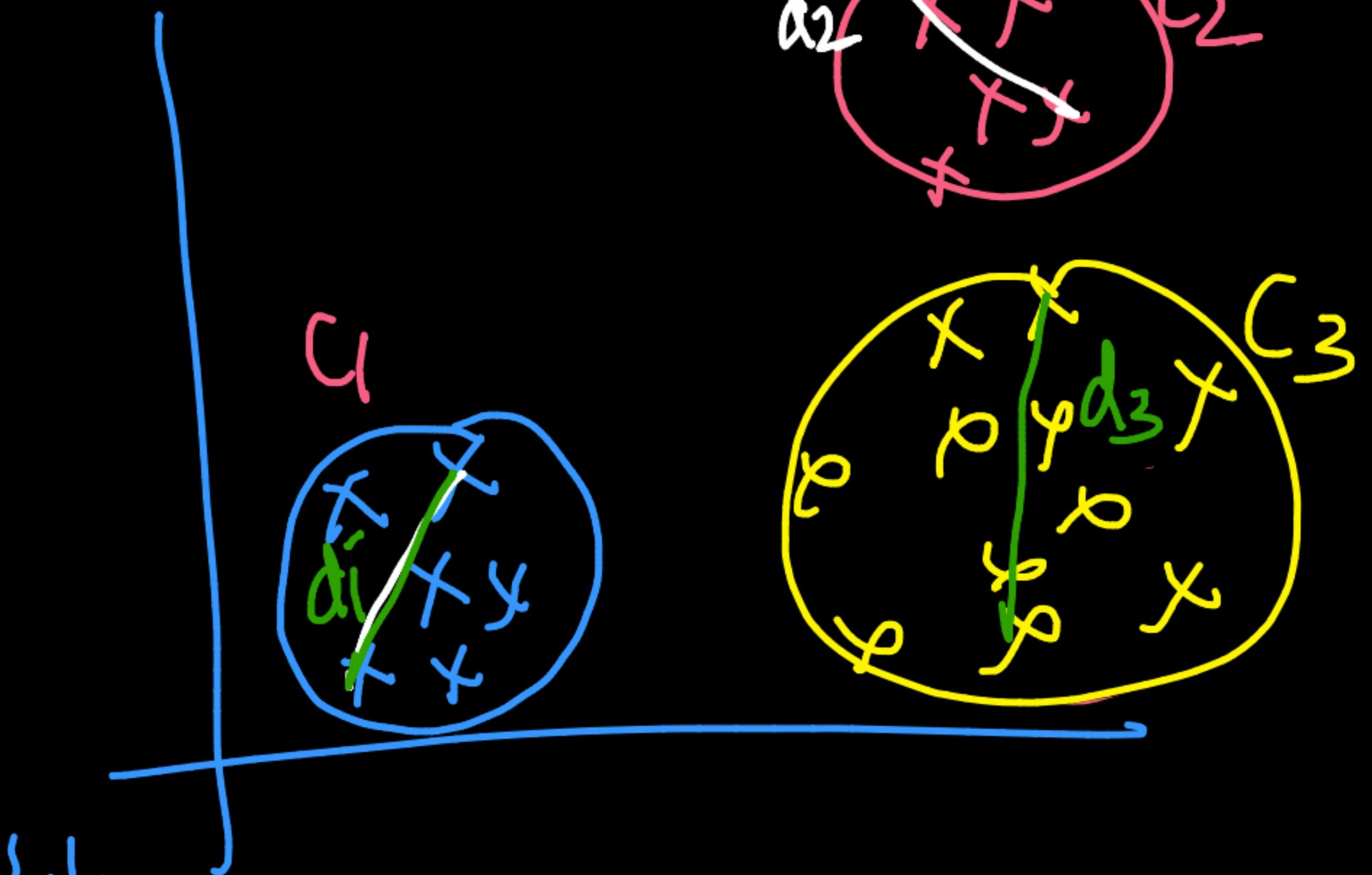


Intera-
clust
dist

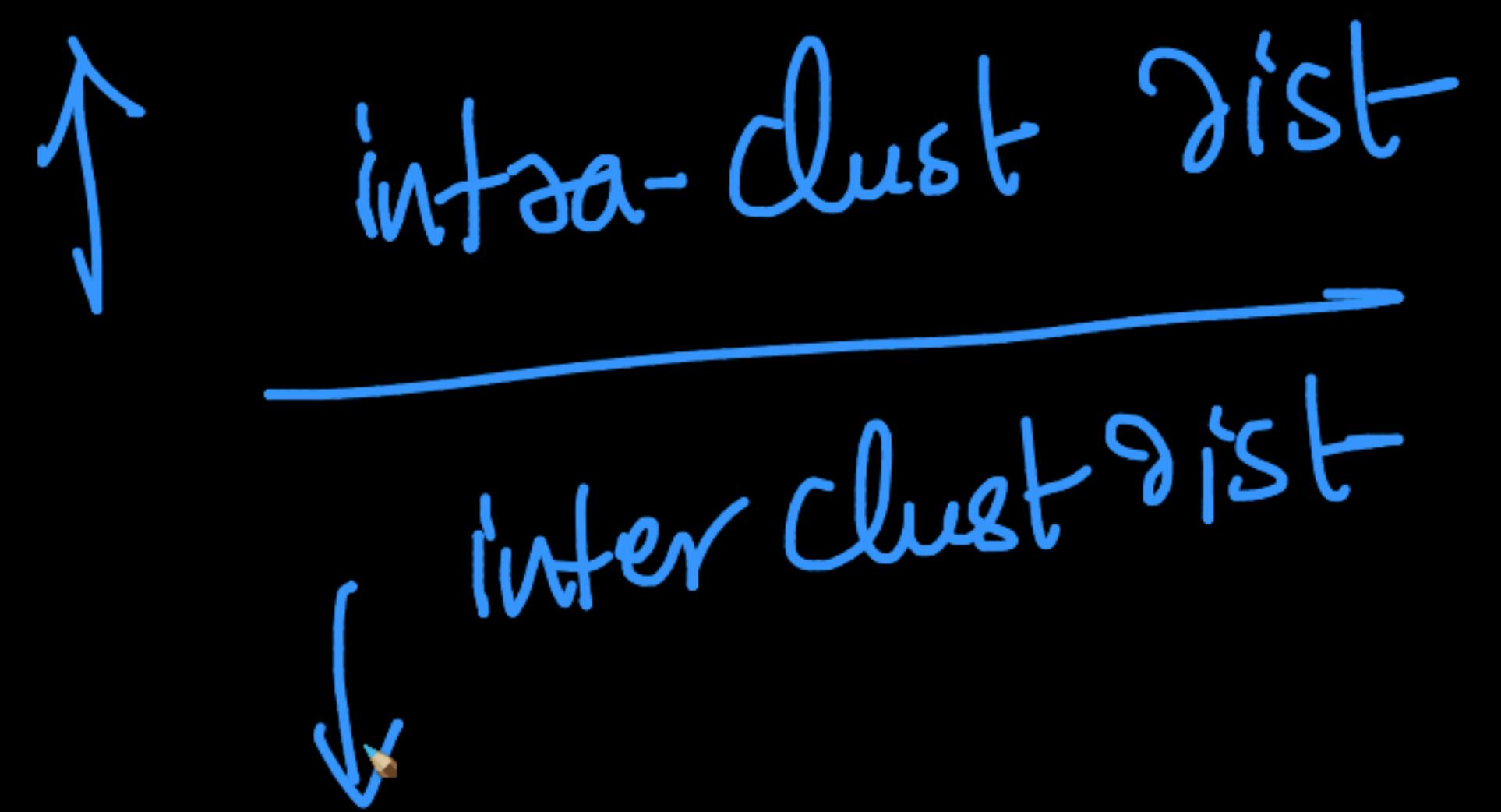


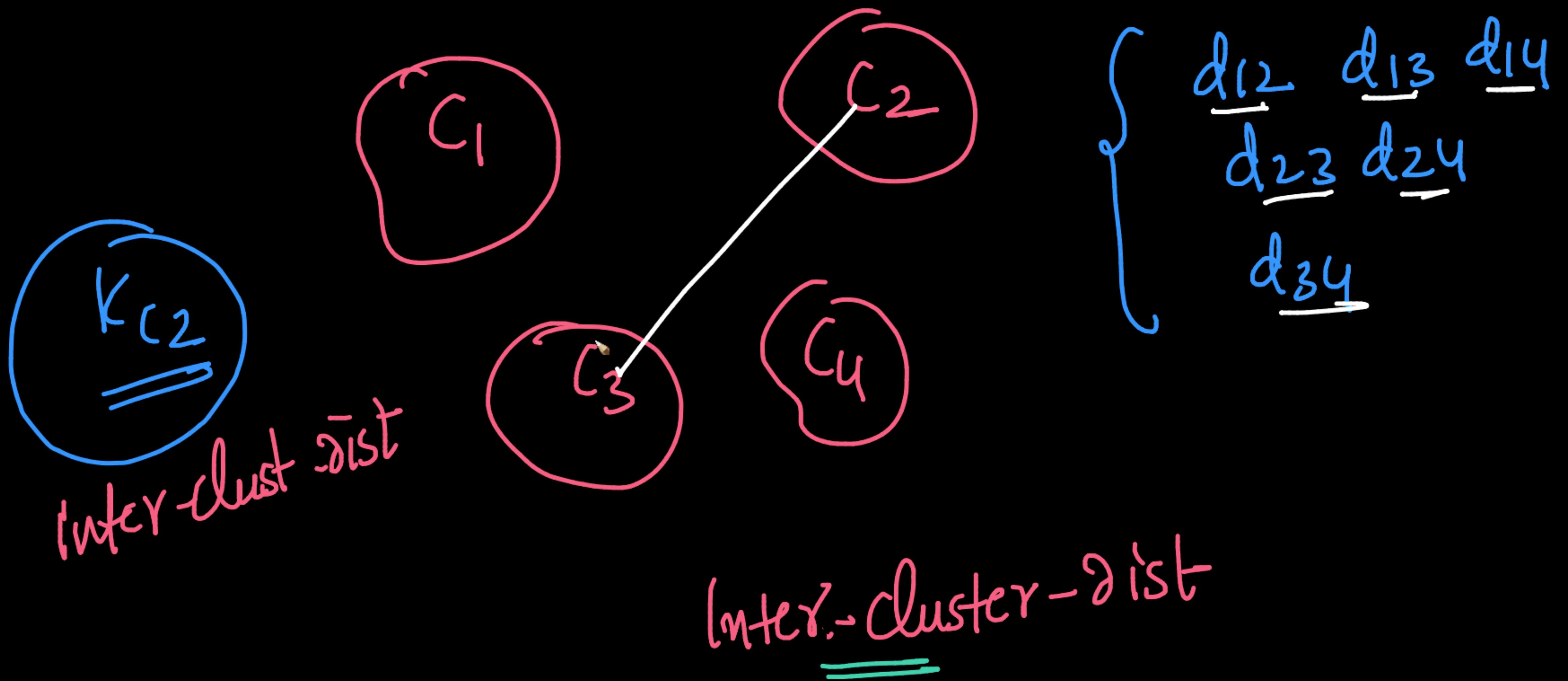
$$\max(d_1, d_2, d_3)$$

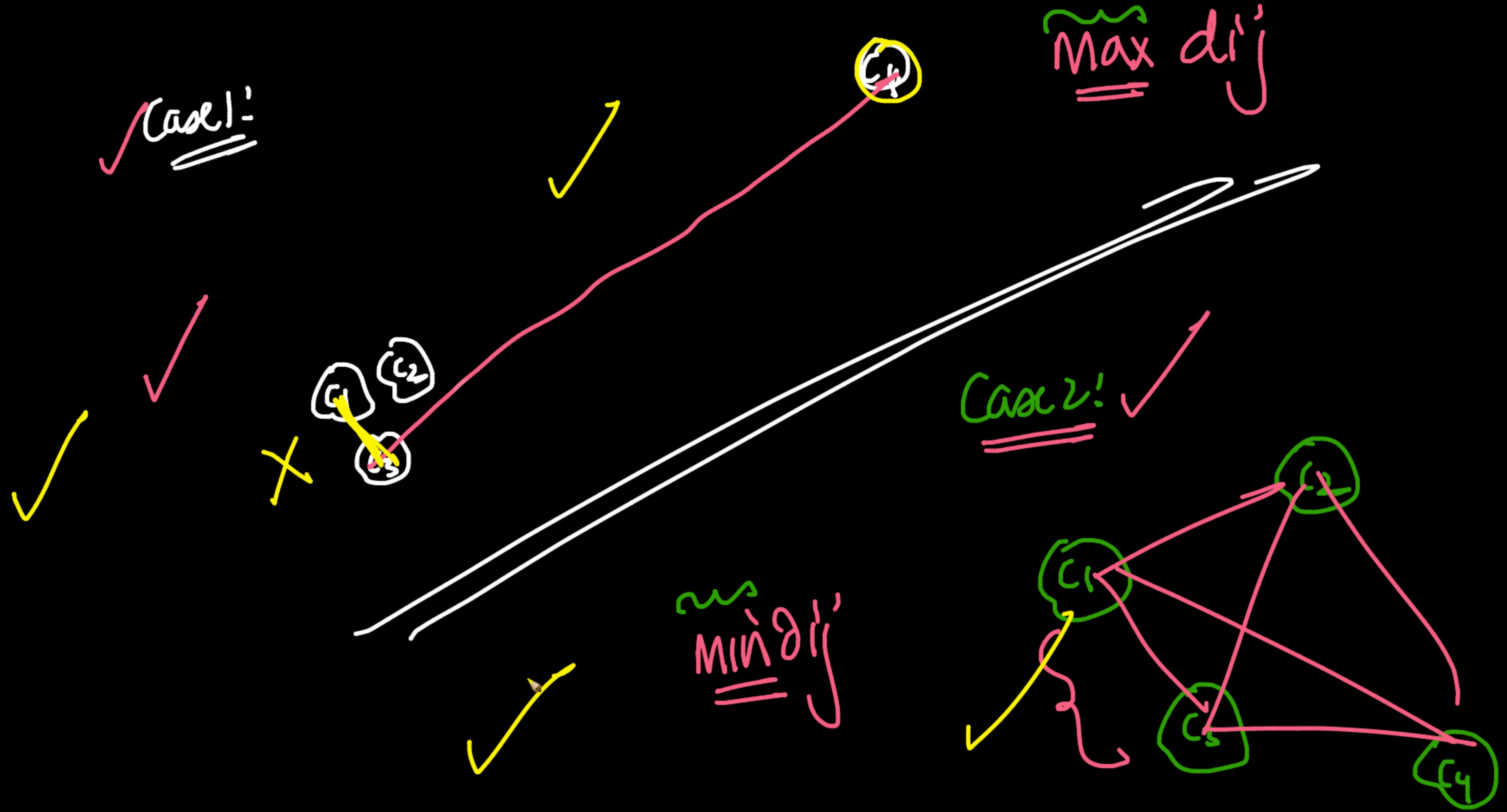
$d'(i) = \max$ dist
b/w 2 points in C_i



Design your,
metric

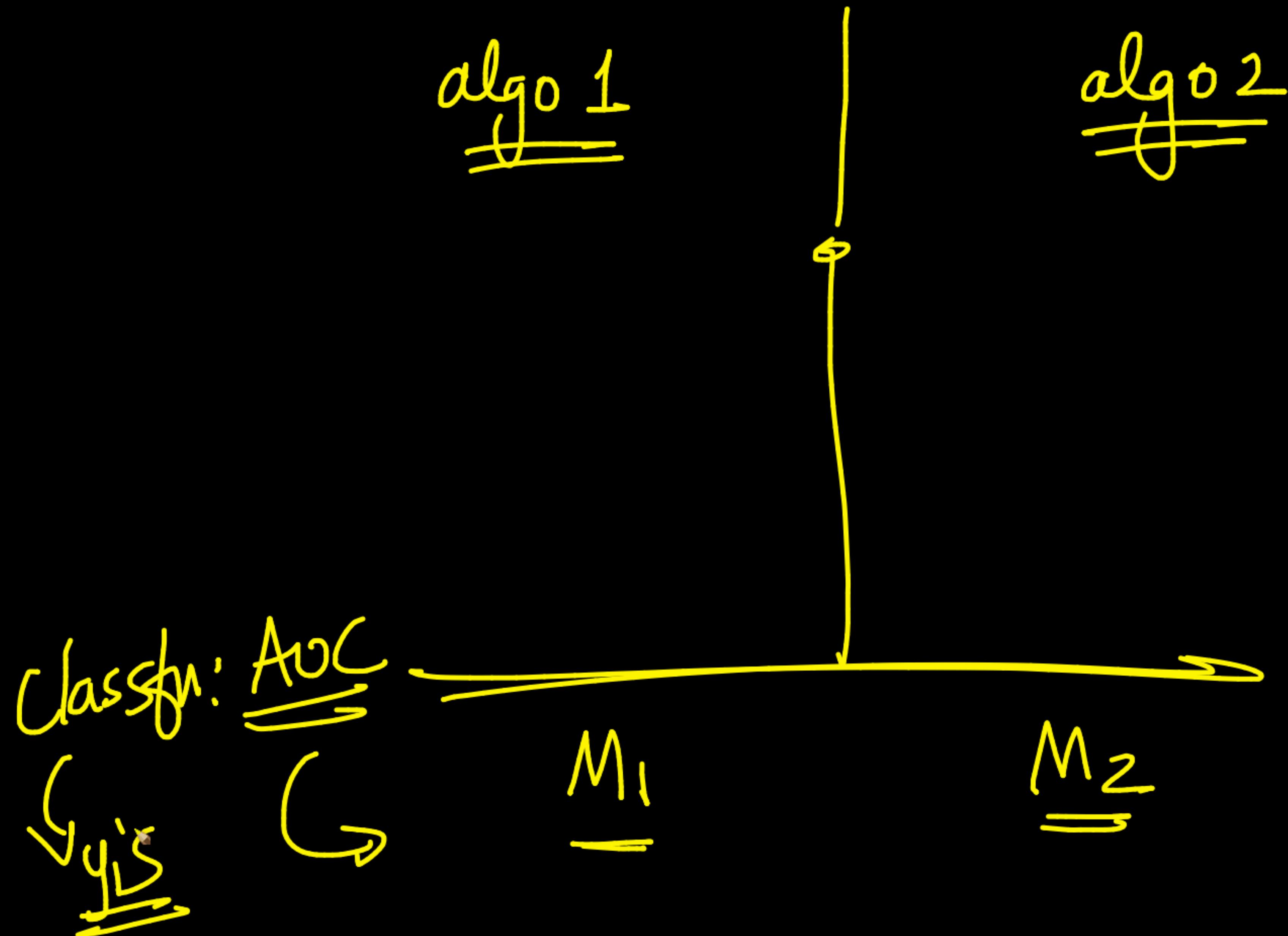


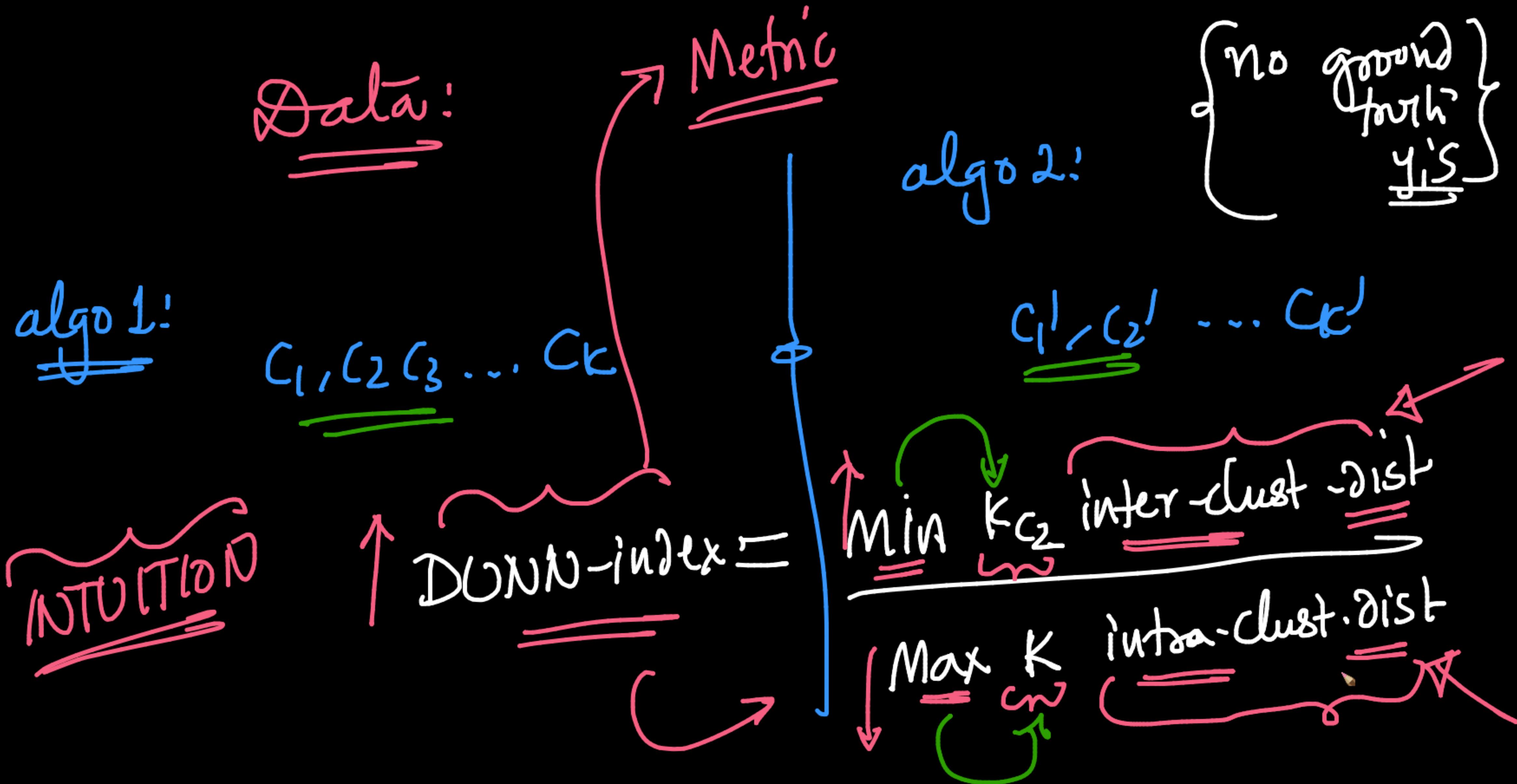




if metric : ~~min max~~ $d(i,j)$

$\max_l d^l(i)$





DUNN - Index → Metric
→ not an optimization problem



$$1 \leq \text{DUNN-Index} = \frac{\min d_{ij}}{\max d_{ik}}$$

↓

$\text{val}_1 = 10$ (let)

$$\mathcal{D} = \{x_i\}_{i=1}^n$$

c_1, c_2, \dots, c_k

(one-clustering)

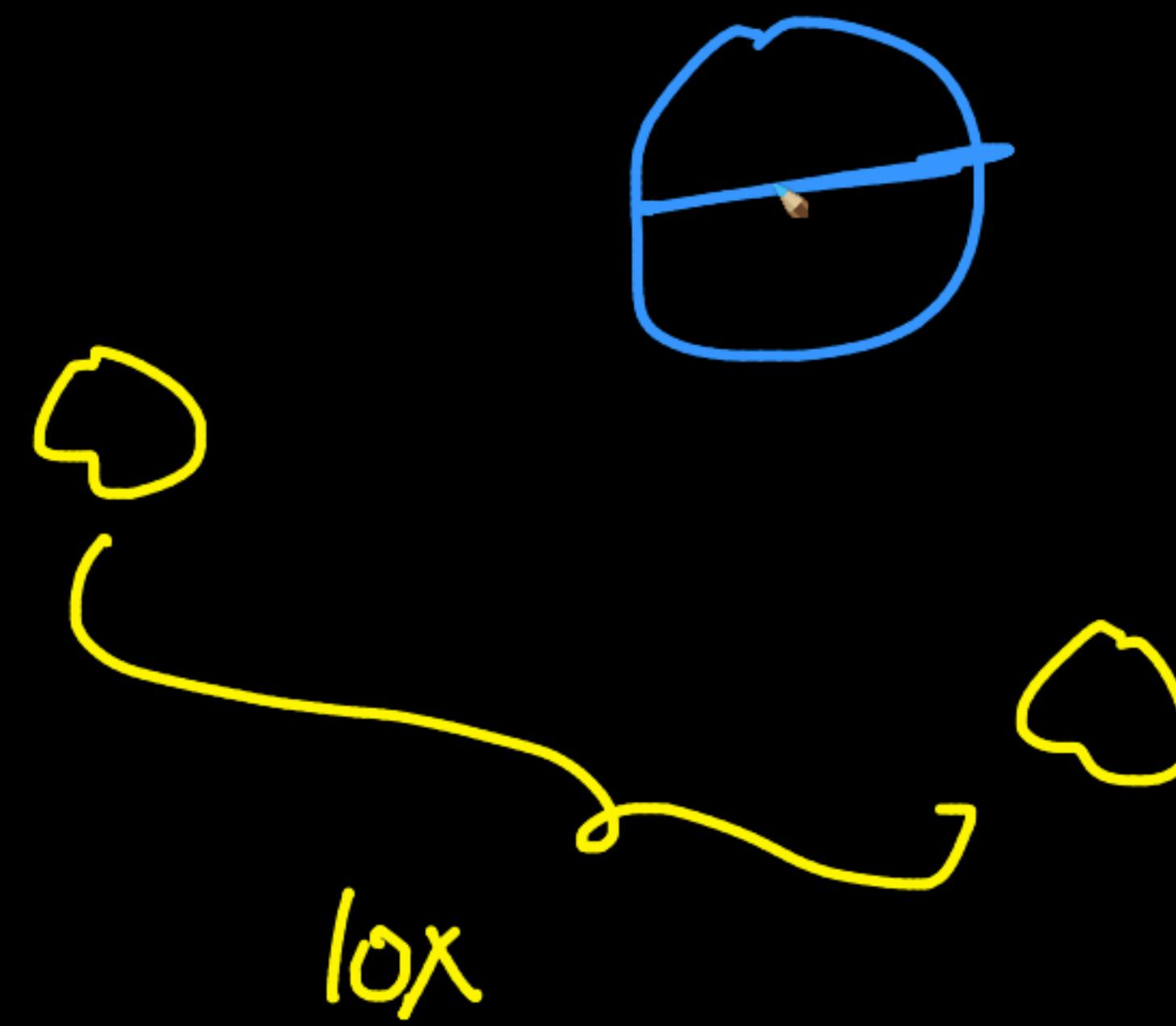
✓

HINT: $AUC = M_1: 0.7$

random: 0.5

ideal - 1

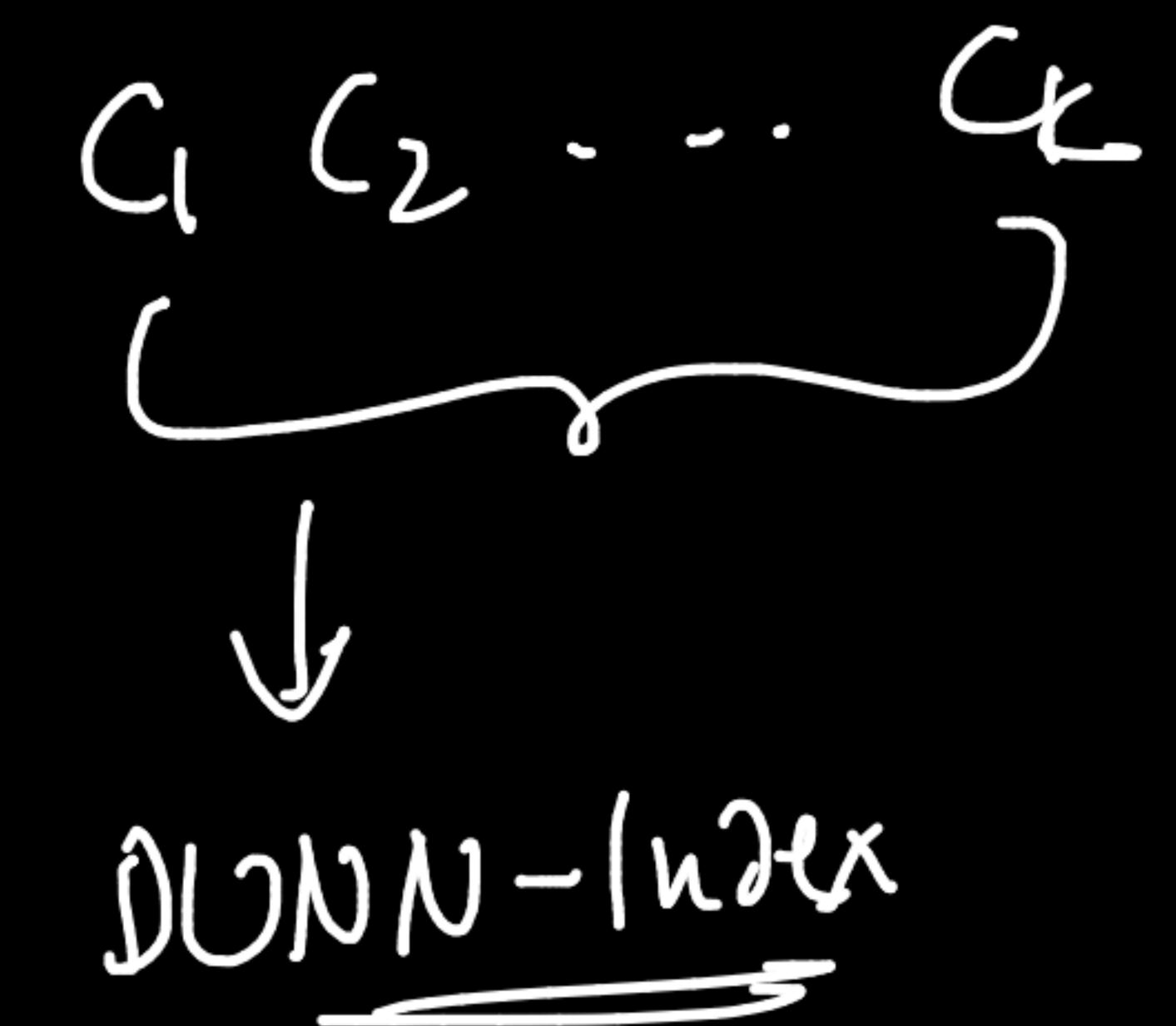
min inter-clust dist
Max intra clust dist



Random clustering

$$\mathcal{D} = \{x_i\}_{i=1}^n$$

baseline clustering
(k-means)



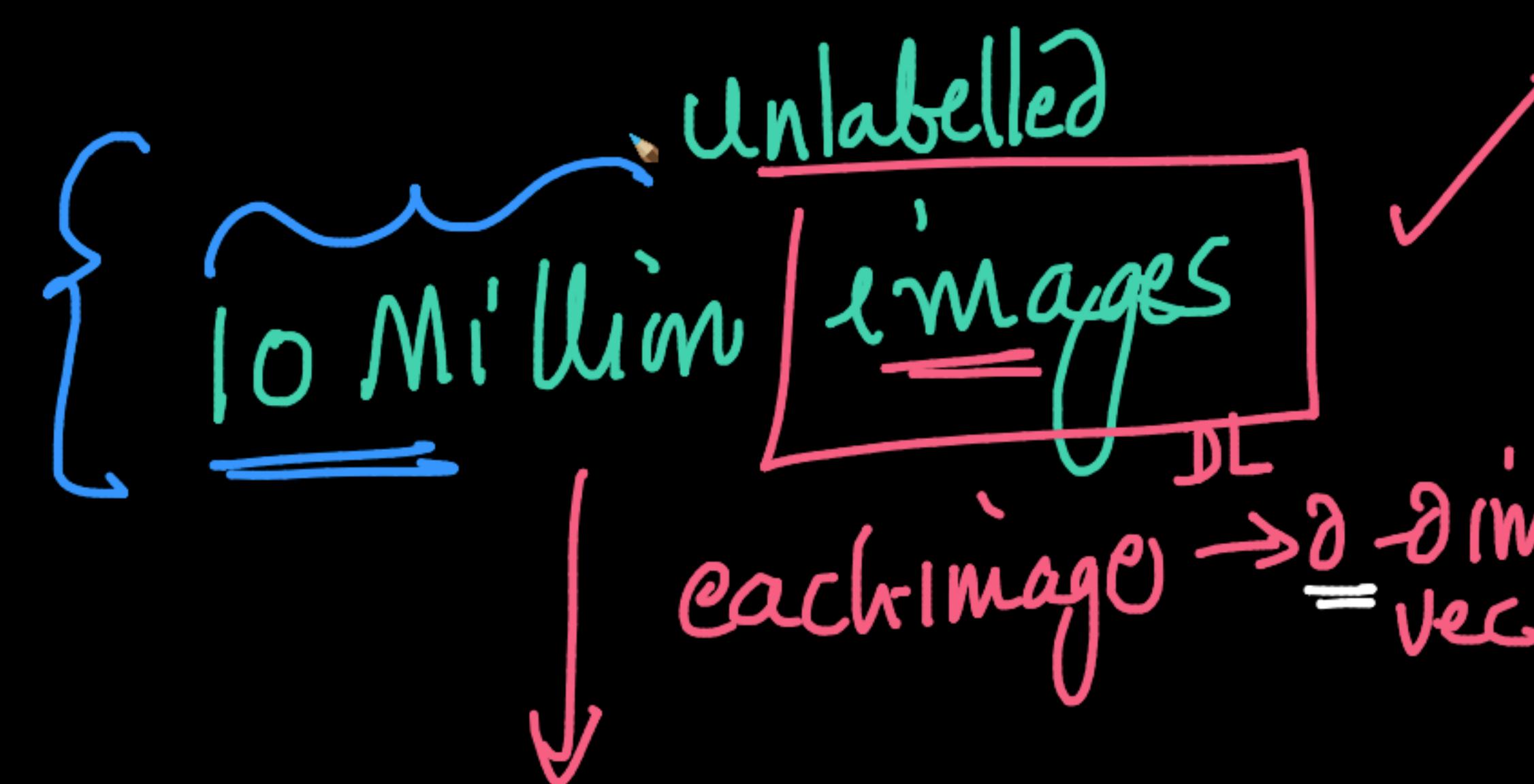
Clustering:-

Input: n 2-dim pts,
 $\# \text{clusters} = k$

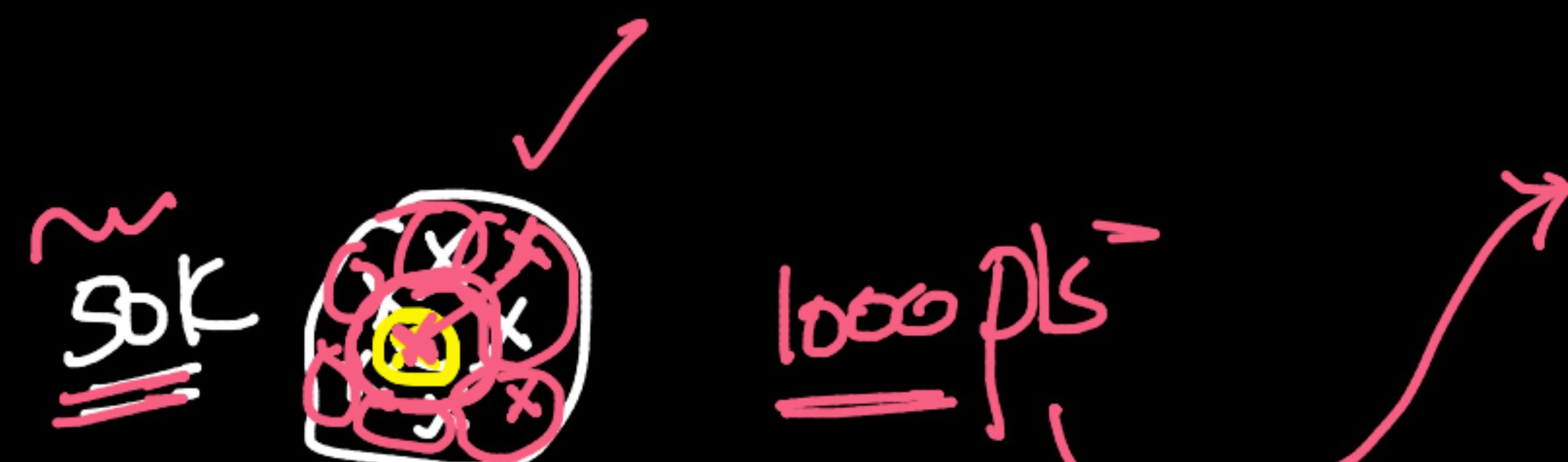
Output: k -clusters
 c_1, c_2, \dots, c_k

x_{new} → NN-like strategies

~~depth
later~~

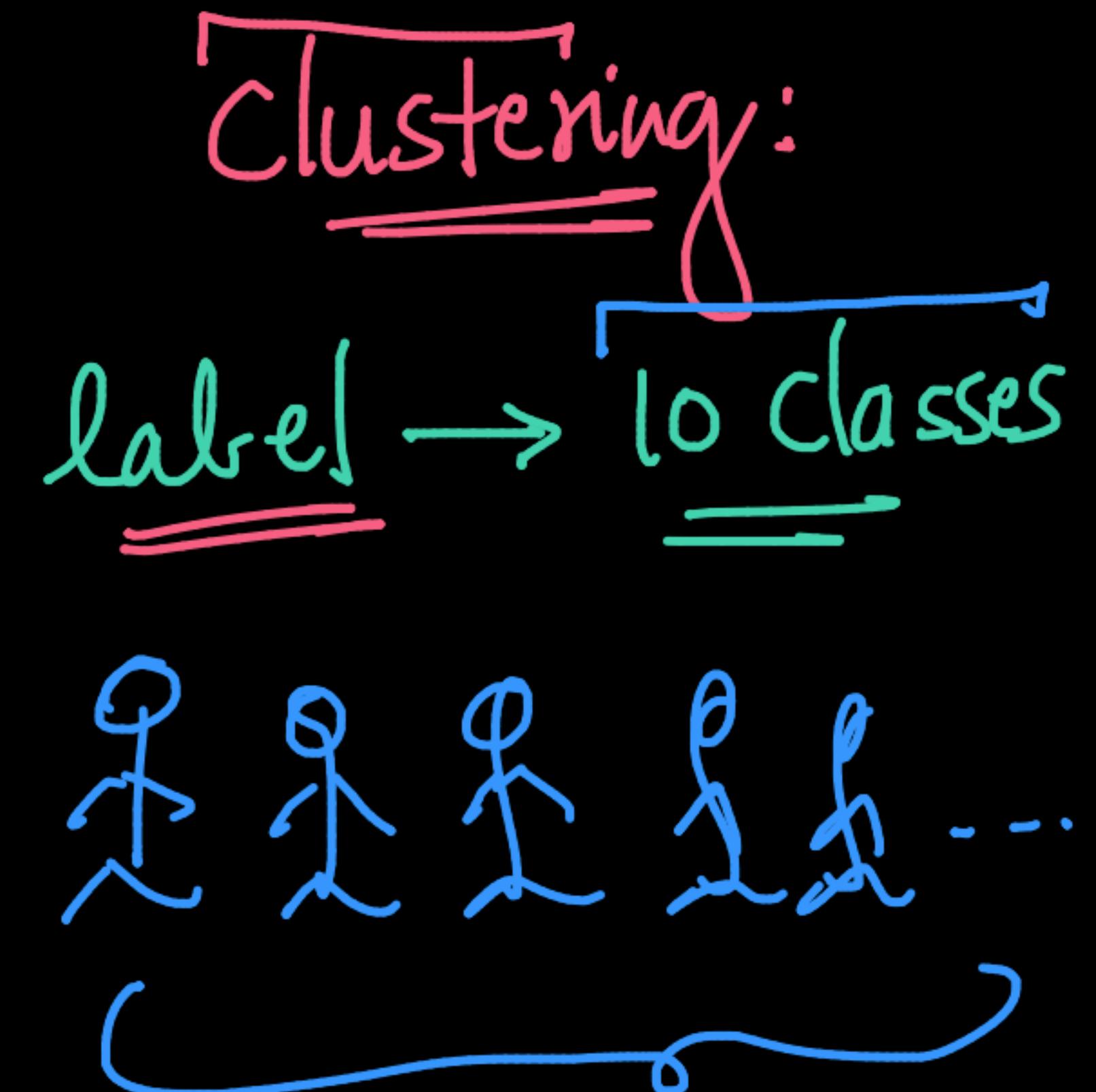


100 - clusters



Mens-jeans

Smart-label



- Simple; popular; baseline
- basic - Variant

Ki:

#clusters

Input:- n 2-dimpts, K

K-means

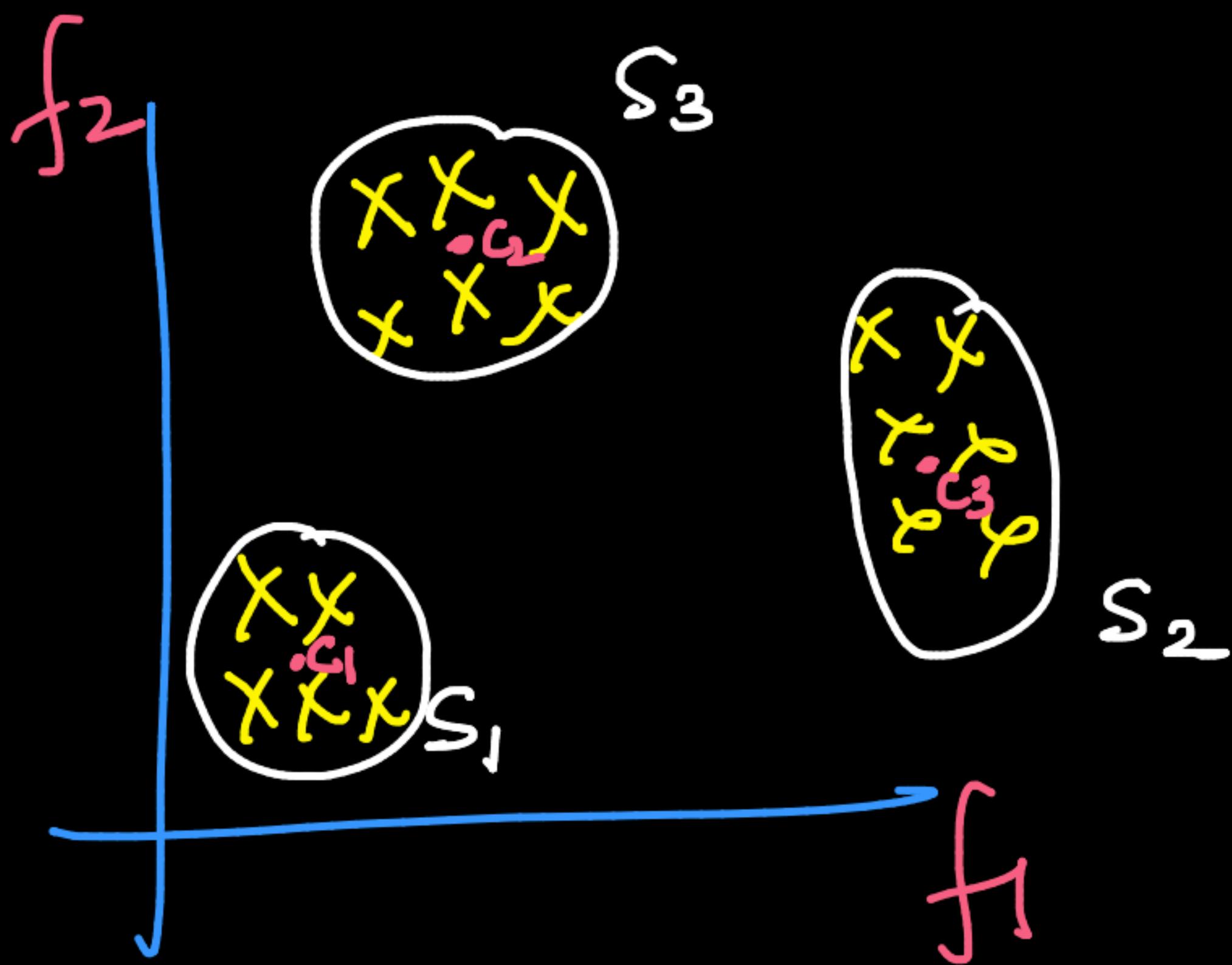
- Geom
- optimization
- Lloyd's algo
- Kmeans ++
- K-Means --

$$S_1 \cup S_2 \cup S_3 = \emptyset$$

$$S_i \rightarrow C_i$$

$$\checkmark S_i \cap S_j = \emptyset$$

each $x_i \in \text{one } S_j$



mean
value { Centroid₁ = $\frac{x_1 + x_2 + \dots + x_m}{m} \in S_1$

Centroid

$$c_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$

= mean vector

Size / cardinality
of set

k-means:
Centroid-based
approach

(later)

Alt!: hierarchy → Agglomerative
density → DBSCAN

Find $s_1, s_2 \dots s_k$

$c_1, c_2 \dots c_k$

optimization

$$\mathcal{D} = \{x_1, x_2, \dots, x_n\}$$
 each $x_i \in \mathbb{R}^d$

Find:- s_1, s_2, \dots, s_k
 c_1, c_2, \dots, c_k

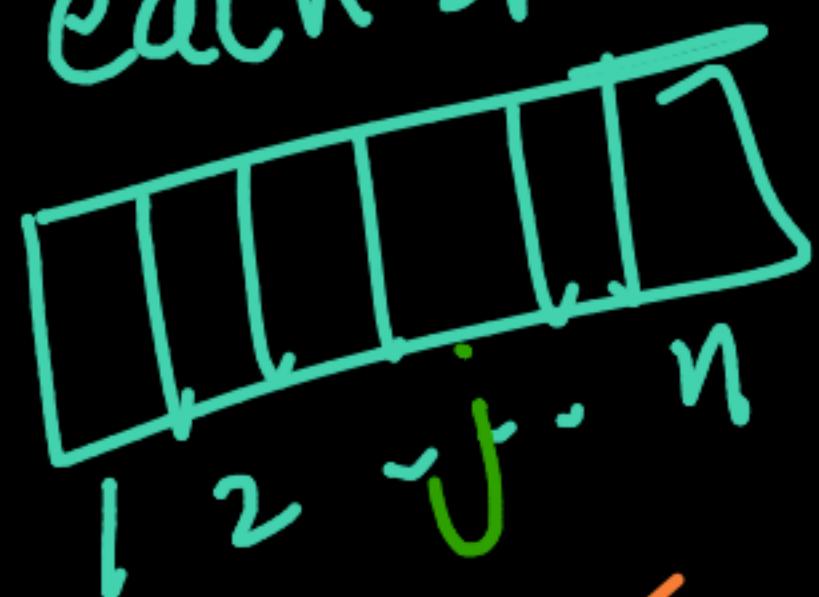
s.t. $\begin{cases} \text{max inter clust dist} \\ \text{and} \\ \text{min intra clust dist} \end{cases}$

between
and
within

GD

$\min_{S_1, S_2, \dots, S_K}$

each S_i



binary vec

$$\sum_{i=1}^K \sum_{x_j \in S_i} \|x_j - c_i\|^2$$

loss

min sum of
within
cluster
distances

s.t. each

$$x_j \in S_i \rightarrow \sum_{j=1}^n s_{ij} = 1$$

$$S_i \cap S_j = \emptyset$$

$$i \neq j$$

$$S_i^T S_j = 0$$

$$\frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$

Centroid

$$c_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$

$s_i = [1 \ 0 \ - \ i \ - \ n]$

if $x_1 \in s_i$
 $x_2 \notin s_i$

→ 0/1 binary vector of size - n

$$S_i \cap S_j = \emptyset \Rightarrow S_i^T S_j = 0$$

~~+ C_{i,j}~~

$S_i \cap S_j = \emptyset$

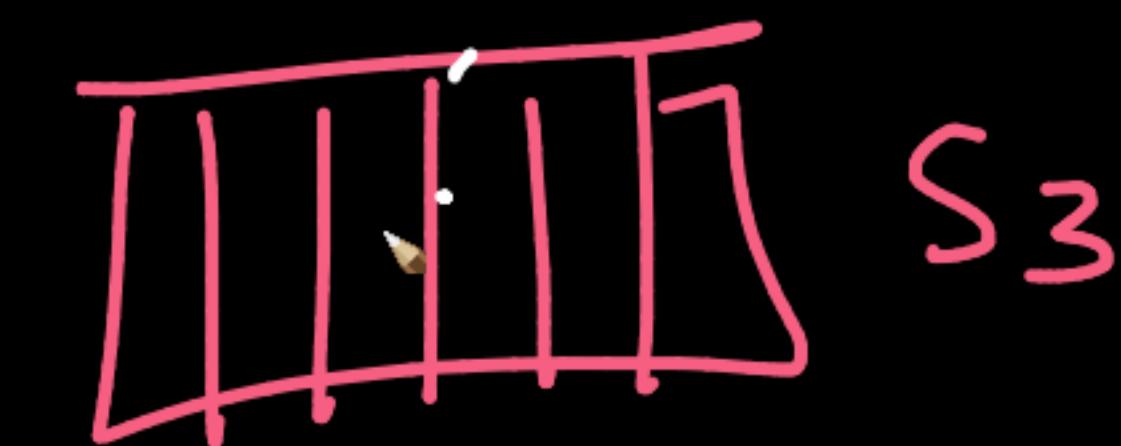
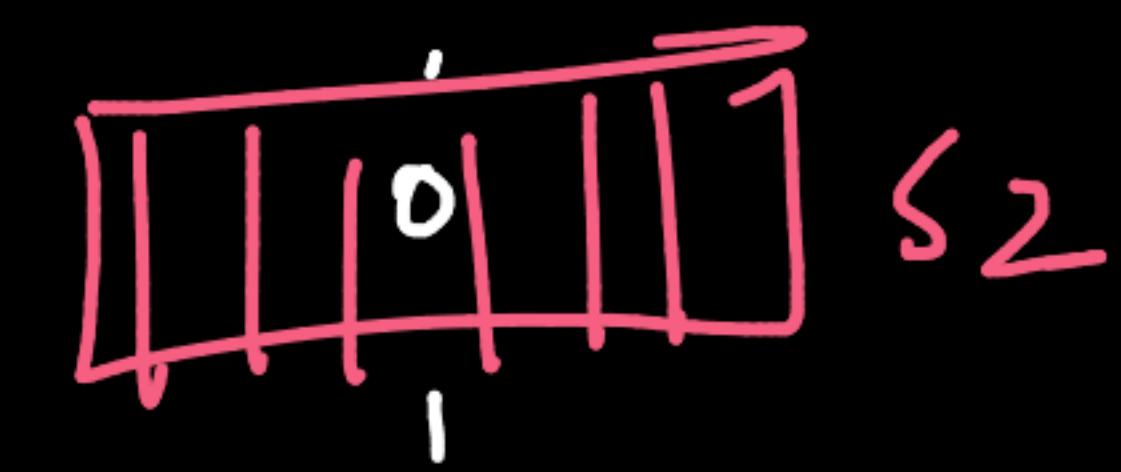
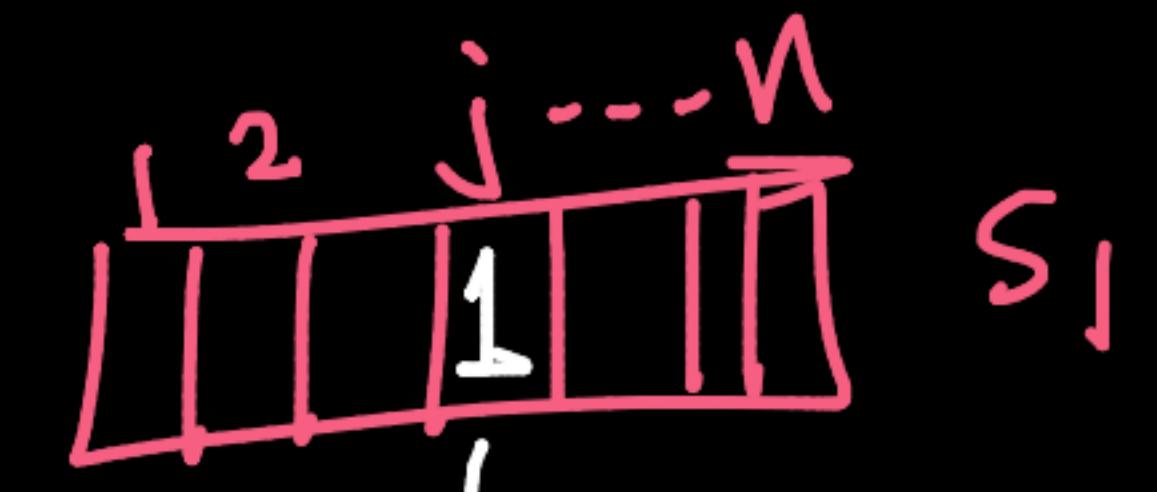
$\begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ | & | & | & | & | & | & | \\ 1 & 0 & & & & & \end{matrix}$

$\begin{matrix} 0 \\ | & | & | & | & | & | & | \\ 2 & 3 & 4 & 5 & 6 & 7 \end{matrix}$

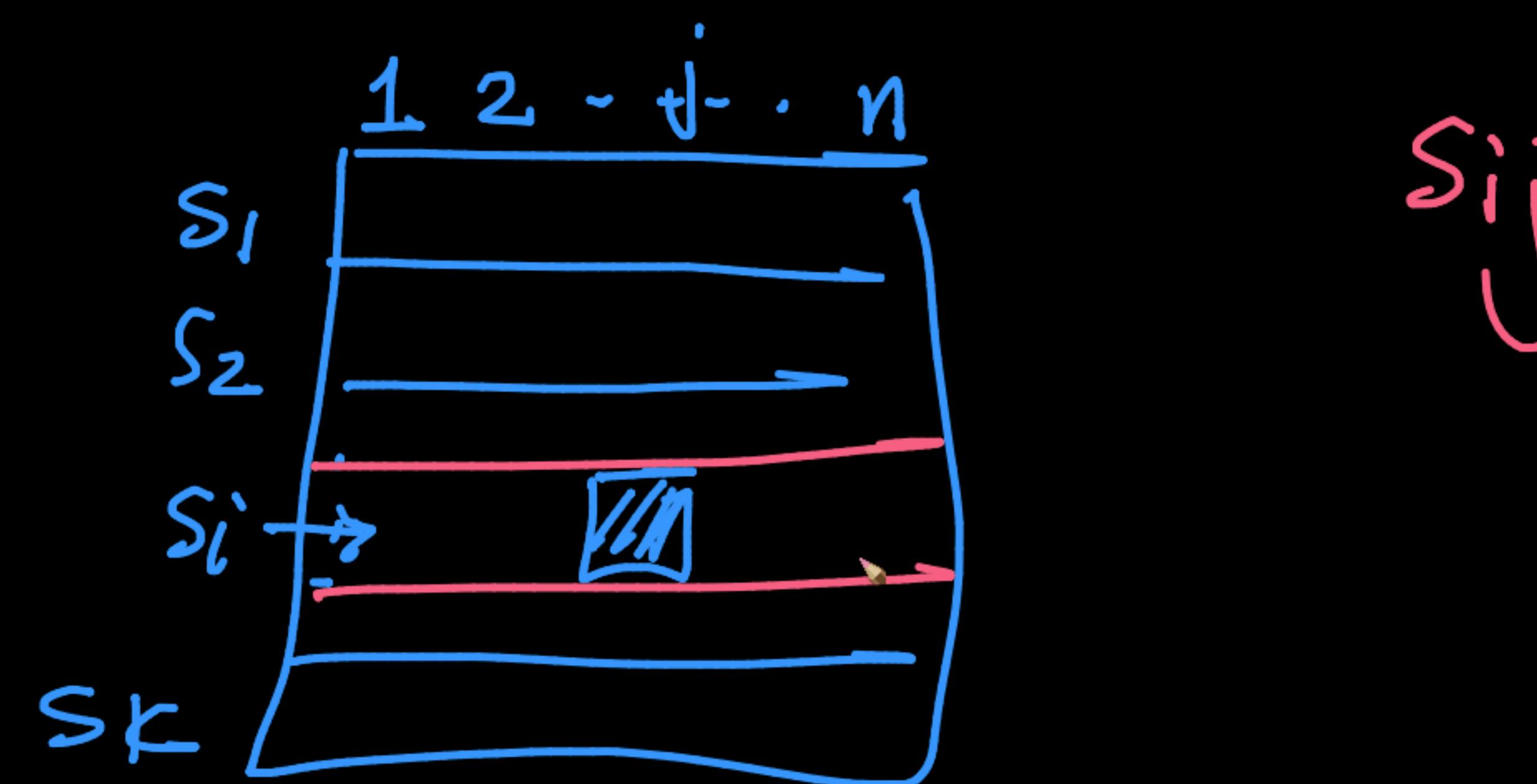
each $x_j \in S_i$ some

$$\sum S_i = [1|1|1|\dots|1]$$

$$\left\{ \sum_{i=1}^k S_{ij} = 1 \quad \forall j : 1 \rightarrow n \right.$$



Two red curved arrows pointing upwards and to the right.



s_{ij}

Opt!
~~not~~

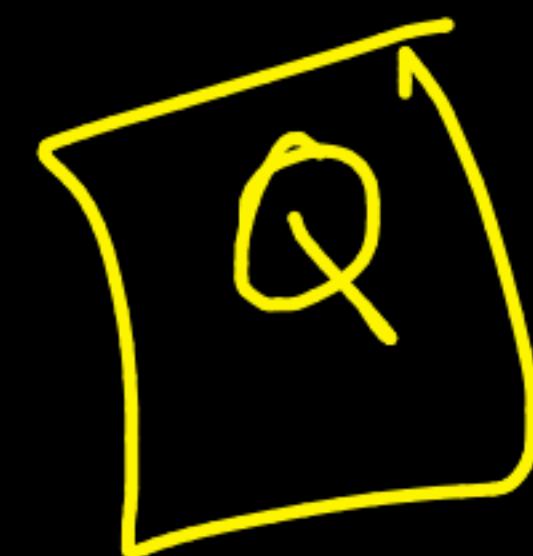
Loss + λ_1
+ λ_2 ...

s_i's

s_i's ...

L

$$\left\{ \frac{\partial L}{\partial s_{ij}} \neq 0 \right. \quad (s_{ij})_{\text{new}} = (s_{ij})^{-n} \frac{\partial L}{\partial s_{ij}}$$



Can we solve k-means using GD
— very clean formulation



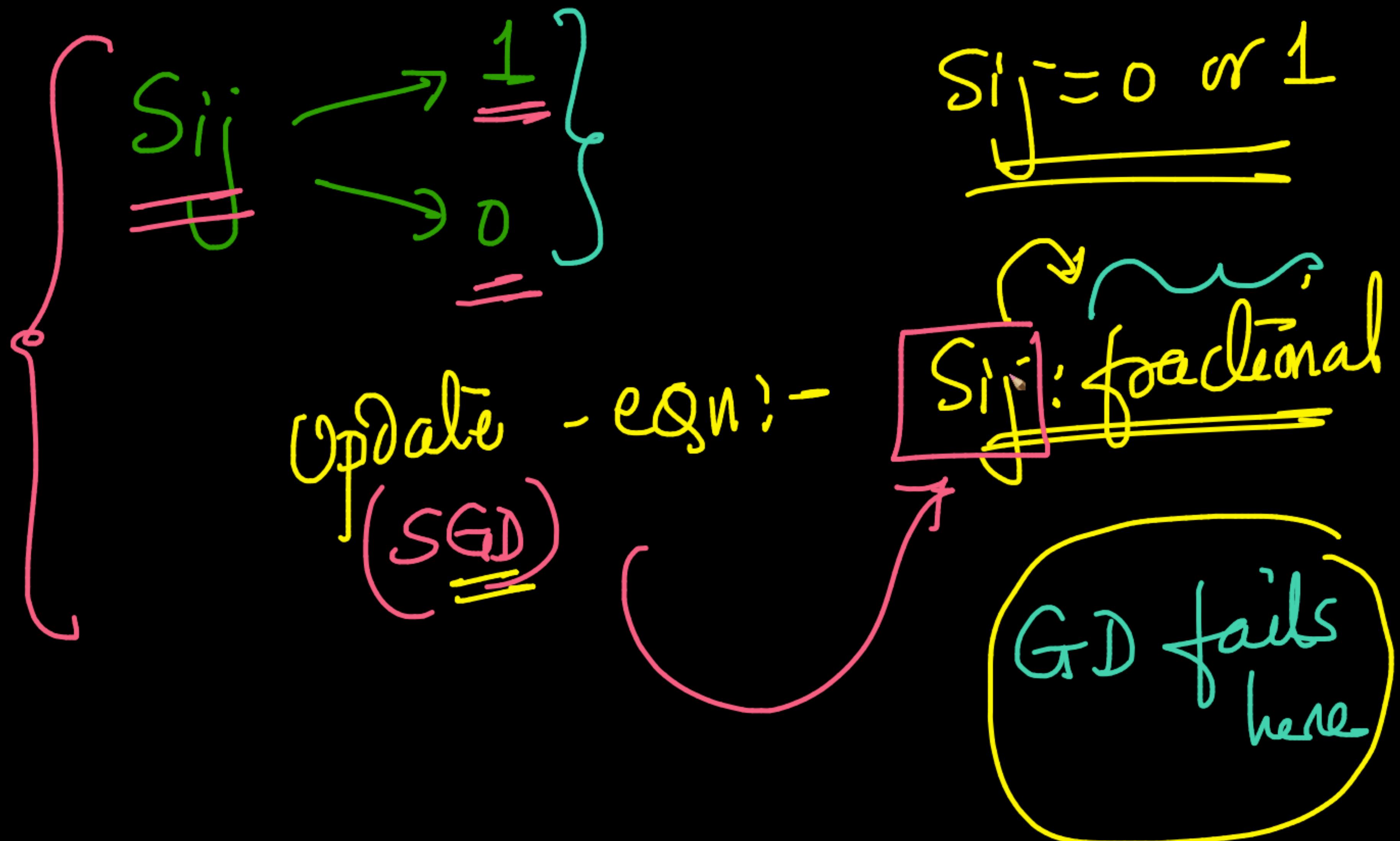
Yes, No

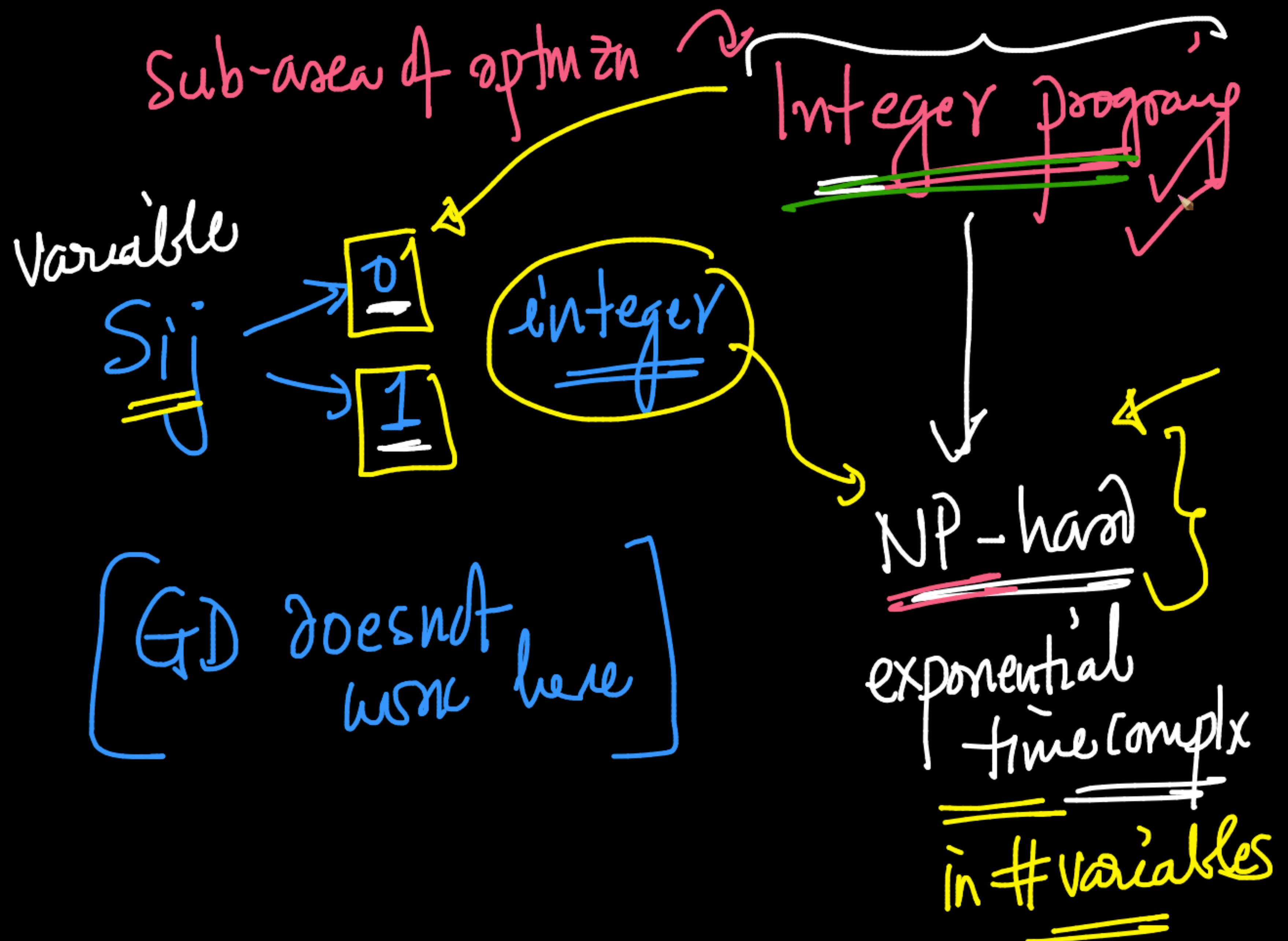
✓ PCA! Where was y_i 's \rightarrow 

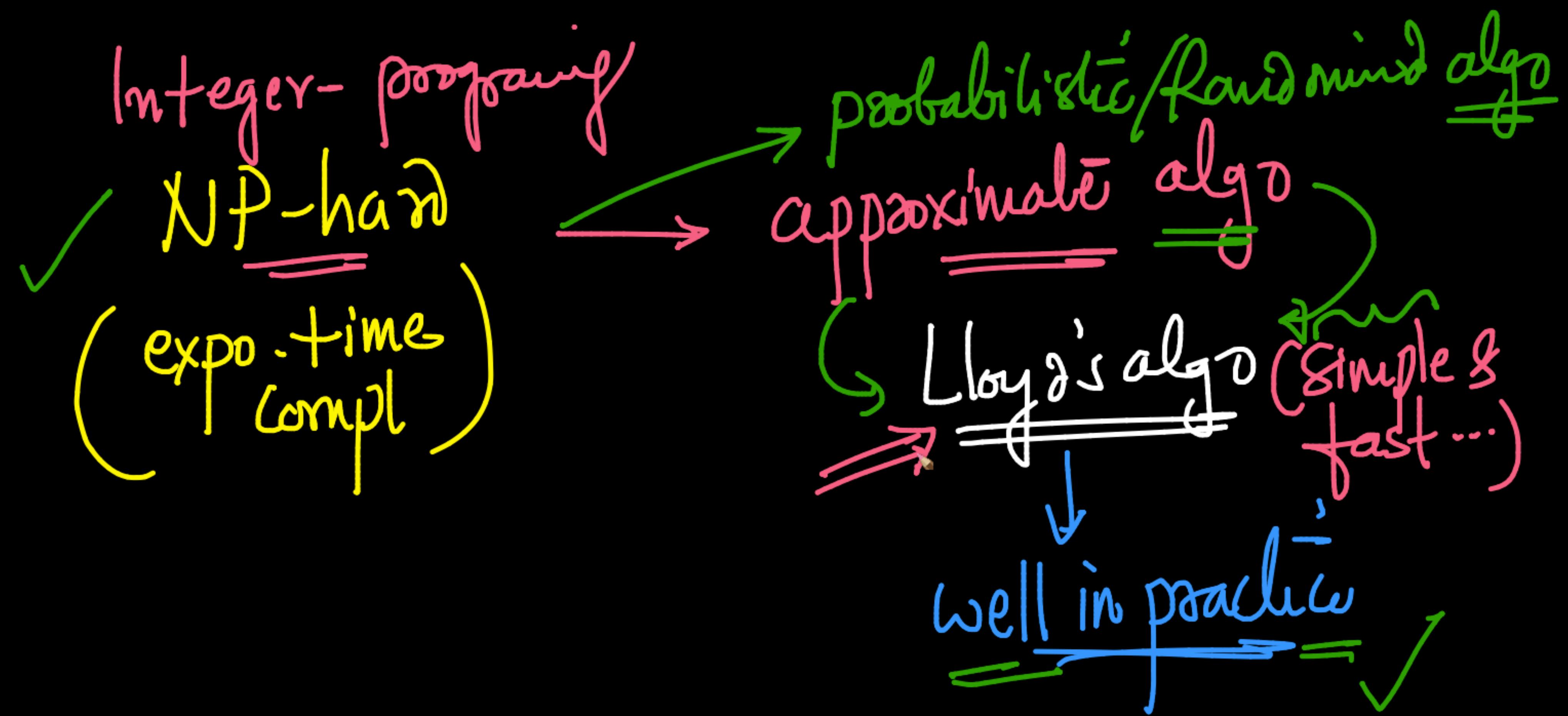
✓ When to stop: $\left| \frac{\partial L}{\partial s_{ij}} \right|$: v-small $\forall i, j$

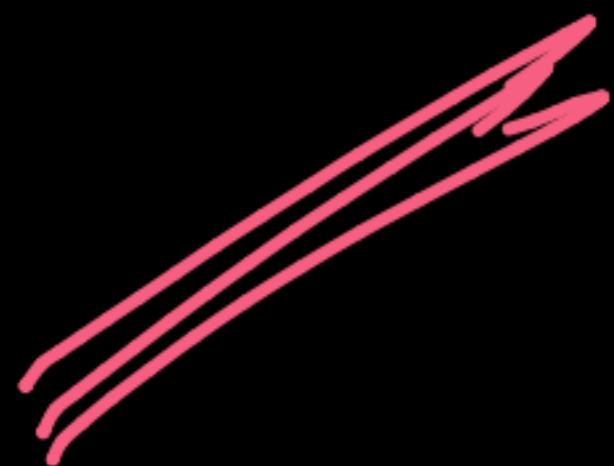
$$\checkmark S_{ij} = \begin{cases} 1 & \text{if } x_j \in S_i \\ 0 & \text{otherwise} \end{cases}$$

J. centroid $C_i \rightarrow S_i$





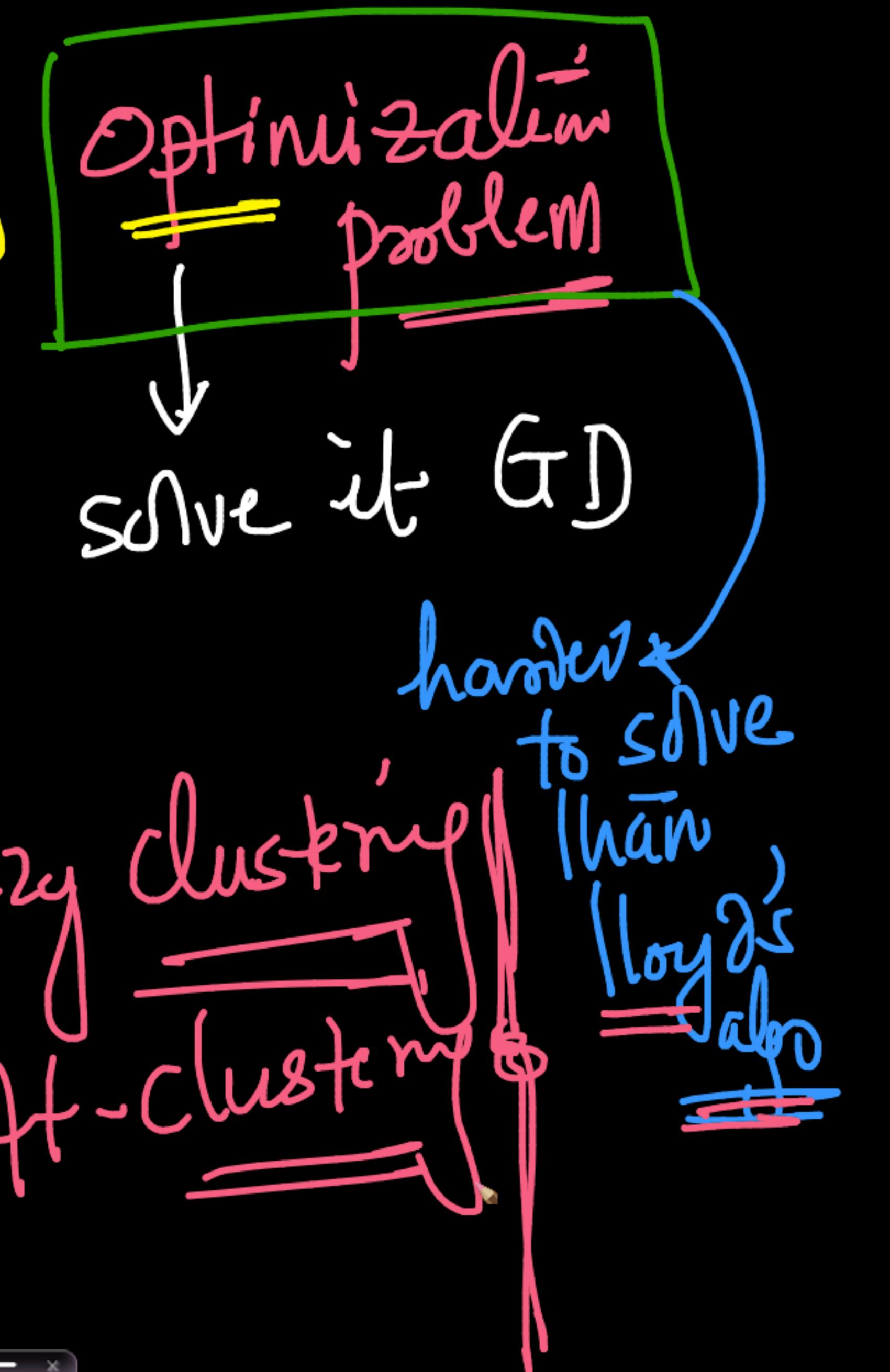




	1	2	...	n
s_1		0.8		
s_2		0.1		
s_3		0.1		
:		0.0		
s_k		0.0		

$$\underline{\text{sum}} = 1$$

fractional belonging
to a cluster





distance-metric

Lloyd's algo

Input: - $\{x_1, x_2, \dots, x_n\} = \mathcal{D}$

$K = \# \text{clusters}$

①

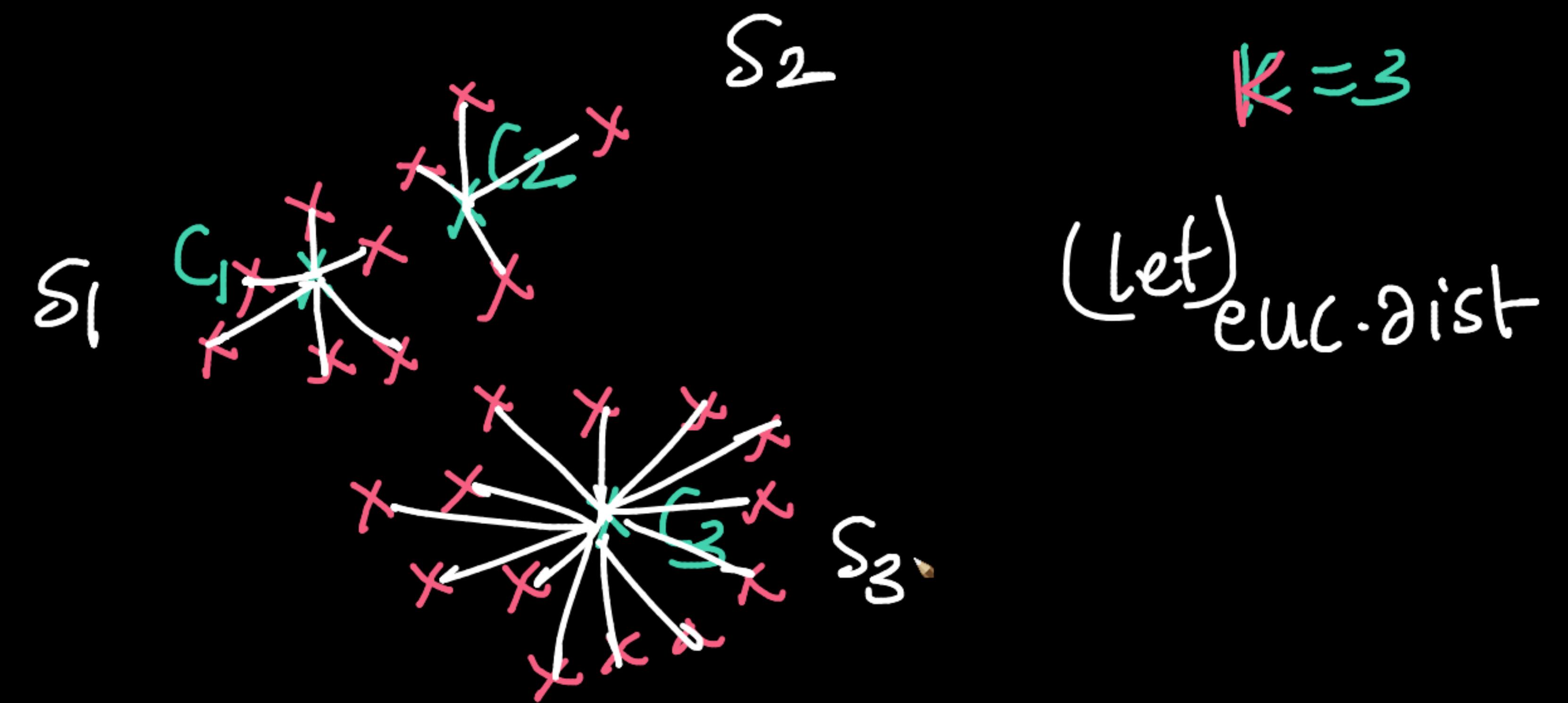
Initialization

: - pick K -pts randomly
from the n -pts
 (c_1, c_2, \dots, c_k)
initial centroids

②

Assignment

- [for each x_i in \mathcal{D}]
- select the nearest centroid: \tilde{c}_j (let)
 - add x_i to \tilde{s}_j
- dist - metric



③

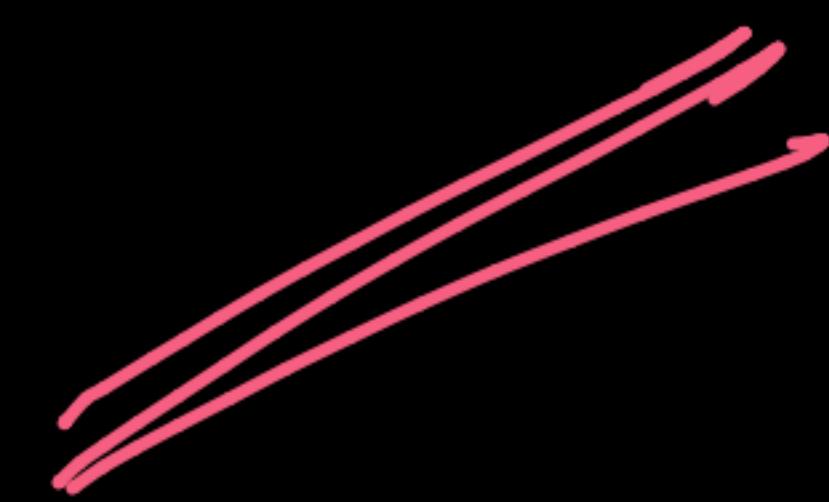
Recompute / update the Centroids

for j in 1 to k

$$\rightarrow c_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$$

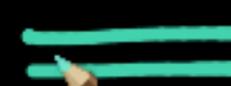
mean-pt

means \Rightarrow prone to outliers



Fix: kMeans++

examples
next class



initialization ↗
sensitive

Lloyd's algo



lots of problems

Fixes

Fix:
(hacks)
kMeans++

examples
(next class)

initialization
sensitive

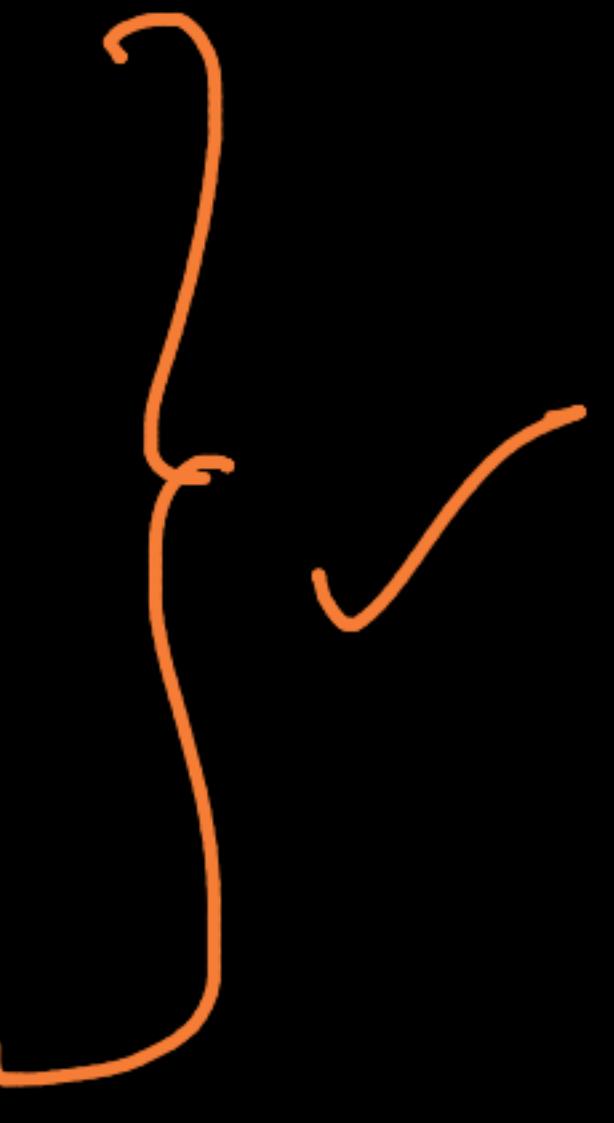
Lloyd's algo

↓
lots of problems

Next - class:

- Init. Sensitivity + k-means ++
- Best K
- (initializations) k-means -- (geom/Math/
real-world)
- K-medoids - PAM
- Code + Case study

K-means

- ↳ real-world ...
 - Math + algo ...
 - code + case-study
- 

Applications in different domains

∞ Curse of Dimensionality.ipynb × W Dunn index - Wikipedia × | W Cluster analysis - Wikipedia × +

colab.research.google.com/drive/1YHHN-_NsZnmlMqclMg0gSzDKiUwU8tTd#scrollTo=9wcHnJwtvj5c

+ Code + Text Reconnect

Reconnect

Up Down Reload Settings Copy Paste Delete More

import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

as dim -> INF; (max_dist - min_dist)/min_dist -> 0

max_min_ratio=[]
for dim in range(1,1000):
 # dimensionality
 d=dim

generate "n" uniformly random points in d-dim space in the unit hypercube of each dimension=1
n = 1000
X = np.random.rand(1000,d) # random numbers in the range [0,1] => unit hypercube
q = np.full([1,d],0.5) # [0.5,0.5,0.5,...] is the central point in the unit hypercube.

dist=[]
compute the pair-wise distances from each other
for i in range(n):
 dist_iq = np.linalg.norm(X[i]-q) # euclidean distance
 dist.append(dist_iq)

max_min_ratio.append((max(dist)-min(dist))/min(dist))