

(Q1)

→ euclidean-dist  
 $\underline{\underline{K\text{-means}}} : \text{optimization problem (course)}$

$$\min_{\underline{\underline{M_1, M_2, \dots, M_K}}} \sum_{i=1}^k \sum_{x_j \in S_i} \|x_j - M_i\|^2$$

min  $\cancel{M_i}$   
 $M_i$ : centroid of  $i$ th cluster  
 $M_1, M_2, \dots, M_K$   
 $d\text{-dim}$   
 $K \times d$

$x_j \in S_i =$  eucl. dist ✓  
 $S_i$ : set of  $x_j$ 's in cluster-i  
 $x_j \in$  one cluster  $S_i$   $\forall j \in$  how to identify clusters  
 Mathematically ✓

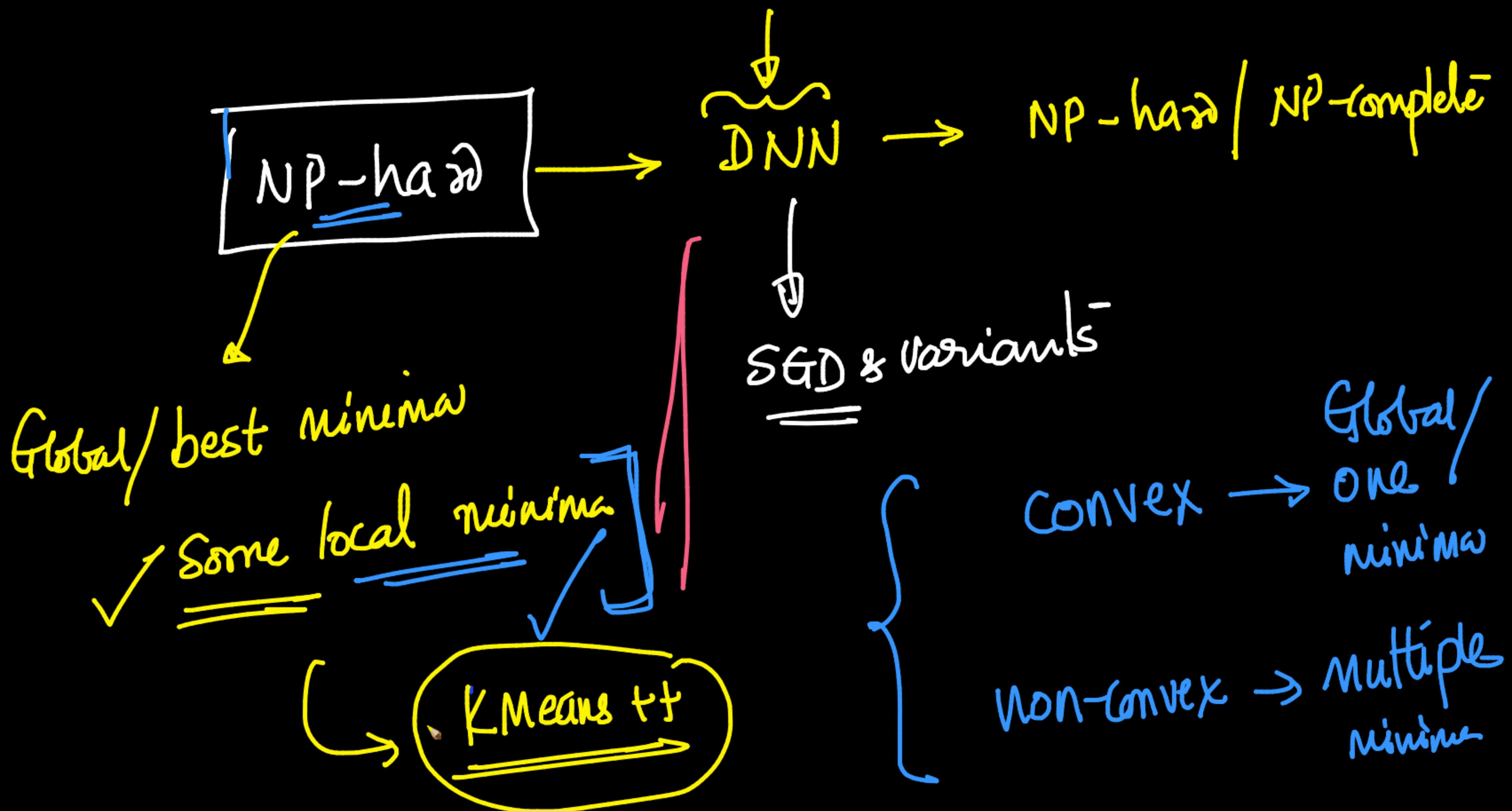
b) Why can't solve it using GD methods?

Log-reg: -  $\min$  log-loss +  $\lambda \|\underline{w}\|^2$

$$\checkmark \text{MF!} - \min_{A, B} \| C - \tilde{A}^T B \|_F^2 \quad \text{no } y_i's$$

- PCA :- optim 2n - prob no  $y_i$ 's
- ~~SVM :-~~

$$\|x_j - \hat{m}_i\|^2$$



$$\min_{m_i} \sum_{j=1}^k \sum_{x_j \in s_i} \|x_j - m_i\|^2 + \lambda_1 \boxed{\quad} + \lambda_2 \boxed{\quad}$$

~~non-differentiable~~ s.t.  $\left[ \begin{array}{l} x_j \in \text{one } s_i \\ s_i \cap s_j = \emptyset \quad i \neq j \end{array} \right] \rightarrow \text{set} = \text{operativ}$

{ C }

Can you modify kMeans-optimization problem  
s.t we can use GD methods [V.H]



$$\min_{m_i, w_{ji}} \sum_{i=1}^K \sum_{j=1}^n \|x_j - m_i\|^2 \cdot w_{ji} + \lambda_1 \boxed{\checkmark} + \lambda_2 \boxed{\checkmark}$$

#clusters  $\rightarrow$  #pts

assignment - mapping  
 $w_{n \times k}$

$$w_{ji} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

$w_{ii}$

s.t.  $x_j \in \text{one } s_i \text{ only}$

~~$s_i \cap s_j = \emptyset \text{ if } i \neq j$~~

$$\sum_{i=1}^K w_{ji} = 1$$

and

$$w_{ji} \in \{0, 1\}$$

GJ

tricky

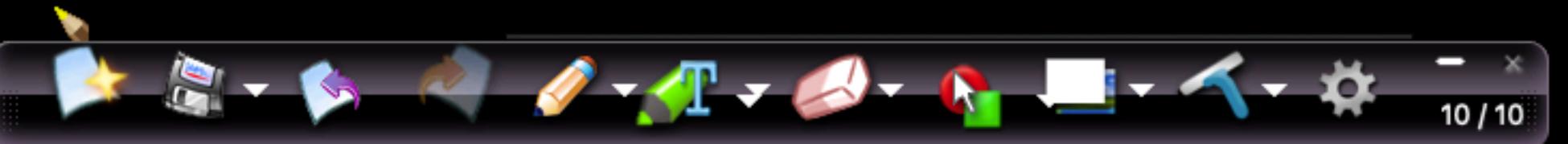
Mathematically  
as a differentiable  
fn. ...

$w_{ij} \in \{0, 1\}$ ,  $h_{ij}$

$w_{ij} = \underline{\underline{0}} \quad \text{or} \quad 1 \quad h_{ij}$

$\checkmark w_{ij} = 0$

$\checkmark w_{ij} = 1$

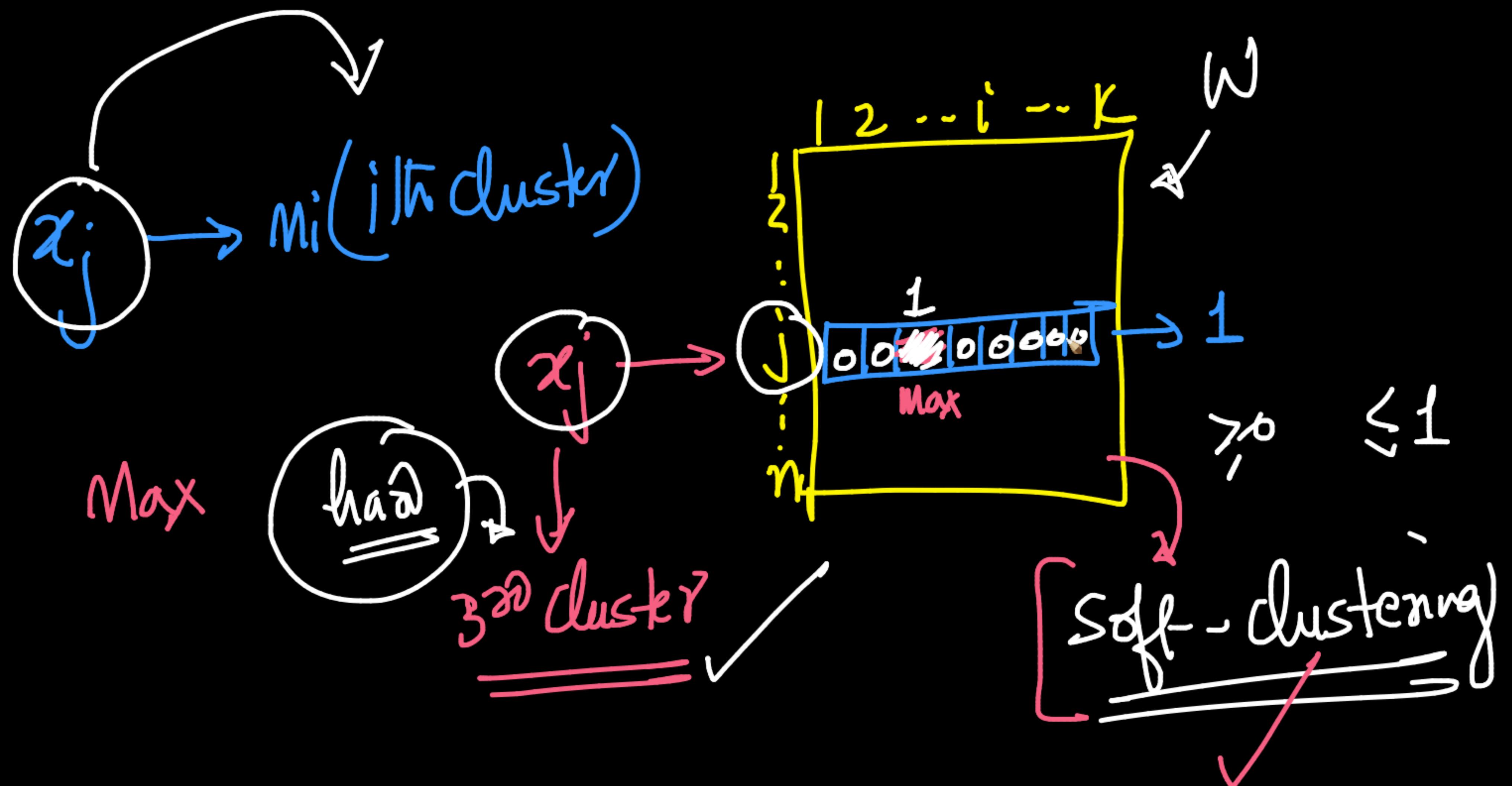


$$\begin{aligned}
 & \min_{\mathbf{w}_{ij}, \mathbf{m}_j} \quad \mathcal{L} \\
 & \text{s.t.} \\
 & \quad \sum_{i=1}^K \sum_{j=1}^n w_{ij} \|x_i - m_j\|^2 + \lambda \|\mathbf{w}\|_1 \\
 & \quad \sum_{i=1}^K w_{ij} = 1 \quad \forall j = 1 \rightarrow n \\
 & \quad w_{ij} \geq 0 \\
 & \quad w_{ij} \leq 1
 \end{aligned}$$

Problem?  
 $w_{ij}$  can take any value  
 $\mathbf{w}$   
 $\mathbf{x}$   
 $\mathbf{m}$   
 $\lambda$   
 $n \times K$

$w_{ij} = P(j \mid i)$  (Weight of cluster)

d-dim



(Q2)

✓ k-means

$$n = \text{~} 1B\text{-points}$$
$$K = 10$$
$$d = 20$$

(EASY) ✓  
practical

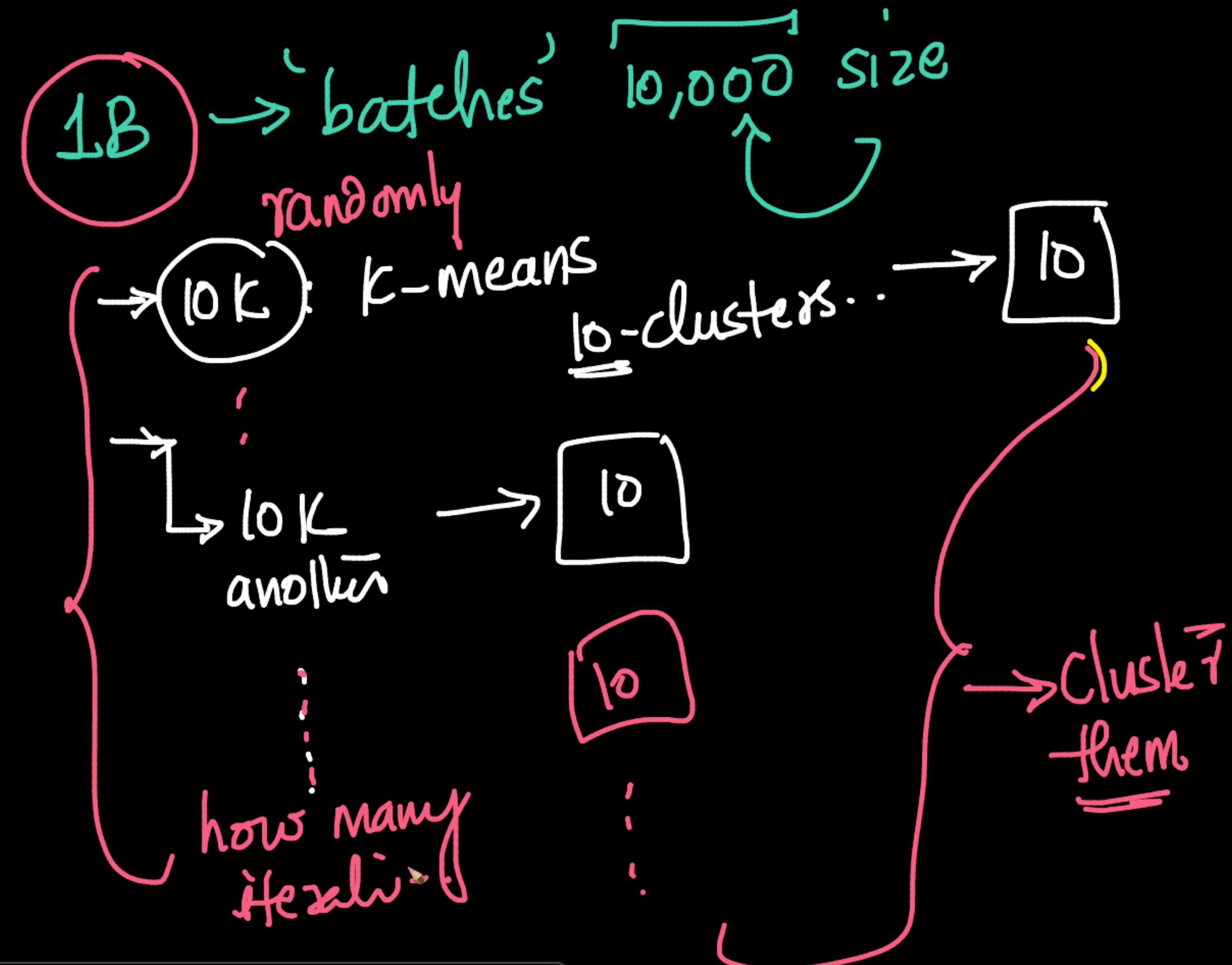
fast - clustering →  
- Single - box  
- < 30 min

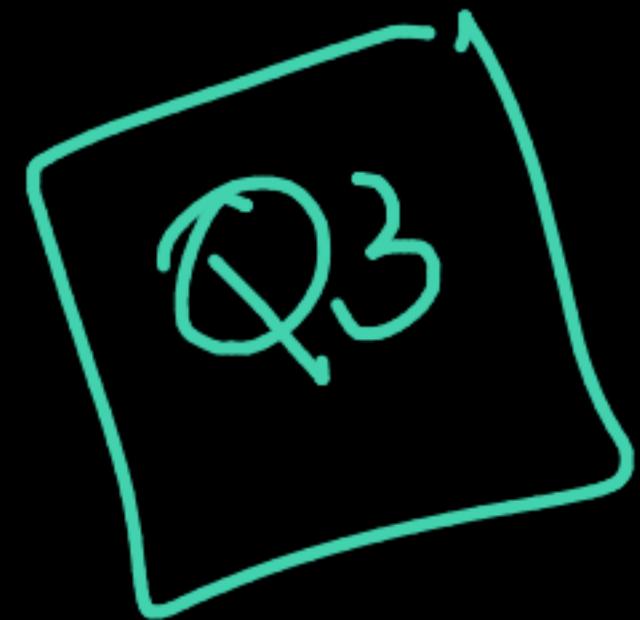
randomly sample  
1M points  
→ generalizes well  
→ KNN like to assign ...

Chi zagg

Soln 2:

[batch  
k-means]





Text-docs: - TF-IDF      fixed  
 $\begin{array}{|c|c|c|c|c|} \hline & 1 & 2 & \dots & d \\ \hline \end{array}$

Practical

a) K-means? → Can be done  
 $\downarrow$   
 $d$  is high → Computation (ignore)

euclidean dist

✓ Alternatives: reduce  $d$  → removing some low utility words  
 PCA SVD NMF  
 Autoencoders...

Alternatives: Cosine - sim / dist ✓



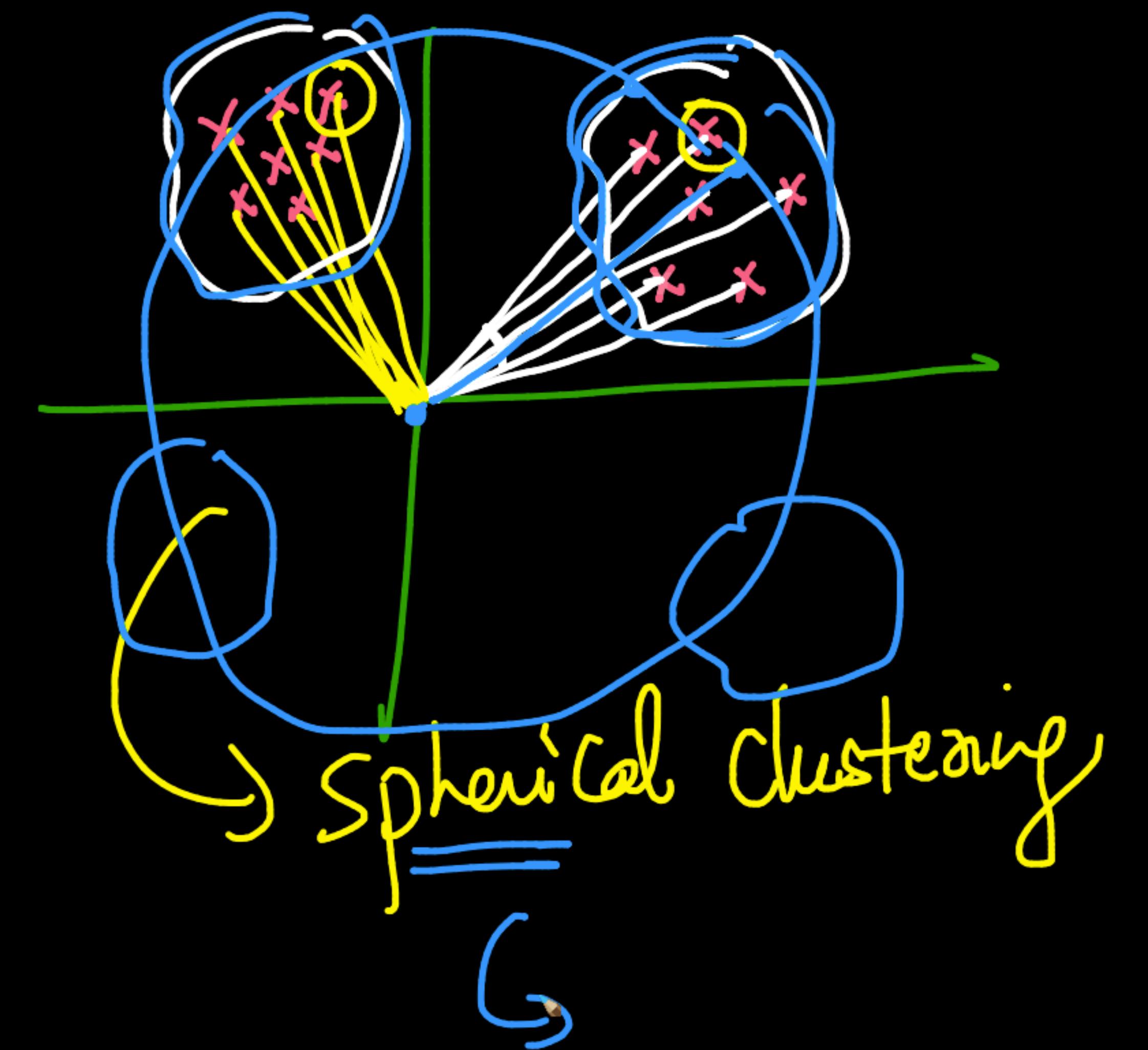
how would you compute centroid

HINT: think  
geometrically...

mean → euclidean-dist



→ cosine-dist

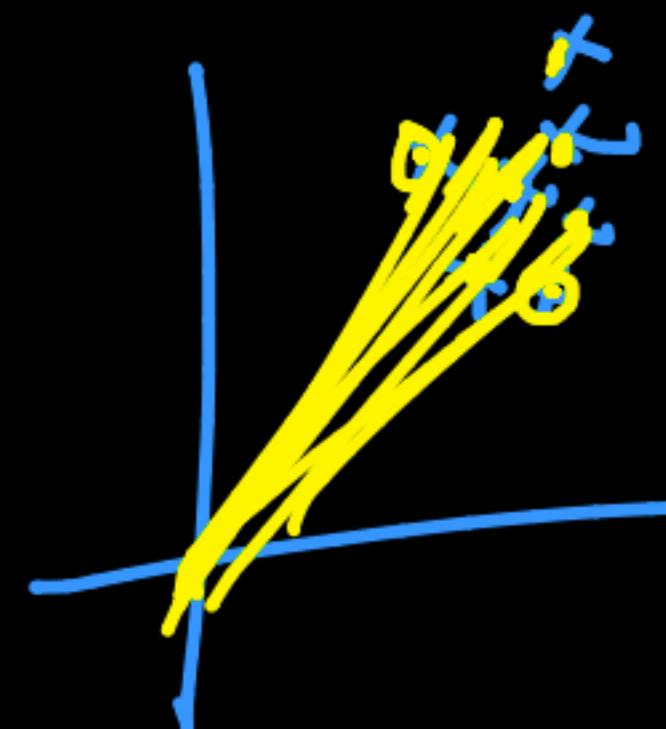


Centroid : Unit vec

Mean-angular-sew

**Q4**

## DB SCAN



range queries  
core pts, border  
outliers

reduce dim

Cosine sim

vary  $\epsilon$  carefully

①

2

dim is large - curse of dim.

$d=1000$

Cat (one-hot encoded)

precautions in practice

Minpts,  $\epsilon$

: params

Small change in

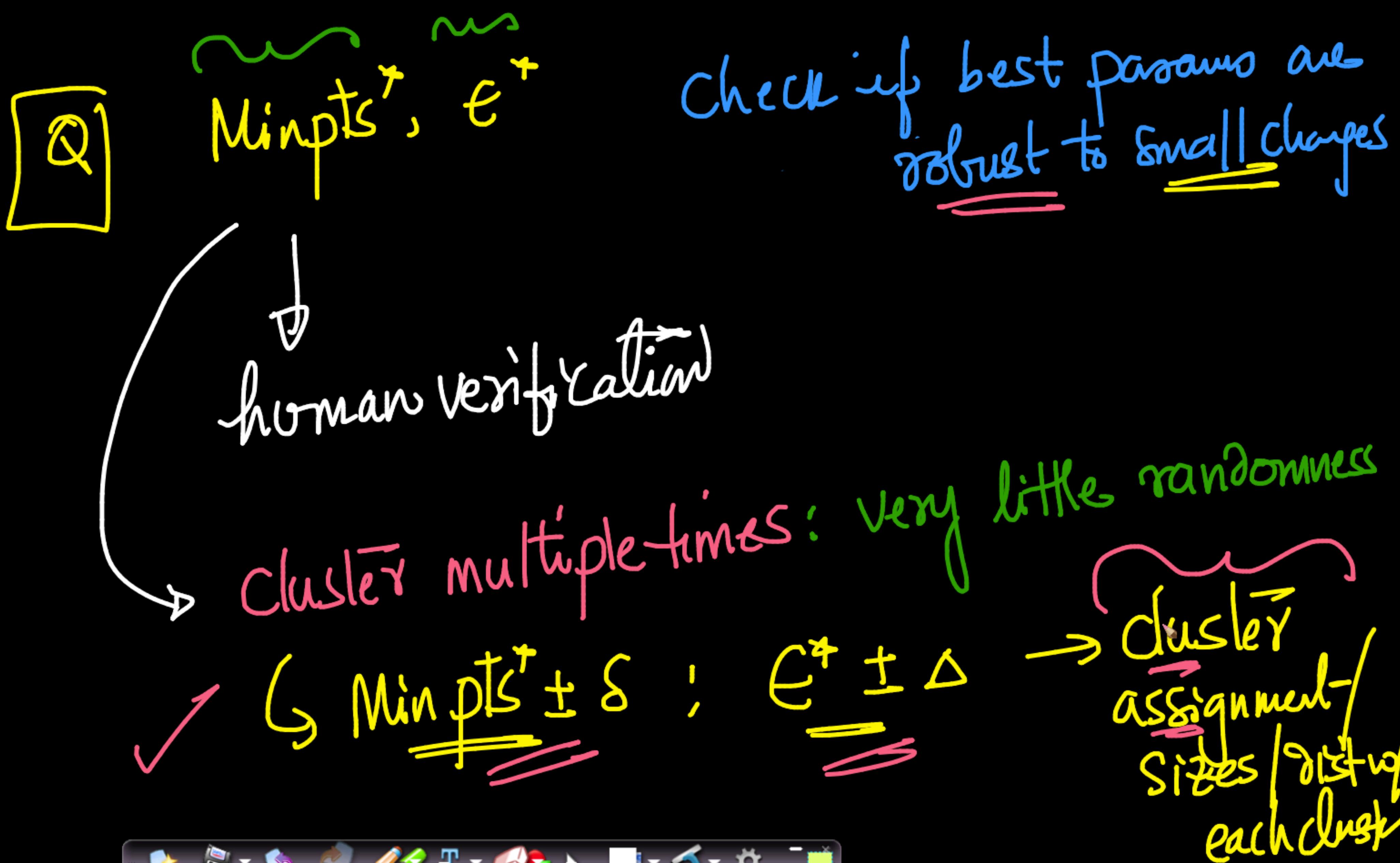
params

different clusters

2

DICE  
 $(Minpts, \epsilon)$   
best

elbow  
method

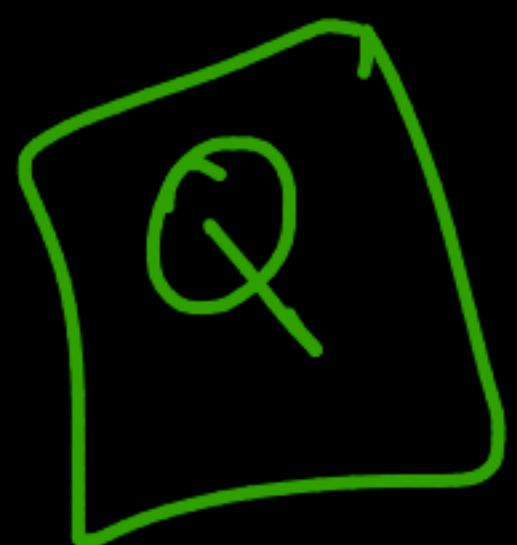


Q

{ Single-linkage: MIN-aggl...  
Complete " : MAX - "

easy

→ cases where they fail (Couse)



KMeans ++

initialization



how to optimize KMeans ++ for speed

↳ parallelizes (not trivially)

①

✓ randomly  
Sample

OK

+



https://www.airmeet.c Airmeet: Interview que Applied AI Course(Scr Unsupervised learning Hierarchical clustering DBSCAN.docx - Goog Microsoft PowerPoint Find Open Datasets a New Tab

Search Google or type a URL

Apps Google Bookmark Google Scholar Share to Google+ Stumble! Live Tv Channels... User generated g... User Generated C... Best implementati... Mod (video gamin... YouTube - Automa... Kinect hack move...

Gmail Images

Google

Search Google or type a URL

Applied Roots Airmeet Google Drive youtube.com GATE CS App...

Google Drive Applied Roots Scaler by Inte... Inbox (1,299) Add shortcut

Customise Chrome

pool\_fn → remove outliers...  
pick farthest pt

