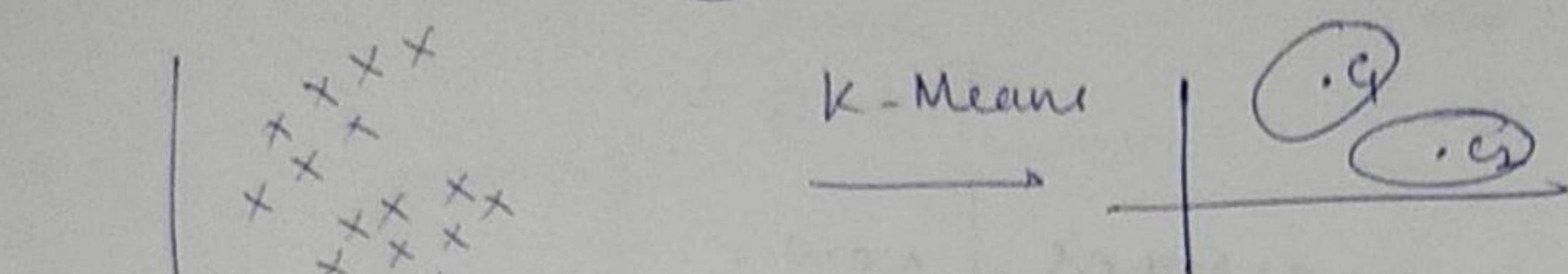
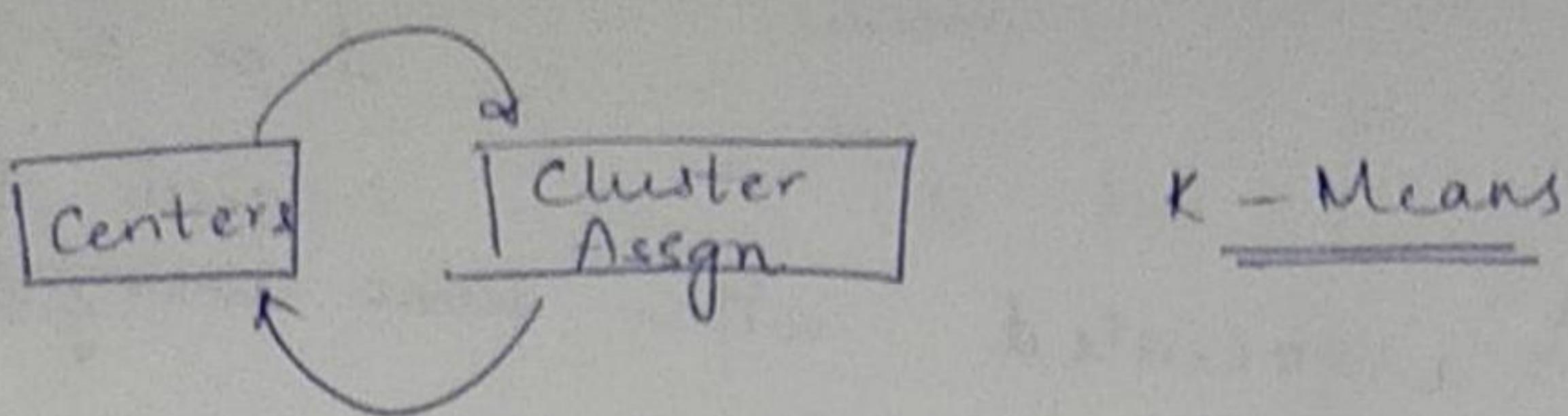


5/11/18

Lec-23

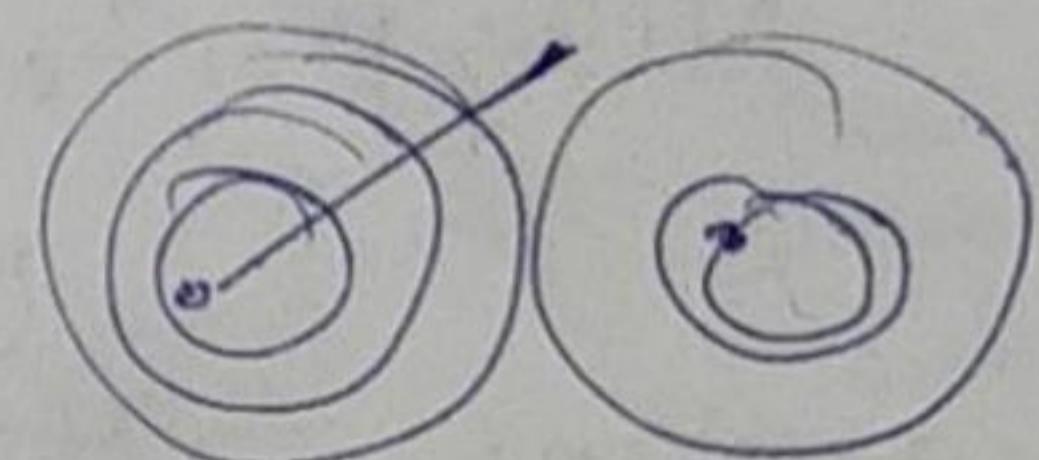
K-Means algorithm

we didn't know the clusters and the centers



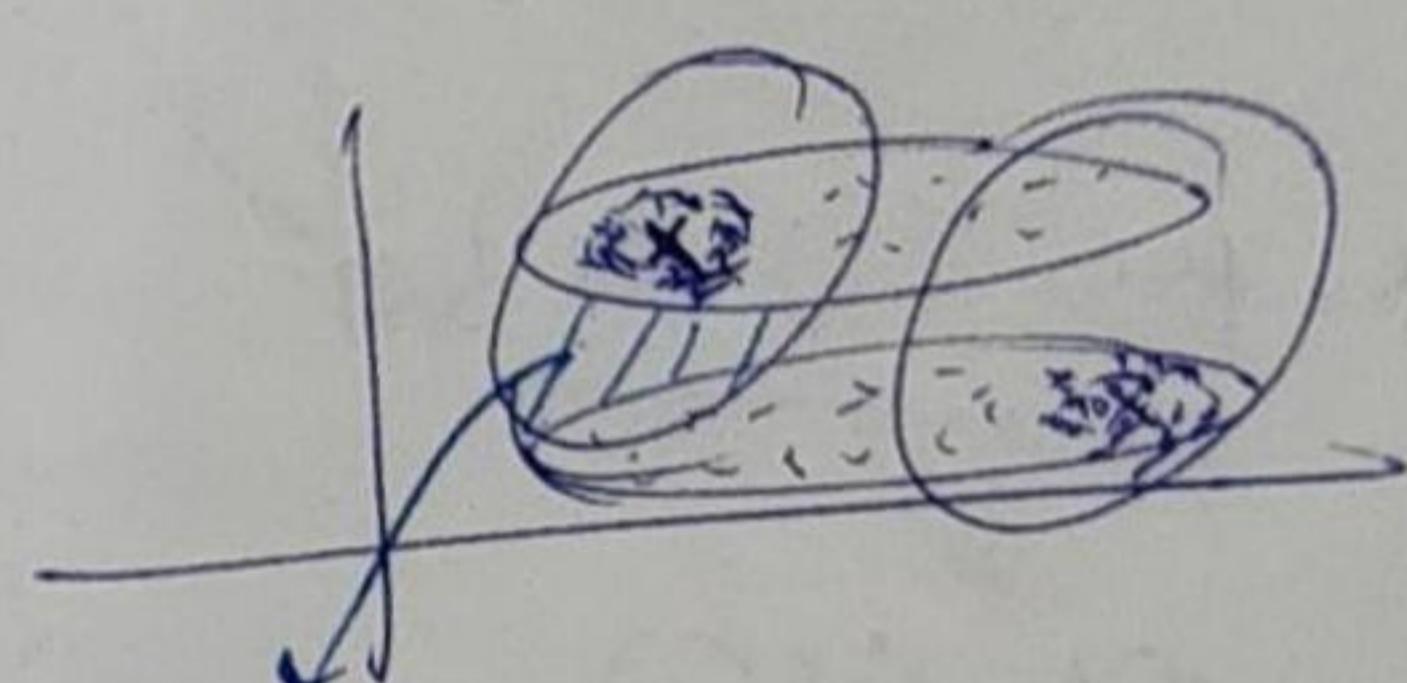
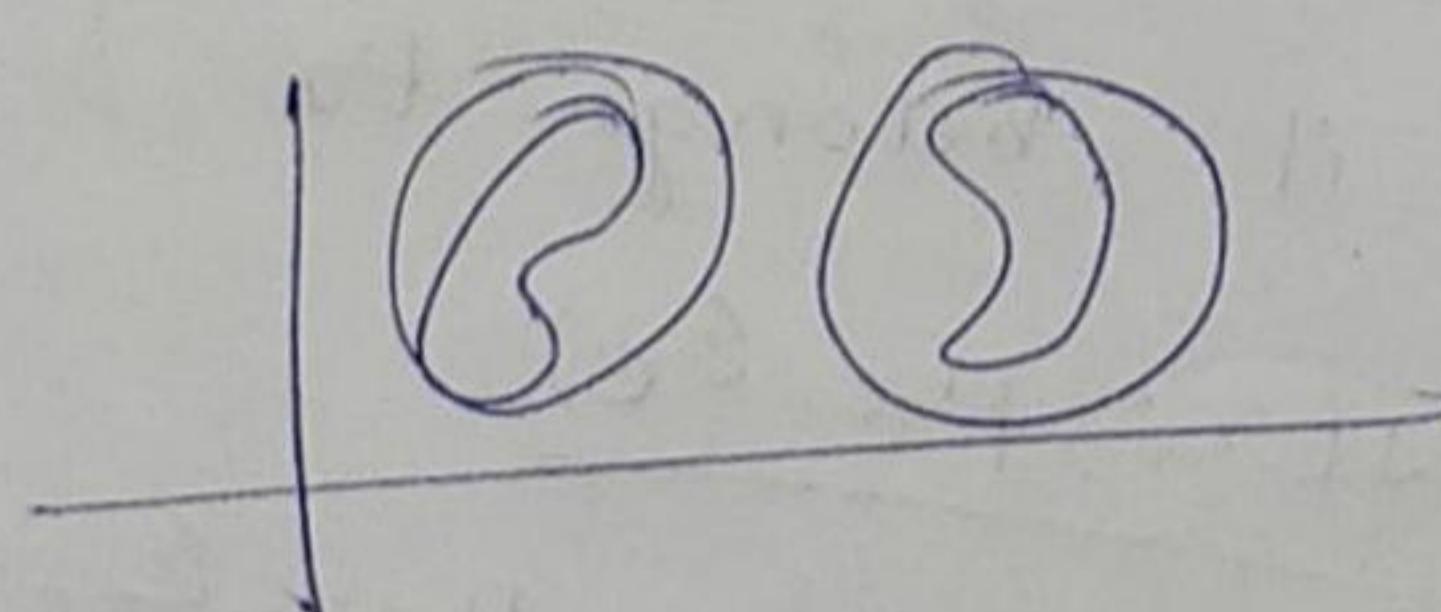
K-Means, how to ensure  $c_1$  and  $c_2$  go to the cluster center

Circles/Spheres are produced, if we consider the Euclidean distance b/w means



So, K-means, won't be good, if points deviate from the cluster // shape's are non spherical.

$\downarrow$   
shape=spherical is not the exact condition  
// Boundary should be separable by spheres



Doesn't think about these gaps //

~~Play~~ → K-means can't handle it.

Soln 1) Make 20 clusters and then combine 'em

2) Mahalanobis distance

$$\sqrt{\sum_i (x_i - \mu)^T (\Sigma^{-1}) \mu} \rightarrow \text{Euclidean}$$

Weighted distances →  
using covariance matrix

~~$\sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$~~

$$(x - \mu)^T \Sigma^{-1} (x - \mu)$$

Covariance matrix of a cluster.

• Euclidean dist  $\rightarrow \Sigma = \sigma^2 \mathbb{I}$

But, now:  $\Sigma$  = generic,  
↳ ellipsoidal Gaussian density

• Ambiguity -  
at intersection?

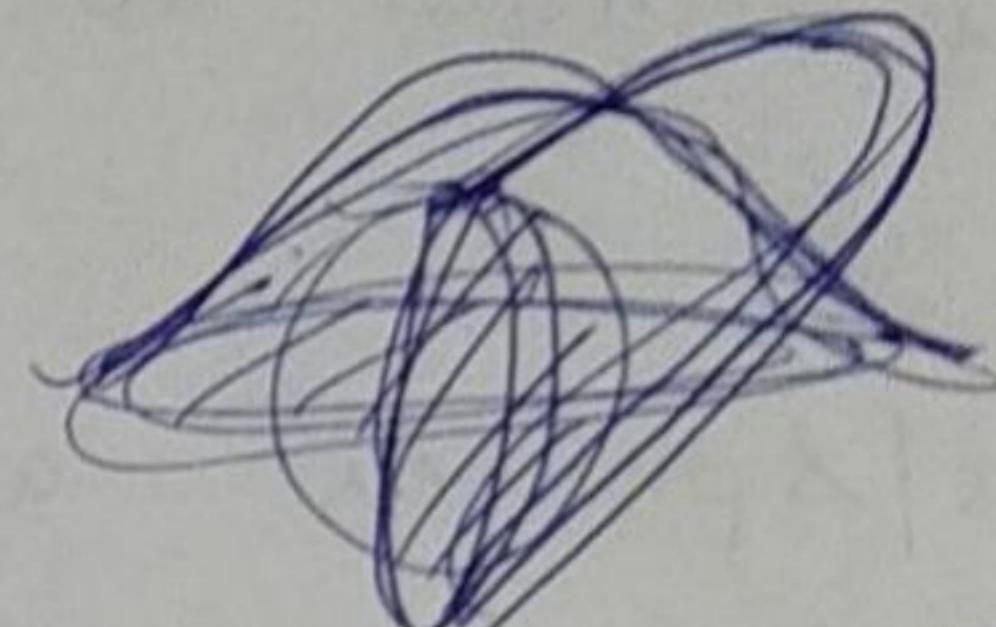
Soft assignment.

Every point is associated with some prob of belonging  
using

So,  $\textcircled{2}$  clusters are Gaussians, now.  
Each cluster  $c$ , has  $\mu_c$  and  $\Sigma_c$  now

$$N(x; \mu_c, \Sigma_c)$$

• Gaussians: weighted a Number of points Belonging to the gaussian



$$\underline{N(x, \mu_c, \Sigma_c, \pi_c)}$$

$\downarrow$   
(Cluster center)  
(Cluster info)

To find which cluster it belongs to

~~Likelihood~~  
~~Prob of class~~

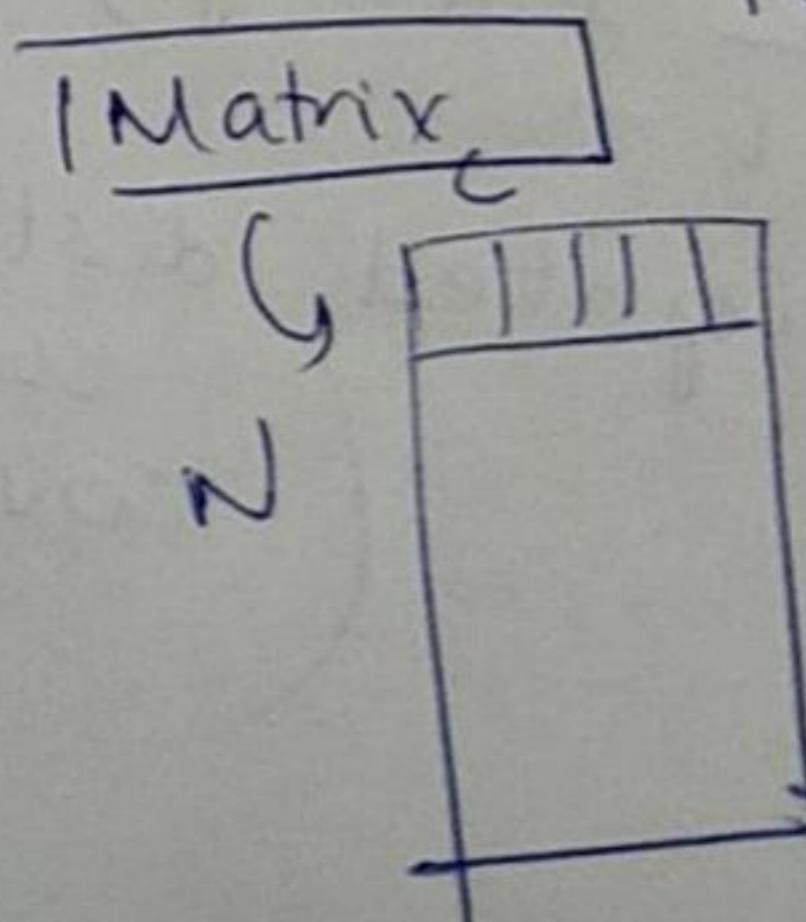
$$\frac{\pi_c N(x_i; \mu_c, \Sigma_c)}{\sum_j \pi_j N(x_i; \mu_j, \Sigma_j)} = r_{ic}$$

$$=$$

responsibility  
of  $\textcircled{c}$ th  
class for  
point i.

The  
cluster assignment  
(permutation)

$\underline{N \times c}$  matrix



one row = 1

$i^{\text{th}}$  row  
has all the likelihoods

Sum of partial assignments  $\eta_i$  = total assignments to a cluster

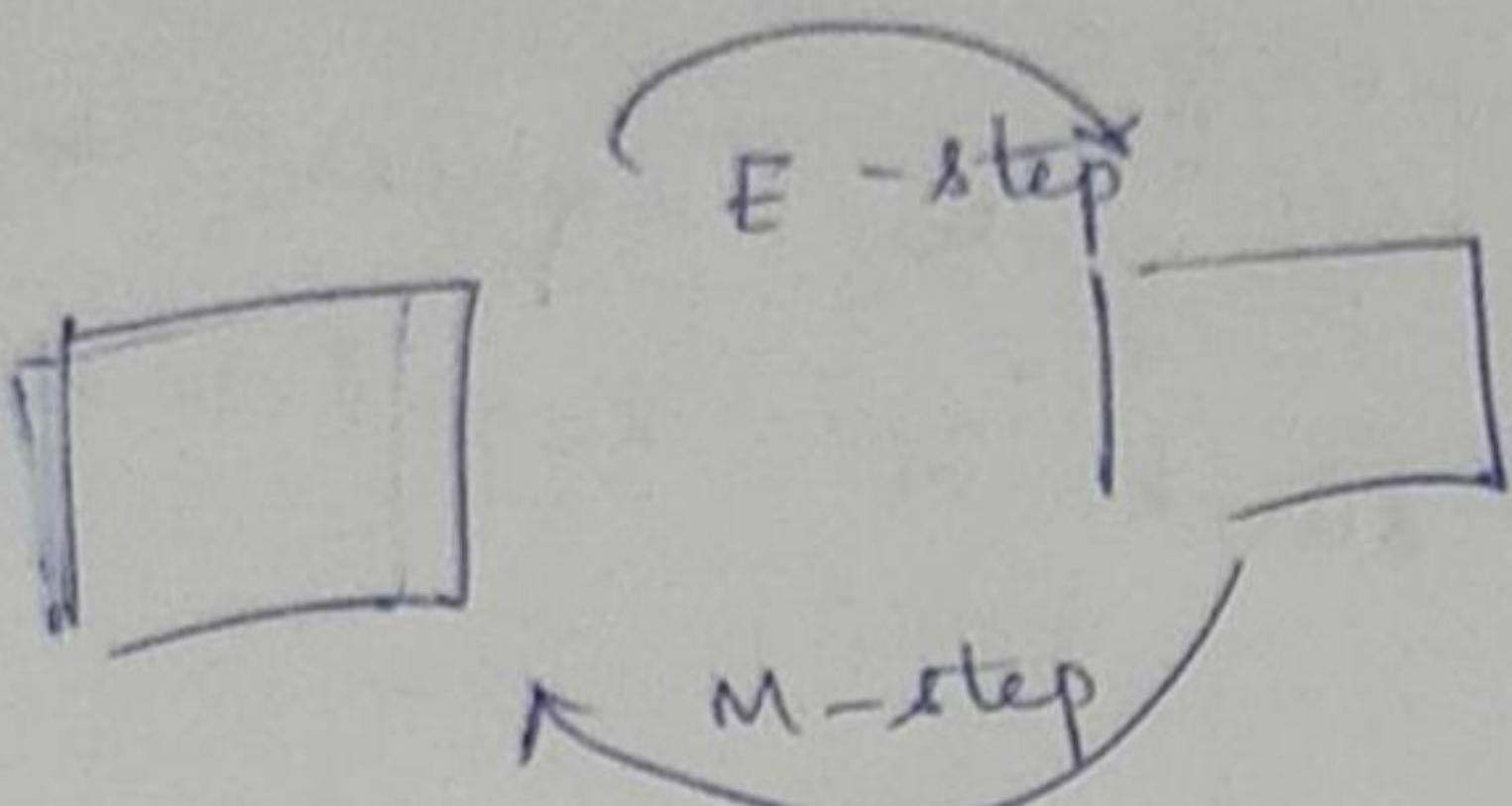
$$m_c = \sum_{i=1}^n \eta_{ic}$$

$$\pi_c = \frac{m_c}{m} = \frac{1}{\sum_{i=1}^n m_i}$$

$$\mu_c' = \frac{\sum_{i=1}^n \eta_{ic} x_i}{\sum_{i=1}^n \eta_{ic}} = \frac{\sum_{i=1}^n \eta_{ic} x_i}{m_c}$$

$$\hat{\Sigma}_c = \frac{1}{m_c} \sum_{i=1}^n (\eta_{ic} x_i - \mu_c')^T (\eta_{ic} x_i - \mu_c')$$

$\hat{\Sigma}, \hat{\mu} \rightarrow$  To denote new  $\mu, \Sigma$



EM-algorithm

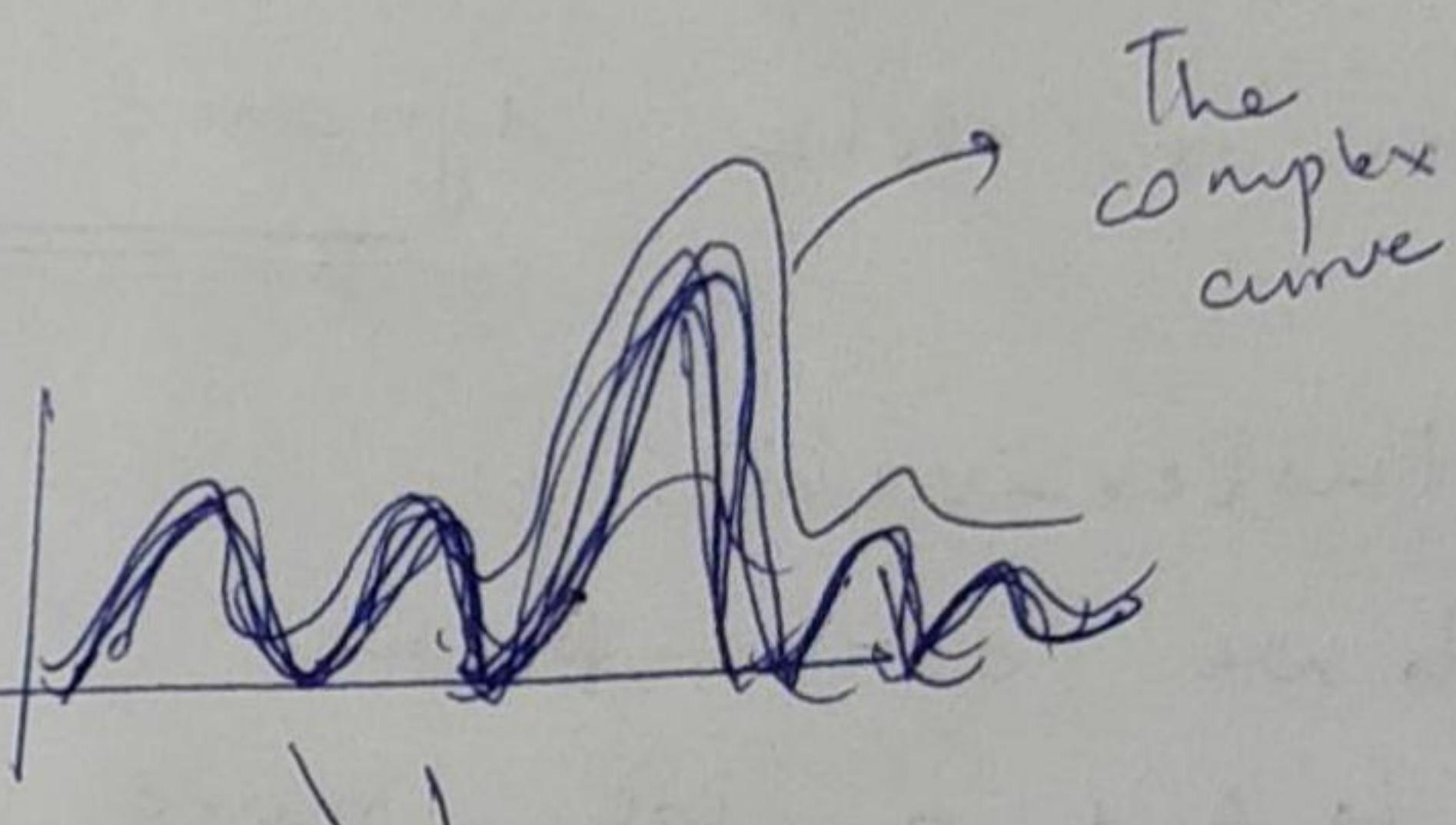
K-means is an EM-algorithm

Gaussian mixture model

MLE parameter estimation

It is a method

$$\text{The } \mu_c' = \frac{1}{m_c} \sum_{i=1}^n \eta_{ic} x_i, \text{ p}(\text{mean of data}) / \text{(no. of data)}$$



The gaussians

The complex curve

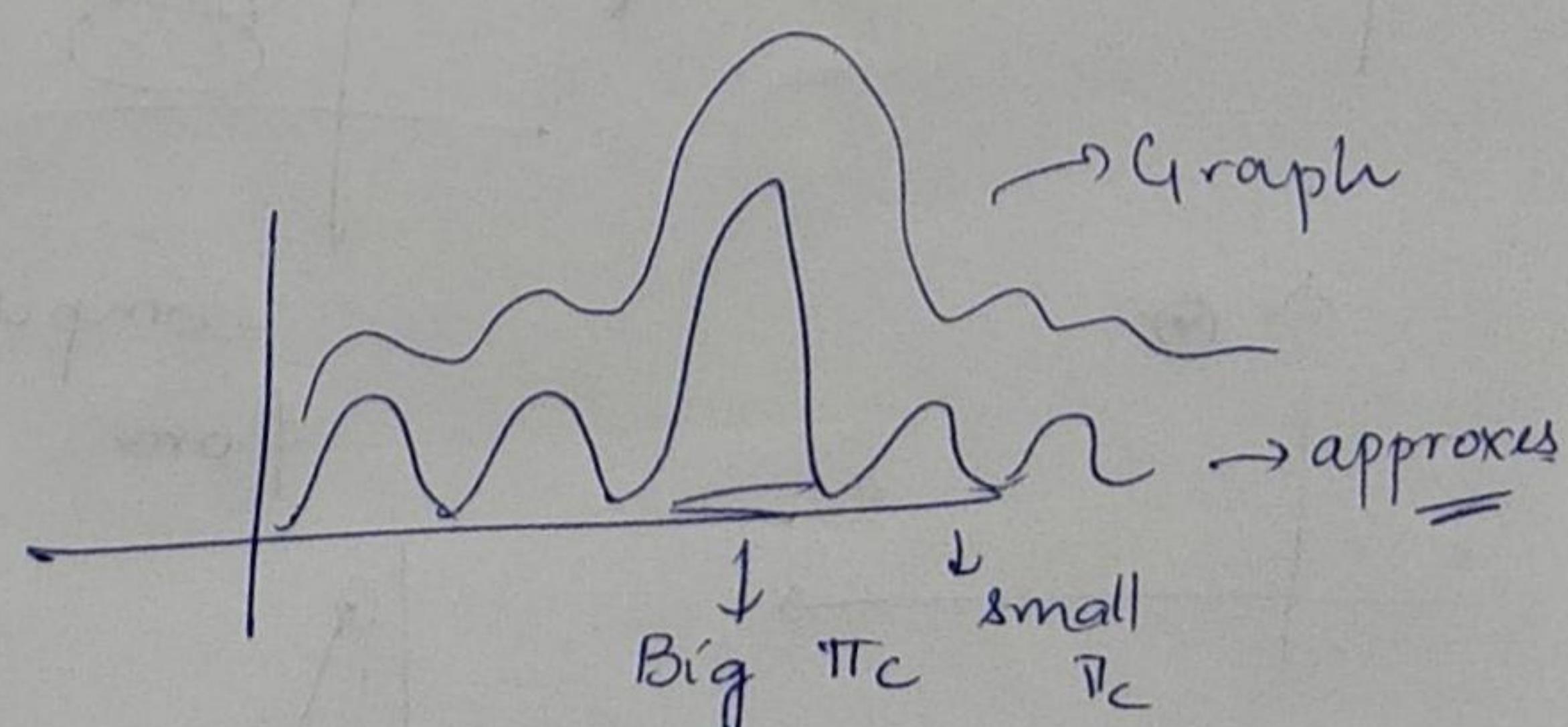
That is why E → is estimation and

M →

EM-Problem:

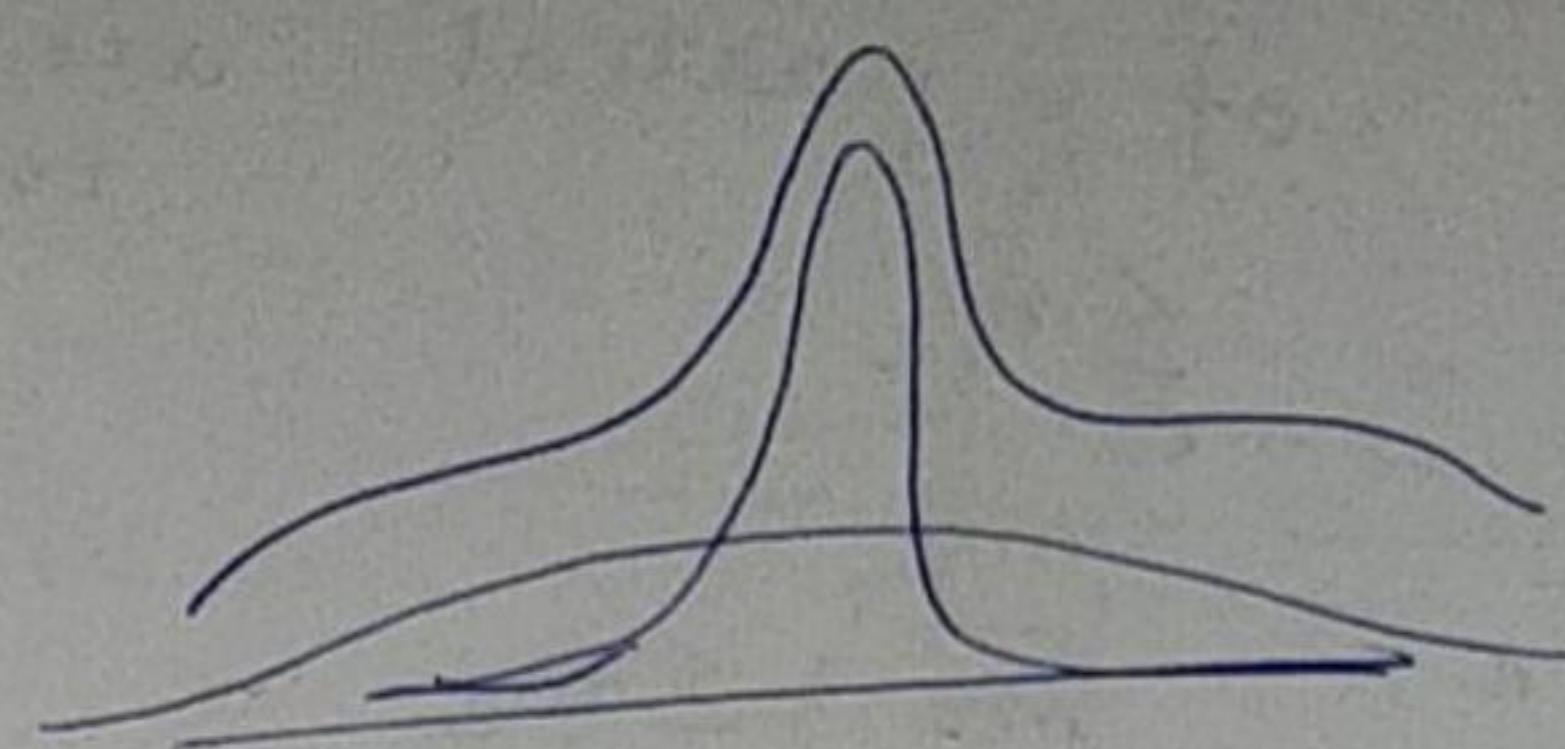
z is ~~unknown~~ the latent variable ( $\eta_{ic}$ ), so that knowing latent variable can help us solve the prob ← nice-versa

Previously, we talked abt a class of probs, A  $\xrightarrow{?}$  B, A, B unknown, now a more general case



$\pi$  is also the number of points described by the gaussian.

Thinly spread  $\rightarrow \pi_c$  less  
Conc  $\rightarrow$  Huge  $\pi_c$



I can generate samples now, using, first a random no for the cluster number (proportional to  $\pi_c$ ) and then random number from the gaussian.

$\frac{\pi_c}{\sum \pi_c}$	$\frac{\pi_c}{\sum \pi_c}$	$\frac{\pi_c}{\sum \pi_c}$
1 to 8	9	10

- EM is guaranteed to converge, it won't stop, but it will reach very small values, so can get stuck in local-minima, but K-means generally doesn't get stuck in E-M can get stuck

// EM is like dynamic programming an idea  $\rightarrow$  Just a class of algos //

### • Cluster-validity:-

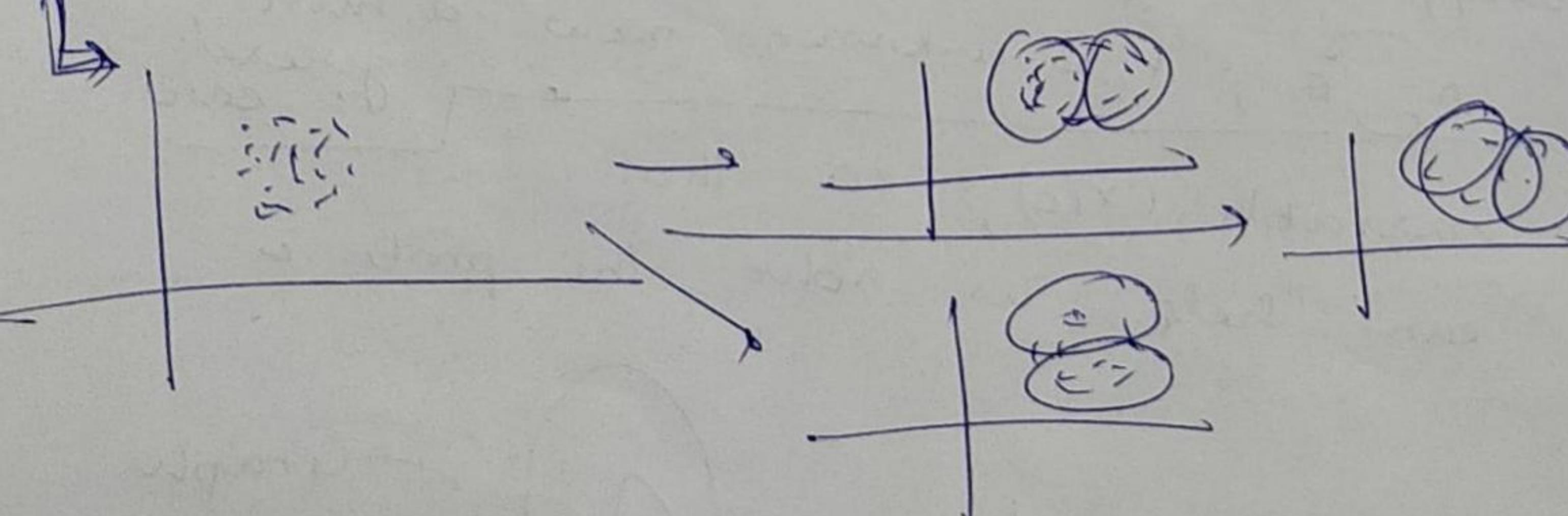
↳ we don't know, how many clusters to choose?

↳ And a lot more (means,  $\Sigma, \dots$ )

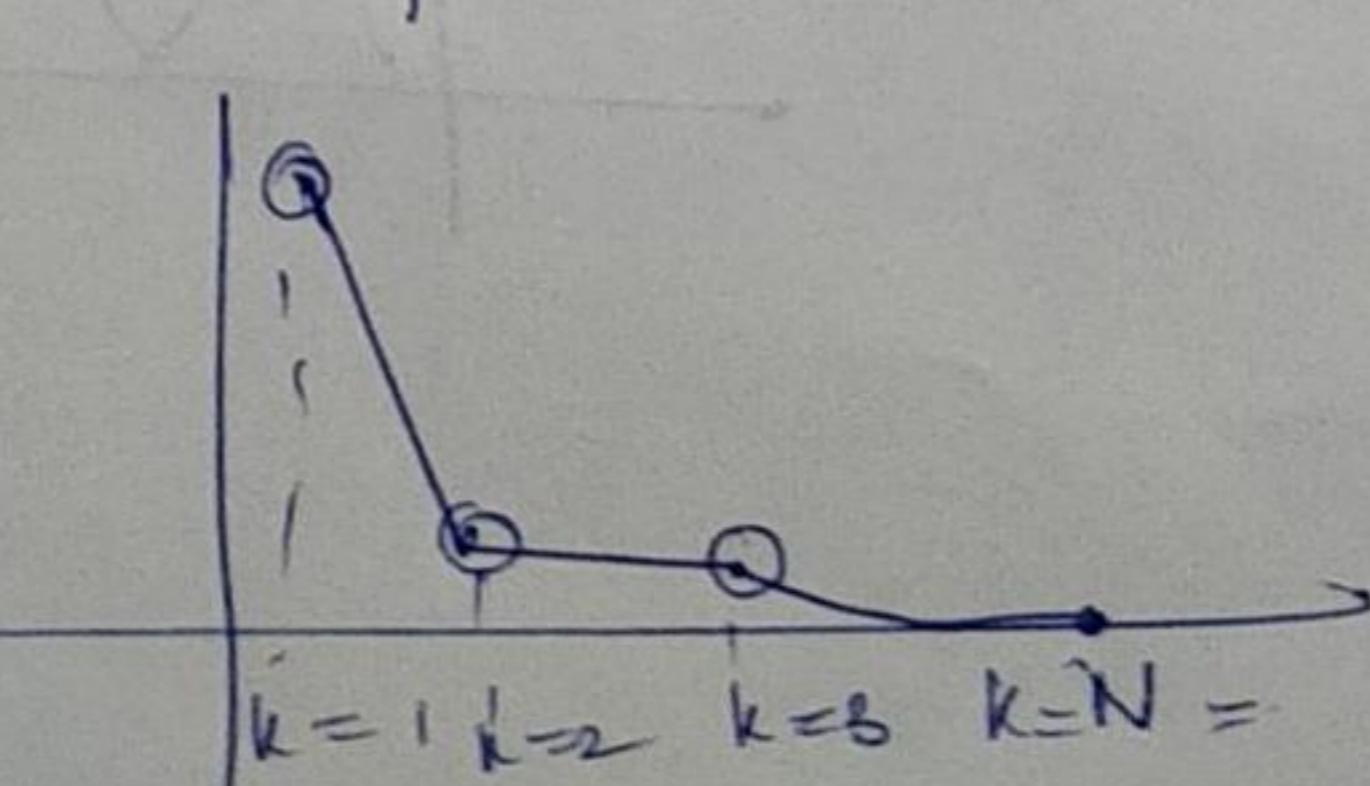
K-means  $\rightarrow$  suppose  $k=3$ , how to we know the right number of clusters.

- ① Stability
  - Ways of measuring validity.
- ② Initialize with different values, do all, lead to the same cluster, How much % give same result

Unstable along different initializations.



Compute averaged distance from cluster centres



② Mean-squared distance to cluster center //

↳ other similar metric too

↳ For GMM, you can use ric. related stuff

Some give good  
graphs, to find  
minima

