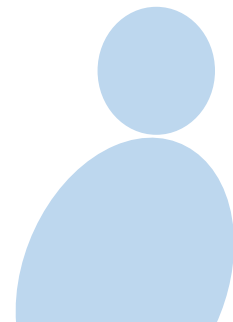


Elective Module:

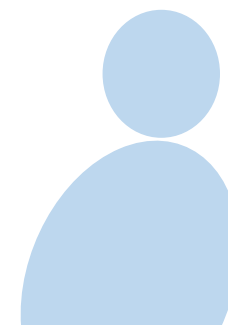
**Advanced ML
Techniques**



Chapter 9

EM Algorithm

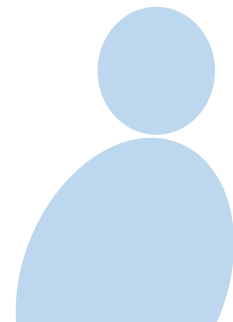
Expectation
Maximization



EM Algorithm

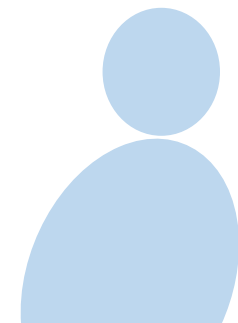
- EM stands for **Expectation-Maximization**.
- This can be used for probabilistic-clustering or soft-clustering

K-Means : Hard clustering



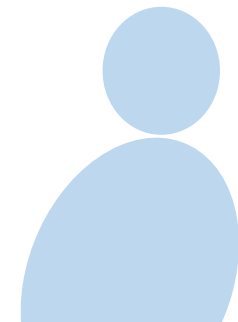
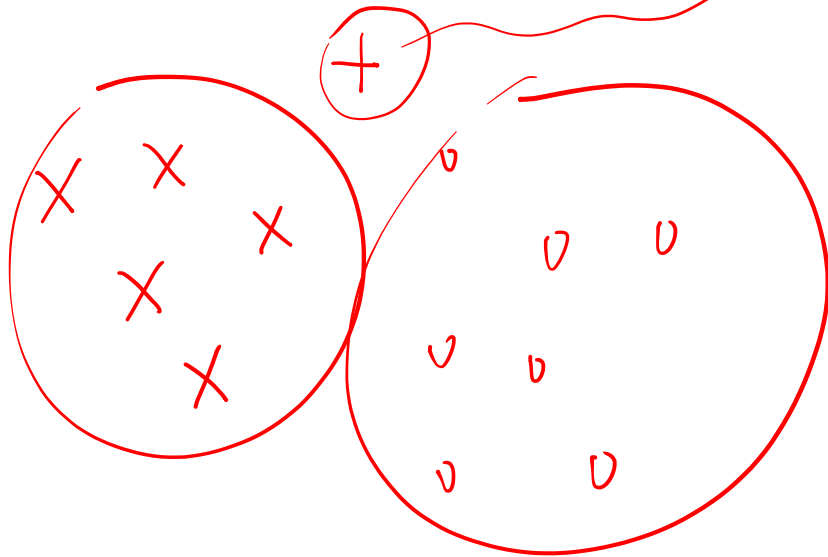
Probabilistic clustering

- What is **probabilistic-clustering**?
-



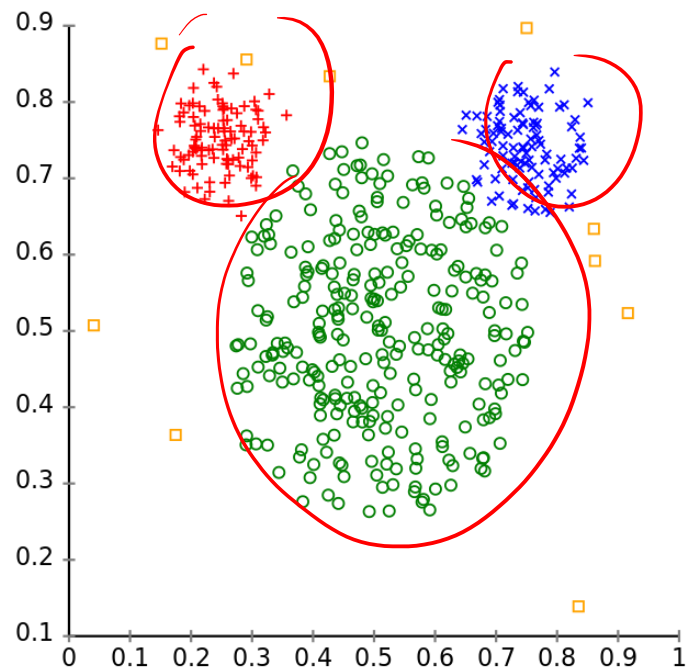
Probabilistic clustering

- Previously we assigned each point to a **unique cluster**.
- Now we calculate the **probability** that a data point belongs to a cluster!

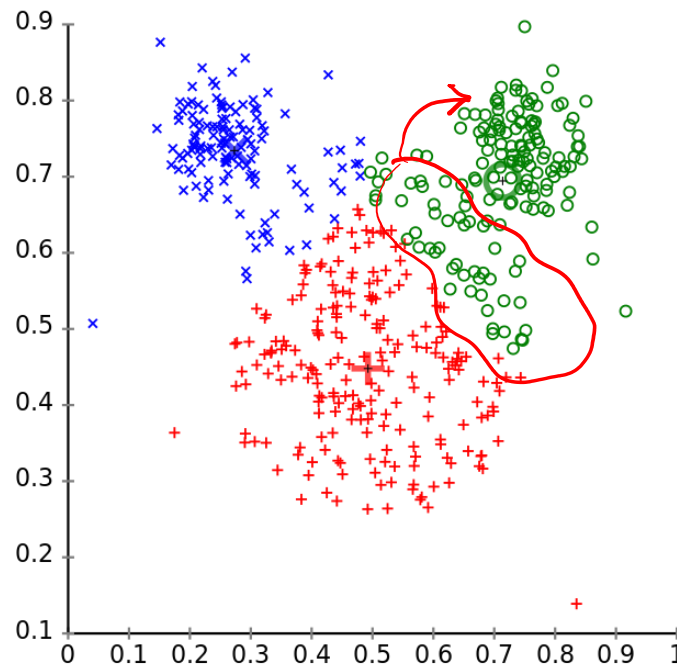


Different cluster analysis results on "mouse" data set:

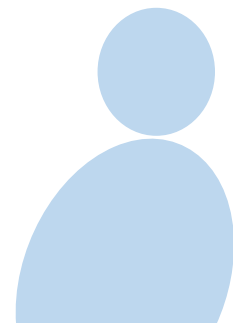
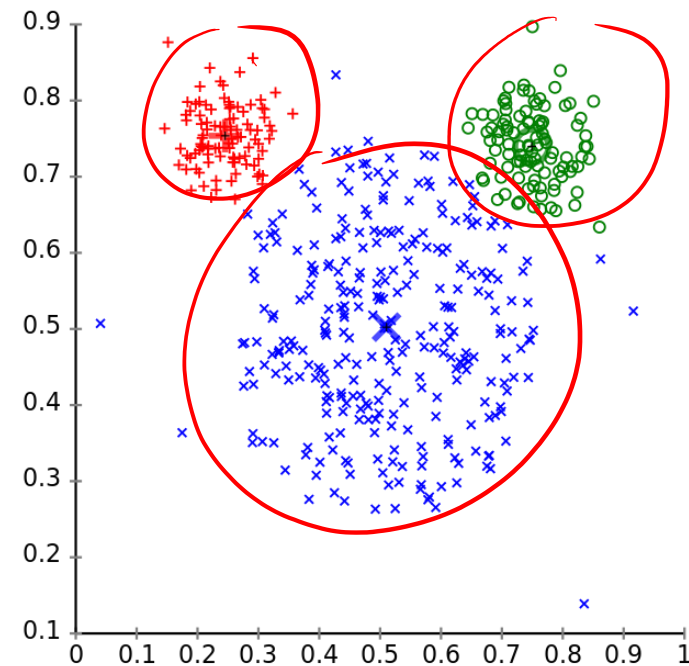
Original Data



k-Means Clustering



EM Clustering



EM Algorithm

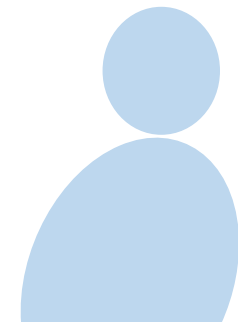
- Consider a **Gaussian cluster model**.
- With probability p_i generate a sample $\bar{\mathbf{x}}$ from **Gaussian** cluster i i.e.

$$\mathcal{N}(\bar{\mu}_i, \sigma^2 I) \sim \Sigma_i$$

Centroid cluster

covariance matrix

Prior probability of i^{th} cluster



EM Algorithm

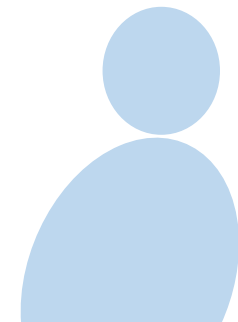
- Consider a **Gaussian cluster model**.
- With probability p_i generate a sample $\bar{\mathbf{x}}$ from **Gaussian cluster i** i.e.

$$N(\bar{\mu}_i, \sigma^2 \mathbf{I})$$

$$N(\bar{\mu}_1, \sigma^2 \mathbf{I}), N(\bar{\mu}_2, \sigma^2 \mathbf{I}), \dots, N(\bar{\mu}_K, \sigma^2 \mathbf{I})$$

K clusters.

K clusters.



EM Algorithm

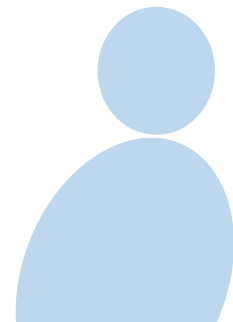
p_1, p_2, \dots, p_k : Prior probabilities of clusters.

- The PDF is given as

$$f(\bar{x}) = \sum_{i=1}^K p_i \left(\frac{1}{2\pi\sigma^2} \right)^{N/2} e^{-\frac{1}{2\sigma^2} \|\bar{x} - \bar{\mu}_i\|^2}$$

Gaussian mixture model.

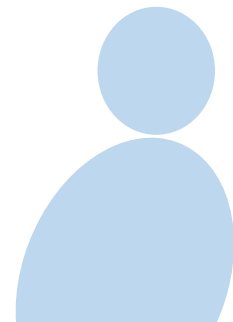
Probability density function



EM Algorithm

- The PDF is given as

$$f_X(\bar{\mathbf{X}}) = \sum_{i=1}^K p_i \times \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{\mathbf{X}} - \bar{\boldsymbol{\mu}}_i\|^2}$$

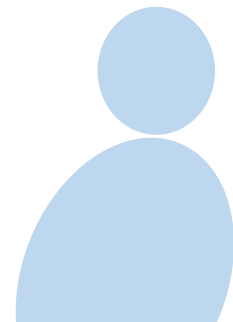


EM Algorithm

$$f_X(\bar{\mathbf{x}}) = \sum_{i=1}^K p_i \times \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{\mathbf{x}} - \bar{\boldsymbol{\mu}}_i\|^2}$$

- This is termed as a Gaussian mixture

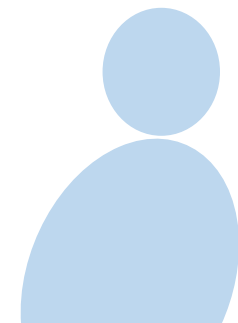
*K Gaussians.
One for each cluster -*



EM Algorithm

$$f_X(\bar{\mathbf{x}}) = \sum_{i=1}^K p_i \times \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{\mathbf{x}} - \bar{\boldsymbol{\mu}}_i\|^2}$$

- This is termed as a **Gaussian mixture**.



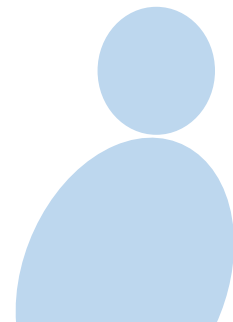
EM Algorithm

- Consider now M data points

$\bar{x}(1), \bar{x}(2), \dots, \bar{x}(M)$

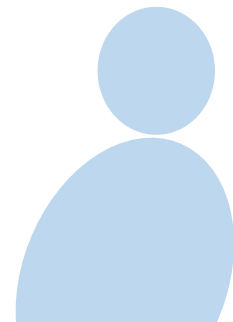
Data

cluster?



EM Algorithm

- Consider now M data points
 $\bar{\mathbf{x}}(1), \bar{\mathbf{x}}(2), \dots, \bar{\mathbf{x}}(M)$



EM Algorithm

- We desire to **estimate** $\bar{\mu}_i$
- As well as the **cluster assignments**
- Performing direct ML estimation is **mathematically intractable**.

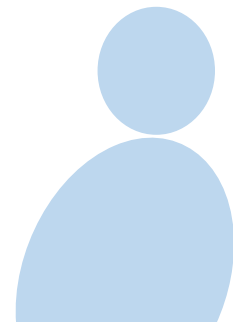
Centroids.

$\alpha_i(j)$.

NOT possible.

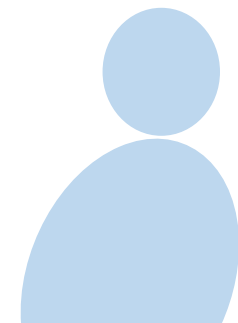
Maximum Likelihood.

cluster assignments.



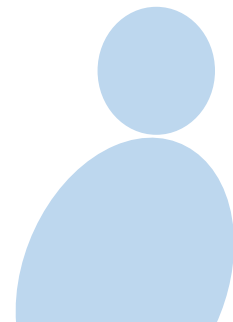
Cluster information

- However, if *cluster assignment* is known, problem is *simple*!



Cluster assignment $M=8$

- For example
- Cluster 1: $\bar{x}(1), \bar{x}(3), \bar{x}(5), \bar{x}(8)$: cluster 1.
- Cluster 2: $\bar{x}(2), \bar{x}(4), \bar{x}(6), \bar{x}(7)$: cluster 2



Cluster assignment

- For example

$$\hat{\mu}_1 = \frac{\bar{x}(1) + \bar{x}(3) + \bar{x}(5) + \bar{x}(8)}{4}$$

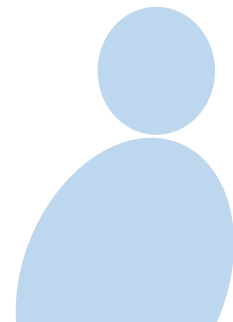
$$\hat{\mu}_2 = \frac{\bar{x}(2) + \bar{x}(4) + \bar{x}(6) + \bar{x}(7)}{4}$$

Average of points in cluster 1.

Average of points in cluster 2

NO Labels!

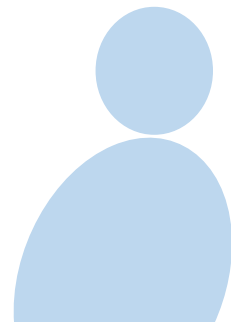
Unsupervised Learning



Cluster assignment

- For example

$$\hat{\mu}_1 = \frac{\bar{\mathbf{x}}(1) + \bar{\mathbf{x}}(3) + \bar{\mathbf{x}}(5) + \bar{\mathbf{x}}(8)}{4}$$
$$\hat{\mu}_2 = \frac{\bar{\mathbf{x}}(2) + \bar{\mathbf{x}}(4) + \bar{\mathbf{x}}(6) + \bar{\mathbf{x}}(7)}{4}$$

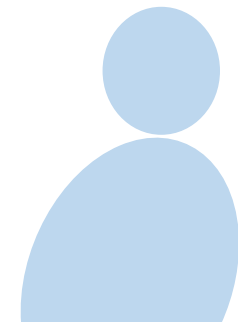


Cluster information

cluster assignment variable.

- We introduce the concept of **missing data** or **latent information**

$$\alpha_i(j) = \begin{cases} 1 & \text{if } \bar{x}(j) \in \mathcal{C}_i \\ 0 & \text{if } \bar{x}(j) \notin \mathcal{C}_i \end{cases}$$

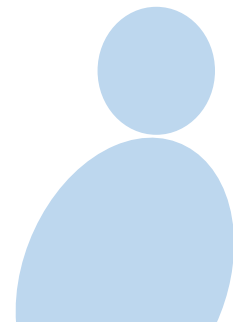


Cluster information

- We introduce the concept of missing data
or latent information

$$\alpha_i(j) = \begin{cases} 1 & \bar{\mathbf{x}}(j) \in \mathcal{C}_i \\ 0 & \bar{\mathbf{x}}(j) \notin \mathcal{C}_i \end{cases}$$

Missing!



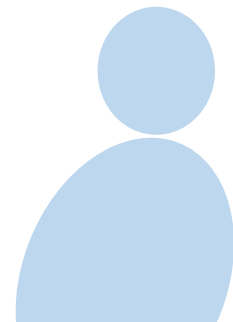
Complete data

$\bar{x}(1), \bar{x}(2), \dots, \bar{x}(M):$

$x_i(j) \quad \begin{matrix} i=1, 2, \dots, K \\ j=1, 2, \dots, M \end{matrix}$

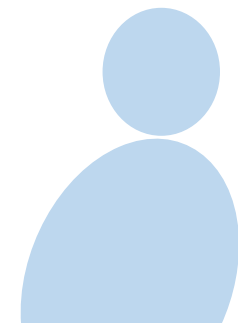
Complete Data.

Missing Data
Latent Information



Complete data

$$\underbrace{\bar{\mathbf{X}}(j), \alpha_i(j)}_{\text{Complete data}}$$



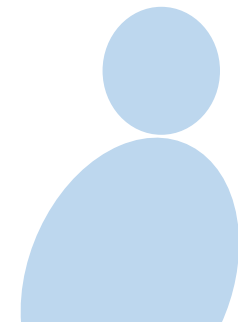
Likelihood

- The likelihood of the complete data

$$\prod_{j=1}^M \prod_{i=1}^K \left(p_i \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{x}(j) - \bar{\mu}_i\|^2} \right)^{\alpha_i(j)} = Q(\bar{\mu})$$

Missing data.

joint PDF of points.
Likelihood function

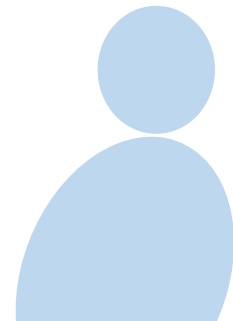


Likelihood

- The likelihood of the *complete data*

$$\prod_{j=1}^M \prod_{i=1}^K \left(p_i \times \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{\mathbf{x}}(j) - \bar{\boldsymbol{\mu}}_i\|^2} \right)^{\alpha_i(j)}$$

Products become sum
for Log.



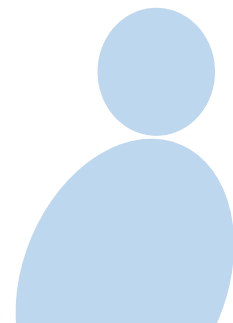
Log-Likelihood

Missing Data!

- The **log-likelihood** of the **complete data**

$$\sum_{j=1}^M \sum_{i=1}^K \alpha_i(j) \left(p_i - \frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\bar{x}(j) - \bar{\mu}_i\|^2 \right)$$

Log likelihood.



Log-Likelihood

- The **log-likelihood** of the **complete data**

$$\sum_{j=1}^M \sum_{i=1}^K \alpha_i(j) \left(\ln p_i - \frac{N}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \|\bar{\mathbf{x}}(j) - \bar{\boldsymbol{\mu}}_i\|^2 \right)$$

Log Likelihood.

Expected value
of $\alpha_i(j)$ in
iteration l .

$$E\{\alpha_i(j)\} = \alpha_i^{(l)}(j)$$

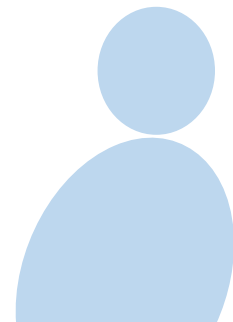
EM Algorithm

E - Step
M - Step.

- EM Algorithm **proceeds iteratively.**
- Consider the $l - 1$ th iteration with centroids

$\bar{\mu}_1^{(l-1)}, \bar{\mu}_2^{(l-1)}, \dots, \bar{\mu}_K^{(l-1)}$: Centroids in $(l-1)^{th}$ iteration

Expectation Maximization



EM Algorithm

- EM Algorithm ***proceeds iteratively***.

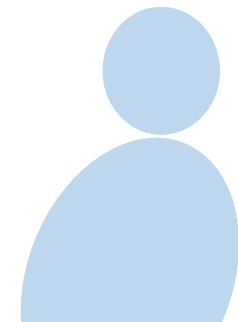
- Consider the $l - 1$ th iteration with centroids

$$\bar{\mu}_0^{(l-1)}, \bar{\mu}_1^{(l-1)}, \dots, \bar{\mu}_K^{(l-1)}$$

$K = \# \text{ clusters}$.

Centroids.
cluster means.

Gaussian mixtures.



EM Algorithm

E - Step
Expectation

- The expected value of the log-likelihood in iteration l is

$$\sum_{j=1}^M \sum_{i=1}^K \alpha_i^{(l)}(j) \left\{ \ln p_i - \frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|x(j) - \mu_i\|^2 \right\}$$

\uparrow
 $E\{\alpha_i(j)\}$

EM Algorithm

- The **expected value** of the log-likelihood in iteration l is

$$\underbrace{\sum_{j=1}^M \sum_{i=1}^K \alpha_i^{(l)}(j) \left(\ln p_i - \frac{N}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \|\bar{\mathbf{x}}(j) - \bar{\boldsymbol{\mu}}_i\|^2 \right)}_{Q(\bar{\boldsymbol{\mu}})}$$

$Q(\bar{\boldsymbol{\mu}})$

EM Algorithm

- How to calculate $\alpha_i^{(l)}(j)$?

Probability!
can take any
value in $[0,1]$

$$\alpha_i^{(l)}(j) = \Pr(C_i | \bar{\mathbf{x}}(j))$$

$$= \frac{\Pr(\bar{\mathbf{x}}(j) | C_i) \cdot P(C_i)}{\sum_k \Pr(\bar{\mathbf{x}}(j) | C_k) P(C_k)}$$

$$= \frac{P_i \left(\frac{1}{2\pi\sigma^2}\right)^{N/2} e^{-\frac{1}{2\sigma^2} \|\bar{\mathbf{x}}(j) - \bar{\mu}_i\|^2}}{\sum_k P_k \cdot \left(\frac{1}{2\pi\sigma^2}\right)^{N/2} e^{-\frac{1}{2\sigma^2} \|\bar{\mathbf{x}}(j) - \bar{\mu}_k\|^2}}$$

$$= 1 \text{ if } \bar{\mathbf{x}}(j) \in C_i$$

$$= \Pr(\bar{\mathbf{x}}(j) \in C_i) \times 1$$

$$+ 0 \times \Pr(\bar{\mathbf{x}}(j) \notin C_i)$$

$$= \Pr(\bar{\mathbf{x}}(j) \in C_i)$$

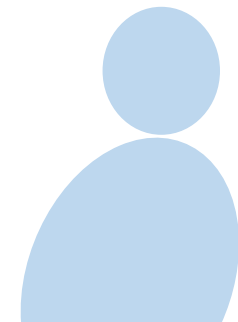
Bayes rule.

$\rightarrow E\{\alpha_i(j)\}$

EM Algorithm

- How to calculate $\alpha_i^{(l)}(j)$?

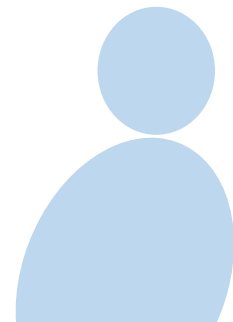
$$\begin{aligned}\alpha_i^{(l)}(j) &= \Pr(\mathcal{C}_i | \bar{\mathbf{x}}(j)) \\ &= \frac{p_i \times \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{\mathbf{x}}(j) - \bar{\boldsymbol{\mu}}_i^{(l-1)}\|^2}}{\sum_{k=1}^K p_k \times \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} \|\bar{\mathbf{x}}(j) - \bar{\boldsymbol{\mu}}_k^{(l-1)}\|^2}}\end{aligned}$$



M-Step

M = Maximization.

- The M-Step is the **maximization step** *Maximize Expected value of log likelihood.*
- To determine $\bar{\mu}_j$, **differentiate** with respect to $\bar{\mu}_j$ and set equal to zero

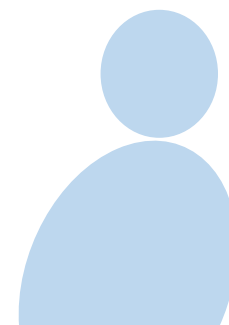


M-Step

$$\nabla_{\bar{\mu}_i} Q(\bar{\mu}) = 0$$

to maximize
take gradient and
set equal to zero.

$$\Rightarrow \bar{\mu}_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{x}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)}$$

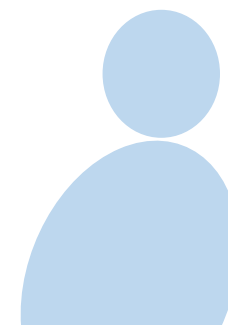


M-Step

$$\nabla_{\bar{\mu}_i} Q(\bar{\mu}) = 0$$

$$\Rightarrow \bar{\mu}_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{\mathbf{x}}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)}$$

Centroid of
ith cluster in
iteration l .

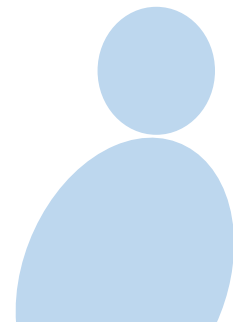


K-Means vs EM

- Compare the expressions in K-Means and EM

$$\text{EM: } \bar{\mu}_i^{(l)} = \frac{\sum_j \alpha_i^{(l)}(j) \bar{x}(j)}{\sum_j \alpha_i^{(l)}(j)} : \text{EM}.$$

$$\text{K - Means: } \bar{\mu}_i^{(l)} = \frac{\sum_j \alpha_i^{(l)}(j) \bar{x}(j)}{\sum_j \alpha_i^{(l)}(j)}$$



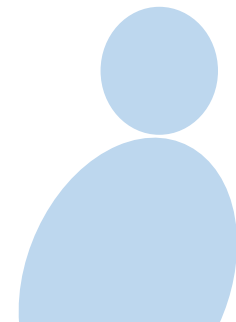
K-Means vs EM

- Compare the expressions in K-Means and EM

$$\underbrace{\bar{\mu}_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{\mathbf{x}}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)}}_{EM}$$

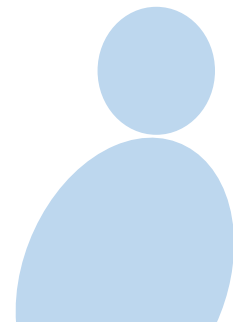
$$\underbrace{\bar{\mu}_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{\mathbf{x}}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)}}_{K-Means}$$

Same Expression!



K-Means vs EM

- Both are **identical!**
- What then is the **difference?**



K-Means vs EM

Weighted
Average

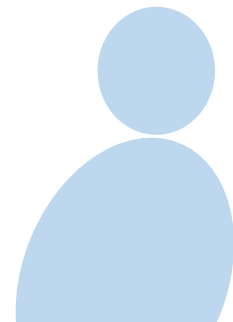
$$\underbrace{\bar{\mu}_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{\mathbf{x}}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)}}_{EM}$$

Soft clustering.
Probabilities.
 $\in [0, 1]$

Average

$$\underbrace{\bar{\mu}_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j) \bar{\mathbf{x}}(j)}{\sum_{j=1}^M \alpha_i^{(l)}(j)}}_{K-Means}$$

$\in \{0, 1\}$
Either 0 or 1
Hard clustering



K-Means vs EM

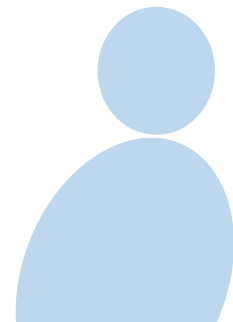
*Can only take
2 possible values.*

- If you observe carefully, in K-Means

$$\alpha_i^{(l)}(j) \in \{0,1\}$$

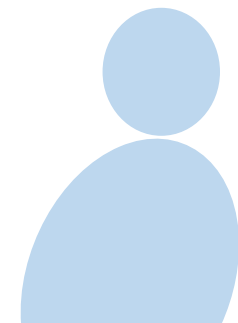
- In EM, $\alpha_i^{(l)}(j) \in [0,1]$

*any value in
interval $[0,1]$.
Ex: 0.95*



K-Means vs EM

- Therefore, in EM, $\alpha_i^{(l)}(j)$ denote probabilities.



Prior probabilities

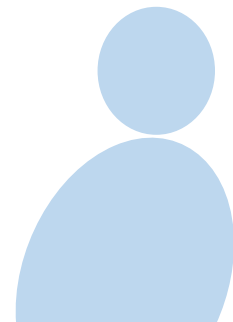
- Finally, the prior probabilities p_i can also be computed as follows

$$p_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j)}{M}$$

Sum of all
posterior probabilities
of points in cluster i

$\Pr(\bar{x}(j) \in \mathcal{C}_i)$

can also
be computed.

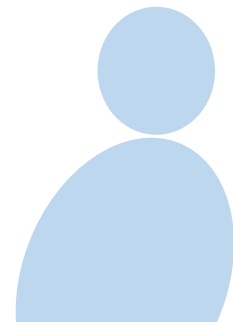


Prior probabilities

- Finally, the prior probabilities p_i can also be computed as follows

$$p_i^{(l)} = \frac{\sum_{j=1}^M \alpha_i^{(l)}(j)}{M}$$

Estimate of
prior probability
in iteration l .



Instructors may use this white area (14.5 cm / 25.4 cm) for the text.
Three options provided below for the font size.

Font: Avenir (Book), Size: 32, Colour: Dark Grey

Font: Avenir (Book), Size: 28, Colour: Dark Grey

Font: Avenir (Book), Size: 24, Colour: Dark Grey

Do not use the space below.

