# Outline
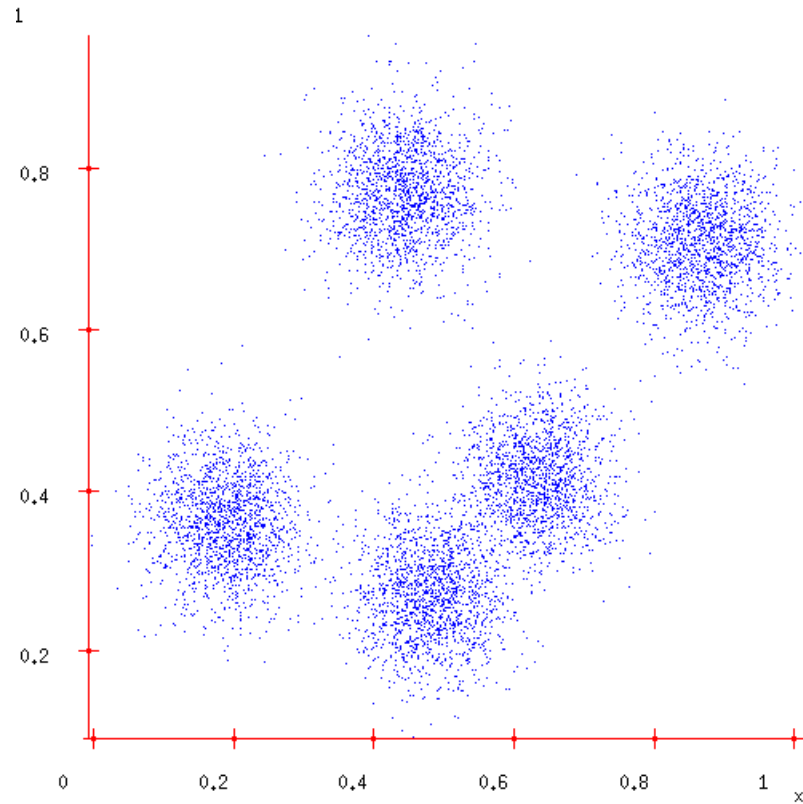
- K-Means

- Hierarchical Clustering

- <span style="color:red">Model-based Clustering (GMM and Expectation Maximization)</span>

- Evaluation of Clustering Algorithms
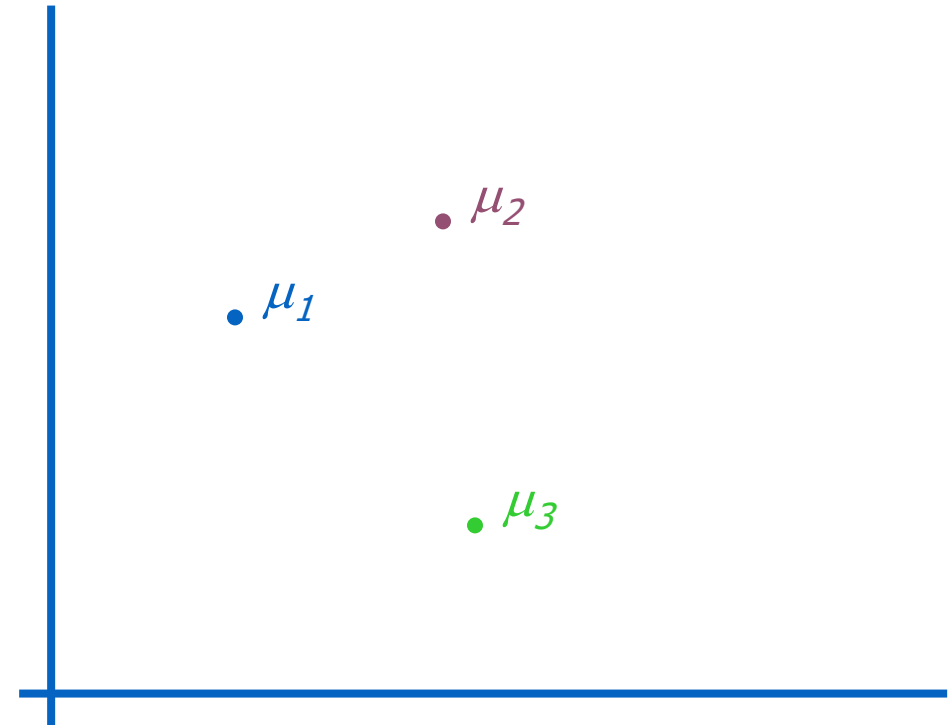
# Model-based Clustering: Gaussian Mixture Model

- Density estimation with multimodal/clumpy data

# Gaussian Mixture Model (GMM)

- The GMM assumption
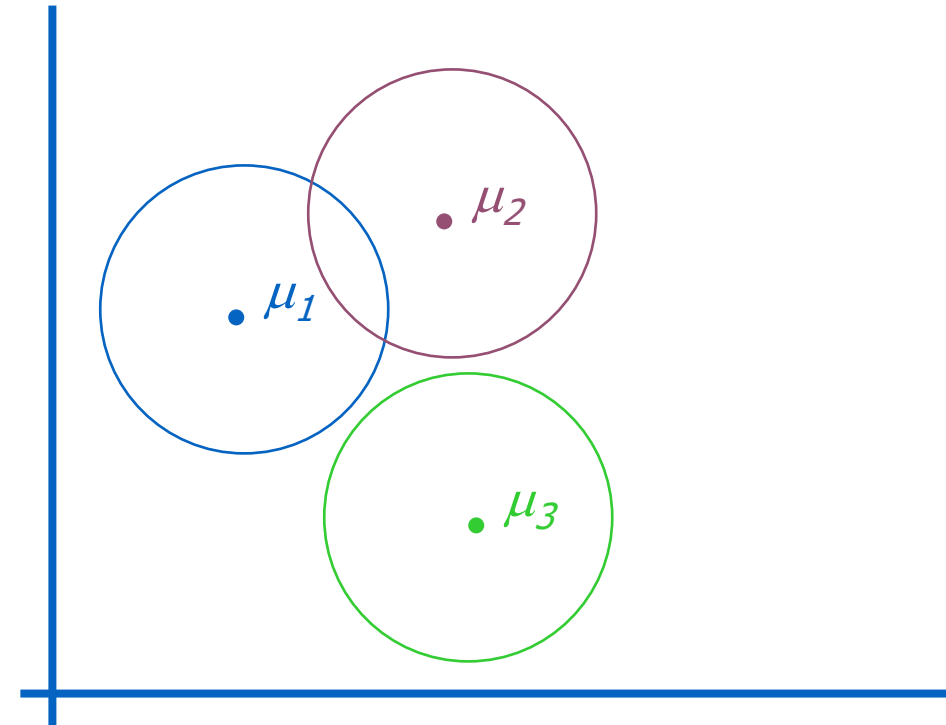- There are k components. The i[th] component is called $\omega_i$
- Component $\omega_i$ has an associated mean vector $\mu_i$

# Gaussian Mixture Model (GMM)

- The GMM assumption
- There are k components. The i[th] component is called $\omega_i$
- Component $\omega_i$ has an associated mean vector $\mu_i$
- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\sigma^2 I$
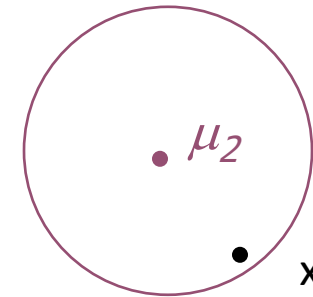
# Gaussian Mixture Model (GMM)

- The GMM assumption
- There are k components. The i$^{th}$ component is called $\omega_i$
- Component $\omega_i$ has an associated mean vector $\mu_i$
- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\sigma^2 I$
- Assume that each datapoint is generated according to the following recipe:

  - Pick a component at random. Choose component i with probability $P(\omega_i)$.

  - Datapoint ~ N($\mu_i$, $\Sigma_i$ )



$\mu_2$

x

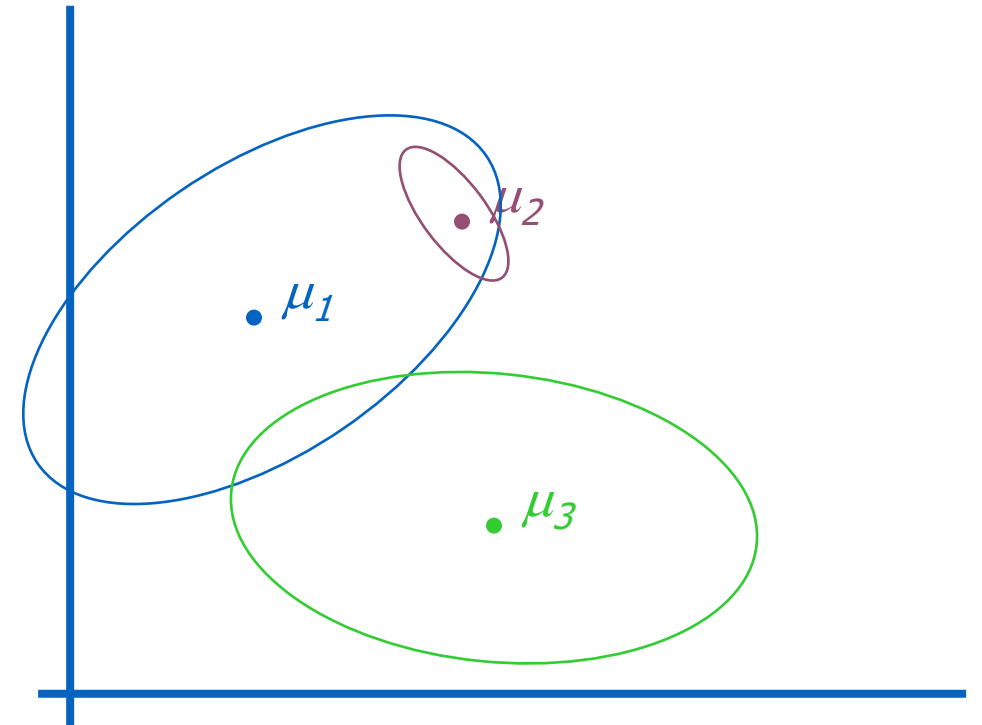आई आई टी हैदराबाद
IIT Hyderabad

# Gaussian Mixture Model (GMM)

- The GMM assumption
- There are k components. The $i^{th}$ component is called $\omega_i$
- Component $\omega_i$ has an associated mean vector $\mu_i$
- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\sigma^2 I$
- Assume that each datapoint is generated according to the following recipe:

  - Pick a component at random. Choose component i with probability $P(\omega_i)$.

  - Datapoint ~ N($\mu_i$, $\Sigma_i$ )

# Gaussian Mixture Model (GMM)

- Given the means and $\sigma^2$, we can compute P(data | $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 .. \boldsymbol{\mu}_k, \sigma^2$). How do we find the $\boldsymbol{\mu}_i$ s and $\sigma^2$ which give max likelihood?

- The normal max likelihood trick:

    Set $\dfrac{d}{d\mu_i}$ log Prob (....) = 0

    and solve for $\mu_i$'s.

- Use gradient descent
    - Slow but doable

- Use a much faster and popular method: EM

आई आई टी हैदराबाद
IIT Hyderabad

# Expectation Maximization (EM)

- We'll get back to unsupervised learning/clustering/GMM soon.
- The EM algorithm was explained and given its name in a classic 1977 paper by Arthur Dempster, Nan Laird, and Donald Rubin.
- They pointed out that the method had been "proposed many times in special circumstances" by earlier authors.
- EM is typically used to compute maximum likelihood estimates given incomplete samples.
  - An excellent way of doing our unsupervised learning problem, as we'll see
  - Many, many other uses, including inference of Hidden Markov Models
- The EM algorithm estimates the parameters of a model iteratively. Starting from some initial guess, each iteration consists of
  - an E step (Expectation step)
  - an M step (Maximization step)

आई आई टी हैदराबाद
IIT Hyderabad

## EM: Trivial Example

Let events be "grades in a class"

| | | |
|---|---|---|
| $w_1$ = Gets an A | | P(A) = ½ |
| $w_2$ = Gets a B | | P(B) = $\mu$ |
| $w_3$ = Gets a C | | P(C) = $2\mu$ |
| $w_4$ = Gets a D | | P(D) = ½-$3\mu$ |

(Note $0 \leq \mu \leq 1/6$)

Assume we want to estimate $\mu$ from data. In a given class, there were

    a  A's
    b  B's
    c  C's
    d  D's

What's the maximum likelihood estimate of $\mu$ given a,b,c,d ?

# EM: Trivial Example

$P(A) = \frac{1}{2}$    $P(B) = \mu$    $P(C) = 2\mu$    $P(D) = \frac{1}{2}-3\mu$

$P(a,b,c,d \mid \mu) = (\frac{1}{2})^a(\mu)^b(2\mu)^c(\frac{1}{2}-3\mu)^d$

$\log P(a,b,c,d \mid \mu) = a\log \frac{1}{2} + b\log \mu + c\log 2\mu + d\log (\frac{1}{2}-3\mu)$

$$\text{FOR MAX LIKE } \mu, \text{ SET } \frac{\partial \mathrm{LogP}}{\partial \mu} = 0$$

$$\frac{\partial \mathrm{LogP}}{\partial \mu} = \frac{b}{\mu} + \frac{2c}{2\mu} - \frac{3d}{1/2 - 3\mu} = 0$$

$$\text{Gives max like } \mu = \frac{b+c}{6(b+c+d)}$$

So if class got

| A | B | C | D |
|---|---|---|---|
| 14 | 6 | 9 | 10 |

$$\text{Max likelihood estimate : } \mu = \frac{1}{10}$$

आई आई टी हैदराबाद
IIT Hyderabad

## EM: Same Example with Hidden Info

Someone tells us that

Number of High grades (A's + B's) = $h$

Number of C's $= c$

Number of D's $= d$

What is the max likelihood estimate of μ now?

Slide Courtesy: Andrew Moore, CMU

## EM: Same Example with Hidden Info

Someone tells us that

Number of High grades (A's + B's) = $h$

Number of C's $= c$

Number of D's $= d$

What is the max likelihood estimate of μ now?
We can answer this circularly as below

**EXPECTATION**

If we know the value of μ we could compute the expected value of $a$ and $b$

Since the ratio a:b should be the same as the ratio ½ : μ

$$a = \frac{\frac{1}{2}}{\frac{1}{2}+\mu} h \qquad b = \frac{\mu}{\frac{1}{2}+\mu} h$$

**MAXIMIZATION**

If we know the expected values of $a$ and $b$ we could compute the maximum likelihood value of μ

$$\mu = \frac{b+c}{6(b+c+d)}$$

IIT Hyderabad

# EM: Solution for Trivial Example

We begin with a guess for μ

We iterate between EXPECTATION and MAXIMIZATION to improve our estimates of μ and *a* and *b*.

Define μ(t) the estimate of μ on the t$^{th}$ iteration

b(t) the estimate of *b* on t$^{th}$ iteration

**E-step**

$$\mu(0) = \text{initial guess}$$

$$b(t) = \frac{\mu(t)h}{\frac{1}{2} + \mu(t)} = \text{E}\big[b \,|\, \mu(t)\big]$$

**M-step**

$$\mu(t+1) = \frac{b(t) + c}{6(b(t) + c + d)}$$

$$= \text{max like est of } \mu \text{ given } b(t)$$

Continue iterating until converged.
Good news: Converging to local optimum is assured.
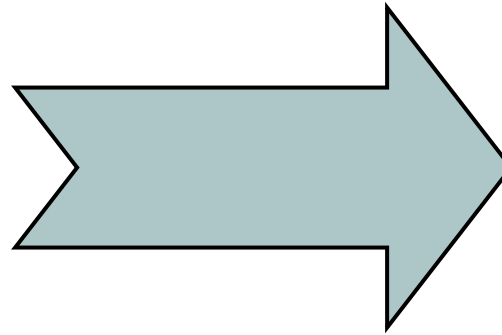Bad news: "local" optimum.

आई आई टी हैदराबाद
IIT Hyderabad

# EM: Convergence

- Convergence proof based on fact that Prob(data | μ) must increase or remain same between each iteration [NOT OBVIOUS]

- But it can never exceed 1 [OBVIOUS]

So it must therefore converge [OBVIOUS]

---

In our example, suppose we had

h = 20
c = 10
d = 10
μ(0) = 0

Convergence is generally <u>linear</u>: error decreases by a constant factor each time step.

| t | μ(t) | b(t) |
|---|------|------|
| 0 | 0 | 0 |
| 1 | 0.0833 | 2.857 |
| 2 | 0.0937 | 3.158 |
| 3 | 0.0947 | 3.185 |
| 4 | 0.0948 | 3.187 |
| 5 | 0.0948 | 3.187 |
| 6 | 0.0948 | 3.187 |

आई आई टी हैदराबाद
IIT Hyderabad

# Back to GMM

Given a training data set: X={x(1),x(2),...,x(n)}
Z={z(1),z(2),...,z(n)}
z(i) is the calss/group label of sample x(i).
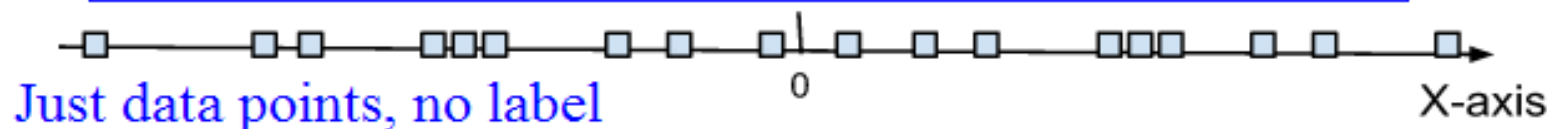As we are in Clustering setting,
X is Given and Z is unknown

Now, we model the data by specifying a joint
distribution p(x(i), z(i))=p(x(i)|z(i))p(z(i))

$$z(i) \sim \text{Multinomial}(\phi)$$
$$\phi_j \geq 0, \sum_{i=1}^{k} \phi_j = 1$$
$$k = \# \text{ of } z(i)\text{'s values}$$
$$\phi_j = p(z(i) = j)$$
$$x(i)|z(i) = j \sim \mathcal{N}(\mu_j, \Sigma_j)$$

$\Rightarrow$ each $x(i)$ was generated by randomly choosing $z(i)$ from $\{1, \ldots, k\}$, and then $x(i)$ was drawn from one of $k$ Gaussians.

The parameters of our model are thus $\phi$, $\mu$ and $\Sigma$.

Just data points, no label

0

X-axis

IIT Hyderabad

# EM for GMM

X={x(1),x(2),...,x(n)} Given — Incomplete Data

Z={z(1),z(2),...,z(n)} unknown

The parameters of our model $\phi, \mu, \Sigma$ — unknown

What is the value of z(i)?

We can answer this question circularly:

**EXPECTATION** — If we know the values of $\phi, \mu, \Sigma$ we could compute the expected values of Z

**MAXIMIZATION** — If we know the expected values of Z we could compute the maximum likelihood value of $\phi, \mu, \Sigma$

We begin with a guess for $\phi, \mu, \Sigma$, and then iterate between EXPECTATION and MAXIMALIZATION to improve our estimates of $\phi, \mu, \Sigma$ and Z

**Continue iterating until converged.**

# EM for GMM

$$\ell(\phi, \mu, \Sigma) = \sum_{i=1}^{m} \log p(x^{(i)} | z^{(i)}; \mu, \Sigma) + \log p(z^{(i)}; \phi)$$

Maximizing this with respect to $\phi$, $\mu$ and $\Sigma$ gives the parameters:

$$
\begin{aligned}
\phi_j &= \frac{1}{m} \sum_{i=1}^{m} 1\{z^{(i)} = j\}, \\
\mu_j &= \frac{\sum_{i=1}^{m} 1\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{m} 1\{z^{(i)} = j\}}, \\
\Sigma_j &= \frac{\sum_{i=1}^{m} 1\{z^{(i)} = j\}(x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^{m} 1\{z^{(i)} = j\}}.
\end{aligned}
$$

Slide Courtesy: Andrew Moore, CMU

आई आई टी हैदराबाद
IIT Hyderabad

# EM for GMM

Repeat until convergence: {

(E-step) For each $i, j$, set

$$w_j^{(i)} := p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$$

(M-step) Update the parameters:

$$\phi_j := \frac{1}{m} \sum_{i=1}^{m} w_j^{(i)},$$

$$\mu_j := \frac{\sum_{i=1}^{m} w_j^{(i)} x^{(i)}}{\sum_{i=1}^{m} w_j^{(i)}},$$

$$\Sigma_j := \frac{\sum_{i=1}^{m} w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^{m} w_j^{(i)}}$$

}

# GMM vs k-Means

Given a training data set: $X=\{x(1),x(2),\dots,x(n)\}$
$Z=\{z(1),z(2),\dots,z(n)\}$
$z(i)$ is the calss/group label of sample $x(i)$.
As we are in Clustering setting,
$X$ is Given and Z is unknown

**Model of EM**

EM model the data by specifying a joint distribution $p(x(i), z(i))=p(x(i)|z(i))p(z(i))$

$$z^{(i)} \sim \text{Multinomial}(\phi)$$
$$\phi_j \geq 0, \sum_{i=1}^{k} \phi_j = 1$$
$$k = \#\ of\ z(i)\text{'s values}$$
$$\phi_j = p(z^{(i)} = j)$$
$$x^{(i)}\big|z^{(i)} = j \sim \mathcal{N}(\mu_j, \Sigma_j)$$

$\Rightarrow$ each $x(i)$ was generated by randomly choosing $z^{(i)}$ from $\{1,\dots,k\}$, and then $x(i)$ was drawn from one of $k$ Gaussians .

K-means is a simplified EM, it assumes that

$$\phi_j = \phi_i = 1/k, \text{ and } \Sigma_j = \Sigma_i \text{ for } i, j=1,2,\dots k$$
k is given by user

$\mu_1, \mu_2 \ \text{------} \ \mu_k$

$\Rightarrow$ are the only unknown parameters of the model (the means of clusters)

Slide Courtesy: Andrew Moore, CMU

23-Nov-20

IIT Hyderabad

19

# GMM: Example

Start: 0<sup>th</sup> iteration

# GMM: Example

After 1ˢᵗ iteration

# GMM: Example

After 2nd iteration



p=0.37
p=0.306
p=0.320

After 3<sup>rd</sup> iteration

GMM: Example

# GMM: Example

After 4th iteration



$p=0.331$

$p=0.288$

# GMM: Example

After 5th iteration



p=0.322

p=0.285

# GMM: Example

After 6<sup>th</sup> iteration

# GMM: Example

After 20th iteration



p=0.234

p=0.334

# More on EM Algorithm

- What are the EM algorithm initialization methods?
  - Random guess.
  - Initialized by k-means. After a few iterations of k-means, using the parameters to initialize EM

- What are the main advantages of parametric methods?
  - You can easily change the model to adapt to different distribution of data sets.
  - Knowledge representation is very compact. Once the model is selected, the model is represented by a specific number of parameters.
  - The number of parameters does not increase with the increasing of training data .