

# 10-701 Introduction to Machine Learning

## The EM Algorithm

Thanks!

---

Spring 2019

Ameet Talwalkar

(slide credit: Virginia Smith)

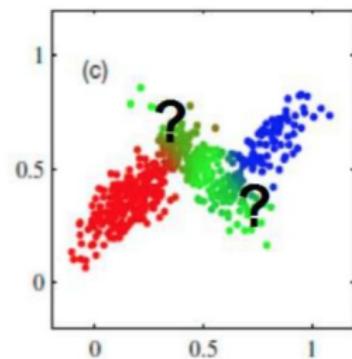
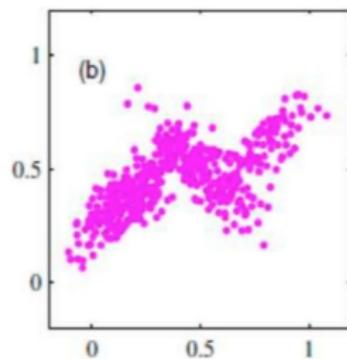
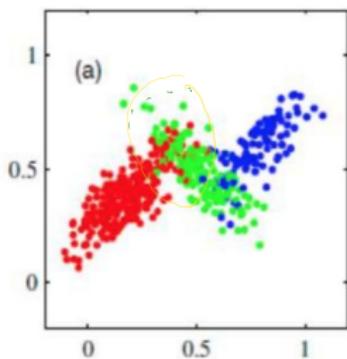
## Gaussian mixture models

---

## Potential issue with $k$ -means ...

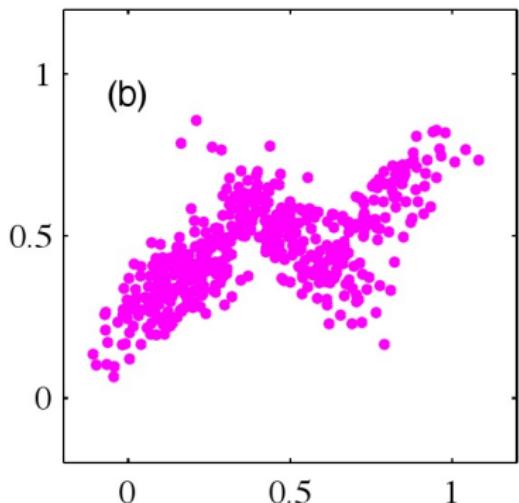
Data points are assigned *deterministically* to one (and only one) cluster

In reality, clusters may overlap, and it may be better to identify the *probability* that a point belongs to each cluster



# Probabilistic interpretation of clustering?

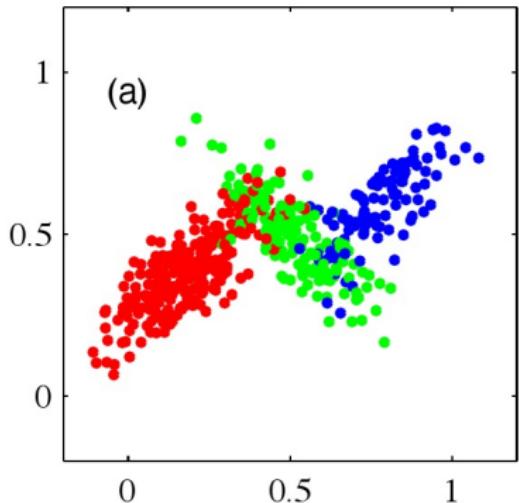
How can we model  $p(\mathbf{x})$  to reflect our intuition that points stay close to their cluster centers?



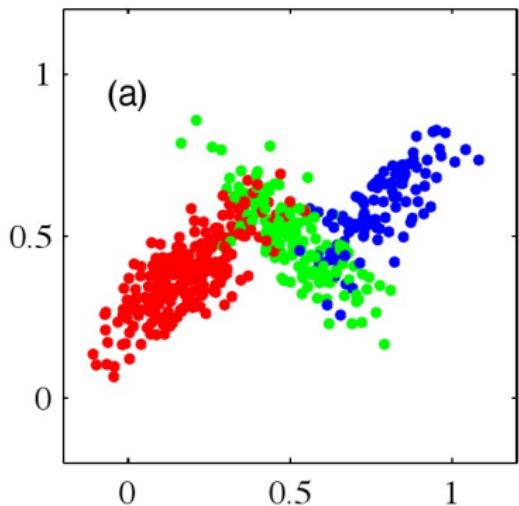
- Points seem to form 3 clusters
- We cannot model  $p(\mathbf{x})$  with simple and known distributions
- E.g., the data is not a Gaussian b/c we have 3 distinct concentrated regions

## Gaussian mixture models: intuition

- **Key idea:** Model each region with a distinct distribution

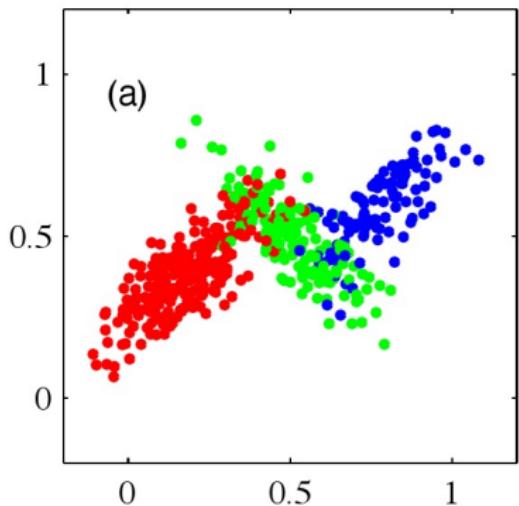


# Gaussian mixture models: intuition



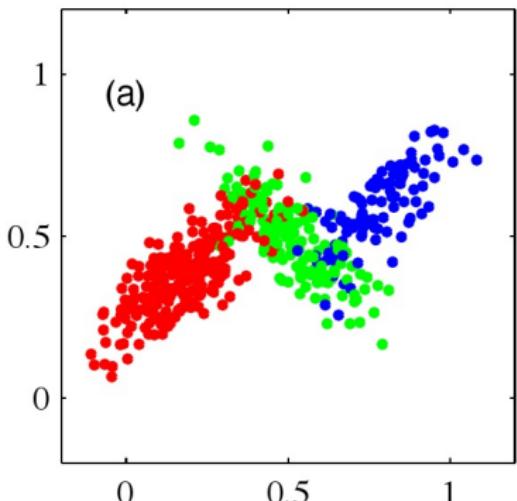
- **Key idea:** Model each region with a distinct distribution
- Can use Gaussians — Gaussian mixture models (GMMs)

# Gaussian mixture models: intuition



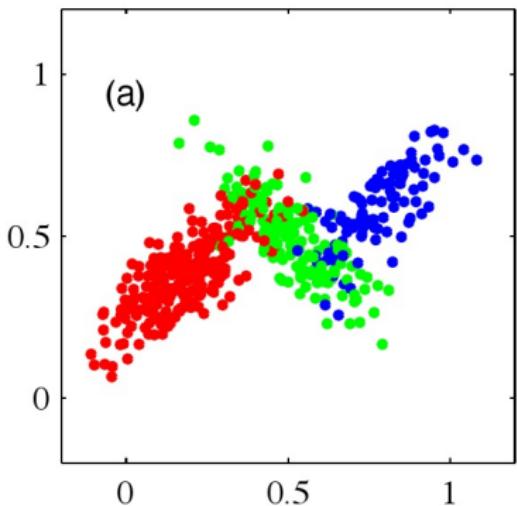
- **Key idea:** Model each region with a distinct distribution
- Can use Gaussians — Gaussian mixture models (GMMs)

# Gaussian mixture models: intuition



- **Key idea:** Model each region with a distinct distribution
- Can use Gaussians — Gaussian mixture models (GMMs)
- \*However\*, we don't know *cluster assignments* (label), *parameters* of Gaussians, or *mixture components*!

# Gaussian mixture models: intuition

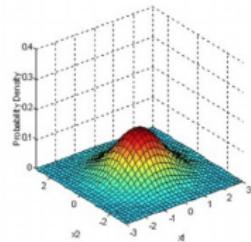


- **Key idea:** Model each region with a distinct distribution
- Can use Gaussians — Gaussian mixture models (GMMs)
- \*However\*, we don't know *cluster assignments* (label), *parameters* of Gaussians, or *mixture components*!
- Must learn from unlabeled data  
 $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$

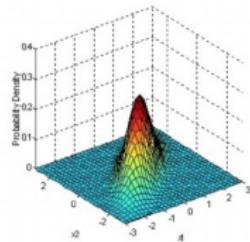
# Recall: Gaussian (normal) distributions

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



$$\boldsymbol{\mu} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix}$$



## Gaussian mixture models: formal definition

GMM has the following density function for  $\mathbf{x}$

$$p(\mathbf{x}) = \sum_{k=1}^K \omega_k N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- $K$ : number of Gaussians — they are called mixture components

## Gaussian mixture models: formal definition

GMM has the following density function for  $\mathbf{x}$

$$p(\mathbf{x}) = \sum_{k=1}^K \omega_k N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- $K$ : number of Gaussians — they are called mixture components
- $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}_k$ : mean and covariance matrix of  $k$ -th component

## Gaussian mixture models: formal definition

GMM has the following density function for  $\mathbf{x}$

$$p(\mathbf{x}) = \sum_{k=1}^K \omega_k N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- $K$ : number of Gaussians — they are called mixture components
- $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}_k$ : mean and covariance matrix of  $k$ -th component
- $\omega_k$ : mixture weights (or priors) represent how much each component contributes to final distribution. They satisfy 2 properties:

$$\forall k, \boxed{\omega_k > 0} \quad \text{and} \quad \boxed{\sum_k \omega_k = 1}$$

These properties ensure  $p(\mathbf{x})$  is in fact a probability density function

## GMM as the marginal distribution of a joint distribution

Consider the following joint distribution

$$\checkmark \quad p(\mathbf{x}, z) = p(z)p(\mathbf{x}|z)$$

where  $z$  is a discrete random variable taking values between 1 and  $K$ .

## GMM as the marginal distribution of a joint distribution

Consider the following joint distribution

$$p(\mathbf{x}, z) = p(z) \underbrace{p(\mathbf{x}|z)}$$

where  $z$  is a discrete random variable taking values between  $1$  and  $K$ .

Denote

$$\checkmark \omega_k = p(z = k)$$

Now, assume the conditional distributions are Gaussian distributions

$$\checkmark p(\mathbf{x}|z = k) = \underbrace{N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

## GMM as the marginal distribution of a joint distribution

Consider the following joint distribution

$$p(\mathbf{x}, z) = p(z)p(\mathbf{x}|z) \Rightarrow p(\mathbf{x}) = \sum_{z=1}^K p(\mathbf{x}, z)$$

where  $z$  is a discrete random variable taking values between 1 and  $K$ .

Denote

$$\omega_k = p(z = k)$$

Now, assume the conditional distributions are Gaussian distributions

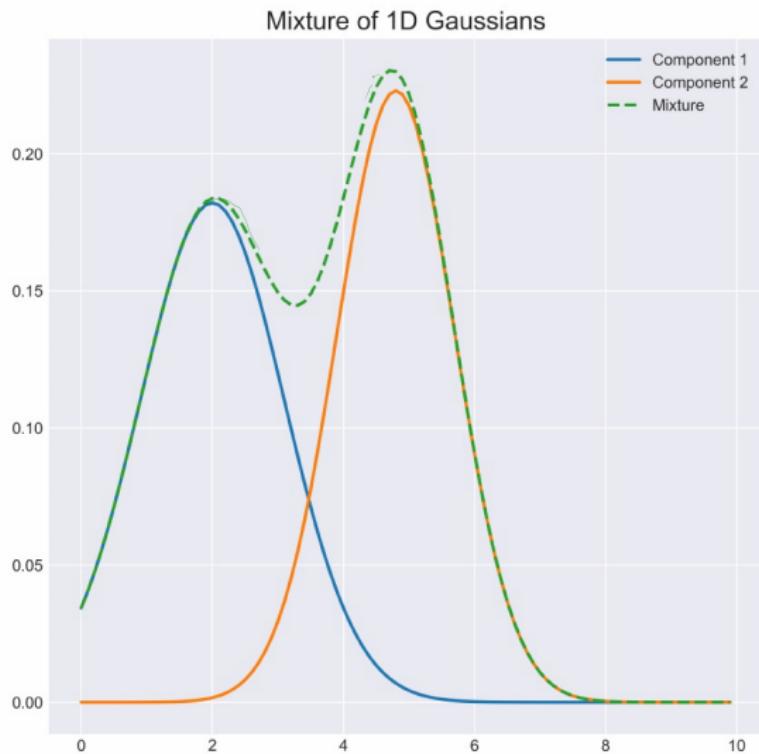
$$p(\mathbf{x}|z = k) = N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Then, the marginal distribution of  $\mathbf{x}$  is

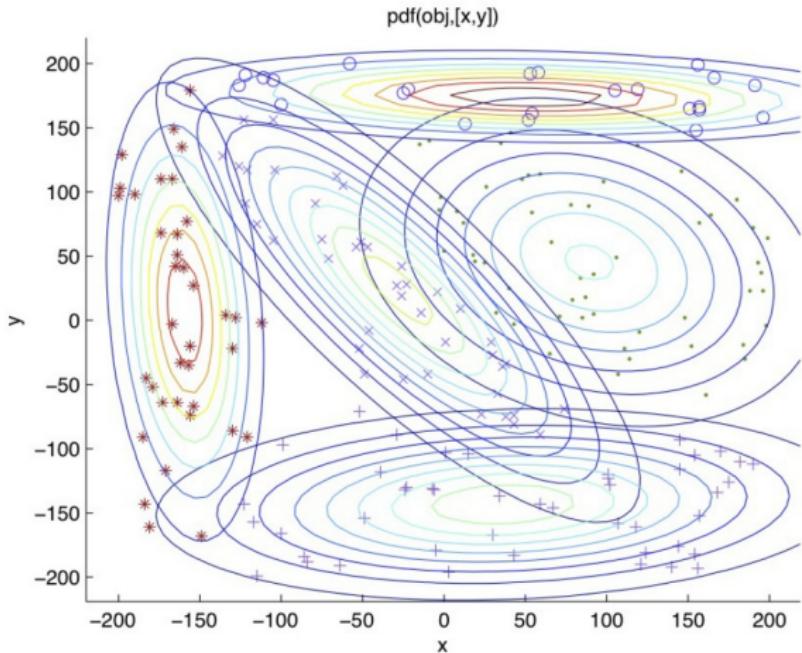
$$\cancel{p(\mathbf{x})} = \sum_{k=1}^K \omega_k N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Namely, the Gaussian mixture model

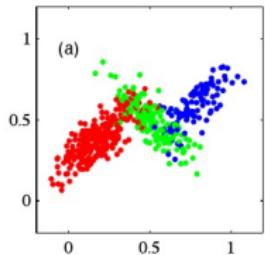
# Gaussian mixtures in 1D



# Gaussian mixture model for clustering



## GMMs: example



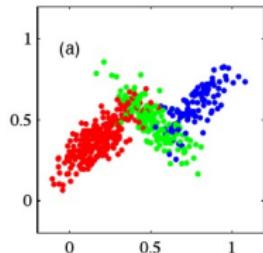
The conditional distribution between  $\mathbf{x}$  and  $z$  (representing color) are

$$\checkmark p(\mathbf{x}|z = \text{red}) = N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \checkmark$$

$$p(\mathbf{x}|z = \text{blue}) = N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \checkmark$$

$$p(\mathbf{x}|z = \text{green}) = N(\mathbf{x}|\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$$

## GMMs: example

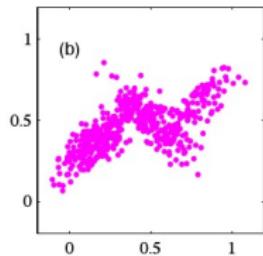


The conditional distribution between  $\mathbf{x}$  and  $z$  (representing color) are

$$p(\mathbf{x}|z = \text{red}) = N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$$

$$p(\mathbf{x}|z = \text{blue}) = N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

$$p(\mathbf{x}|z = \text{green}) = N(\mathbf{x}|\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$$



The marginal distribution is thus

✓  $p(\mathbf{x}) = \underbrace{p(\text{red})N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)}_{+ p(\text{blue})N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)} + \underbrace{p(\text{green})N(\mathbf{x}|\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)}$

## Parameter estimation for Gaussian mixture models

The parameters in GMMs are:

# Parameter estimation for Gaussian mixture models

The parameters in GMMs are:

$$\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$$

Let's first consider the simple/unrealistic case where we have labels  $z$

Define  $\boxed{\mathcal{D}' = \{\mathbf{x}_n, z_n\}_{n=1}^N}$ ,  $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$

- $\mathcal{D}'$  is the **complete** data
- $\mathcal{D}$  the **incomplete** data

How can we learn our parameters?

# Parameter estimation for Gaussian mixture models

The parameters in GMMs are:

$$\checkmark \theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$$

Let's first consider the simple/unrealistic case where we have labels  $z$

$$\text{Define } \mathcal{D}' = \{\mathbf{x}_n, z_n\}_{n=1}^N, \mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$$

- $\mathcal{D}'$  is the **complete** data
- $\mathcal{D}$  the **incomplete** data

How can we learn our parameters?

Given  $\mathcal{D}'$ , the maximum likelihood estimation of the  $\theta$  is given by

$$\theta = \arg \max \log \mathcal{D}' = \sum_n \log p(\mathbf{x}_n, z_n)$$

$\log p(\mathcal{D}'|\theta) = \log \prod_{n=1}^N p(\mathbf{x}_n, z_n; \theta)$

## Parameter estimation for GMMs: complete data

The complete likelihood is decomposable

$$\sum_n \log p(\mathbf{x}_n, z_n) = \underbrace{\sum_n \log p(z_n)p(\mathbf{x}_n|z_n)}_{\text{decomposable}} = \sum_k \left[ \sum_{n:z_n=k} \log p(z_n)p(\mathbf{x}_n|z_n) \right]$$

where we have grouped data by cluster labels  $z_n$ .

## Parameter estimation for GMMs: complete data

The complete likelihood is decomposable

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_n \log p(z_n) p(\mathbf{x}_n | z_n) = \sum_k \sum_{n:z_n=k} \log p(z_n) p(\mathbf{x}_n | z_n)$$

where we have grouped data by cluster labels  $z_n$ .

Let  $\gamma_{nk} \in \{0, 1\}$  be a binary variable that indicates whether  $z_n = k$ :

## Parameter estimation for GMMs: complete data

The complete likelihood is decomposable

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_n \log p(z_n)p(\mathbf{x}_n|z_n) = \sum_k \left[ \sum_{n:z_n=k} \log p(z_n)p(\mathbf{x}_n|z_n) \right]$$

where we have grouped data by cluster labels  $z_n$ .

Let  $\gamma_{nk} \in \{0, 1\}$  be a binary variable that indicates whether  $z_n = k$ :

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \underbrace{\gamma_{nk}}_{\text{indicates } z_n = k} \log p(z = k)p(\mathbf{x}_n|z = k)$$

## Parameter estimation for GMMs: complete data

The complete likelihood is decomposable

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_n \log p(z_n) p(\mathbf{x}_n | z_n) = \sum_k \sum_{n:z_n=k} \log p(z_n) p(\mathbf{x}_n | z_n)$$

where we have grouped data by cluster labels  $z_n$ .

Let  $\gamma_{nk} \in \{0, 1\}$  be a binary variable that indicates whether  $z_n = k$ :

$$\begin{aligned}\sum_n \log p(\mathbf{x}_n, z_n) &= \sum_k \sum_n \gamma_{nk} \log p(z = k) p(\mathbf{x}_n | z = k) \\ &= \sum_k \sum_n \gamma_{nk} [\log \omega_k + \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]\end{aligned}$$

Note: in the complete setting the  $\gamma_{nk}$  just add to the notation, but later we will 'relax' these variables and allow them to take on fractional values

## Parameter estimation for GMMs: complete data

From our previous discussion, we have

$$\sum_n \log p(x_n, z_n) = \sum_k \sum_n \gamma_{nk} [\log \omega_k + \log N(x_n | \mu_k, \Sigma_k)]$$

## Parameter estimation for GMMs: complete data

From our previous discussion, we have

$$\sum_n \log p(x_n, z_n) = \sum_k \sum_n \gamma_{nk} [\log \omega_k + \log N(x_n | \mu_k, \Sigma_k)]$$

Regrouping, we have

$$\sum_n \log p(x_n, z_n) = \sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log N(x_n | \mu_k, \Sigma_k) \right\}$$

## Parameter estimation for GMMs: complete data

From our previous discussion, we have

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} [\log \omega_k + \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

Regrouping, we have

$$l(\theta) = \sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \underbrace{\gamma_{nk} \log \omega_k}_{\text{depends on } \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k} + \sum_k \left\{ \sum_n \gamma_{nk} \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} + \lambda \left( \sum_{k=1}^K \omega_k - 1 \right)$$

$$\frac{\partial l}{\partial \omega_k} =$$

The term inside the braces depends on  $k$ -th component's parameters. It is now easy to show that (left as an exercise) the MLE is:

HW: Check!

How to compute  $\omega_k$ ?

$$\cancel{\omega_k} = \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}}, \quad \boldsymbol{\mu}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

What's the intuition?  $\frac{\partial l}{\partial \omega_k} = \sum_n \gamma_{nk} \frac{1}{\omega_k} + \lambda = 0 \Rightarrow \sum_k \sum_n \gamma_{nk} + \sum_k \lambda \omega_k = 0$

$$\Rightarrow \omega_k = \frac{1}{N} \sum_n \gamma_{nk} \Rightarrow \lambda = -\frac{1}{N}$$

Since  $\gamma_{nk}$  is binary, the previous solution is nothing but:

- $\omega_k$ : fraction of total data points whose cluster label  $z_n$  is  $k$ 
  - note that  $\sum_k \sum_n \gamma_{nk} = N$
- $\mu_k$ : mean of all data points whose  $z_n$  is  $k$
- $\Sigma_k$ : covariance of all data points whose  $z_n$  is  $k$

Since  $\gamma_{nk}$  is binary, the previous solution is nothing but:

- $\omega_k$ : fraction of total data points whose cluster label  $z_n$  is  $k$ 
  - note that  $\sum_k \sum_n \gamma_{nk} = N$
- $\mu_k$ : mean of all data points whose  $z_n$  is  $k$
- $\Sigma_k$ : covariance of all data points whose  $z_n$  is  $k$

Recall that this depends on us knowing the true cluster labels  $z_n$

This intuition will help us develop an algorithm for estimating  $\theta$  when we \*do not\* know  $z_n$  (incomplete data)

## GMMs and Incomplete Data

---

# Parameter estimation for GMMs: Incomplete data

## GMM Parameters

$$\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$$

## Incomplete Data

Our data contains observed and unobserved data, and hence is incomplete

- Observed:  $\mathcal{D} = \{\mathbf{x}_n\}$
- Unobserved (hidden):  $\{\mathbf{z}_n\}$

# Parameter estimation for GMMs: Incomplete data

## GMM Parameters

$$\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$$

## Incomplete Data

Our data contains observed and unobserved data, and hence is incomplete

- Observed:  $\mathcal{D} = \{x_n\}$  ✓
- Unobserved (hidden):  $\{z_n\}$

**Goal** Obtain the maximum likelihood estimate of  $\theta$ :

$$\theta = \arg \max \ell(\theta) = \arg \max \log \mathcal{D} = \arg \max \sum_n \log p(x_n | \theta)$$

$$\arg \max \sum_{n=1}^N \log p(x_n | \theta)$$

# Parameter estimation for GMMs: Incomplete data

## GMM Parameters

$$\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$$

## Incomplete Data

Our data contains observed and unobserved data, and hence is incomplete

- Observed:  $\mathcal{D} = \{x_n\}$
- Unobserved (hidden):  $\{z_n\}$

**Goal** Obtain the maximum likelihood estimate of  $\theta$ :

$$\begin{aligned}\theta &= \arg \max \ell(\theta) = \arg \max \log \mathcal{D} = \boxed{\arg \max \sum_n \log p(x_n | \theta)} \\ &= \arg \max \sum_n \log \sum_{z_n} p(x_n, z_n | \theta)\end{aligned}$$

The objective function  $\ell(\theta)$  is called the *incomplete* log-likelihood.

## Issue with Incomplete log-likelihood

No simple way to optimize the incomplete log-likelihood (exercise: try to take derivative with respect to parameters, set it to zero and solve)

## Issue with Incomplete log-likelihood

No simple way to optimize the incomplete log-likelihood (exercise: try to take derivative with respect to parameters, set it to zero and solve)

**EM algorithm** provides a strategy for iteratively optimizing this function

## Issue with Incomplete log-likelihood

No simple way to optimize the incomplete log-likelihood (exercise: try to take derivative with respect to parameters, set it to zero and solve)

**EM algorithm** provides a strategy for iteratively optimizing this function

Two steps as they apply to GMM:

- E-step: 'guess' values of the  $z_n$  using existing values of  $\theta$
- M-step: solve for new values of  $\theta$  given imputed values for  $z_n$  (i.e., maximize complete likelihood!)

## E-step: Soft cluster assignments

We define  $\underline{\gamma_{nk}}$  as  $p(z_n = k | \mathbf{x}_n, \theta)$

- This is the posterior distribution of  $z_n$  given  $\mathbf{x}_n$  and  $\theta$

## E-step: Soft cluster assignments

We define  $\gamma_{nk}$  as  $p(z_n = k | \mathbf{x}_n, \theta)$

- This is the posterior distribution of  $z_n$  given  $\mathbf{x}_n$  and  $\theta$
- Recall that in complete data setting  $\gamma_{nk}$  was binary

## E-step: Soft cluster assignments

We define  $\gamma_{nk}$  as  $p(z_n = k | \mathbf{x}_n, \theta)$

- This is the posterior distribution of  $z_n$  given  $\mathbf{x}_n$  and  $\theta$
- Recall that in complete data setting  $\gamma_{nk}$  was binary
- Now it's a "soft" assignment of  $\mathbf{x}_n$  to  $k$ -th component, with  $\mathbf{x}_n$  assigned to each component with some probability

## E-step: Soft cluster assignments

We define  $\gamma_{nk}$  as  $p(z_n = k | \mathbf{x}_n, \theta)$

- This is the posterior distribution of  $z_n$  given  $\mathbf{x}_n$  and  $\theta$
- Recall that in complete data setting  $\gamma_{nk}$  was binary
- Now it's a "soft" assignment of  $\mathbf{x}_n$  to  $k$ -th component, with  $\mathbf{x}_n$  assigned to each component with some probability

## E-step: Soft cluster assignments

We define  $\gamma_{nk}$  as  $p(z_n = k | \mathbf{x}_n, \theta)$

- This is the posterior distribution of  $z_n$  given  $\mathbf{x}_n$  and  $\theta$
- Recall that in complete data setting  $\gamma_{nk}$  was binary
- Now it's a "soft" assignment of  $\mathbf{x}_n$  to  $k$ -th component, with  $\mathbf{x}_n$  assigned to each component with some probability

Given an estimate of  $\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$ , we can compute  $\gamma_{nk}$  as follows:

$$\begin{aligned}\gamma_{nk} &= p(z_n = k | \mathbf{x}_n) \\ &= \frac{p(\mathbf{x}_n | z_n = k) p(z_n = k)}{p(\mathbf{x}_n)}\end{aligned}$$

$$P(A|B) = P(B|A) * P(A) / P(B)$$

## E-step: Soft cluster assignments

We define  $\gamma_{nk}$  as  $p(z_n = k | \mathbf{x}_n, \theta)$

- This is the posterior distribution of  $z_n$  given  $\mathbf{x}_n$  and  $\theta$
- Recall that in complete data setting  $\gamma_{nk}$  was binary
- Now it's a "soft" assignment of  $\mathbf{x}_n$  to  $k$ -th component, with  $\mathbf{x}_n$  assigned to each component with some probability

Given an estimate of  $\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$ , we can compute  $\gamma_{nk}$  as follows:

$$\begin{aligned}\gamma_{nk} &= p(z_n = k | \mathbf{x}_n) \\ &= \frac{p(\mathbf{x}_n | z_n = k)p(z_n = k)}{p(\mathbf{x}_n)} \\ &= \frac{p(\mathbf{x}_n | z_n = k)p(z_n = k)}{\sum_{k'=1}^K p(\mathbf{x}_n | z_n = k')p(z_n = k')}\end{aligned}$$

## M-step: Maximize complete likelihood

Recall definition of complete likelihood from earlier:

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Previously  $\gamma_{nk}$  was binary, but now we define  $\gamma_{nk} = p(z_n = k | \mathbf{x}_n)$  (E-step)

## M-step: Maximize complete likelihood

Recall definition of complete likelihood from earlier:

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Previously  $\gamma_{nk}$  was binary, but now we define  $\gamma_{nk} = p(z_n = k | \mathbf{x}_n)$  (E-step)

We get the same simple expression for the MLE as before!

$$\omega_k = \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}}, \quad \boldsymbol{\mu}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

## M-step: Maximize complete likelihood

Recall definition of complete likelihood from earlier:

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Previously  $\gamma_{nk}$  was binary, but now we define  $\gamma_{nk} = p(z_n = k | \mathbf{x}_n)$  (E-step)

We get the same simple expression for the MLE as before!

$$\omega_k = \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}}, \quad \boldsymbol{\mu}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

Intuition: Each point now contributes some fractional component to each of the parameters, with weights determined by  $\gamma_{nk}$

## EM procedure for GMM

Alternate between estimating  $\gamma_{nk}$  and estimating  $\theta$

- Initialize  $\theta$  with some values (random or otherwise)
- Repeat
  - ✓ • E-Step: Compute  $\gamma_{nk}$  using the current  $\theta$
  - ✓ • M-Step: Update  $\theta$  using the  $\gamma_{nk}$  we just computed
- Until Convergence

## EM procedure for GMM

Alternate between estimating  $\gamma_{nk}$  and estimating  $\theta$

- Initialize  $\theta$  with some values (random or otherwise)
- Repeat
  - E-Step: Compute  $\gamma_{nk}$  using the current  $\theta$
  - M-Step: Update  $\theta$  using the  $\gamma_{nk}$  we just computed
- Until Convergence

Questions to be answered next

- How does GMM relate to  $K$ -means?
- Is this procedure reasonable, i.e., are we optimizing a sensible criterion?
- Will this procedure converge?

## GMMs and K-means

GMMs provide probabilistic interpretation for K-means

GMMs provide probabilistic interpretation for K-means

HW:

GMMs reduce to K-means under the following assumptions (in which case EM for GMM parameter estimation simplifies to K-means):

- Assume all Gaussians have  $\sigma^2 \mathbf{I}$  covariance matrices
- Further assume  $\sigma \rightarrow 0$ , so we only need to estimate  $\mu_k$ , i.e., means

K-means is often called “hard” GMM or GMMs is called “soft” K-means

The posterior  $\gamma_{nk}$  provides a probabilistic assignment for  $x_n$  to cluster  $k$

$$\begin{aligned}
 V_{nk} &= p(z_n=k | x_n; \theta) \\
 &= \frac{p(x_n | z_n=k) p(z_n=k)}{p(x_n)} \\
 &= \frac{\omega_k N(x_n; \mu_k, \Sigma_k)}{\sum_{j=1}^K \omega_j N(x_n; \mu_j, \Sigma_j)} \\
 &= \frac{\omega_k \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) \right\}}{\sum_{j=1}^K \omega_j \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_n - \mu_j)^T \Sigma_j^{-1} (x_n - \mu_j) \right\}}
 \end{aligned}$$

*Substitution  $\Sigma_k = \sigma^2 I$*

$$V_{nk} = \frac{\omega_k \exp \left\{ -\frac{1}{2\sigma^2} \|x_n - \mu_k\|^2 \right\}}{\sum_{j=1}^K \omega_j \exp \left\{ -\frac{1}{2\sigma^2} \|x_n - \mu_j\|^2 \right\}} = \frac{\omega_k}{\sum_{j=1}^K \omega_j \exp \left\{ -\frac{1}{2\sigma^2} (\|x_n - \mu_j\|^2 - \|x_n - \mu_k\|^2) \right\}}$$

if  $\sigma \rightarrow 0$ ,  $V_{nk} \rightarrow \begin{cases} 1 & \text{if } x_n \text{ is closest to } \mu_k \\ 0 & \text{if } x_n \text{ is not closest to } \mu_k \end{cases}$  | K-mean  
labels -

## GMMs vs. $k$ -means

### Pros/Cons

- $k$ -means is a simpler, more straightforward method, but might not be as accurate because of deterministic clustering

## GMMs vs. $k$ -means

### Pros/Cons

- $k$ -means is a simpler, more straightforward method, but might not be as accurate because of deterministic clustering
- GMMs can be more accurate, as they model more information (soft clustering, variance), but can be more expensive to compute

# GMMs vs. $k$ -means

## Pros/Cons

- $k$ -means is a simpler, more straightforward method, but might not be as accurate because of deterministic clustering
- GMMs can be more accurate, as they model more information (soft clustering, variance), but can be more expensive to compute
- Both methods have a similar set of practical issues (having to select  $k$ , the distance, and the initialization)

