



UNIVERSITÀ
di **VERONA**

Machine Learning

PARAMETRIC UNSUPERVISED LEARNING

Cigdem Beyan
A. Y. 2024/2025



UNSUPERVISED LEARNING

- **Non-Parametric**

- **Clustering:** no assumptions are made about the underlying densities, instead we seek a partition of the data into clusters

- **Parametric approach**

- Assume the parametric distribution of data
- Model the underlying **class-conditional densities** with a mixture of parametric densities
- The objective is to find the model parameters
 - Estimate parameters of the distribution assumed



PARAMETRIC UNSUPERVISED LEARNING

- We assume the data was generated by a **model** with **known shape** but **unknown parameters**
- What is good about having a **model**? We can:
 - Adjust the parameters of the model to maximize the probability that our model produced for the observed data
 - Compute the **likelihood** of data
 - Compare multiple algorithms' performance: whichever gives the higher likelihood for the observed data

RECAP: PARAMETRIC SUPERVISED LEARNING

- Given m classes
- Samples $\mathbf{x}_1, \dots, \mathbf{x}_n$ each of class $1, 2, \dots, m$
- Suppose \mathbf{D}_i has samples from class i
- The probability distribution for class i is $p_i(\mathbf{x} | \theta_i)$
 - For class 1: $p_1(\mathbf{x} | \theta_1)$
 - For class 2: $p_2(\mathbf{x} | \theta_2)$
- Use **Maximum likelihood** method to estimate parameters θ_i



PARAMETRIC UNSUPERVISED LEARNING

- We do not have the true classes for samples. But we still know that
 - We have **m** classes
 - Have samples $\mathbf{x}_1, \dots, \mathbf{x}_n$ each of **unknown** class
 - Probability distribution for class **i** is $\mathbf{p}_i(\mathbf{x} | \theta_i)$
- The task is in a way to determine the classes and parameters simultaneously

EXAMPLE APPLICATION

- MRI brain image segmentation
 - Different brain tissues have different intensities
 - There are 6 major tissue types
 - Each type of tissue can be modelled by a Gaussian distribution
 - Parameters: **mean** and **standard deviation** → **unknown**
- Segmenting (classifying) the brain image into different tissue classes is our task. But we do not know:
 - which image pixel corresponds to which tissue (class)
 - parameters for **Gaussian distribution**

GAUSSIAN MIXTURE MODEL (GMM)

- Gaussian mixtures of **D-dimensional variables** with **K Gaussian components** is written as:

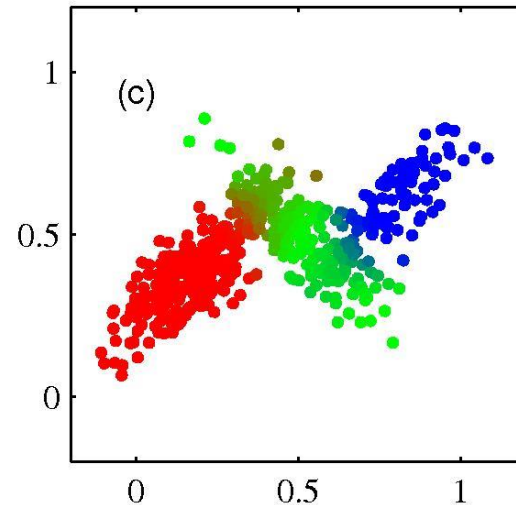
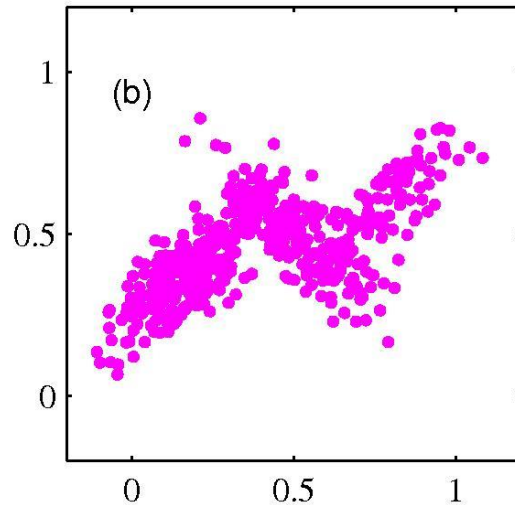
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)$$

- GMM introduces a D-dimensional **binary** variable called **latent variable** shown as **z** (which is not observed, but inferred during modelling process) to indicate which Gaussian component a data point belongs to.
 - Let $\mathbf{z} = z_1, \dots, z_K$ $z_k \in \{0, 1\}$, $\sum_k z_k = 1$: Only one component of **z** can be active (equal to 1) at a time
 - We have K possible states of **z** corresponding to **K components**
 - One of K representation is: $[0 \dots 0 \ 1 \ 0 \ 0 \dots 0]$ (one-hot vector)
 - E.g., $\mathbf{z} = [0, 0, 1, 0, 0]$ indicates that the data point belongs to the third Gaussian component.

GAUSSIAN MIXTURE MODEL (GMM)

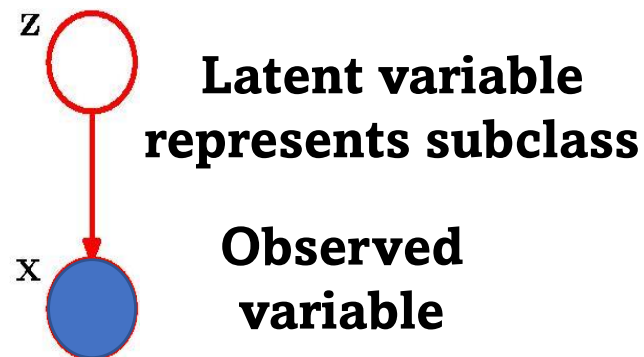
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)$$

- The aim is to find **maximum likelihood** parameters: $\pi_k, \boldsymbol{\mu}_k, \Sigma_k$
- Each data point is associated with a subclass (sub-group) k with **probability** π_k



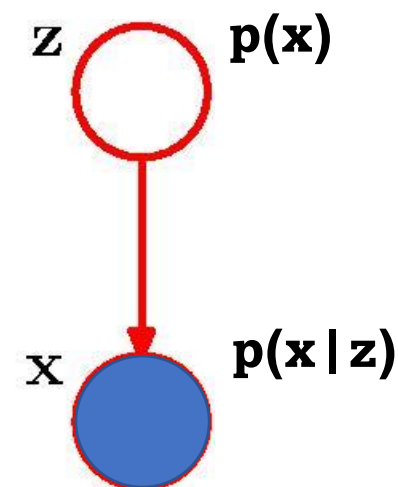
GAUSSIAN MIXTURE MODEL (GMM)

- The joint probability of **latent variable (\mathbf{z})** and **observed variable (\mathbf{x})**
 - $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})$
 - \mathbf{x} is the observed variable
 - \mathbf{z} is the hidden or missing variable: latent variable
 - **Marginal distribution or prior probability $p(\mathbf{z})$:** Representing the likelihood of \mathbf{z} independently of \mathbf{x} .
 - **Conditional distribution $p(\mathbf{x} | \mathbf{z})$:** describes how the observed variable \mathbf{x} depends on the state of the latent variable \mathbf{z} .
 - Graphical representation:



GAUSSIAN MIXTURE MODEL (GMM)

- By using $\mathbf{p}(\mathbf{z})$ and $\mathbf{p}(\mathbf{x} | \mathbf{z})$
 - we specify $\mathbf{p}(\mathbf{x})$ in terms of \mathbf{x} and \mathbf{z}
- Each component \mathbf{z}_k is assigned a probability
 - $p(z_k = 1) = \pi_k$, where parameters π_k satisfy $0 \leq \pi_k \leq 1$ and $\sum_k \pi_k = 1$
- Because \mathbf{z} uses **1-of-K** it follows that $p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$ since $\mathbf{z}_k \in \{0, 1\}$ and components of \mathbf{z} are **mutually exclusive** and **therefore independent**.
With one mixture component: $p(z_1) = \pi_1^{z_1}$
With two components: $p(z_1, z_2) = \pi_1^{z_1} \pi_2^{z_2}$



GAUSSIAN MIXTURE MODEL (GMM)

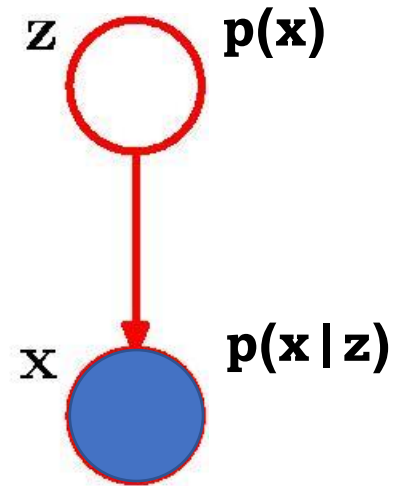
- For a particular component (value of \mathbf{z})

$$p(\mathbf{x} | z_k = 1) = N(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k)$$

- Thus $\mathbf{p}(\mathbf{x} | \mathbf{z})$ can be written in the form

$$p(\mathbf{x} | \mathbf{z}) = \prod_{k=1}^K N(\mathbf{x} | \boldsymbol{\mu}_k, \Sigma_k)^{z_k}$$

- Using the one-hot property of \mathbf{z} , only the term corresponding to the active component ($\mathbf{z}_{\mathbf{k}} = 1$) contributes to the product.



GAUSSIAN MIXTURE MODEL (GMM)

- The joint distribution $\mathbf{p}(\mathbf{x}, \mathbf{z})$ is equal to $\mathbf{p}(\mathbf{z})\mathbf{p}(\mathbf{x} | \mathbf{z})$
- Based on law of total probability, marginal distribution of \mathbf{x} is obtained by summing over all possible states of \mathbf{z} to give

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) = \sum_{\mathbf{z}} \prod_{k=1}^K \pi_k^{z_k} N(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} = \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$$

- Since $\mathbf{z}_k \in \{0, 1\}$
- This is the standard form of a Gaussian Mixture

GAUSSIAN MIXTURE MODEL (GMM)

- If we have observations $\mathbf{x}_1, \dots, \mathbf{x}_N$
- Because marginal distribution is in the form:
$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$$
 - Every observed data point \mathbf{x}_n there is a corresponding latent vector \mathbf{z}_n , i.e., its sub-class
- Thus, we have found a formulation of GMM involving a latent variable meaning that we can directly incorporate the relationships between the \mathbf{x}_n and \mathbf{z}_n (i.e., component memberships).
 - This approach simplifies parameter estimation by working with $p(\mathbf{x}, \mathbf{z})$ rather than $p(\mathbf{x})$

GAUSSIAN MIXTURE MODEL (GMM)

- Another conditional probability is called **Responsibility**, which is the **posterior probability** $p(\mathbf{z}|\mathbf{x})$ of a data point \mathbf{x} being associated with a particular component k . $p(\mathbf{z}_k = 1|\mathbf{x})$ is denoted $\gamma(\mathbf{z}_k)$

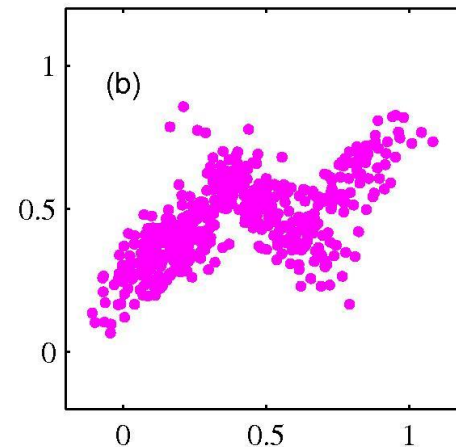
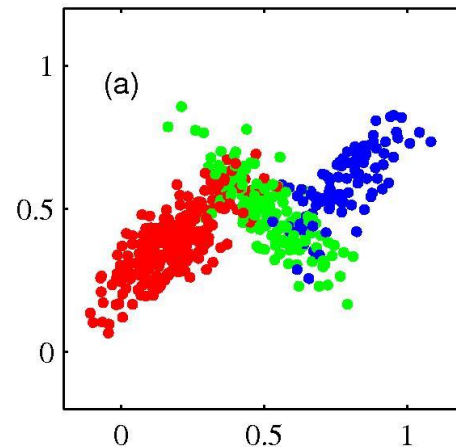
- From Bayes theorem
$$p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1) p(\mathbf{x}|z_k = 1)}{p(\mathbf{x})}.$$

$$\begin{aligned}\gamma(\mathbf{z}_k) &\equiv p(z_k = 1|\mathbf{x}) = \\ &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} = \frac{\pi_k \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\boldsymbol{\mu}_j, \Sigma_j)}\end{aligned}$$

- π_k can be viewed as the **prior probability** of $\mathbf{z}_k = 1$
- $\gamma(\mathbf{z}_k)$ can be interpreted as the **responsibility** that component k takes for 'explaining' the observation \mathbf{x} .

GETTING DATA FROM GMM SYNTHETICALLY

- By generating synthetic data, you can compare model predictions against the true data distribution, which is useful for evaluating the model's performance, diagnosing issues, and fine-tuning parameters.
- Two ways:
 - We can generate a value of \mathbf{z}^* from $p(\mathbf{z})$, and then, generate a value for \mathbf{x} from $p(\mathbf{x} | \mathbf{z}) \rightarrow$ plot (a) in the figure
 - Generate from $p(\mathbf{x})$ by ignoring the values of $\mathbf{z} \rightarrow$ plot (b) in the figure



MAXIMUM LIKELIHOOD FOR GMM

- We would like to model **data set** $\{x_1, \dots, x_N\}$ using a mixture of Gaussians (N items each of dimension D).

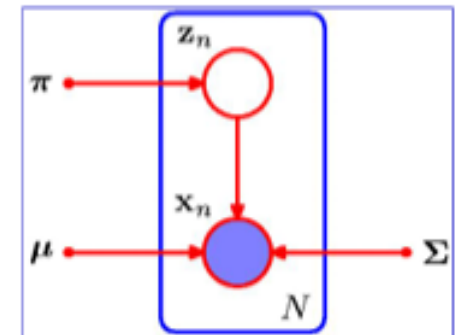
- \mathbf{X} is represented by a $N \times D$ matrix
 - n^{th} row is given by x_n^T

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

- \mathbf{Z} represents N latent variables with $N \times K$ matrix
 - n^{th} row is given by z_n^T

$$Z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_N \end{bmatrix}$$

- The goal is to estimate **means**, **covariances**, and **mixing coefficients** by maximizing the likelihood of observing the dataset X . Maximizing the likelihood involves adjusting the parameters so that the probability of observing the data under the model is as high as possible.



LIKELIHOOD FUNCTION FOR GMM

- Given mixture density function

$$p(\mathbf{x}) = \sum_z p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) = \sum_z \prod_{k=1}^K \pi_k^{z_k} N(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} = \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$$

- The joint likelihood for the dataset \mathbf{X} is

$$p(\mathbf{X} | \pi, \mu, \Sigma) = \prod_{n=1}^N \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

The product is over the N
Independent and identically
distributed samples

Since the data points are independently drawn.

- Log-likelihood function is

$$\ln p(\mathbf{X} | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

We need to maximize
this

LOG-LIKELIHOOD FUNCTION FOR GMM

$$\ln p(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

- The goal is to estimate the three sets of parameters:
 π_k, μ_k, Σ_k
 - We need to take derivatives w.r.t. each while keeping the others constant
 - However, there are no closed-form solutions (check the formulas to see why)
 - Task is not straightforward since summation appears in Gaussian and logarithm does not operate on Gaussian
- A gradient-based optimization is possible,
 - we consider the iterative EM algorithm

- Before proceeding with the Expectation Maximization (Maximum Likelihood Estimation) we need to mention two technical issues:
 - Problem of singularities with Gaussian mixtures
 - Singularities occur when one or more of the components of the mixture model become overly narrow, essentially collapsing to a single point or becoming degenerate (i.e., zero variance).
 - Problem of Identifiability of mixtures
 - A model is said to be identifiable if, for a given set of observed data, there is a unique set of model parameters that can explain the data.
 - If a model is non-identifiable, it means that there are multiple different sets of parameters that could explain the same data, leading to ambiguity in the model's interpretation.

SINGULARITIES PROBLEM

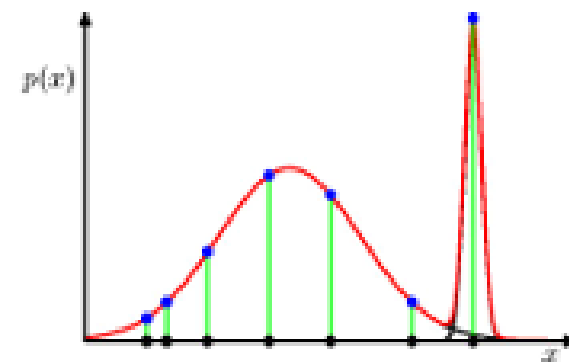
- Consider a Gaussian mixture
 - Having components with covariance matrices $\Sigma_k = \sigma_k^2 I$
- Data point that falls on a mean $\mu_j = x_n$ contribute to the likelihood function

$$N(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 I) = \frac{1}{(2\pi)^{1/2}} \frac{1}{\sigma_j}$$

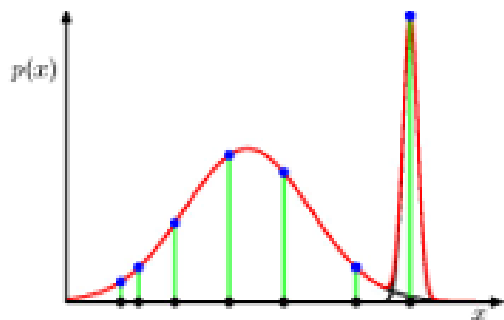
since $\exp(\mathbf{x}_n - \mu_j)^2 = 1$ (see Appendix – a for the proof)

- As $\sigma_k \rightarrow 0$ terms goes to infinity
- Thus, maximization of log-likelihood is not well-posed

$$\ln p(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$



SINGULARITIES PROBLEM



- One component assigns finite values because its variance is almost zero.
- Instead, the other component has a large variance.
- The solution in the case of $\sigma_k \rightarrow 0$ is resetting the mean or covariance to a **small but nonzero value**.

IDENTIFIABILITY PROBLEM

- A density $p(x | \theta)$ is identifiable if $\theta \neq \theta'$ then there is an x for which $p(x | \theta) \neq p(x | \theta')$
- A **K-component** mixture will have a total of **K! equivalent solutions**
 - Corresponding to K! ways of assigning K sets of parameters to K components
 - E.g., for K=3 K!=6: 123, 132, 213, 231, 312, 321
 - For any given point in the space of parameter values there will be a further K!-1 additional points all giving the same distribution
- However, any of the equivalent solutions are as good as the other
- The EM algorithm can help address the problem of identifiability by iteratively refining the estimates of the parameters, although it does not necessarily solve the issue of multiple equivalent solutions.
 - To avoid non-identifiable solutions, GMM parameters can be initialized using K-means, providing a better starting point for EM.

EXPECTATION-MAXIMIZATION (EM) FOR GMM

- EM is a method for finding **maximum likelihood** solutions for models with **latent variables**
- We begin with Log-likelihood function

$$\ln p(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

- We want to find **π_k, μ_k, Σ_k** that maximizes this function
- This task is not straightforward since **summation** appears in Gaussian and **logarithm does not operate on Gaussian**

EXPECTATION-MAXIMIZATION (EM) FOR GMM

■

$$\ln p(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

- Take derivatives one by one
 - Means μ_k and set to zero
 - Covariance matrices Σ_k and set to zero
 - mixing coefficients π_k , and set to zero

EXPECTATION-MAXIMIZATION (EM) FOR GMM

$$\ln p(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

- Take derivative w.r.t. the means μ_k and set to zero
 - When you use the exponential form of Gaussian and then use the formulas $\frac{d}{dx} \ln u = \frac{u'}{u}$ and $\frac{d}{dx} e^u = e^u u'$

- We obtain:

$$0 = \sum_{n=1}^N \underbrace{\frac{\pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(\mathbf{x}_n | \mu_j, \Sigma_j)}}_{\gamma(z_{nk})} \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)$$

Inverse of covariance matrix

$\gamma(z_{nk})$ the posterior probabilities

EXPECTATION-MAXIMIZATION (EM) FOR GMM

$$0 = \sum_{n=1}^N \underbrace{\frac{\pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(\mathbf{x}_n | \mu_j, \Sigma_j)}}_{\gamma(z_{nk})} \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)$$

Inverse of covariance matrix

$\gamma(z_{nk})$ the posterior probabilities

- Mean of k^{th} Gaussian component is the weighted mean of all the points in the data set:
 - where data point \mathbf{x}_n is weighted by the posterior probability that component k was responsible for generating \mathbf{x}_n
- Consequently, rearranging the above formula and multiplying Σ_k

EXPECTATION-MAXIMIZATION (EM) FOR GMM

- We obtain the maximum likelihood estimation for **means** as:
(refer to Appendix –B for mathematical derivations)

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

- The effective number of points assigned to *cluster k (mixture k)*:

$$N_k = \sum_{n=1}^N \gamma(z_{nk})$$

EXPECTATION-MAXIMIZATION (EM) FOR GMM

- Maximum likelihood solution for **covariance**:
 - Take derivative w.r.t. the means Σ_k and set to zero

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T$$

- Notice that each data point is **weighted** by the corresponding **posterior probability**!!!!
- The denominator is the same as in the **mean**, i.e., the number of data points in a component.

EXPECTATION-MAXIMIZATION (EM) FOR GMM

- Maximize $\ln p(X | \pi, \mu, \Sigma)$ w.r.t. π_k
 - Must consider that mixing coefficients sum to one
 - To solve this, we need to use Lagrange multiplier and maximizing

$$\ln p(X | \pi, \mu, \Sigma) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

- And then set derivative wrt. π_k to zero, giving us:

$$\pi_k = \frac{N_k}{N}$$

EXPECTATION-MAXIMIZATION (EM) FOR GMM

- GMM maximum likelihood parameter estimates:

Means:

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

Covariance matrices:

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T$$

Mixing Coefficients:

$$\pi_k = \frac{N_k}{N}$$

$$N_k = \sum_{n=1}^N \gamma(z_{nk})$$

- All three are in terms of responsibilities therefore we have not completely solved the problem.

EM FORMULATION FOR GMM

- The results for π_k, μ_k, Σ_k are not closed form solutions for the parameters since γ_{nk} (responsibilities) depend on those parameters

$$\gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}$$

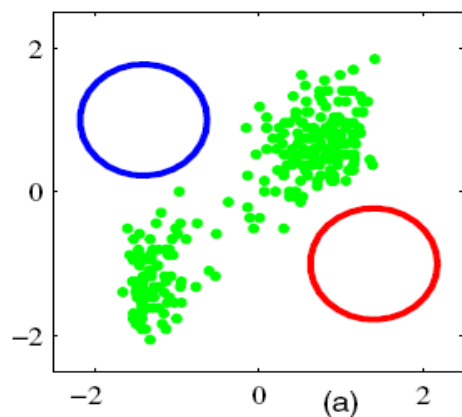
- This results in building an iterative solution, which is an instance of EM algorithm for GMM.

EM FORMULATION FOR GMM

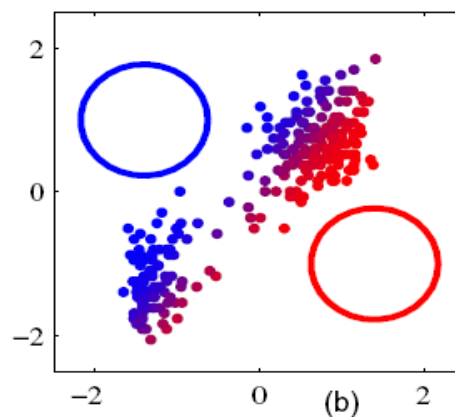
- Initialize the parameters of GMM: means, covariances and mixing coefficients
- Apply two updated: **E step** & **M step**
 - **E step**: use the current value of parameters to evaluate posterior probabilities or responsibilities
 - **M step**: use these posterior probabilities to re-estimate means, covariances, and mixing coefficients

EM FORMULATION FOR GMM

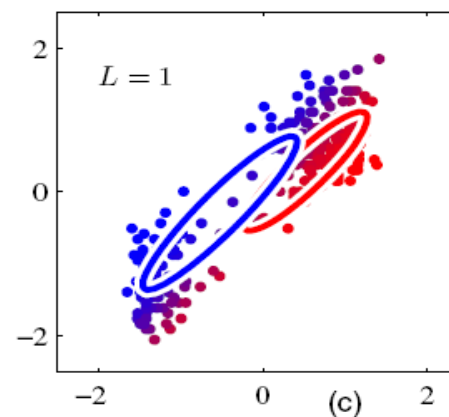
**Data points and
initial mixture model**



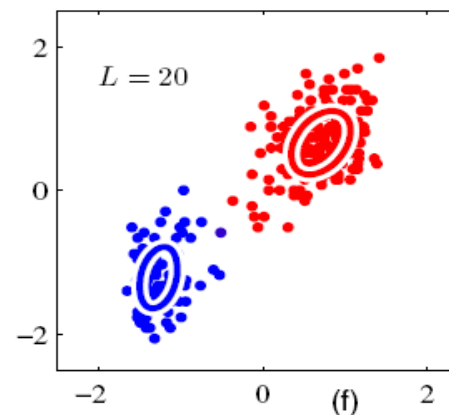
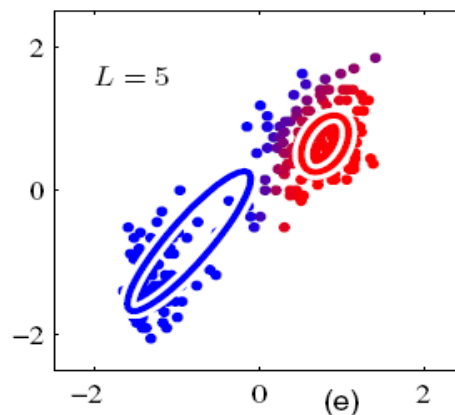
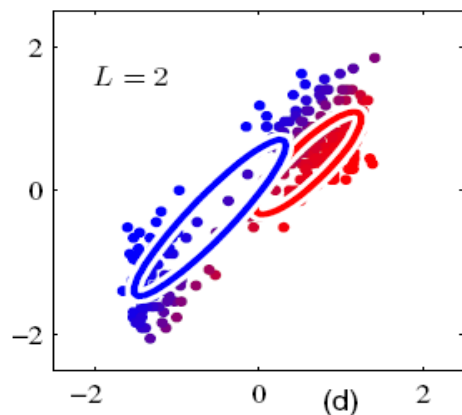
**Initial E step:
determine
responsibilities**



**After the first M step
re-evaluate
parameters**



**After 2, 5 and 20
cycles →**



EM FORMULATION FOR GMM

- Given a Gaussian mixture model
- Goal is to maximize the likelihood function w.r.t. the parameters (means, covariances and mixing coefficients)
- **Step 1:** Initialize the π_k, μ_k, Σ_k and evaluate the initial value of log-likelihood

$$\ln p(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

- **Step 2: E-step:** Evaluate responsibilities using current parameter values

$$\gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}$$

EM FORMULATION FOR GMM

- **Step 3: M-step:** Re-estimate parameters using current responsibilities

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T$$

$$\pi_k = \frac{N_k}{N}$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk})$$

EM FORMULATION FOR GMM

- **Step 4:** Evaluate the **log-likelihood**

$$\ln p(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

- And check the convergence of either **parameters** or **log-likelihood**
- If the convergence is not satisfied GOTO Step 2



PROPERTIES OF EM FOR GMM

- Takes many **more iterations** than **K-means**
 - Each cycle requires significantly more comparison
- Common to run **K-means** first to find **suitable initialization**
- Covariance matrices can be initialized to **covariances** of **clusters** found by **K-means**
- **EM** is not guaranteed to find the **global maximum** of the log-likelihood function



EXAMPLE

<https://colab.research.google.com/drive/1hn0UuluAxRPzNAMZmEyLPsJQ6tT9vRHz?usp=sharing>



GMM SUMMARY

- A probabilistic view of unsupervised learning (even clustering).
- Each cluster corresponds to a **different Gaussian**.
- Model using **latent variables**.
- General approach, can eventually replace Gaussian with other distributions (continuous or discrete) .
- More generally, mixture model are very powerful models, universal approximator.
- Optimization is done using the **EM algorithm**.

APPENDIX -A

- In a Gaussian Mixture Model (GMM), when a data point x_n falls exactly on the mean μ_j of a particular component j , it contributes to the likelihood function through the probability density function of the corresponding Gaussian distribution.
- The pdf of a Gaussian distribution at a specific point x is given by:

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- For a GMM component j with mean μ_j and covariance matrix Σ_j , the pdf at x_n is:

$$\mathcal{N}(x_n|\mu_j, \Sigma_j) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_j|}} e^{-\frac{1}{2}(x_n - \mu_j)^T \Sigma_j^{-1} (x_n - \mu_j)}$$

APPENDIX - A

$$\mathcal{N}(x_n | \mu_j, \Sigma_j) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_j|}} e^{-\frac{1}{2} (x_n - \mu_j)^T \Sigma_j^{-1} (x_n - \mu_j)}$$

- where $|\Sigma_j|$ denotes the determinant of the covariance matrix Σ_j .
 Σ_j^{-1} is the inverse of the covariance matrix.
 d is the dimensionality of the data space.
- When x_n falls exactly on the mean μ_j , the term $(x_n - \mu_j)$ becomes zero in the exponent. This means that the exponent reduces to zero, and $e^0 = 1$. Consequently, the likelihood contribution from x_n to the Gaussian component j becomes:

$$\mathcal{N}(x_n | \mu_j, \Sigma_j) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_j|}}$$

APPENDIX - B

- The log-likelihood function for a Gaussian Mixture Model (GMM) is given by:

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \ln \sum_{n=1}^N \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$

- First, let's differentiate the log-likelihood function w.r.t $\boldsymbol{\mu}_k$:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln \sum_{n=1}^N \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \\ &= \frac{1}{\sum_{n=1}^N \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)} \sum_{n=1}^N \left(\frac{\partial}{\partial \boldsymbol{\mu}_k} \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \right) \end{aligned}$$

- Next, we differentiate the term $\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ with respect to $\boldsymbol{\mu}_k$:

$$\frac{\partial}{\partial \boldsymbol{\mu}_k} \left(\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) = \pi_k \frac{\partial}{\partial \boldsymbol{\mu}_k} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

APPENDIX - B

- The derivative of the Gaussian distribution with respect to the mean vector $\boldsymbol{\mu}_k$ is:

$$\frac{\partial}{\partial \boldsymbol{\mu}_k} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = -\frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Substituting this back into our expression, we get:

$$\frac{\partial}{\partial \boldsymbol{\mu}_k} \ln p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sum_{n=1}^N \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)} \sum_{n=1}^N \left(-\pi_k \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$

- Now, we set this derivative equal to zero and solve for $\boldsymbol{\mu}_k$ to find the stationary points:

$$\frac{1}{\sum_{n=1}^N \left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)} \sum_{n=1}^N \left(-\pi_k \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) = 0$$

APPENDIX - B

- Solving this equation for $\boldsymbol{\mu}_k$ generally requires numerical optimization methods due to the complexity of the expression and the non-linearity introduced by the Gaussian density. The resulting solution will be the MLE estimate for $\boldsymbol{\mu}_k$.
- The resulting solution for $\boldsymbol{\mu}_k$ in the context of the EM algorithm for GMMs involves a weighted average of the data points assigned to the **k -th** component, where the weights are determined by the **responsibilities** of each component for each data point.
- Specifically, the solution for $\boldsymbol{\mu}_k$ is given by:
- γ_{nk} represents the responsibility of the k -th component for the n -th data point.

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}}$$



That's all

Cigdem Beyan
cigdem.beyan@univr.it
<https://cbeyan.github.io/>