DDA3020 Machine Learning: Lecture 15 Gaussian Mixture Models

JIA, Kui School of Data Science, CUHK-SZ

November 14/16, 2023



2 The Latent Variable Perspective for Gaussian Mixture Model

3 Expectation Maximization for Fitting Gaussian Mixture Model

4 Relation to k-Means

• We model the joint distribution over (\mathbf{x}, z) as follows

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

where \mathbf{x} denotes a feature variable, and z denotes the class label variable.

- \bullet However, we do not have the class labels z in unsupervised clustering.
- $\bullet\,$ In this case, we can model the marginal distribution over ${\bf x}$ as follows

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

• This is called as mixture models.

Gaussian Mixture Model (GMM)

The most common mixture model is called as Gaussian mixture model (GMM).

• A GMM represents a **distribution** as

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

with π_k the mixing coefficients, where: $\sum_{k=1}^{K} \pi_k = 1$ and $\pi_k \ge 0, \forall k$. And,

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right),$$

with $|\Sigma_k| = \det(\Sigma_k)$ denotes the determinant of Σ_k , d indicates the dimension of \mathbf{x} .

- GMM is a density estimator. If given enough Gaussian components, GMM is universal approximator of densities.
- In general, mixture models are very powerful, but difficult to optimize

Visualizing a Mixture of Gaussians – 1D Gaussians

• If you fit a Gaussian to data:



• Now, we are trying to fit a GMM (with K = 2 in this example):



[Slide credit: K. Kutulakos]

JIA, Kui School of Data Science, CUHKDDA3020 Machine Learning: Lect

Visualizing a Mixture of Gaussians – 2D Gaussians



Fitting GMMs: Maximum Likelihood

• The log likelihood is

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

where
$$\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}, \mathbf{\Theta} = \{\pi, \mu, \Sigma\}, \pi = \{\pi_1, \dots, \pi_K\}, \mu = \{\mu_1, \dots, \mu_K\}, \Sigma = \{\Sigma_1, \dots, \Sigma_K\}, \text{and}$$

- We aim to learn the parameters Θ by maximizing the above log likelihood.
- Due to the log-sum-exp operation, we cannot obtain a closed-form solution by setting the derivative to zero.
- Of course you can choose gradient based method. However, in the following we will introduce a more elegant optimization method.

D Gaussian Mixture Model for Density Estimation

2 The Latent Variable Perspective for Gaussian Mixture Model

3 Expectation Maximization for Fitting Gaussian Mixture Model

4 Relation to k-Means



- We introduce a hidden (latent) variable z, indicating which Gaussian component generates the observation **x**, with some probability
- Let $z \sim \text{Categorical}(\boldsymbol{\pi})$, where $\boldsymbol{\pi} \geq 0$, $\sum_k \pi_k = 1$
- Then:

$$p(\mathbf{x}) = \sum_{k=1}^{K} p(\mathbf{x}, z = k) = \sum_{k=1}^{K} \underbrace{p(z = k)}_{\pi_k} \underbrace{p(\mathbf{x} \mid z = k)}_{\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

• This breaks a complicated distribution into simple components - the price is the hidden variable.

Latent variable model (LVM):

- Definition: A latent variable model is a statistical model that relates a set of observable variables to a set of latent variables.
- Some model variables may be unobserved, either at training or at testing time, or both. Variables which are always unobserved are called **latent** variables, or sometimes hidden variables.
- We may want to intentionally introduce latent variables to model complex dependencies between variables this can actually simplify the model
- According to the type of latent variables, there are two types of LVMs,
 - LVM with continuous latent variables, e.g., factor analysis
 - LVM with discrete latent variables, e.g., mixture models

Reference:

https://en.wikipedia.org/wiki/Latent_variable_model

Back to GMM

• A Gaussian mixture distribution:

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

- We had: $z \sim \text{Categorical}(\pi)$, *i.e.*, $p(z = k \mid \pi) = \pi_k$, where $\pi \geq 0$, $\sum_k \pi_k = 1$
- Joint distribution: $p(\mathbf{x}, z) = p(z)p(\mathbf{x} \mid z)$
- Log-likelihood:

$$\ell(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \ln p(\mathbf{X} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln p\left(\mathbf{x}^{(n)} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$$

$$=\sum_{n=1}^{N}\ln\sum_{k=1}^{K}p\left(\mathbf{x}^{(n)}, z^{(n)}=k \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$$

$$=\sum_{n=1}^{N}\ln\sum_{k=1}^{K}p\left(\mathbf{x}^{(n)} \mid z^{(n)}=k; \boldsymbol{\mu}, \boldsymbol{\Sigma}\right)p(z^{(n)}=k \mid \boldsymbol{\pi})$$

• Note: we have a hidden variable $z^{(n)}$ for every observation $\mathbf{x}^{(n)}$

• How can we optimize this problem, since there is sum inside the log?

• If we knew $z^{(n)}$ for every $\mathbf{x}^{(n)}$, the maximum log likelihood problem is easy:

$$\max \ell(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \sum_{k=1}^{K} \mathbb{1}_{[z^{(n)}=k]} \cdot p\left(\mathbf{x}^{(n)}, z^{(n)} = k \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$$
$$= \sum_{n=1}^{N} \left[\ln \left(\mathbb{1}_{[z^{(n)}=k]} \cdot p\left(\mathbf{x}^{(n)} \mid z^{(n)} = k; \boldsymbol{\mu}, \boldsymbol{\Sigma}\right) \right) + \ln \left(\mathbb{1}_{[z^{(n)}=k]} \cdot p\left(z^{(n)} = k \mid \boldsymbol{\pi}\right) \right) \right]$$

with the constraint $1 - \sum_{k=1}^{K} \pi_k = 0$.

• For the above constrained optimization problem, we also resort to KKT conditions based on Lagrangian function, as follows:

$$\mathcal{L}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \lambda) = -\ell(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda(1 - \sum_{k=1}^{K} \pi_k).$$

Note that the original primal problem is maximizing $\ell(\pi, \mu, \Sigma)$, thus it should be $-\ell(\pi, \mu, \Sigma)$ in the Lagrangian function.

• Lagrangian function:

$$\mathcal{L}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \lambda) = -\ell(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \lambda(1 - \sum_{k=1}^{K} \pi_k).$$

• Since there is only one equality constraint, there are only stationary and primal feasibility constraints in KKT conditions. Specifically, $\forall k$, we have

$$\frac{\partial \mathcal{L}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\lambda})}{\partial \boldsymbol{\mu}_{k}} = \frac{-\partial \sum_{n=1}^{N} \mathbb{1}_{[z^{(n)}=k]} \ln p(\mathbf{x}^{(n)}; \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})}{\partial \boldsymbol{\mu}_{k}} = \mathbf{0}, \quad (1)$$

$$\frac{\partial \mathcal{L}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\lambda})}{\partial \boldsymbol{\Sigma}_{k}} = \frac{-\partial \sum_{n=1}^{N} \mathbb{1}_{[z^{(n)}=k]} \ln p(\mathbf{x}^{(n)}; \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})}{\partial \boldsymbol{\Sigma}_{k}} = \mathbf{0}, \quad (2)$$

$$\frac{\partial \mathcal{L}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\lambda})}{\partial \boldsymbol{\pi}_{k}} = \frac{\partial \sum_{n=1}^{N} \mathbb{1}_{[z^{(n)}=k]} \ln \boldsymbol{\pi}_{k}}{\partial \boldsymbol{\pi}_{k}} - \boldsymbol{\lambda} = 0, \quad (3)$$

$$\mathbb{1} - \sum_{k=1}^{K} \boldsymbol{\pi}_{k} = 0. \quad (4)$$

• The solution to μ_k :

37

Ν

$$\frac{\partial \mathcal{L}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \lambda)}{\partial \boldsymbol{\mu}_{k}} = \frac{-\partial \left[\sum_{n=1}^{N} \mathbb{1}_{[z^{(n)}=k]} \ln p(\mathbf{x}^{(n)}; \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})\right]}{\partial \boldsymbol{\mu}_{k}} = \mathbf{0}, \quad (5)$$
$$\Rightarrow \frac{\partial \left[\sum_{n=1}^{N} \mathbb{1}_{[z^{(n)}=k]} (-\frac{1}{2}) (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})^{\top} \boldsymbol{\Sigma}_{k}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})\right]}{\partial \boldsymbol{\mu}_{k}} = \mathbf{0}, \quad (6)$$

$$\Rightarrow \sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \boldsymbol{\Sigma}_{k}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k}) = \mathbf{0},$$
(7)

$$\Rightarrow \sum_{n=1}^{\infty} \mathbf{1}_{\left[z^{(n)}=k\right]}(\mathbf{x}^{(n)}-\boldsymbol{\mu}_k) = \mathbf{0},\tag{8}$$

$$\Rightarrow \boldsymbol{\mu}_{k} = \frac{\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \mathbf{x}^{(n)}}{\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]}}.$$
(9)

• The solution to Σ_k :

$$\frac{\partial \mathcal{L}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \lambda)}{\partial \boldsymbol{\Sigma}_{k}} = \frac{-\partial \left[\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \ln p(\mathbf{x}^{(n)}; \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})\right]}{\partial \boldsymbol{\Sigma}_{k}} = \mathbf{0}, \quad (10)$$
$$\Rightarrow \frac{\partial \left[\frac{1}{2} \sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \left(\ln |\boldsymbol{\Sigma}_{k}| + (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})^{\top} \boldsymbol{\Sigma}_{k}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})\right)\right]}{\partial \boldsymbol{\Sigma}_{k}} = \mathbf{0}. \quad (11)$$

• We define $\Lambda_k = \Sigma_k^{-1}$, which is called precision matrix. Then, the above equation is equivalent to:

$$\frac{\partial \left[\frac{1}{2}\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \left(\ln |\mathbf{\Lambda}_{k}^{-1}| + (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})^{\top} \mathbf{\Lambda}_{k} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})\right)\right]}{\partial \mathbf{\Lambda}_{k}} = \mathbf{0}, \quad (12)$$

$$= \frac{\partial \left[\frac{1}{2}\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \left(-\ln |\mathbf{\Lambda}_{k}| + (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})^{\top} \mathbf{\Lambda}_{k} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})\right)\right]}{\partial \mathbf{\Lambda}_{k}} = \mathbf{0}, \quad (13)$$

where we utilize that $|\mathbf{\Lambda}_k^{-1}| = \frac{1}{|\mathbf{\Lambda}_k|}$.

The solution to Σ_k :

• Its solution is derived as follows:

$$\frac{\partial \left[\frac{1}{2}\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \left(-\ln |\mathbf{\Lambda}_{k}| + (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})^{\top} \mathbf{\Lambda}_{k} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})\right)\right]}{\partial \mathbf{\Lambda}_{k}} = \mathbf{0}, \quad (14)$$

$$\Rightarrow \frac{1}{2}\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \left(-\mathbf{\Lambda}_{k}^{-1} + (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k}) (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})^{\top}\right) = \mathbf{0}, \quad (15)$$

$$\Rightarrow \mathbf{\Lambda}_{k}^{-1} = \mathbf{\Sigma}_{k} = \frac{\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k}) (\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k})^{\top}}{\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]}}, \quad (16)$$

where we utilize that $\frac{d \ln |\Lambda_k|}{d\Lambda_k} = \Lambda_k^{-1}$. Reference of matrix derivatives: https://en.wikipedia.org/wiki/Matrix_calculus

• The solution to π :

$$\frac{\partial \mathcal{L}(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \lambda)}{\partial \pi_k} = \frac{\sum_{n=1}^N \mathbf{1}_{[z^{(n)}=k]} \ln \pi_k}{\partial \pi_k} - \lambda = 0, \quad (17)$$
$$\Rightarrow \sum_{n=1}^N \frac{\mathbf{1}_{[z^{(n)}=k]}}{\pi_k} = \lambda \Rightarrow \sum_{n=1}^N \mathbf{1}_{[z^{(n)}=k]} = \pi_k \lambda. \quad (18)$$

• Combining with $1 - \sum_{k=1}^{K} \pi_k = 0$, we can obtain $\lambda = N$, which is replaced back to obtain

$$\pi_k = \frac{1}{N} \sum_{n=1}^N \mathbf{1}_{[z^{(n)}=k]}.$$

Solution summary:

• If we knew $z^{(n)}$ for every $\mathbf{x}^{(n)}$, the maximum log likelihood problem is easy:

$$\max \ell(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln p\left(\mathbf{x}^{(n)}, z^{(n)} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\right),$$

with the constraint $1 - \sum_{k=1}^{K} \pi_k = 0.$

• Its solution is

$$\mu_{k} = \frac{\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \mathbf{x}^{(n)}}{\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]}}$$
$$\boldsymbol{\Sigma}_{k} = \frac{\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]} \left(\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k} \right) \left(\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k} \right)^{T}}{\sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]}}$$
$$\pi_{k} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{1}_{[z^{(n)}=k]}$$

However, how to learn the parameters when we don't know $z^{(n)}$ for every $\mathbf{x}^{(n)}$? We have to resort to the expectation maximization algorithm. D Gaussian Mixture Model for Density Estimation

2 The Latent Variable Perspective for Gaussian Mixture Model

(3) Expectation Maximization for Fitting Gaussian Mixture Model

4 Relation to k-Means

Intuitively, How Can We Fit a Mixture of Gaussians?

Optimization uses the **Expectation Maximization algorithm**, which alternates between two steps:

- E-step: Compute the posterior probability over z given the current model, *i.e.*, $p(z|\mathbf{x}; \Theta)$, which tells how much do we think each Gaussian generates each data point.
- M-step: Assuming that the data was really generated this way, update the parameters of each Gaussian component to maximize the probability that it would generate the data it is currently responsible for.



Elegant and powerful method for finding maximum likelihood solutions for models with latent variables

- E-step:
 - In order to adjust the parameters, we must first solve the inference problem: which Gaussian component generated each datapoint?
 - We cannot ensure, so it's a distribution over all possibilities.

$$\gamma_k^{(n)} = p\left(z^{(n)} = k \mid \mathbf{x}^{(n)}; \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\right).$$

- M-step:
 - Each Gaussian gets a certain amount of posterior probability for each datapoint.
 - We fit each Gaussian to the weighted datapoints.
 - We can derive closed form updates for all parameters

Lets see how EM works on GMM:

Conditional probability (using Bayes rule) of z given \mathbf{x}

$$\gamma_k = p(z = k \mid \mathbf{x}) = \frac{p(z = k)p(\mathbf{x} \mid z = k)}{p(\mathbf{x})}$$
$$= \frac{p(z = k)p(\mathbf{x} \mid z = k)}{\sum_{j=1}^{K} p(z = j)p(\mathbf{x} \mid z = j)}$$
$$= \frac{\pi_k \mathcal{N} \left(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\right)}{\sum_{j=1}^{K} \pi_j \mathcal{N} \left(\mathbf{x} \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\right)}.$$

 γ_k can be viewed as the **responsibility** of cluster k towards **x**

GMM E-Step: expected log likelihood

Once we computed $\gamma_k^{(n)} = p(z^{(n)} = k | \mathbf{x}^{(n)})$, we can compute the expected log likelihood, as follows:

$$\sum_{n} \mathbb{E}_{P(z^{(n)}|\mathbf{x}^{(n)})} \left[\ln \left(P(\mathbf{x}^{(n)}, z^{(n)} \mid \Theta) \right) \right]$$

= $\sum_{n} \sum_{k} \gamma_{k}^{(n)} \left(\ln \left(P(z^{(n)} = k \mid \Theta) \right) + \ln \left(P(\mathbf{x}^{(n)} \mid z^{(n)} = k, \Theta) \right) \right)$
= $\sum_{n} \sum_{k} \gamma_{k}^{(n)} \left(\ln (\pi_{k}) + \ln \left(\mathcal{N}(\mathbf{x}^{(n)} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right) \right)$
= $\sum_{n} \sum_{k} \gamma_{k}^{(n)} \ln (\pi_{k}) + \sum_{n} \sum_{k} \gamma_{k}^{(n)} \ln \left(\mathcal{N}(\mathbf{x}^{(n)} \mid \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right),$

where $\Theta = \{\mu, \Sigma, \pi\}$. Note that the above expectation is fully decomposed to each data *n* and each cluster *k*, which will facilitate the parameter learning in the following maximization step.

GMM M-Step

Maximization step:

• Given the posterior probability $\gamma_k^{(n)} = p\left(z^{(n)} = k \mid \mathbf{x}^{(n)}\right)$, we want to update the model parameters $\Theta = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}\}$ by maximizing the expected log likelihood, *i.e.*,

$$\max_{\Theta} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{k}^{(n)} \ln\left(\pi_{k}\right) + \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{k}^{(n)} \ln\left(\mathcal{N}(\mathbf{x}^{(n)}|\boldsymbol{\mu}_{k},\boldsymbol{\Sigma}_{k})\right), \text{ s.t. } \sum_{k=1}^{K} \pi_{k} = 1.$$

• Following the derivations introduced in previous slides (see page 12-17), it is easy to obtain the following solutions:

$$\boldsymbol{\mu}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{k}^{(n)} \mathbf{x}^{(n)}$$
$$\boldsymbol{\Sigma}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{k}^{(n)} \left(\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k} \right) \left(\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k} \right)^{\top}$$
$$\pi_{k} = \frac{N_{k}}{N}, \text{ with } N_{k} = \sum_{n=1}^{N} \gamma_{k}^{(n)}$$

EM for fitting GMM

Summary of EM: Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k

Iterate until convergence:

• E-step: Evaluate the responsibilities given current parameters

$$\gamma_k^{(n)} = p\left(z^{(n)} = k \mid \mathbf{x}^{(n)}; \Theta\right) = \frac{\pi_k \mathcal{N}\left(\mathbf{x}^{(n)} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\right)}{\sum_{j=1}^K \pi_j \mathcal{N}\left(\mathbf{x}^{(n)} \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\right)}.$$

• M-step: Re-estimate the parameters given current responsibilities

$$\boldsymbol{\mu}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{k}^{(n)} \mathbf{x}^{(n)},$$
$$\boldsymbol{\Sigma}_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma_{k}^{(n)} \left(\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k} \right) \left(\mathbf{x}^{(n)} - \boldsymbol{\mu}_{k} \right)^{\top},$$

$$\pi_k = \frac{N_k}{N}$$
, with $N_k = \sum_{n=1}^N \gamma_k^{(n)}$.

• Evaluate log likelihood and check for convergence

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

EM for fitting GMM



D Gaussian Mixture Model for Density Estimation

2 The Latent Variable Perspective for Gaussian Mixture Model

3 Expectation Maximization for Fitting Gaussian Mixture Model

4 Relation to k-Means

The K-Means Algorithm:

- Assignment step: Assign each data point to the closest cluster, *i.e.*, hard assignment
- **Refitting step**: Move each cluster center to the center of gravity of the data assigned to it

The EM Algorithm:

- **E-step**: Given the current model, compute the posterior probability over z for each data point, like soft assignment
- **M-step**: Given the posterior probability, update the model parameters by maximizing the expected log-likelihood

- If fixing the covariance matrices Σ as the identity matrix I for all Gaussian components, EM for GMMs is reduced to a soft version of K-means
- Instead of hard assignments in the E-step, EM does **soft assignments** based on the softmax of the squared Euclidean distance from each point to each cluster.
- Each center moved by **weighted means** of the data, with weights given by soft assignments
- In K-means, weights are 0 or 1
- GMM provides a probabilistic view of clustering Each cluster corresponds to a different Gaussian component.

Where does EM come from?

The final issue:

• Let's recall the original objective function (*i.e.*, log-likelihood) of fitting GMM, as follows

$$\ell(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \ln p(\mathbf{X} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \sum_{z^{(n)}=1}^{K} p\left(\mathbf{x}^{(n)}, z^{(n)} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$$

• However, in the EM algorithm introduced above, this objective function never occurs. Instead, we maximize the following expected log-likelihood in the M-step,

$$\sum_{n=1}^{N} \mathbb{E}_{p(z^{(n)}|\mathbf{x}^{(n)})} \bigg[\ln \left(p(\mathbf{x}^{(n)}, z^{(n)} \mid \Theta) \right) \bigg].$$

- So, why EM is a good algorithm for fitting GMM? What is relationship between above two objective functions?
- In the next lecture, we will provide a principled justification of the EM algorithm to answer these questions.

Examples of Clustering Results



Reference

Further readings:

- Chapter 9 in the book "Pattern Recognition and Machine Learning". Link
- Demo with code: https://scikit-learn.org/stable/modules/generated sklearn.mixture.GaussianMixture.html#sklearn.mixture.GaussianMix