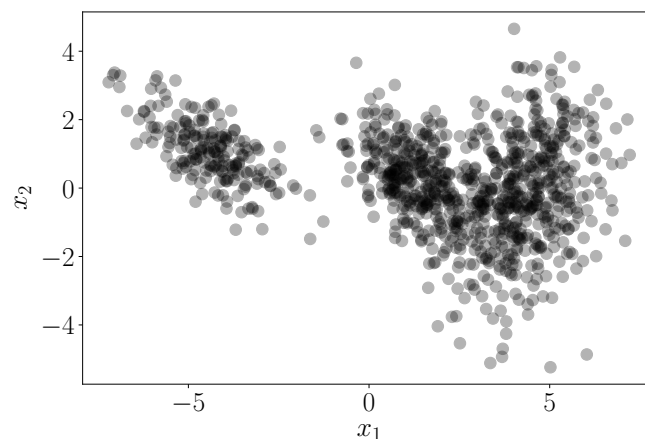


Density Estimation with Gaussian Mixture Models

In earlier chapters, we covered already two fundamental problems in machine learning: regression (Chapter 9) and dimensionality reduction (Chapter 10). In this chapter, we will have a look at a third pillar of machine learning: density estimation. On our journey, we introduce important concepts, such as the Expectation Maximization (EM) algorithm and a latent variable perspective of density estimation with mixture models.

When we apply machine learning to data we often aim to represent data in some way. A straightforward way is to take the data points themselves as the representation of the data, see Figure 11.1 for an example. However, this approach may be unhelpful if the dataset is huge or if we are interested in representing characteristics of the data. In density estimation, we represent the data compactly using a density from a parametric family, e.g., a Gaussian or Beta distribution. For example, we may be looking for the mean and variance of a dataset in order to represent the data compactly using a Gaussian distribution. The mean and variance can be found using tools we discussed in Section 8.2: maximum likelihood or maximum-a-posteriori estimation. We can then use the mean and variance of this Gaussian to represent the distribution underlying the data, i.e., we think of the dataset to be a typical realization from this distribution if we were to sample from it.

Figure 11.1
Two-dimensional dataset that cannot be meaningfully represented by a Gaussian.



In practice, the Gaussian (or similarly all other distributions we encountered so far) have limited modeling capabilities. For example, a Gaussian approximation of the density that generated the data in Figure 11.1 would be a poor approximation. In the following, we will look at a more expressive family of distributions, which we can use for density estimation: *mixture models*.

mixture models

Mixture models can be used to describe a distribution $p(\mathbf{x})$ by a convex combination of K simple (base) distributions

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}) \quad (11.1)$$

$$0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1, \quad (11.2)$$

where the components p_k are members of a family of basic distributions, e.g., Gaussians, Bernoullis or Gammas, and the π_k are *mixture weights*. Mixture models are more expressive than the corresponding base distributions because they allow for multimodal data representations, i.e., they can describe datasets with multiple “clusters”, such as the example in Figure 11.1.

mixture weights

We will focus on Gaussian mixture models (GMMs), where the basic distributions are Gaussians. For a given dataset, we aim to maximize the likelihood of the model parameters to train the GMM. For this purpose we will use results from Chapter 5, Section 7.2 and Chapter 6. However, unlike other application we discussed earlier (linear regression or PCA), we will not find a closed-form maximum likelihood solution. Instead, we will arrive at a set of dependent simultaneous equations, which we can only solve iteratively.

11.1 Gaussian Mixture Model

A *Gaussian mixture model* is a density model where we combine a finite number of K Gaussian distributions $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ so that

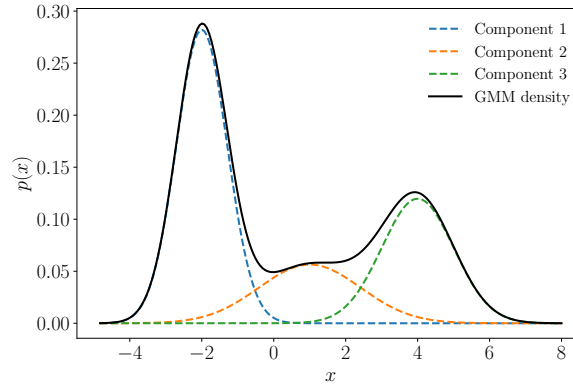
Gaussian mixture model

$$p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (11.3)$$

$$0 \leq \pi_k \leq 1, \quad \sum_{k=1}^K \pi_k = 1, \quad (11.4)$$

where we defined $\boldsymbol{\theta} := \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k : k = 1, \dots, K\}$ as the collection of all parameters of the model. This convex combination of Gaussian distribution gives us significantly more flexibility for modeling complex densities than a simple Gaussian distribution (which we recover from (11.3) for $K = 1$). An illustration is given in Figure 11.2, displaying the weighted

Figure 11.2
Gaussian mixture model. The Gaussian mixture distribution (black) is composed of a convex combination of Gaussian distributions and is more expressive than any individual component. Dashed lines represent the weighted Gaussian components.



components and the mixture density, which is given as

$$p(x | \theta) = 0.5\mathcal{N}(x | -2, \frac{1}{2}) + 0.2\mathcal{N}(x | 1, 2) + 0.3\mathcal{N}(x | 4, 1). \quad (11.5)$$

6307

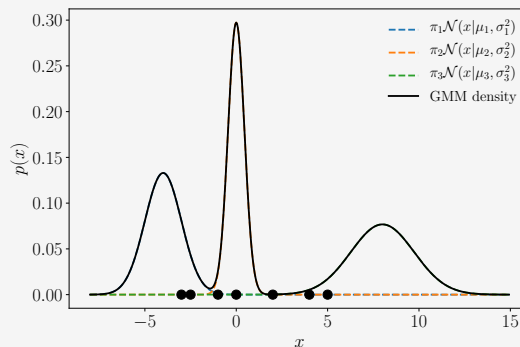
6308

11.2 Parameter Learning via Maximum Likelihood

6309 Assume we are given a dataset $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ where \mathbf{x}_n , $n = 1, \dots, N$
 6310 are drawn i.i.d. from an unknown distribution $p(\mathbf{x})$. Our objective is to
 6311 find a good approximation/representation of this unknown distribution
 6312 $p(\mathbf{x})$ by means of a Gaussian mixture model (GMM) with K mixture com-
 6313 ponents. The parameters of the GMM are the K means μ_k , the covariances
 6314 Σ_k and mixture weights π_k . We summarize all these free parameters in
 6315 $\theta := \{\pi_k, \mu_k, \Sigma_k : k = 1, \dots, K\}$.

Example 11.1 (Initial setting)

Figure 11.3 Initial setting: GMM (black) with mixture three mixture components (dashed) and seven data points (discs).



Throughout this chapter, we will have a simple running example that helps us illustrate and visualize important concepts.

We consider a one-dimensional dataset $\mathcal{X} = \{-3, -2.5, -1, 0, 2, 4, 5\}$ consisting of seven data points and wish to find a GMM with $K = 3$ components that models the density of the data. We initialize the mixture components as

$$p_1(x) = \mathcal{N}(x \mid -4, 1) \quad (11.6)$$

$$p_2(x) = \mathcal{N}(x \mid 0, 0.2) \quad (11.7)$$

$$p_3(x) = \mathcal{N}(x \mid 8, 3) \quad (11.8)$$

and assign them equal weights $\pi_1 = \pi_2 = \pi_3 = \frac{1}{3}$. The corresponding model (and the data points) are shown in Figure 11.3.

In the following, we detail how to obtain a maximum likelihood estimate θ_{ML} of the model parameters θ . We start by writing down the likelihood, i.e., the predictive distribution of the training data given the parameters. We exploit our i.i.d. assumption, which leads to the factorized likelihood

$$p(\mathcal{X} \mid \theta) = \prod_{n=1}^N p(\mathbf{x}_n \mid \theta), \quad p(\mathbf{x}_n \mid \theta) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (11.9)$$

where every individual likelihood term $p(\mathbf{x}_n \mid \theta)$ is a Gaussian mixture density. Then, we obtain the log-likelihood as

$$\log p(\mathcal{X} \mid \theta) = \sum_{n=1}^N \log p(\mathbf{x}_n \mid \theta) = \underbrace{\sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{=: \mathcal{L}}. \quad (11.10)$$

We aim to find parameters θ_{ML}^* that maximize the log-likelihood \mathcal{L} defined in (11.10). Our “normal” procedure would be to compute the gradient $d\mathcal{L}/d\theta$ of the log-likelihood with respect to the model parameters θ , set it to $\mathbf{0}$ and solve for θ . However, unlike our previous examples for maximum likelihood estimation (e.g., when we discussed linear regression in Section 9.2), we cannot obtain a closed-form solution. However, we can exploit an iterative scheme to find good model parameters θ_{ML} , which will out to be the EM algorithm for Gaussian mixture models. The key idea is to update one model parameter at a time while keeping the others fixed.

Remark. If we were to consider a single Gaussian as the desired density, the sum over k in (11.10) vanishes, and the log can be applied directly to the Gaussian component, such that we get

$$\log \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log \det(\boldsymbol{\Sigma}) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}). \quad (11.11)$$

This simple form allows us find closed-form maximum likelihood estimates of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, as discussed in Chapter 8. In (11.10), we cannot move

the log into the sum over k so that we cannot obtain a simple closed-form maximum likelihood solution. \diamond

Any local optimum of a function exhibits the property that its gradient with respect to the parameters must vanish (necessary condition), see Chapter 7. In our case, we obtain the following necessary conditions when we optimize the log-likelihood in (11.10) with respect to the GMM parameters $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$:

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_k} = \mathbf{0}^\top \iff \sum_{n=1}^N \frac{\partial \log p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \mathbf{0}^\top, \quad (11.12)$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\Sigma}_k} = \mathbf{0} \iff \sum_{n=1}^N \frac{\partial \log p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_k} = \mathbf{0}, \quad (11.13)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = 0 \iff \sum_{n=1}^N \frac{\partial \log p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \pi_k} = 0. \quad (11.14)$$

For all three necessary conditions, by applying the chain rule (see Section 5.2.2), we require partial derivatives of the form

$$\frac{\partial \log p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{p(\mathbf{x}_n | \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \quad (11.15)$$

where $\boldsymbol{\theta} = \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k, k = 1, \dots, K\}$ are the model parameters and

$$\frac{1}{p(\mathbf{x}_n | \boldsymbol{\theta})} = \frac{1}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \quad (11.16)$$

In the following, we will compute the partial derivatives (11.12)–(11.14). But before we do this, we introduce a quantity that will play a central role in the remainder of this chapter: responsibilities.

11.2.1 Responsibilities

We define the quantity

$$r_{nk} := \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (11.17)$$

responsibility

as the *responsibility* of the k th mixture component for the n th data point. The responsibility r_{nk} of the k th mixture component for data point \mathbf{x}_n is proportional to the likelihood

$$p(\mathbf{x}_n | \pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (11.18)$$

of the mixture component given the data point. Therefore, mixture components have a high responsibility for a data point when the data point could be a plausible sample from that mixture component. Note that $\mathbf{r}_n := [r_{n1}, \dots, r_{nK}]^\top \in \mathbb{R}^K$ is a (normalized) probability vector, i.e.,

\mathbf{r}_n follows a Boltzmann/Gibbs distribution.

$\sum_k r_{nk} = 1$ with $r_{nk} \geq 0$. This probability vector distributes probability mass among the K mixture components, and we can think of \mathbf{r}_n as a “soft assignment” of \mathbf{x}_n to the K mixture components. Therefore, the responsibility r_{nk} from (11.17) represents the probability that \mathbf{x}_n has been generated by the k th mixture component.

The responsibility r_{nk} is the probability that the k th mixture component generated the n th data point.

Example 11.2 (Responsibilities)

For our example from Figure 11.3 we compute the responsibilities r_{nk}

$$\begin{bmatrix} 1.0 & 0.0 & 0.0 \\ 1.0 & 0.0 & 0.0 \\ 0.057 & 0.943 & 0.0 \\ 0.001 & 0.999 & 0.0 \\ 0.0 & 0.066 & 0.934 \\ 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 \end{bmatrix} \in \mathbb{R}^{N \times K}. \quad (11.19)$$

Here, the n th row tells us the responsibilities of all mixture components for \mathbf{x}_n . The sum of all K responsibilities for a data point (sum of every row) is 1. The k th column gives us an overview of the responsibility of the k th mixture component. We can see that the third mixture component (third column) is not responsible for any of the first four data points, but takes much responsibility of the remaining data points. The sum of all entries of a column gives us the values N_k , i.e., the total responsibility of the k th mixture component. In our example, we get $N_1 = 2.058$, $N_2 = 2.008$, $N_3 = 2.934$.

In the following, we determine the updates of the model parameters $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$ for given responsibilities. We will see that the update equations all depend on the responsibilities, which makes a closed-form solution to the maximum likelihood estimation problem impossible. However, for given responsibilities we will be updating one model parameter at a time, while keeping the others fixed. After this, we will re-compute the responsibilities. Iterating these two steps will eventually converge to a local optimum and is a specific instantiation of the EM algorithm. We will discuss this in some more detail in Section 11.3.

11.2.2 Updating the Means

Theorem 11.1 (Update of the GMM Means). *The update of the mean parameters $\boldsymbol{\mu}_k$, $k = 1, \dots, K$, of the GMM is given by*

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{\sum_{n=1}^N r_{nk} \mathbf{x}_n}{\sum_{n=1}^N r_{nk}}, \quad (11.20)$$

where the responsibilities r_{nk} are defined in (11.17).

Remark. The update of the means μ_k of the individual mixture components in (11.20) depends on all means, covariance matrices Σ_k and mixture weights π_k via r_{nk} given in (11.17). Therefore, we cannot obtain a closed-form solution for all μ_k at once. \diamond

Proof From (11.15), we see that the gradient of the log-likelihood with respect to the mean parameters μ_k , $k = 1, \dots, K$, requires us to compute the partial derivative

$$\frac{\partial p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \mu_k} = \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}{\partial \mu_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\partial \mu_k} \quad (11.21)$$

$$= \pi_k (\mathbf{x}_n - \mu_k)^\top \Sigma_k^{-1} \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \quad (11.22)$$

where we exploited that only the k th mixture component depends on μ_k .

We use our result from (11.22) in (11.15) and put everything together so that the desired partial derivative of \mathcal{L} with respect to μ_k is given as

$$\frac{\partial \mathcal{L}}{\partial \mu_k} = \sum_{n=1}^N \frac{\partial \log p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \mu_k} = \sum_{n=1}^N \frac{1}{p(\mathbf{x}_n | \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \mu_k} \quad (11.23a)$$

$$= \sum_{n=1}^N (\mathbf{x}_n - \mu_k)^\top \Sigma_k^{-1} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}}_{=r_{nk}} \quad (11.23b)$$

$$= \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \mu_k)^\top \Sigma_k^{-1}. \quad (11.23c)$$

Here, we used the identity from (11.16) and the result of the partial derivative in (11.22) to get to (11.23b). The values r_{nk} are the responsibilities we defined in (11.17).

We now solve (11.23c) for μ_k^{new} so that $\frac{\partial \mathcal{L}(\mu_k^{\text{new}})}{\partial \mu_k} = \mathbf{0}^\top$ and obtain

$$\sum_{n=1}^N r_{nk} \mathbf{x}_n = \sum_{n=1}^N r_{nk} \mu_k^{\text{new}} \iff \mu_k^{\text{new}} = \frac{\sum_{n=1}^N r_{nk} \mathbf{x}_n}{\sum_{n=1}^N r_{nk}} = \frac{1}{N_k} \sum_{n=1}^N r_{nk} \mathbf{x}_n, \quad (11.24)$$

where we defined

$$N_k := \sum_{n=1}^N r_{nk} \quad (11.25)$$

as the total responsibility of the k th mixture component for the entire dataset. This concludes the proof of Theorem 11.1. \square

Intuitively, (11.20) can be interpreted as an importance-weighted Monte-Carlo estimate of the mean, where the importance weights of data point \mathbf{x}_n are the responsibilities r_{nk} of the k th cluster for \mathbf{x}_n , $k = 1, \dots, K$.

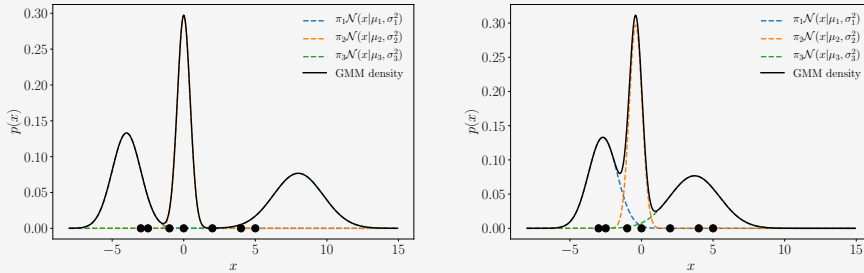
Therefore, the mean μ_k is pulled toward a data point x_n with strength given by r_{nk} . The means are pulled stronger toward data points for which the corresponding mixture component has a high responsibility, i.e., a high likelihood. Figure 11.4 illustrates this. We can also interpret the mean update in (11.20) as the expected value of all data points under the distribution given by

$$\mathbf{r}_k := [r_{1k}, \dots, r_{Nk}]^\top / N_k, \quad (11.26)$$

which is a normalized probability vector, i.e.,

$$\mu_k \leftarrow \mathbb{E}_{\mathbf{r}_k}[\mathcal{X}]. \quad (11.27)$$

Example 11.3 (Mean Updates)



(a) GMM density and individual components prior to updating the mean values.

(b) GMM density and individual components after updating the mean values.

In our example from Figure 11.3, the mean values are updated as follows:

$$\mu_1 : -4 \rightarrow -2.7 \quad (11.28)$$

$$\mu_2 : 0 \rightarrow -0.4 \quad (11.29)$$

$$\mu_3 : 8 \rightarrow 3.7 \quad (11.30)$$

Here, we see that the means of the first and third mixture component move toward the regime of the data, whereas the mean of the second component does not change so dramatically. Figure 11.5 illustrates this change, where Figure 11.5(a) shows the GMM density prior to updating the means and Figure 11.5(b) shows the GMM density after updating the mean values μ_k .

Figure 11.4 Update of the mean parameter of mixture component in a GMM. The mean μ is being pulled toward individual data points with the weights given by the corresponding responsibilities.

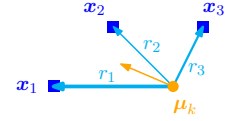


Figure 11.5 Effect of updating the mean values in a GMM. (a) GMM before updating the mean values; (b) GMM after updating the mean values μ_k while retaining the variances and mixture weights.

The update of the mean parameters in (11.20) look fairly straightforward. However, note that the responsibilities r_{nk} are a function of π_j, μ_j, Σ_j for all $j = 1, \dots, K$, such that the updates in (11.20) depend on all parameters of the GMM, and a closed-form solution, which we obtained for linear regression in Section 9.2 or PCA in Chapter 10, cannot be obtained.

11.2.3 Updating the Covariances

Theorem 11.2 (Updates of the GMM Covariances). *The update of the covariance parameters Σ_k , $k = 1, \dots, K$ of the GMM is given by*

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top, \quad (11.31)$$

where r_{nk} and N_k are defined in (11.17) and (11.25), respectively.

Proof To prove Theorem 11.2 our approach is to compute the partial derivatives of the log-likelihood \mathcal{L} with respect to the covariances Σ_k , set them to $\mathbf{0}$ and solve for Σ_k . We start with our general approach

$$\frac{\partial \mathcal{L}}{\partial \Sigma_k} = \sum_{n=1}^N \frac{\partial \log p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \Sigma_k} = \sum_{n=1}^N \frac{1}{p(\mathbf{x}_n | \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \Sigma_k}. \quad (11.32)$$

We already know $1/p(\mathbf{x}_n | \boldsymbol{\theta})$ from (11.16). To obtain the remaining partial derivative $\partial p(\mathbf{x}_n | \boldsymbol{\theta}) / \partial \Sigma_k$, we write down the definition of the Gaussian distribution $p(\mathbf{x}_n | \boldsymbol{\theta})$, see (11.9), and drop all terms but the k th. We then obtain

$$\frac{\partial p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \Sigma_k} \quad (11.33a)$$

$$= \frac{\partial}{\partial \Sigma_k} \left((2\pi)^{-\frac{D}{2}} \det(\Sigma_k)^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) \right) \quad (11.33b)$$

$$= \pi_k (2\pi)^{-\frac{D}{2}} \left[\frac{\partial}{\partial \Sigma_k} \det(\Sigma_k)^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) + \det(\Sigma_k)^{-\frac{1}{2}} \frac{\partial}{\partial \Sigma_k} \exp \left(-\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) \right]. \quad (11.33c)$$

We now use the identities

$$\frac{\partial}{\partial \Sigma_k} \det(\Sigma_k)^{-\frac{1}{2}} \stackrel{(5.100)}{=} -\frac{1}{2} \det(\Sigma_k)^{-\frac{1}{2}} \Sigma_k^{-1}, \quad (11.34)$$

$$\frac{\partial}{\partial \Sigma_k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \stackrel{(5.105)}{=} -\Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} \quad (11.35)$$

and obtain (after some re-arranging) the desired partial derivative required in (11.32) as

$$\begin{aligned} \frac{\partial p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \Sigma_k} &= \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k) \\ &\quad \cdot \left[-\frac{1}{2} (\Sigma_k^{-1} - \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1}) \right]. \end{aligned} \quad (11.36)$$

Putting everything together, the partial derivative of the log-likelihood

with respect to Σ_k is given by

$$\frac{\partial \mathcal{L}}{\partial \Sigma_k} = \sum_{n=1}^N \frac{\partial \log p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \Sigma_k} = \sum_{n=1}^N \frac{1}{p(\mathbf{x}_n | \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_n | \boldsymbol{\theta})}{\partial \Sigma_k} \quad (11.37a)$$

$$= \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Sigma_k)}{\underbrace{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma_j)}_{=r_{nk}}} \cdot \left[-\frac{1}{2} (\Sigma_k^{-1} - \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1}) \right] \quad (11.37b)$$

$$= -\frac{1}{2} \sum_{n=1}^N r_{nk} (\Sigma_k^{-1} - \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1}) \quad (11.37c)$$

$$= -\frac{1}{2} \Sigma_k^{-1} \underbrace{\sum_{n=1}^N r_{nk}}_{=N_k} + \frac{1}{2} \Sigma_k^{-1} \left(\sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \right) \Sigma_k^{-1}. \quad (11.37d)$$

We see that the responsibilities r_{nk} also appear in this partial derivative. Setting this partial derivative to $\mathbf{0}$, we obtain the necessary optimality condition

$$N_k \Sigma_k^{-1} = \Sigma_k^{-1} \left(\sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \right) \Sigma_k^{-1} \quad (11.38a)$$

$$\iff N_k \mathbf{I} = \left(\sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \right) \Sigma_k^{-1} \quad (11.38b)$$

By solving for Σ_k , we obtain

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top, \quad (11.39)$$

where \mathbf{r}_k is the probability vector defined in (11.26). This gives us a simple update rule for Σ_k for $k = 1, \dots, K$ and proves Theorem 11.2. \square

Similar to the update of $\boldsymbol{\mu}_k$ in (11.20), we can interpret the update of the covariance in (11.31) as an importance-weighted expected value of the square of the centered data $\tilde{\mathcal{X}}_k := \{\mathbf{x}_1 - \boldsymbol{\mu}_k, \dots, \mathbf{x}_N - \boldsymbol{\mu}_k\}$.

Example 11.4 (Variance Updates)

In our example from Figure 11.3, the variances are updated as follows:

$$\sigma_1^2 : 1 \rightarrow 0.14 \quad (11.40)$$

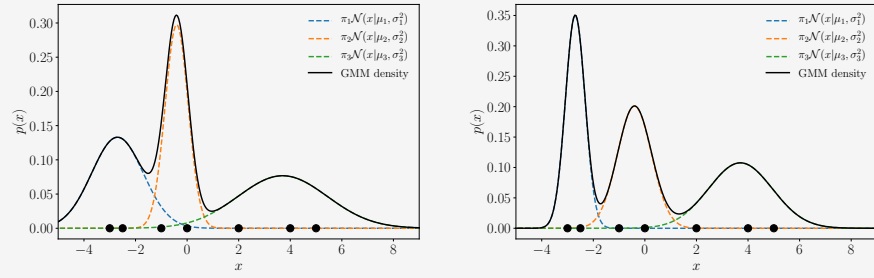
$$\sigma_2^2 : 0.2 \rightarrow 0.44 \quad (11.41)$$

$$\sigma_3^2 : 3 \rightarrow 1.53 \quad (11.42)$$

Here, we see that the variances of the first and third component shrink significantly, the variance of the second component increases slightly.

Figure 11.6 illustrates this setting. Figure 11.6(a) is identical (but zoomed in) to Figure 11.5(b) and shows the GMM density and its individual components prior to updating the variances. Figure 11.6(b) shows the GMM density after updating the variances.

Figure 11.6 Effect of updating the variances in a GMM. (a) GMM before updating the variances; (b) GMM after updating the variances while retaining the means and mixture weights.



(a) GMM density and individual components prior to updating the variances.

(b) GMM density and individual components after updating the variances.

Similar to the update of the mean parameters, we can interpret (11.31) as a Monte-Carlo estimate of the weighted covariance of data points x_n associated with the k th mixture component, where the weights are the responsibilities r_{nk} . As with the updates of the mean parameters, this update depends on all $\pi_j, \mu_j, \Sigma_j, j = 1, \dots, K$, through the responsibilities r_{nk} , which prohibits a closed-form solution.

11.2.4 Updating the Mixture Weights

Theorem 11.3 (Update of the GMM Mixture Weights). *The mixture weights of the GMM are updated as*

$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad (11.43)$$

for $k = 1, \dots, K$, where N is the number of data points and N_k is defined in (11.25).

Proof To find the partial derivative of the log-likelihood with respect to the weight parameters $\pi_k, k = 1, \dots, K$, we take the constraint $\sum_k \pi_k = 1$ into account by using Lagrange multipliers (see Section 7.2). The Lagrangian is

$$\mathcal{L} = \mathcal{L} + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \quad (11.44a)$$

$$= \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right), \quad (11.44b)$$

where \mathcal{L} is the log-likelihood from (11.10) and the second term encodes for the equality constraint that all the mixture weights need to sum up to 1. We obtain the partial derivative with respect to π_k as

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \sum_{n=1}^N \frac{\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} + \lambda \quad (11.45a)$$

$$= \frac{1}{\pi_k} \underbrace{\sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{=N_k} + \lambda = \frac{N_k}{\pi_k} + \lambda, \quad (11.45b)$$

and the partial derivative with respect to the Lagrange multiplier λ as

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sum_{k=1}^K \pi_k - 1. \quad (11.46)$$

Setting both partial derivatives to 0 (necessary condition for optimum) yields the system of equations

$$\pi_k = -\frac{N_k}{\lambda}, \quad (11.47)$$

$$1 = \sum_{k=1}^K \pi_k. \quad (11.48)$$

Using (11.47) in (11.48) and solving for π_k , we obtain

$$\sum_{k=1}^K \pi_k = 1 \iff -\sum_{k=1}^K \frac{N_k}{\lambda} = 1 \iff -\frac{N}{\lambda} = 1 \iff \lambda = -N. \quad (11.49)$$

This allows us to substitute $-N$ for λ in (11.47) to obtain

$$\pi_k^{\text{new}} = \frac{N_k}{N}, \quad (11.50)$$

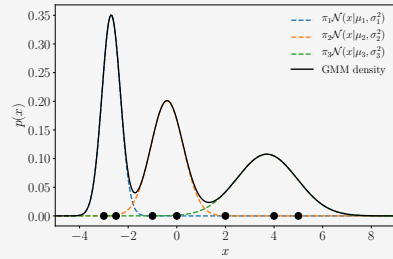
which gives us the update for the weight parameters π_k and proves Theorem 11.3. \square

We can identify the mixture weight in (11.43) as the ratio of the total responsibility of the k th cluster and the number of data points. Since $N = \sum_k N_k$ the number of data points can also be interpreted as the total responsibility of all mixture components together, such that π_k is the relative importance of the k th mixture component for the dataset.

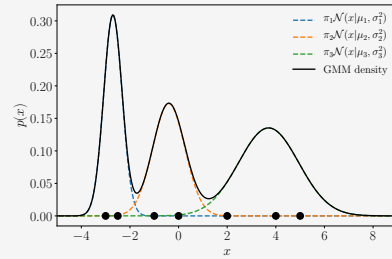
Remark. Since $N_k = \sum_{i=1}^N r_{nk}$, the update equation (11.43) for the mixture weights π_k also depends on all $\pi_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j, j = 1, \dots, K$ via the responsibilities r_{nk} . \diamond

Example 11.5 (Weight Parameter Updates)

Figure 11.7 Effect of updating the mixture weights in a GMM. (a) GMM before updating the mixture weights; (b) GMM after updating the mixture weights while retaining the means and variances. Note the different scalings of the vertical axes.



(a) GMM density and individual components prior to updating the mixture weights.



(b) GMM density and individual components after updating the mixture weights.

In our running example from Figure 11.3, the mixture weights are updated as follows:

$$\pi_1 : \frac{1}{3} \rightarrow 0.29 \quad (11.51)$$

$$\pi_2 : \frac{1}{3} \rightarrow 0.29 \quad (11.52)$$

$$\pi_3 : \frac{1}{3} \rightarrow 0.42 \quad (11.53)$$

Here we see that the third component gets more weight/importance, while the other components become slightly less important. Figure 11.7 illustrates the effect of updating the mixture weights. Figure 11.7(a) is identical to Figure 11.6(b) and shows the GMM density and its individual components prior to updating the mixture weights. Figure 11.7(b) shows the GMM density after updating the mixture weights.

Overall, having updated the means, the variances and the weights once, we obtain the GMM shown in Figure 11.7(b). Compared with the initialization shown in Figure 11.3, we can see that the parameter updates caused the GMM density to shift some of its mass toward the data points.

After updating the means, variances and weights once, the GMM fit in Figure 11.7(b) is already remarkably better than its initialization from Figure 11.3. This is also evidenced by the log-likelihood values, which increased from 28.3 (initialization) to 14.4 after one complete update cycle.

11.3 EM Algorithm

Unfortunately, the updates in (11.20), (11.31), and (11.43) do not constitute a closed-form solution for the updates of the parameters μ_k , Σ_k , π_k of the mixture model because the responsibilities r_{nk} depend on those parameters in a complex way. However, the results suggest a simple *iterative scheme* for finding a solution to the parameters estimation problem via maximum likelihood. The Expectation Maximization algorithm (*EM algo-*

EM algorithm

rithm) was proposed by Dempster et al. (1977) and is a general iterative scheme for learning parameters (maximum likelihood or MAP) in mixture models and, more generally, latent-variable models.

In our example of the Gaussian mixture model, we choose initial values for μ_k, Σ_k, π_k and alternate until convergence between

- *E-step*: Evaluate the responsibilities r_{nk} (posterior probability of data point n belonging to mixture component k).
- *M-step*: Use the updated responsibilities to re-estimate the parameters μ_k, Σ_k, π_k .

Every step in the EM algorithm increases the log-likelihood function (Neal and Hinton, 1999). For convergence, we can check the log-likelihood or the parameters directly. A concrete instantiation of the EM algorithm for estimating the parameters of a GMM is as follows:

- 1 Initialize μ_k, Σ_k, π_k
- 2 *E-step*: Evaluate responsibilities r_{nk} for every data point x_n using current parameters π_k, μ_k, Σ_k :

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_j \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)}. \quad (11.54)$$

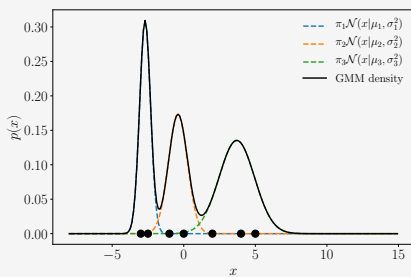
- 3 *M-step*: Re-estimate parameters π_k, μ_k, Σ_k using the current responsibilities r_{nk} (from E-step):

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n, \quad (11.55)$$

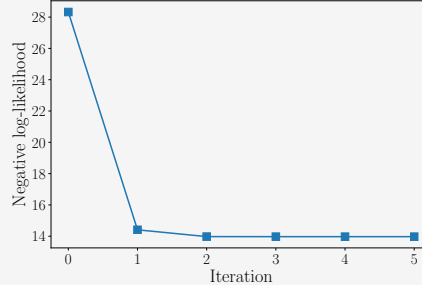
$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^\top, \quad (11.56)$$

$$\pi_k = \frac{N_k}{N}. \quad (11.57)$$

Example 11.6 (GMM Fit)



(a) Final GMM fit. After five iterations, the EM algorithm converges and returns this mixture density.

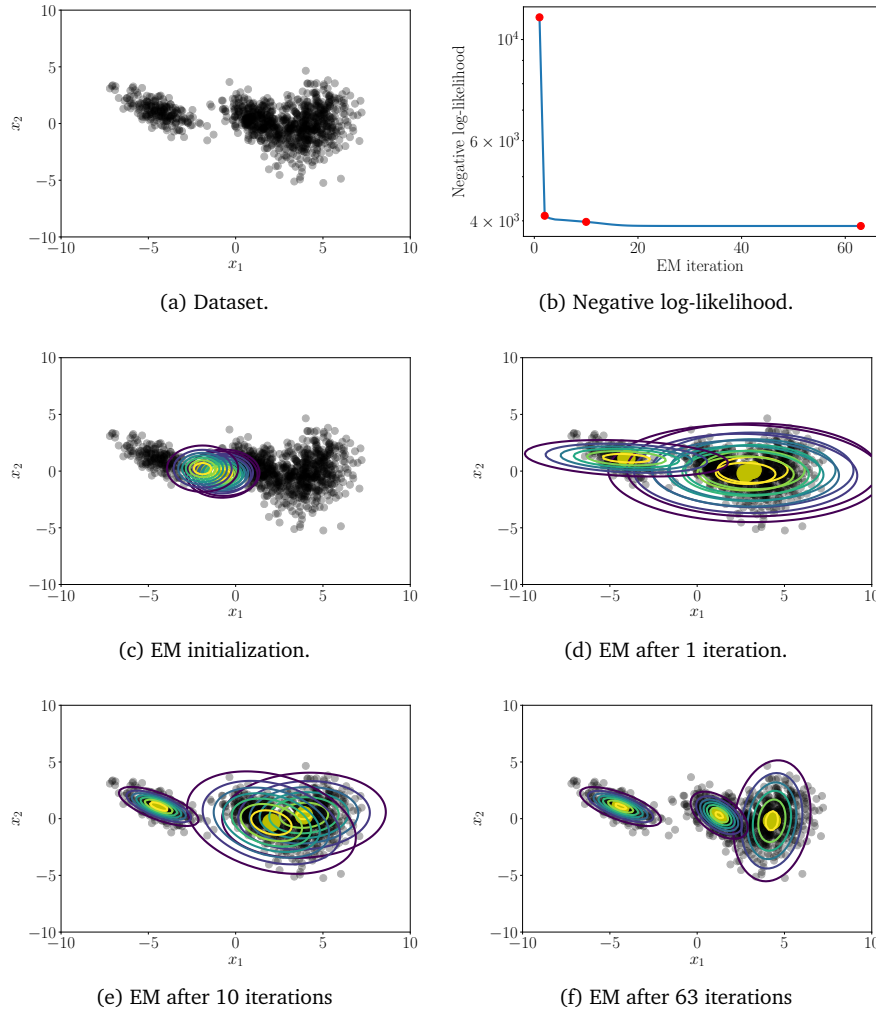


(b) Negative log-likelihood as a function of the EM iterations.

Figure 11.8 EM algorithm applied to the GMM from Figure 11.2.

Figure 11.9

Illustration of the EM algorithm for fitting a Gaussian mixture model with three components to a two-dimensional dataset. (a) Dataset; (b) Negative log-likelihood (lower is better) as a function of the EM iterations. The red dots indicate the iterations for which the corresponding GMM fits are shown in (c)–(f). The yellow discs indicate the mean of the Gaussian distribution.



When we run EM on our example from Figure 11.3, we obtain the final result shown in Figure 11.8(a) after five iterations, and Figure 11.8(b) shows how the negative log-likelihood evolves as a function of the EM iterations. The final GMM is given as

$$p(x) = 0.29\mathcal{N}(x \mid -2.75, 0.06) + 0.28\mathcal{N}(x \mid -0.50, 0.25) + 0.43\mathcal{N}(x \mid 3.64, 1.63). \quad (11.58)$$

We applied the EM algorithm to the two-dimensional dataset shown in Figure 11.1 with $K = 3$ mixture components. Figure 11.9 illustrates a some steps of the EM algorithm and the negative log-likelihood as a function of the EM iteration (Figure 11.9(b)).

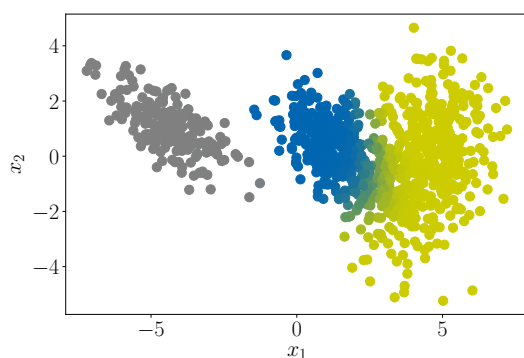


Figure 11.10
Dataset colored according to the responsibilities of the mixture components when EM converges. While a single mixture component is clearly responsible for the data on the left, the overlap of the two data clusters on the right could have been generated by two mixture components.

Figure 11.10 visualizes the final responsibilities of the mixture components for the data points. It becomes clear that there are data points that cannot be uniquely assigned to a single (either blue or yellow) component, such that the responsibilities of these two clusters for those points are around 0.5.

11.4 Latent Variable Perspective

We can look at the GMM from the perspective of a discrete latent variable model, i.e., where the latent variable z can attain only a finite set of values. This is in contrast to PCA where the latent variables were continuous-valued numbers in \mathbb{R}^M .

The advantages of the probabilistic perspective are that (i) it will justify some ad-hoc decisions we made in the previous sections, (ii) it allows for a concrete interpretation of the responsibilities as posterior probabilities, (iii) the iterative algorithm for updating the model parameters can be derived in a principled manner as the EM algorithm for maximum likelihood parameter estimation in latent variable models.

11.4.1 Generative Process and Probabilistic Model

To derive the probabilistic model for GMMs, it is useful to think about the generative process, i.e., the process that allows us to generate data, using a probabilistic model.

We assume a mixture model with K components and that a data point \mathbf{x} can be generated by exactly one mixture component. We introduce a binary indicator variable $z_k \in \{0, 1\}$ with two states (see Section 6.2) that indicates whether the k th mixture component generated that data point so that

$$p(\mathbf{x} | z_k = 1) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad (11.59)$$

We define $\mathbf{z} := [z_1, \dots, z_K]^\top \in \mathbb{R}^K$ as a probability vector consisting of

Sometimes this kind of probability distribution is called “multinoulli”, a generalization of the Bernoulli distribution to more than two values (Murphy, 2012). one-hot encoding 1-of- K representation

$K - 1$ many 0s and exactly one 1. For example, for $K = 3$, a valid \mathbf{z} would be $\mathbf{z} = [z_1, z_2, z_3]^\top = [0, 1, 0]^\top$, which would select the second mixture component since $z_2 = 1$. The properties of \mathbf{z} imply that $\sum_{k=1}^K z_k = 1$. Therefore, \mathbf{z} is a *one-hot encoding* (also: *1-of- K representation*).

Thus far, we assumed that the indicator variables z_k are known. However, in practice, this is not the case, and we place a prior distribution on \mathbf{z} . Given that we do not know which mixture component generated the data point, we treat the indicator vector \mathbf{z} as a latent variable and place a prior

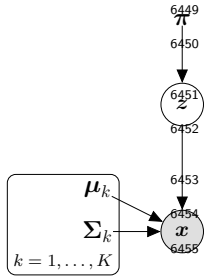
$$p(\mathbf{z}) = \boldsymbol{\pi} = [\pi_1, \dots, \pi_K]^\top, \quad \sum_{k=1}^K \pi_k = 1, \quad (11.60)$$

on it. Then the k th entry

$$\pi_k = p(z_k = 1) \quad (11.61)$$

of this probability vector describes the probability that the k th mixture component generated data point \mathbf{x} .

Figure 11.11 Graphical model for a GMM with a single data point.



Remark (Sampling from a GMM). The construction of this latent variable model (see the corresponding graphical model in Figure 11.11) lends itself to a very simple sampling procedure (generative process) to generate data:

- 1 Sample $z^{(i)} \sim p(\mathbf{z})$
- 2 Sample $\mathbf{x}^{(i)} \sim p(\mathbf{x} | z^{(i)} = 1)$

In the first step, we select a mixture component i (via the one-hot encoding \mathbf{z}) at random according to $p(\mathbf{z}) = \boldsymbol{\pi}$; in the second step we draw a sample from the corresponding mixture component. When we discard the samples of the latent variable so that we are left with the $\mathbf{x}^{(i)}$, we have valid samples from the GMM.

This kind of sampling, where samples of random variables depend on samples from the variable’s parents in the graphical model, is called *ancestral sampling*. \diamond

Generally, a probabilistic model is defined by the joint distribution of the data and the latent variables (see Section 8.3). With the prior $p(\mathbf{z})$ defined in (11.60)–(11.61) and the conditional $p(\mathbf{x} | \mathbf{z})$ from (11.59) we obtain all K components of this joint distribution via

$$p(\mathbf{x}, z_k = 1) = p(\mathbf{x} | z_k = 1)p(z_k = 1) = \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (11.62)$$

for $k = 1, \dots, K$, so that

$$p(\mathbf{x}, \mathbf{z}) = \begin{bmatrix} p(\mathbf{x}, z_1 = 1) \\ \vdots \\ p(\mathbf{x}, z_K = 1) \end{bmatrix} = \begin{bmatrix} \pi_1 \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \\ \vdots \\ \pi_K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K) \end{bmatrix}, \quad (11.63)$$

which fully specifies the probabilistic model.

11.4.2 Likelihood

To obtain the likelihood $p(\mathbf{x} | \boldsymbol{\theta})$ in a latent variable model, we need to marginalize out the latent variables. In our case, this can be done by summing out all latent variables from the joint $p(\mathbf{x}, \mathbf{z})$ in (11.63) so that

$$p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{\mathbf{z}} p(\mathbf{x} | \boldsymbol{\theta}, \mathbf{z}) p(\mathbf{z} | \boldsymbol{\theta}), \quad \boldsymbol{\theta} := \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k : k = 1, \dots, K\}. \quad (11.64)$$

We now explicitly condition on the parameters $\boldsymbol{\theta}$ of the probabilistic model, which we previously omitted. In (11.64), we sum over all K possible one-hot encodings of \mathbf{z} , which is denoted by $\sum_{\mathbf{z}}$. Since there is only a single non-zero single entry in each \mathbf{z} there are only K possible configurations/settings of \mathbf{z} . For example, if $K = 3$ then \mathbf{z} can have the configurations

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}. \quad (11.65)$$

Summing over all possible configurations of \mathbf{z} in (11.64) is equivalent to looking at the non-zero entry of the \mathbf{z} -vector and write

$$p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{\mathbf{z}} p(\mathbf{x} | \boldsymbol{\theta}, \mathbf{z}) p(\mathbf{z} | \boldsymbol{\theta}) \quad (11.66a)$$

$$= \sum_{k=1}^K p(\mathbf{x} | \boldsymbol{\theta}, z_k = 1) p(z_k = 1 | \boldsymbol{\theta}) \quad (11.66b)$$

so that the desired marginal distribution is given as

$$p(\mathbf{x} | \boldsymbol{\theta}) \stackrel{(11.66b)}{=} \sum_{k=1}^K p(\mathbf{x} | \boldsymbol{\theta}, z_k = 1) p(z_k = 1 | \boldsymbol{\theta}) \quad (11.67a)$$

$$= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (11.67b)$$

which we identify as the GMM model from (11.3). Given a dataset \mathcal{X} we immediately obtain the likelihood

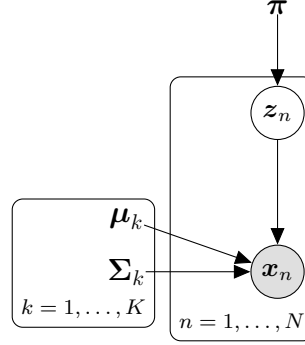
$$p(\mathcal{X} | \boldsymbol{\theta}) = \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\theta}) \stackrel{(11.67b)}{=} \prod_{n=1}^N \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (11.68)$$

which is exactly the GMM likelihood from (11.9). Therefore, the latent variable model with latent indicators z_k is an equivalent way of thinking about a Gaussian mixture model.

11.4.3 Posterior Distribution

Let us have a brief look at the posterior distribution on the latent variable \mathbf{z} . According to Bayes' theorem, the posterior of the k th component having

Figure 11.12
Graphical model for
a GMM with N data
points.



generated data point \mathbf{x}

$$p(z_k = 1 | \mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x} | z_k = 1)}{p(\mathbf{x})}, \quad (11.69)$$

where the marginal $p(\mathbf{x})$ is given in (11.67b). This yields the posterior distribution for the k th indicator variable z_k

$$p(z_k = 1 | \mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x} | z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x} | z_j = 1)} = \frac{\pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \mu_j, \Sigma_j)}, \quad (11.70)$$

which we identify as the responsibility of the k th mixture component for data point \mathbf{x} . Note that we omitted the explicit conditioning on the GMM parameters π_k, μ_k, Σ_k where $k = 1, \dots, K$.

11.4.4 Extension to a Full Dataset

Thus far, we only discussed the case where the dataset consists only of a single data point \mathbf{x} . However, the concepts of the prior and the posterior can be directly extended to the case of N data points $\mathcal{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$.

In the probabilistic interpretation of the GMM, every data point \mathbf{x}_n possesses its own latent variable

$$\mathbf{z}_n = [z_{n1}, \dots, z_{nK}]^\top \in \mathbb{R}^K. \quad (11.71)$$

Previously (when we only considered a single data point \mathbf{x}) we omitted the index n , but now this becomes important.

We share the same prior distribution π across all latent variables \mathbf{z}_n . The corresponding graphical model is shown in Figure 11.12, where we use the plate notation.

The conditional distribution $p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{z}_1, \dots, \mathbf{z}_N)$ factorizes over the data points and is given as

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{z}_1, \dots, \mathbf{z}_N) = \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n). \quad (11.72)$$

To obtain the posterior distribution $p(z_{nk} = 1 | \mathbf{x}_n)$ we follow the same reasoning as in Section 11.4.3 and apply Bayes' theorem to obtain

$$p(z_{nk} = 1 | \mathbf{x}_n) = \frac{p(\mathbf{x}_n | z_{nk} = 1)p(z_{nk} = 1)}{\sum_{j=1}^K p(\mathbf{x}_n | z_{nj} = 1)p(z_{nj} = 1)} \quad (11.73a)$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = r_{nk}. \quad (11.73b)$$

This means that $p(z_k = 1 | \mathbf{x}_n)$ is the (posterior) probability that the k th mixture component generated data point \mathbf{x}_n and corresponds to the responsibility r_{nk} we introduced in (11.17). Now, the responsibilities also have not only an intuitive but also a mathematically justified interpretation as posterior probabilities.

11.4.5 EM Algorithm Revisited

The EM algorithm that we introduced as an iterative scheme for maximum likelihood estimation can be derived in a principled way from the latent variable perspective. Given a current setting $\boldsymbol{\theta}^{(t)}$ of model parameters, the E-step calculates the expected log-likelihood

$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(t)}) = \mathbb{E}_{\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}^{(t)}} [\log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})] \quad (11.74a)$$

$$= \int \log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}^{(t)}) d\mathbf{z}, \quad (11.74b)$$

where the expectation of $\log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$ is taken with respect to the posterior $p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}^{(t)})$ of the latent variables. The M-step selects an updated set of model parameters $\boldsymbol{\theta}^{(t+1)}$ by maximizing (11.74b).

Although an EM iteration does increase the log-likelihood, there are no guarantees that EM converges to the maximum likelihood solution. It is possible that the EM algorithm converges to a local maximum of the log-likelihood. Different initializations of the parameters $\boldsymbol{\theta}$ could be used in multiple EM runs to reduce the risk of ending up in a bad local optimum. We do not go into further details here, but refer to the excellent expositions by Rogers and Girolami (2016) and Bishop (2006).

11.5 Further Reading

The GMM can be considered a generative model in the sense that it is straightforward to generate new data using ancestral sampling (Bishop, 2006). For given GMM parameters $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$, $k = 1, \dots, K$, we sample an index k from the probability vector $[\pi_1, \dots, \pi_K]^\top$ and then sample a data point $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$. If we repeat this N times, we obtain a dataset that has been generated by a GMM. Figure 11.1 was generated using this procedure.

Throughout this chapter, we assumed that the number of components K is known. In practice, this is often not the case. However, we could use nested cross-validation, as discussed in Section 8.5.1, to find good models.

Gaussian mixture models are closely related to the K -means clustering algorithm. K -means also uses the EM algorithm to assign data points to clusters. If we treat the means in the GMM as cluster centers and ignore the covariances, we arrive at K -means. As also nicely described by MacKay (2003a), K -means makes a “hard” assignments of data points to cluster centers μ_k , whereas a GMM makes a “soft” assignment via the responsibilities.

We only touched upon the latent variable perspective of GMMs and the EM algorithm. Note that EM can be used for parameter learning in general latent variable models, e.g., nonlinear state-space models (Ghahramani and Roweis, 1999; Roweis and Ghahramani, 1999) and for reinforcement learning as discussed by Barber (2012). Therefore, the latent variable perspective of a GMM is useful to derive the corresponding EM algorithm in a principled way (Bishop, 2006; Barber, 2012; Murphy, 2012).

We only discussed maximum likelihood estimation (via the EM algorithm) for finding GMM parameters. The standard criticisms of maximum likelihood also apply here:

- As in linear regression, maximum likelihood can suffer from severe overfitting. In the GMM case, this happens when the mean of a mixture component is identical to a data point and the covariance tends to $\mathbf{0}$. Then, the likelihood approaches infinity. Bishop (2006) and Barber (2012) discuss this issue in detail.
- We only obtain a point estimate of the parameters π_k, μ_k, Σ_k for $k = 1, \dots, K$, which does not give any indication of uncertainty in the parameter values. A Bayesian approach would place a prior on the parameters, which can be used to obtain a posterior distribution on the parameters. This posterior allows us to compute the model evidence (marginal likelihood), which can be used for model comparison, which gives us a principled way to determine the number of mixture components. Unfortunately, closed-form inference is not possible in this setting because there is no conjugate prior for this model. However, approximations, such as variational inference, can be used to obtain an approximate posterior (Bishop, 2006).

In this chapter, we discussed mixture models for density estimation. There is a plethora of density estimation techniques available. In practice, we often use histograms and kernel density estimation.

Histograms provide a non-parametric way to represent continuous densities and have been proposed by Pearson (1895). A histogram is constructed by “binning” the data space and count how many data points fall into each bin. Then a bar is drawn at the center of each bin, and the height of the bar is proportional to the number of data points within that bin. The

histograms

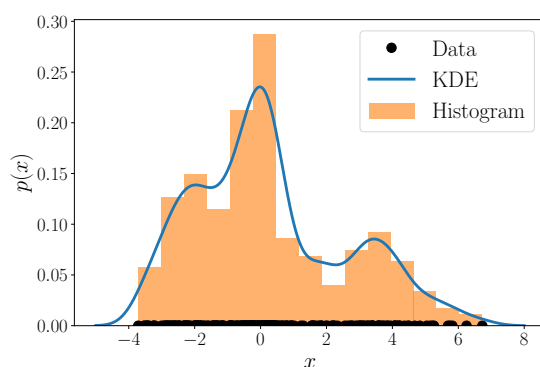


Figure 11.13
Histogram (orange bars) and kernel density estimation (blue line). The kernel density estimator produces a smooth estimate of the underlying density, whereas the histogram is an unsmoothed count measure of how many data points (black) fall into a single bin.

bin size is a critical hyper-parameter, and a bad choice can lead to overfitting and underfitting. Cross-validation, as discussed in Section 8.1.4, can be used to determine a good bin size.

Kernel density estimation, independently proposed by Rosenblatt (1956) and Parzen (1962), is a nonparametric way for density estimation. Given N i.i.d. samples, the kernel density estimator represents the underlying distribution as

$$p(\mathbf{x}) = \frac{1}{Nh} \sum_{n=1}^N k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right), \quad (11.75)$$

where k is a kernel function, i.e., a non-negative function that integrates to 1 and $h > 0$ is a smoothing/bandwidth parameter, which plays a similar role as the bin size in histograms. Note that we place a kernel on every single data point \mathbf{x}_n in the dataset. Commonly used kernel functions are the uniform distribution and the Gaussian distribution. Kernel density estimates are closely related to histograms, but by choosing a suitable kernel, we can guarantee smoothness of the density estimate. Figure 11.13 illustrates the difference between a histogram and a kernel density estimator (with a Gaussian-shaped kernel) for a given dataset of 250 data points.

kernel density estimation