

# 10

## Approximate Inference

**目次:** A central task in the application of probabilistic models is the evaluation of the posterior distribution  $p(\mathbf{Z}|\mathbf{X})$  of the latent variables  $\mathbf{Z}$  given the observed (visible) data variables  $\mathbf{X}$ , and the evaluation of expectations computed with respect to this distribution. The model might also contain some deterministic parameters, which we will leave implicit for the moment, or it may be a fully Bayesian model in which any unknown parameters are given prior distributions and are absorbed into the set of latent variables denoted by the vector  $\mathbf{Z}$ . For instance, in the EM algorithm we need to evaluate the expectation of the complete-data log likelihood with respect to the posterior distribution of the latent variables. For many models of practical interest, it will be infeasible to evaluate the posterior distribution or indeed to compute expectations with respect to this distribution. This could be because the dimensionality of the latent space is too high to work with directly or because the posterior distribution has a highly complex form for which expectations are not analytically tractable. In the case of continuous variables, the required integrations may not have closed-form

对隐藏量/参数的处理，  
没讲得太清楚了！

精确计算通常很困难

analytical solutions, while the dimensionality of the space and the complexity of the integrand may prohibit numerical integration. For discrete variables, the marginalizations involve summing over all possible configurations of the hidden variables, and though this is always possible in principle, we often find in practice that there may be exponentially many hidden states so that exact calculation is prohibitively expensive.]

两种近似计算的方法  
随机近似：MCMC等  
确定性近似：  
Laplace近似，变分推断等

In such situations, we need to resort to approximation schemes, and these fall broadly into two classes, according to whether they rely on stochastic or deterministic approximations. Stochastic techniques such as Markov chain Monte Carlo, described in Chapter 11, have enabled the widespread use of Bayesian methods across many domains. They generally have the property that given infinite computational resource, they can generate exact results, and the approximation arises from the use of a finite amount of processor time. In practice, sampling methods can be computationally demanding, often limiting their use to small-scale problems. Also, it can be difficult to know whether a sampling scheme is generating independent samples from the required distribution.

// In this chapter, we introduce a range of deterministic approximation schemes, some of which scale well to large applications. These are based on analytical approximations to the posterior distribution, for example by assuming that it factorizes in a particular way or that it has a specific parametric form such as a Gaussian. As such, they can never generate exact results, and so their strengths and weaknesses are complementary to those of sampling methods.

In Section 4.4, we discussed the Laplace approximation, which is based on a local Gaussian approximation to a mode (i.e., a maximum) of the distribution. Here we turn to a family of approximation techniques called *variational inference* or *variational Bayes*, which use more global criteria and which have been widely applied. We conclude with a brief introduction to an alternative variational framework known as *expectation propagation*.]

## 10.1. Variational Inference

**变分法**

Variational methods have their origins in the 18<sup>th</sup> century with the work of Euler, Lagrange, and others on the *calculus of variations*. Standard calculus is concerned with finding derivatives of functions. We can think of a function as a mapping that takes the value of a variable as the input and returns the value of the function as the output. The derivative of the function then describes how the output value varies as we make infinitesimal changes to the input value. Similarly, we can define a *functional* as a mapping that takes a function as the input and that returns the value of the functional as the output. An example would be the entropy  $H[p]$ , which takes a probability distribution  $p(x)$  as the input and returns the quantity

$$H[p] = -\int p(x) \ln p(x) dx \quad (10.1)$$

as the output. We can then introduce the concept of a *functional derivative*, which expresses how the value of the functional changes in response to infinitesimal changes to the input function (Feynman *et al.*, 1964). The rules for the calculus of variations mirror those of standard calculus and are discussed in Appendix D. Many problems can be expressed in terms of an optimization problem in which the quantity being optimized is a functional. The solution is obtained by exploring all possible input functions to find the one that maximizes, or minimizes, the functional. Variational methods have broad applicability and include such areas as finite element methods (Kapur, 1989) and maximum entropy (Schwarz, 1988). 有限元

变分法产生近似结果  
通常是由限制函数的取值范围  
引起的

Although there is nothing intrinsically approximate about variational methods, they do naturally lend themselves to finding approximate solutions. This is done by restricting the range of functions over which the optimization is performed, for instance by considering only quadratic functions or by considering functions composed of a linear combination of fixed basis functions in which only the coefficients of the linear combination can vary. In the case of applications to probabilistic inference, the restriction may for example take the form of factorization assumptions (Jordan *et al.*, 1999; Jaakkola, 2001).

举例：根据模型推断时  
问题分化

Now let us consider in more detail how the concept of variational optimization can be applied to the inference problem. Suppose we have a fully Bayesian model in which all parameters are given prior distributions. The model may also have latent variables as well as parameters, and we shall denote the set of all latent variables and parameters by  $\mathbf{Z}$ . Similarly, we denote the set of all observed variables by  $\mathbf{X}$ . For example, we might have a set of  $N$  independent, identically distributed data, for which  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and  $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$ . Our probabilistic model specifies the joint distribution  $p(\mathbf{X}, \mathbf{Z})$ , and our goal is to find an approximation for the posterior distribution  $p(\mathbf{Z}|\mathbf{X})$  as well as for the model evidence  $p(\mathbf{X})$ . As in our discussion of EM, we can decompose the log marginal probability using

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p) \quad (10.2)$$

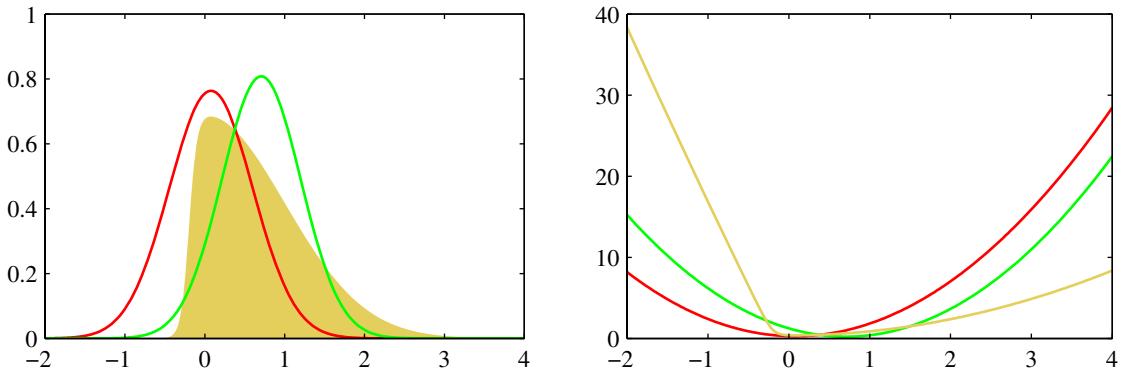
where we have defined

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right\} d\mathbf{Z} \quad (10.3)$$

$$\text{KL}(q||p) = - \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} \right\} d\mathbf{Z}. \quad (10.4)$$

注意：与EM不同的是，这里将参数也放入隐变量中！

This differs from our discussion of EM only in that the parameter vector  $\theta$  no longer appears, because the parameters are now stochastic variables and are absorbed into  $\mathbf{Z}$ . Since in this chapter we will mainly be interested in continuous variables we have used integrations rather than summations in formulating this decomposition. However, the analysis goes through unchanged if some or all of the variables are discrete simply by replacing the integrations with summations as required. As before, we can maximize the lower bound  $\mathcal{L}(q)$  by optimization with respect to the distribution  $q(\mathbf{Z})$ , which is equivalent to minimizing the KL divergence. If we allow any possible choice for  $q(\mathbf{Z})$ , then the maximum of the lower bound occurs when the KL divergence vanishes, which occurs when  $q(\mathbf{Z})$  equals the posterior distribution  $p(\mathbf{Z}|\mathbf{X})$ .



**Figure 10.1** Illustration of the variational approximation for the example considered earlier in Figure 4.14. The left-hand plot shows the original distribution (yellow) along with the Laplace (red) and variational (green) approximations, and the right-hand plot shows the negative logarithms of the corresponding curves.

However, we shall suppose the model is such that working with the true posterior distribution is intractable.

We therefore consider instead a restricted family of distributions  $q(\mathbf{Z})$  and then seek the member of this family for which the KL divergence is minimized. Our goal is to restrict the family sufficiently that they comprise only tractable distributions, while at the same time allowing the family to be sufficiently rich and flexible that it can provide a good approximation to the true posterior distribution. It is important to emphasize that the restriction is imposed purely to achieve tractability, and that subject to this requirement we should use as rich a family of approximating distributions as possible. In particular, there is no ‘over-fitting’ associated with highly flexible distributions. Using more flexible approximations simply allows us to approach the true posterior distribution more closely.

One way to restrict the family of approximating distributions is to use a parametric distribution  $q(\mathbf{Z}|\boldsymbol{\omega})$  governed by a set of parameters  $\boldsymbol{\omega}$ . The lower bound  $\mathcal{L}(q)$  then becomes a function of  $\boldsymbol{\omega}$ , and we can exploit standard nonlinear optimization techniques to determine the optimal values for the parameters. An example of this approach, in which the variational distribution is a Gaussian and we have optimized with respect to its mean and variance, is shown in Figure 10.1. ]

### 10.1.1 Factorized distributions

Here we consider an alternative way in which to restrict the family of distributions  $q(\mathbf{Z})$ . Suppose we partition the elements of  $\mathbf{Z}$  into disjoint groups that we denote by  $\mathbf{Z}_i$  where  $i = 1, \dots, M$ . We then assume that the  $q$  distribution factorizes with respect to these groups, so that

$$q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i). \quad (10.5)$$

$\text{除(10.5)之外其他假設}$

It should be emphasized that we are making no further assumptions about the distribution. In particular, we place no restriction on the functional forms of the individual factors  $q_i(\mathbf{Z}_i)$ . This factorized form of variational inference corresponds to an approximation framework developed in physics called *mean field theory* (Parisi, 1988).

平均场模型泛函优化的求解

Amongst all distributions  $q(\mathbf{Z})$  having the form (10.5), we now seek that distribution for which the lower bound  $\mathcal{L}(q)$  is largest. We therefore wish to make a free form (variational) optimization of  $\mathcal{L}(q)$  with respect to all of the distributions  $q_i(\mathbf{Z}_i)$ , which we do by optimizing with respect to each of the factors in turn. To achieve this, we first substitute (10.5) into (10.3) and then dissect out the dependence on one of the factors  $q_j(\mathbf{Z}_j)$ . Denoting  $q_j(\mathbf{Z}_j)$  by simply  $q_j$  to keep the notation uncluttered, we then obtain

$$\begin{aligned}\mathcal{L}(q) &= \int \prod_i q_i \left\{ \ln p(\mathbf{X}, \mathbf{Z}) - \sum_i \ln q_i \right\} d\mathbf{Z} \\ &= \int q_j \left\{ \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i d\mathbf{Z}_i \right\} d\mathbf{Z}_j - \int q_j \ln q_j d\mathbf{Z}_j + \text{const} \\ &= \int q_j \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j - \int q_j \ln q_j d\mathbf{Z}_j + \text{const}\end{aligned}\quad (10.6)$$

where we have defined a new distribution  $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$  by the relation

$$\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{const}. \quad (10.7)$$

Here the notation  $\mathbb{E}_{i \neq j} [\dots]$  denotes an expectation with respect to the  $q$  distributions over all variables  $\mathbf{z}_i$  for  $i \neq j$ , so that

$$\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i d\mathbf{Z}_i. \quad (10.8)$$

Now suppose we keep the  $\{q_{i \neq j}\}$  fixed and maximize  $\mathcal{L}(q)$  in (10.6) with respect to all possible forms for the distribution  $q_j(\mathbf{Z}_j)$ . This is easily done by recognizing that (10.6) is a negative Kullback-Leibler divergence between  $q_j(\mathbf{Z}_j)$  and  $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ . Thus maximizing (10.6) is equivalent to minimizing the Kullback-Leibler



Leonhard Euler  
1707–1783

Euler was a Swiss mathematician and physicist who worked in St. Petersburg and Berlin and who is widely considered to be one of the greatest mathematicians of all time. He is certainly the most prolific, and his collected works fill 75 volumes. Amongst his many

contributions, he formulated the modern theory of the function, he developed (together with Lagrange) the calculus of variations, and he discovered the formula  $e^{i\pi} = -1$ , which relates four of the most important numbers in mathematics. During the last 17 years of his life, he was almost totally blind, and yet he produced nearly half of his results during this period.

divergence, and the minimum occurs when  $q_j(\mathbf{Z}_j) = \tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ . Thus we obtain a general expression for the optimal solution  $q_j^*(\mathbf{Z}_j)$  given by

*注意, 这里的 X 是 observed*

*data, 而不是变量*

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{const.} \quad (10.9)$$

It is worth taking a few moments to study the form of this solution as it provides the basis for applications of variational methods. It says that the log of the optimal solution for factor  $q_j$  is obtained simply by considering the log of the joint distribution over all hidden and visible variables and then taking the expectation with respect to all of the other factors  $\{q_i\}$  for  $i \neq j$ .

The additive constant in (10.9) is set by normalizing the distribution  $q_j^*(\mathbf{Z}_j)$ . Thus if we take the exponential of both sides and normalize, we have

$$q_j^*(\mathbf{Z}_j) = \frac{\exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})])}{\int \exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})]) d\mathbf{Z}_j}.$$

In practice, we shall find it more convenient to work with the form (10.9) and then reinstate the normalization constant (where required) by inspection. This will become clear from subsequent examples.

*解法这个计算过程 //* The set of equations given by (10.9) for  $j = 1, \dots, M$  represent a set of consistency conditions for the maximum of the lower bound subject to the factorization constraint. However, they do not represent an explicit solution because the expression on the right-hand side of (10.9) for the optimum  $q_j^*(\mathbf{Z}_j)$  depends on expectations computed with respect to the other factors  $q_i(\mathbf{Z}_i)$  for  $i \neq j$ . We will therefore seek a consistent solution by first initializing all of the factors  $q_i(\mathbf{Z}_i)$  appropriately and then cycling through the factors and replacing each in turn with a revised estimate given by the right-hand side of (10.9) evaluated using the current estimates for all of the other factors. Convergence is guaranteed because bound is convex with respect to each of the factors  $q_i(\mathbf{Z}_i)$  (Boyd and Vandenberghe, 2004). }]

### 10.1.2 Properties of factorized approximations

*以高斯分布为例*

[ Our approach to variational inference is based on a factorized approximation to the true posterior distribution. Let us consider for a moment the problem of approximating a general distribution by a factorized distribution. To begin with, we discuss the problem of approximating a Gaussian distribution using a factorized Gaussian, which will provide useful insight into the types of inaccuracy introduced in using factorized approximations. Consider a Gaussian distribution  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$  over two correlated variables  $\mathbf{z} = (z_1, z_2)$  in which the mean and precision have elements

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \boldsymbol{\Lambda} = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix} \quad (10.10)$$

and  $\Lambda_{21} = \Lambda_{12}$  due to the symmetry of the precision matrix. Now suppose we wish to approximate this distribution using a factorized Gaussian of the form  $q(\mathbf{z}) = q_1(z_1)q_2(z_2)$ . We first apply the general result (10.9) to find an expression for the

optimal factor  $q_1^*(z_1)$ . In doing so it is useful to note that on the right-hand side we only need to retain those terms that have some functional dependence on  $z_1$  because all other terms can be absorbed into the normalization constant. Thus we have

$$\begin{aligned}\ln q_1^*(z_1) &= \mathbb{E}_{z_2}[\ln p(\mathbf{z})] + \text{const} \quad \text{对应式(10.9)中以 } X \text{ 为空} \\ &= \mathbb{E}_{z_2} \left[ -\frac{1}{2}(z_1 - \mu_1)^2 \Lambda_{11} - (z_1 - \mu_1) \Lambda_{12} (z_2 - \mu_2) \right] + \text{const} \\ &= -\frac{1}{2} z_1^2 \Lambda_{11} + z_1 \mu_1 \Lambda_{11} - z_1 \Lambda_{12} (\mathbb{E}[z_2] - \mu_2) + \text{const}. \quad (10.11)\end{aligned}$$

Next we observe that the right-hand side of this expression is a quadratic function of  $z_1$ , and so we can identify  $q^*(z_1)$  as a Gaussian distribution. It is worth emphasizing that we did not assume that  $q(z_i)$  is Gaussian, but rather we derived this result by variational optimization of the KL divergence over all possible distributions  $q(z_i)$ . Note also that we do not need to consider the additive constant in (10.9) explicitly because it represents the normalization constant that can be found at the end by inspection if required. Using the technique of completing the square, we can identify the mean and precision of this Gaussian, giving

$$q_1^*(z_1) = \mathcal{N}(z_1 | m_1, \Lambda_{11}^{-1}) \quad (10.12)$$

where

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (\mathbb{E}[z_2] - \mu_2). \quad (10.13)$$

By symmetry,  $q_2^*(z_2)$  is also Gaussian and can be written as

$$q_2^*(z_2) = \mathcal{N}(z_2 | m_2, \Lambda_{22}^{-1}) \quad (10.14)$$

in which

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21} (\mathbb{E}[z_1] - \mu_1). \quad (10.15)$$

Note that these solutions are coupled, so that  $q^*(z_1)$  depends on expectations computed with respect to  $q^*(z_2)$  and vice versa. In general, we address this by treating the variational solutions as re-estimation equations and cycling through the variables in turn updating them until some convergence criterion is satisfied. We shall see an example of this shortly. Here, however, we note that the problem is sufficiently simple that a closed form solution can be found. In particular, because  $\mathbb{E}[z_1] = m_1$  and  $\mathbb{E}[z_2] = m_2$ , we see that the two equations are satisfied if we take  $\mathbb{E}[z_1] = \mu_1$  and  $\mathbb{E}[z_2] = \mu_2$ , and it is easily shown that this is the only solution provided the distribution is nonsingular. This result is illustrated in Figure 10.2(a). We see that the mean is correctly captured but that the variance of  $q(\mathbf{z})$  is controlled by the direction of smallest variance of  $p(\mathbf{z})$ , and that the variance along the orthogonal direction is significantly under-estimated. It is a general result that a factorized variational approximation tends to give approximations to the posterior distribution that are too compact.

// By way of comparison, suppose instead that we had been minimizing the reverse Kullback-Leibler divergence  $\text{KL}(p\|q)$ . As we shall see, this form of KL divergence

### Section 2.3.1

迭代更新

但这个例子比较简单，直接得理论解

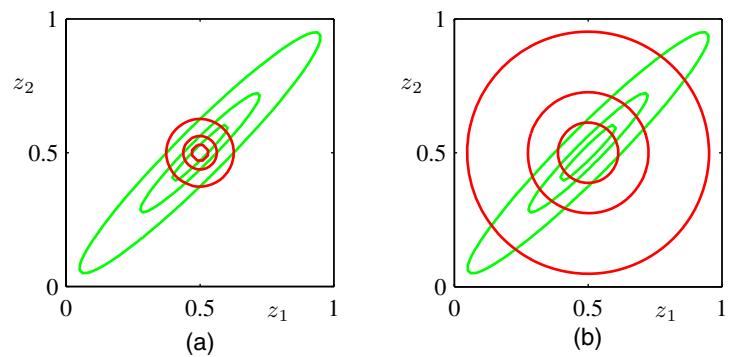
### Exercise 10.2

近似结果的特征

前后分割介绍的基于  $\text{KL}(q\|p)$  和  $\text{KL}(p\|q)$  的

近似，并以高斯分布为例  
给出了近似结果的不同特点，见 Figure 10.2

**Figure 10.2** Comparison of the two alternative forms for the Kullback-Leibler divergence. The green contours corresponding to 1, 2, and 3 standard deviations for a correlated Gaussian distribution  $p(\mathbf{z})$  over two variables  $z_1$  and  $z_2$ , and the red contours represent the corresponding levels for an approximating distribution  $q(\mathbf{z})$  over the same variables given by the product of two independent univariate Gaussian distributions whose parameters are obtained by minimization of (a) the Kullback-Leibler divergence  $\text{KL}(q\|p)$ , and (b) the reverse Kullback-Leibler divergence  $\text{KL}(p\|q)$ .



$\left\{ \begin{array}{l} \text{KL}(q\|p): \text{Variational approximation} \\ \text{KL}(p\|q): \text{expectation propagation} \end{array} \right.$

### Section 10.7

is used in an alternative approximate inference framework called (*expectation propagation*). We therefore consider the general problem of minimizing  $\text{KL}(p\|q)$  when  $q(\mathbf{Z})$  is a factorized approximation of the form (10.5). The KL divergence can then be written in the form

$$\text{KL}(p\|q) = - \int p(\mathbf{Z}) \left[ \sum_{i=1}^M \ln q_i(\mathbf{Z}_i) \right] d\mathbf{Z} + \text{const} \quad (10.16)$$

where the constant term is simply the entropy of  $p(\mathbf{Z})$  and so does not depend on  $q(\mathbf{Z})$ . We can now optimize with respect to each of the factors  $q_j(\mathbf{Z}_j)$ , which is easily done using a Lagrange multiplier to give

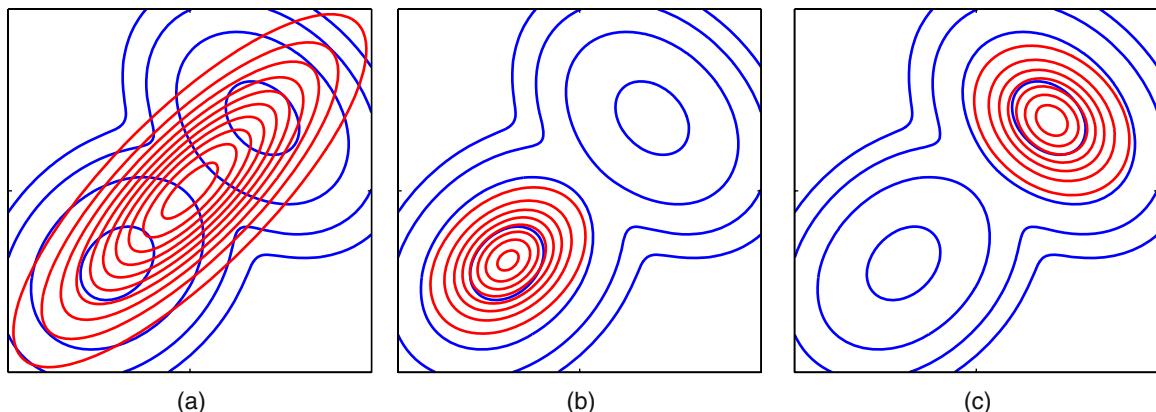
$$q_j^*(\mathbf{Z}_j) = \int p(\mathbf{Z}) \prod_{i \neq j} d\mathbf{Z}_i = p(\mathbf{Z}_j). \quad (10.17)$$

In this case, we find that the optimal solution for  $q_j(\mathbf{Z}_j)$  is just given by the corresponding marginal distribution of  $p(\mathbf{Z})$ . Note that this is a closed-form solution and so does not require iteration.

**高斯分布** To apply this result to the illustrative example of a Gaussian distribution  $p(\mathbf{z})$  over a vector  $\mathbf{z}$  we can use (2.98), which gives the result shown in Figure 10.2(b). **优势结果的特征** We see that once again the mean of the approximation is correct, but that it places significant probability mass in regions of variable space that have very low probability.

**分析和比较** The difference between these two results can be understood by noting that there is a large positive contribution to the Kullback-Leibler divergence

$$\text{KL}(q\|p) = - \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z})}{q(\mathbf{Z})} \right\} d\mathbf{Z} \quad (10.18)$$



**Figure 10.3** Another comparison of the two alternative forms for the Kullback-Leibler divergence. (a) The blue contours show a bimodal distribution  $p(\mathbf{Z})$  given by a mixture of two Gaussians, and the red contours correspond to the single Gaussian distribution  $q(\mathbf{Z})$  that best approximates  $p(\mathbf{Z})$  in the sense of minimizing the Kullback-Leibler divergence  $\text{KL}(p\|q)$ . (b) As in (a) but now the red contours correspond to a Gaussian distribution  $q(\mathbf{Z})$  found by numerical minimization of the Kullback-Leibler divergence  $\text{KL}(q\|p)$ . (c) As in (b) but showing a different local minimum of the Kullback-Leibler divergence.

这里是说，为了  $\min \text{KL}(q\|p)$

对  $-\int q \ln p$  这一项，  
P 非常  
处 q 一定也零，因此若 P  
是 unimodal，则基于  $\text{KL}(q\|p)$   
from regions of  $\mathbf{Z}$  space in which  $p(\mathbf{Z})$  is near zero unless  $q(\mathbf{Z})$  is also close to zero. Thus minimizing this form of KL divergence leads to distributions  $q(\mathbf{Z})$  that avoid regions in which  $p(\mathbf{Z})$  is small. Conversely, the Kullback-Leibler divergence  $\text{KL}(p\|q)$  is minimized by distributions  $q(\mathbf{Z})$  that are nonzero in regions where  $p(\mathbf{Z})$  is nonzero.

靠近 P，q 只能去抄一个  
峰，否则各个峰之间极片  
较小的区域就会存在较大 q，也就是这里所说的  
q 应“避免” P 直接小的  
区域。对  $\text{KL}(p\|q)$  而言  
则也类似。

We can gain further insight into the different behaviour of the two KL divergences if we consider approximating a multimodal distribution by a unimodal one, as illustrated in Figure 10.3. In practical applications, the true posterior distribution will often be multimodal, with most of the posterior mass concentrated in some number of relatively small regions of parameter space. These multiple modes may arise through nonidentifiability in the latent space or through complex nonlinear dependence on the parameters. Both types of multimodality were encountered in Chapter 9 in the context of Gaussian mixtures, where they manifested themselves as multiple maxima in the likelihood function, and a variational treatment based on the minimization of  $\text{KL}(q\|p)$  will tend to find one of these modes. By contrast, if we were to minimize  $\text{KL}(p\|q)$ , the resulting approximations would average across all of the modes and, in the context of the mixture model, would lead to poor predictive distributions (because the average of two good parameter values is typically itself not a good parameter value). It is possible to make use of  $\text{KL}(p\|q)$  to define a useful inference procedure, but this requires a rather different approach to the one discussed here, and will be considered in detail when we discuss expectation propagation. ]

### Section 10.7

$\text{KL}(q\|p), \text{KL}(p\|q)$  The two forms of Kullback-Leibler divergence are members of the *alpha family*

in 简一形式

of divergences (Ali and Silvey, 1966; Amari, 1985; Minka, 2005) defined by

$$D_\alpha(p\|q) = \frac{4}{1-\alpha^2} \left( 1 - \int p(x)^{(1+\alpha)/2} q(x)^{(1-\alpha)/2} dx \right) \quad (10.19)$$

where  $-\infty < \alpha < \infty$  is a continuous parameter. The Kullback-Leibler divergence  $\text{KL}(p\|q)$  corresponds to the limit  $\alpha \rightarrow 1$ , whereas  $\text{KL}(q\|p)$  corresponds to the limit  $\alpha \rightarrow -1$ . For all values of  $\alpha$  we have  $D_\alpha(p\|q) \geq 0$ , with equality if, and only if,  $p(x) = q(x)$ . Suppose  $p(x)$  is a fixed distribution, and we minimize  $D_\alpha(p\|q)$  with respect to some set of distributions  $q(x)$ . Then for  $\alpha \leq -1$  the divergence is zero forcing, so that any values of  $x$  for which  $p(x) = 0$  will have  $q(x) = 0$ , and typically  $q(x)$  will under-estimate the support of  $p(x)$  and will tend to seek the mode with the largest mass. Conversely for  $\alpha \geq 1$  the divergence is zero-avoiding, so that values of  $x$  for which  $p(x) > 0$  will have  $q(x) > 0$ , and typically  $q(x)$  will stretch to cover all of  $p(x)$ , and will over-estimate the support of  $p(x)$ . When  $\alpha = 0$  we obtain a symmetric divergence that is linearly related to the Hellinger distance given by

$$D_H(p\|q) = \int (p(x)^{1/2} - q(x)^{1/2})^2 dx. \quad (10.20)$$

The square root of the Hellinger distance is a valid distance metric.

### 10.1.3 Example: The univariate Gaussian

We now illustrate the factorized variational approximation using a Gaussian distribution over a single variable  $x$  (MacKay, 2003). Our goal is to infer the posterior distribution for the mean  $\mu$  and precision  $\tau$ , given a data set  $\mathcal{D} = \{x_1, \dots, x_N\}$  of observed values of  $x$  which are assumed to be drawn independently from the Gaussian. The likelihood function is given by

$$p(\mathcal{D}|\mu, \tau) = \left( \frac{\tau}{2\pi} \right)^{N/2} \exp \left\{ -\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}. \quad (10.21)$$

We now introduce conjugate prior distributions for  $\mu$  and  $\tau$  given by

$$p(\mu|\tau) = \mathcal{N}(\mu|\mu_0, (\lambda_0 \tau)^{-1}) \quad (10.22)$$

$$p(\tau) = \text{Gam}(\tau|a_0, b_0) \quad (10.23)$$

where  $\text{Gam}(\tau|a_0, b_0)$  is the gamma distribution defined by (2.146). Together these distributions constitute a Gaussian-Gamma conjugate prior distribution.

For this simple problem the posterior distribution can be found exactly, and again takes the form of a Gaussian-gamma distribution. However, for tutorial purposes we will consider a factorized variational approximation to the posterior distribution given by

前面的式子在这里初是  
参数  $\mu, \tau$

$$q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau). \quad (10.24)$$

### Section 2.3.6

### Exercise 2.44

Note that the true posterior distribution does not factorize in this way. The optimum factors  $q_\mu(\mu)$  and  $q_\tau(\tau)$  can be obtained from the general result (10.9) as follows. For  $q_\mu(\mu)$  we have

$$\begin{aligned}\ln q_\mu^*(\mu) &= \mathbb{E}_\tau [\ln p(\mathcal{D}|\mu, \tau) + \ln p(\mu|\tau)] + \text{const} \\ &= -\frac{\mathbb{E}[\tau]}{2} \left\{ \lambda_0(\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right\} + \text{const.}\end{aligned}\quad (10.25)$$

Completing the square over  $\mu$  we see that  $q_\mu(\mu)$  is a Gaussian  $\mathcal{N}(\mu|\mu_N, \lambda_N^{-1})$  with mean and precision given by

$$\mu_N = \frac{\lambda_0\mu_0 + N\bar{x}}{\lambda_0 + N} \quad (10.26)$$

$$\lambda_N = (\lambda_0 + N)\mathbb{E}[\tau]. \quad (10.27)$$

Note that for  $N \rightarrow \infty$  this gives the maximum likelihood result in which  $\mu_N = \bar{x}$  and the precision is infinite.

Similarly, the optimal solution for the factor  $q_\tau(\tau)$  is given by

$$\begin{aligned}\ln q_\tau^*(\tau) &= \mathbb{E}_\mu [\ln p(\mathcal{D}|\mu, \tau) + \ln p(\mu|\tau)] + \ln p(\tau) + \text{const} \\ &= (a_0 - 1) \ln \tau - b_0 \tau + \frac{N}{2} \ln \tau + \frac{1}{2} \ln \tau \\ &\quad - \frac{\tau}{2} \mathbb{E}_\mu \left[ \sum_{n=1}^N (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2 \right] + \text{const}\end{aligned}\quad (10.28)$$

and hence  $q_\tau(\tau)$  is a gamma distribution  $\text{Gam}(\tau|a_N, b_N)$  with parameters

$$a_N = a_0 + \frac{N+1}{2} \quad (10.29)$$

$$b_N = b_0 + \frac{1}{2} \mathbb{E}_\mu \left[ \sum_{n=1}^N (x_n - \mu)^2 + \lambda_0(\mu - \mu_0)^2 \right]. \quad (10.30)$$

### Exercise 10.8

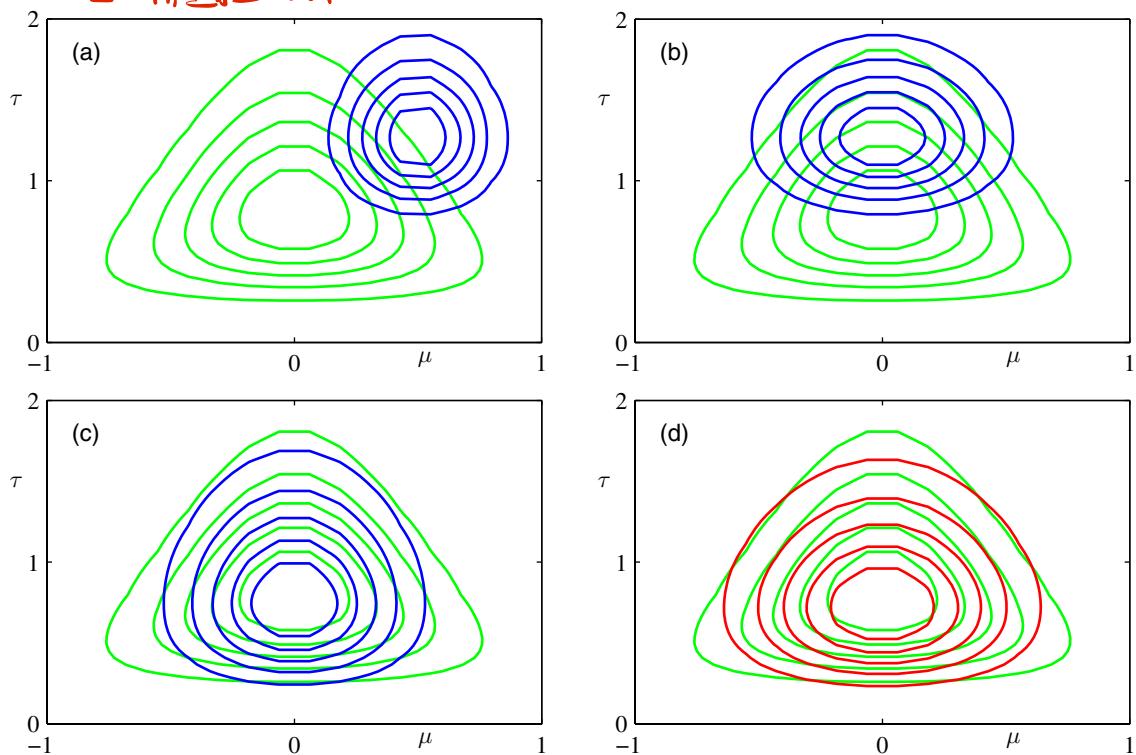
### Section 10.4.1

Again this exhibits the expected behaviour when  $N \rightarrow \infty$ .

It should be emphasized that we did not assume these specific functional forms for the optimal distributions  $q_\mu(\mu)$  and  $q_\tau(\tau)$ . They arose naturally from the structure of the likelihood function and the corresponding conjugate priors.

**迭代计算** Thus we have expressions for the optimal distributions  $q_\mu(\mu)$  and  $q_\tau(\tau)$  each of which depends on moments evaluated with respect to the other distribution. One approach to finding a solution is therefore to make an initial guess for, say, the moment  $\mathbb{E}[\tau]$  and use this to re-compute the distribution  $q_\mu(\mu)$ . Given this revised distribution we can then extract the required moments  $\mathbb{E}[\mu]$  and  $\mathbb{E}[\mu^2]$ , and use these to recompute the distribution  $q_\tau(\tau)$ , and so on. Since the space of hidden variables for this example is only two dimensional, we can illustrate the variational approximation to the posterior distribution by plotting contours of both the true posterior and the factorized approximation, as illustrated in Figure 10.4.

迭代计算过程的图示



**Figure 10.4** Illustration of variational inference for the mean  $\mu$  and precision  $\tau$  of a univariate Gaussian distribution. Contours of the true posterior distribution  $p(\mu, \tau|D)$  are shown in green. (a) Contours of the initial factorized approximation  $q_\mu(\mu)q_\tau(\tau)$  are shown in blue. (b) After re-estimating the factor  $q_\mu(\mu)$ . (c) After re-estimating the factor  $q_\tau(\tau)$ . (d) Contours of the optimal factorized approximation, to which the iterative scheme converges, are shown in red.

这里问题很简单，可以直解  
得到理论解：取特定  
的初始值（对应非informative prior），迭代-矩阵即  
可收敛。

In general, we will need to use an iterative approach such as this in order to solve for the optimal factorized posterior distribution. For the very simple example we are considering here, however, we can find an explicit solution by solving the simultaneous equations for the optimal factors  $q_\mu(\mu)$  and  $q_\tau(\tau)$ . Before doing this, we can simplify these expressions by considering broad, noninformative priors in which  $\mu_0 = a_0 = b_0 = \lambda_0 = 0$ . Although these parameter settings correspond to improper priors, we see that the posterior distribution is still well defined. Using the standard result  $\mathbb{E}[\tau] = a_N/b_N$  for the mean of a gamma distribution, together with (10.29) and (10.30), we have

$$\frac{1}{\mathbb{E}[\tau]} = \mathbb{E} \left[ \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2 \right] = (\bar{x}^2 - 2\bar{x}\mathbb{E}[\mu] + \mathbb{E}[\mu^2]) \frac{N}{N+1} \quad (10.31)$$

Then, using (10.26) and (10.27), we obtain the first and second order moments of

$q_\mu(\mu)$  in the form

$$\mathbb{E}[\mu] = \bar{x}, \quad \mathbb{E}[\mu^2] = \bar{x}^2 + \frac{1}{N\mathbb{E}[\tau]}. \quad (10.32)$$

### Exercise 10.9

We can now substitute these moments into (10.31) and then solve for  $\mathbb{E}[\tau]$  to give

$$\frac{1}{\mathbb{E}[\tau]} = \frac{1}{N} (\bar{x}^2 - \bar{x}^2)$$

For a comprehensive treatment of Bayesian inference for the Gaussian distribution, including a discussion of the advantages compared to maximum likelihood, see Minka (1998).

We recognize the right-hand side as the familiar unbiased estimator for the variance of a univariate Gaussian distribution, and so we see that the use of a Bayesian approach has avoided the bias of the maximum likelihood solution.

### Section 1.2.4

#### 10.1.4 Model comparison

贝叶斯模型选择  
该部分就是要 find  
approximate posterior  
distributions over models  
using variational inference

As well as performing inference over the hidden variables  $\mathbf{Z}$ , we may also wish to compare a set of candidate models, labelled by the index  $m$ , and having prior probabilities  $p(m)$ . Our goal is then to approximate the posterior probabilities  $p(\mathbf{m}|\mathbf{X})$ , where  $\mathbf{X}$  is the observed data. This is a slightly more complex situation than that considered so far because different models may have different structure and indeed different dimensionality for the hidden variables  $\mathbf{Z}$ . We cannot therefore simply consider a factorized approximation  $q(\mathbf{Z})q(m)$ , but must instead recognize that the posterior over  $\mathbf{Z}$  must be conditioned on  $m$ , and so we must consider

### Exercise 10.10

近似分布形式  $q(\mathbf{Z}, m) = q(\mathbf{Z}|m)q(m)$ . We can readily verify the following decomposition based on this variational distribution

步骤和前面一样：① 豪分下界

不能像前面一样没为  
 $q(\mathbf{Z}, m) = q(\mathbf{Z})q(m)$

$$\ln p(\mathbf{X}) = \mathcal{L}_M - \sum_m \sum_{\mathbf{Z}} q(\mathbf{Z}|m)q(m) \ln \left\{ \frac{p(\mathbf{Z}, m|\mathbf{X})}{q(\mathbf{Z}|m)q(m)} \right\} \quad (10.34)$$

where the  $\mathcal{L}_M$  is a lower bound on  $\ln p(\mathbf{X})$  and is given by

$$\mathcal{L}_M = \sum_m \sum_{\mathbf{Z}} q(\mathbf{Z}|m)q(m) \ln \left\{ \frac{p(\mathbf{Z}, \mathbf{X}, m)}{q(\mathbf{Z}|m)q(m)} \right\}. \quad (10.35)$$

### Exercise 10.11

Here we are assuming discrete  $\mathbf{Z}$ , but the same analysis applies to continuous latent variables provided the summations are replaced with integrations. We can maximize  $\mathcal{L}_M$  with respect to the distribution  $q(m)$  using a Lagrange multiplier, with the result

$$q(m) \propto p(m) \exp\{\mathcal{L}_M\} \quad \text{where} \quad \mathcal{L}_M = \sum_m q(\mathbf{Z}|m) \ln \frac{p(\mathbf{Z}, \mathbf{X}|m)}{q(\mathbf{Z}|m)}. \quad (10.36)$$

③ 迭代计算

However, if we maximize  $\mathcal{L}_M$  with respect to the  $q(\mathbf{Z}|m)$ , we find that the solutions for different  $m$  are coupled, as we expect because they are conditioned on  $m$ . We proceed instead by first optimizing each of the  $q(\mathbf{Z}|m)$  individually by optimization

$$\mathcal{L}_m = \sum_{\mathbf{Z}} q(\mathbf{Z}|m) \ln \frac{p(\mathbf{Z}, \mathbf{X}|m)}{q(\mathbf{Z}|m)}.$$

*or equivalently by optimization of  $L_m$*

of (10.35), and then subsequently determining the  $q(m)$  using (10.36). After normalization the resulting values for  $q(m)$  can be used for model selection or model averaging in the usual way.

## 变分法用在高斯混合模型中

### 10.2. Illustration: Variational Mixture of Gaussians

回顾本节的  
内容及意义

We now return to our discussion of the Gaussian mixture model and apply the variational inference machinery developed in the previous section. This will provide a good illustration of the application of variational methods and will also demonstrate how a Bayesian treatment elegantly resolves many of the difficulties associated with the maximum likelihood approach (Attias, 1999b). The reader is encouraged to work through this example in detail as it provides many insights into the practical application of variational methods. Many Bayesian models, corresponding to much more sophisticated distributions, can be solved by straightforward extensions and generalizations of this analysis.

高斯混合模型的介绍 Our starting point is the likelihood function for the Gaussian mixture model, illustrated by the graphical model in Figure 9.6. For each observation  $\mathbf{x}_n$  we have a corresponding latent variable  $\mathbf{z}_n$  comprising a 1-of- $K$  binary vector with elements  $z_{nk}$  for  $k = 1, \dots, K$ . As before we denote the observed data set by  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , and similarly we denote the latent variables by  $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$ . From (9.10) we can write down the conditional distribution of  $\mathbf{Z}$ , given the mixing coefficients  $\boldsymbol{\pi}$ , in the form

$$p(\mathbf{Z}|\boldsymbol{\pi}) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}}. \quad (10.37)$$

Similarly, from (9.11), we can write down the conditional distribution of the observed data vectors, given the latent variables and the component parameters

$$p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})^{z_{nk}} \quad (10.38)$$

where  $\boldsymbol{\mu} = \{\boldsymbol{\mu}_k\}$  and  $\boldsymbol{\Lambda} = \{\boldsymbol{\Lambda}_k\}$ . Note that we are working in terms of precision matrices rather than covariance matrices as this somewhat simplifies the mathematics.

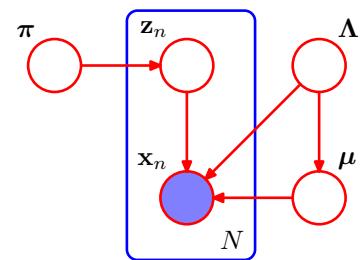
Next we introduce priors over the parameters  $\boldsymbol{\mu}$ ,  $\boldsymbol{\Lambda}$  and  $\boldsymbol{\pi}$ . The analysis is considerably simplified if we use conjugate prior distributions. We therefore choose a Dirichlet distribution over the mixing coefficients  $\boldsymbol{\pi}$

$$p(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}_0) = C(\boldsymbol{\alpha}_0) \prod_{k=1}^K \pi_k^{\alpha_0-1} \quad (10.39)$$

where by symmetry we have chosen the same parameter  $\alpha_0$  for each of the components, and  $C(\boldsymbol{\alpha}_0)$  is the normalization constant for the Dirichlet distribution defined

#### Section 10.4.1

**Figure 10.5** Directed acyclic graph representing the Bayesian mixture of Gaussians model, in which the box (plate) denotes a set of  $N$  i.i.d. observations. Here  $\mu$  denotes  $\{\mu_k\}$  and  $\Lambda$  denotes  $\{\Lambda_k\}$ .



### Section 2.2.1

by (B.23). As we have seen, the parameter  $\alpha_0$  can be interpreted as the effective prior number of observations associated with each component of the mixture. If the value of  $\alpha_0$  is small, then the posterior distribution will be influenced primarily by the data rather than by the prior.

Similarly, we introduce an independent Gaussian-Wishart prior governing the mean and precision of each Gaussian component, given by

$$\begin{aligned} p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) &= p(\boldsymbol{\mu}|\boldsymbol{\Lambda})p(\boldsymbol{\Lambda}) \\ &= \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, (\beta_0 \boldsymbol{\Lambda}_k)^{-1}) \mathcal{W}(\boldsymbol{\Lambda}_k | \mathbf{W}_0, \nu_0) \end{aligned} \quad (10.40)$$

### Section 2.3.6

because this represents the conjugate prior distribution when both the mean and precision are unknown. Typically we would choose  $\mathbf{m}_0 = \mathbf{0}$  by symmetry.

The resulting model can be represented as a directed graph as shown in Figure 10.5. Note that there is a link from  $\boldsymbol{\Lambda}$  to  $\boldsymbol{\mu}$  since the variance of the distribution over  $\boldsymbol{\mu}$  in (10.40) is a function of  $\boldsymbol{\Lambda}$ . [

隐藏量与参数  
辨证关系

This example provides a nice illustration of the distinction between latent variables and parameters. Variables such as  $\mathbf{z}_n$  that appear inside the plate are regarded as latent variables because the number of such variables grows with the size of the data set. By contrast, variables such as  $\boldsymbol{\mu}$  that are outside the plate are fixed in number independently of the size of the data set, and so are regarded as parameters.

From the perspective of graphical models, however, there is really no fundamental difference between them.]

#### 10.2.1 Variational distribution

In order to formulate a variational treatment of this model, we next write down the joint distribution of all of the random variables, which is given by

$$p(\mathbf{X}, \mathbf{Z}, \pi, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})p(\mathbf{Z}|\pi)p(\pi)p(\boldsymbol{\mu}|\boldsymbol{\Lambda})p(\boldsymbol{\Lambda}) \quad (10.41)$$

in which the various factors are defined above. The reader should take a moment to verify that this decomposition does indeed correspond to the probabilistic graphical model shown in Figure 10.5. Note that only the variables  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  are observed.

注意，EM算法推导时总是condition on参数，而这里的步骤是将参数和隐藏变量放在一样进行考虑，概率的条件项是空集。

We now consider a variational distribution which factorizes between the latent variables and the parameters so that

$$q(\mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = q(\mathbf{Z})q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}). \quad (10.42)$$

式(10.4)是引例唯一  
假设

这样得真清楚。  
真心！

It is remarkable that this is the *only* assumption that we need to make in order to obtain a tractable practical solution to our Bayesian mixture model. In particular, the functional form of the factors  $q(\mathbf{Z})$  and  $q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$  will be determined automatically by optimization of the variational distribution. Note that we are omitting the subscripts on the  $q$  distributions, much as we do with the  $p$  distributions in (10.41), and are relying on the arguments to distinguish the different distributions.

The corresponding sequential update equations for these factors can be easily derived by making use of the general result (10.9). Let us consider the derivation of the update equation for the factor  $q(\mathbf{Z})$ . The log of the optimized factor is given by

$$\ln q^*(\mathbf{Z}) = \mathbb{E}_{\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}} [\ln p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] + \text{const.} \quad (10.43)$$

We now make use of the decomposition (10.41). Note that we are only interested in the functional dependence of the right-hand side on the variable  $\mathbf{Z}$ . Thus any terms that do not depend on  $\mathbf{Z}$  can be absorbed into the additive normalization constant, giving

$$\ln q^*(\mathbf{Z}) = \mathbb{E}_{\boldsymbol{\pi}} [\ln p(\mathbf{Z} | \boldsymbol{\pi})] + \mathbb{E}_{\boldsymbol{\mu}, \boldsymbol{\Lambda}} [\ln p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] + \text{const.} \quad (10.44)$$

Substituting for the two conditional distributions on the right-hand side, and again absorbing any terms that are independent of  $\mathbf{Z}$  into the additive constant, we have

步。其中M步的优化过程也是一个迭代优化过程，  
分为计算 $q(\mathbf{Z})$ 和 $q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$ 两个stage，而从计算 $q(\mathbf{Z})$ 开始。  
程来看， $q(\mathbf{Z})$ 的计算可以与EM算法的E步对应， $q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$ 的计算可以与EM算法的M步对应。  
where we have defined  $\mathbf{z}$ 定义为“隐变量”，则 E步： $q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$ 的计算，M步： $q(\mathbf{Z})$ 的计算。

$$\begin{aligned} \ln \rho_{nk} &= \mathbb{E}[\ln \pi_k] + \frac{1}{2} \mathbb{E}[\ln |\boldsymbol{\Lambda}_k|] - \frac{D}{2} \ln(2\pi) \\ &\quad - \frac{1}{2} \mathbb{E}_{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k} [(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Lambda}_k (\mathbf{x}_n - \boldsymbol{\mu}_k)] \end{aligned} \quad (10.46)$$

where  $D$  is the dimensionality of the data variable  $\mathbf{x}$ . Taking the exponential of both sides of (10.45) we obtain

$$q^*(\mathbf{Z}) \propto \prod_{n=1}^N \prod_{k=1}^K \rho_{nk}^{z_{nk}}. \quad (10.47)$$

Requiring that this distribution be normalized, and noting that for each value of  $n$  the quantities  $z_{nk}$  are binary and sum to 1 over all values of  $k$ , we obtain

$$q^*(\mathbf{Z}) = \prod_{n=1}^N \prod_{k=1}^K r_{nk}^{z_{nk}} \quad (10.48)$$

我们先假设式(10.46)中的期望  
能够被计算出来以便推导能进  
续推导。事实上，这些期望是对

$q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$ 进行的而 $q(\mathbf{Z})$ ，  
还未推导出来。下面，我们将先  
推导出 $q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$ 的计算公式

再回过头来给出式(10.46)

Exercise 10.12

中期望的具体计算公式。

where

$$r_{nk} = \frac{\rho_{nk}}{\sum_{j=1}^K \rho_{nj}}. \quad (10.49)$$

可以看出来：

基础高斯分布得到的

近似后验分布和它的函数形式  
式子先验分布  $p(\pi|\Lambda)$  相同

(有种类别的意思)。由此，我们自然地推得了  $q(\pi|Z)$   
的函数形式，而没有直接  
假设  $q(\pi|Z)$  的函数形式。

We see that the optimal solution for the factor  $q(Z)$  takes the same functional form as the prior  $p(Z|\pi)$ . Note that because  $\rho_{nk}$  is given by the exponential of a real quantity, the quantities  $r_{nk}$  will be nonnegative and will sum to one, as required.

For the discrete distribution  $q^*(Z)$  we have the standard result

$$\mathbb{E}[z_{nk}] = r_{nk} \quad (10.50)$$

from which we see that the quantities  $r_{nk}$  are playing the role of responsibilities. Note that the optimal solution for  $q^*(Z)$  depends on moments evaluated with respect to the distributions of other variables, and so again the variational update equations are coupled and must be solved iteratively. 这个计算

At this point, we shall find it convenient to define three statistics of the observed data set evaluated with respect to the responsibilities, given by

$$N_k = \sum_{n=1}^N r_{nk} \quad (10.51)$$

$$\bar{x}_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} \mathbf{x}_n \quad (10.52)$$

$$\mathbf{S}_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \bar{x}_k)(\mathbf{x}_n - \bar{x}_k)^T. \quad (10.53)$$

Note that these are analogous to quantities evaluated in the maximum likelihood EM algorithm for the Gaussian mixture model.

// Now let us consider the factor  $q(\pi, \mu, \Lambda)$  in the variational posterior distribution. Again using the general result (10.9) we have

$$\begin{aligned} \ln q^*(\pi, \mu, \Lambda) &= \ln p(\pi) + \sum_{k=1}^K \ln p(\mu_k, \Lambda_k) + \mathbb{E}_Z [\ln p(Z|\pi)] \\ &\quad + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}[z_{nk}] \ln \mathcal{N}(\mathbf{x}_n | \mu_k, \Lambda_k^{-1}) + \text{const.} \end{aligned} \quad (10.54)$$

We observe that the right-hand side of this expression decomposes into a sum of terms involving only  $\pi$  together with terms only involving  $\mu$  and  $\Lambda$ , which implies that the variational posterior  $q(\pi, \mu, \Lambda)$  factorizes to give  $q(\pi)q(\mu, \Lambda)$ . Furthermore, the terms involving  $\mu$  and  $\Lambda$  themselves comprise a sum over  $k$  of terms involving  $\mu_k$  and  $\Lambda_k$  leading to the further factorization

$$q(\pi, \mu, \Lambda) = q(\pi) \prod_{k=1}^K q(\mu_k, \Lambda_k). \quad (10.55)$$

基础高斯分布的  
解  $q^*(\pi, \mu, \Lambda)$  可以进一步  
分解，即其形式为式  
(10.55)

{ ①

Identifying the terms on the right-hand side of (10.54) that depend on  $\pi$ , we have

$$\ln q^*(\boldsymbol{\pi}) = (\alpha_0 - 1) \sum_{k=1}^K \ln \pi_k + \sum_{k=1}^K \sum_{n=1}^N r_{nk} \ln \pi_k + \text{const} \quad (10.56)$$

where we have used (10.50). Taking the exponential of both sides, we recognize  $q^*(\boldsymbol{\pi})$  as a Dirichlet distribution

$$q^*(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}) \quad (10.57)$$

where  $\boldsymbol{\alpha}$  has components  $\alpha_k$  given by

{ ②

$$\alpha_k = \alpha_0 + N_k. \quad (10.58)$$

// Finally, the variational posterior distribution  $q^*(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$  does not factorize into the product of the marginals, but we can always use the product rule to write it in the form  $q^*(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) = q^*(\boldsymbol{\mu}_k | \boldsymbol{\Lambda}_k)q^*(\boldsymbol{\Lambda}_k)$ . The two factors can be found by inspecting (10.54) and reading off those terms that involve  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Lambda}_k$ . The result, as expected, is a Gaussian-Wishart distribution and is given by

$$q^*(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) = \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_k, (\beta_k \boldsymbol{\Lambda}_k)^{-1}) \mathcal{W}(\boldsymbol{\Lambda}_k | \mathbf{W}_k, \nu_k) \quad (10.59)$$

where we have defined

$$\beta_k = \beta_0 + N_k \quad (10.60)$$

$$\mathbf{m}_k = \frac{1}{\beta_k} (\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k) \quad (10.61)$$

$$\mathbf{W}_k^{-1} = \mathbf{W}_0^{-1} + N_k \mathbf{S}_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^T \quad (10.62)$$

$$\nu_k = \nu_0 + N_k. \quad (10.63)$$

从逻辑上来说上述(10.52),  $q(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})$

更新公式对应EM的M步，而  
实际情况已退化不见。

需要计算的量在MLE中也需要  
要计算。

现在回过头来给出  $p_{nk}$  的  
计算公式，进一步可得  $r_{nk}$

**Exercise 10.14**

计算公式。从计算上讲，  
这部分的计算又相当于  
EM算法的E步。

These update equations are analogous to the M-step equations of the EM algorithm for the maximum likelihood solution of the mixture of Gaussians. We see that the computations that must be performed in order to update the variational posterior distribution over the model parameters involve evaluation of the same sums over the data set, as arose in the maximum likelihood treatment.

In order to perform this variational M step, we need the expectations  $\mathbb{E}[z_{nk}] = r_{nk}$  representing the responsibilities. These are obtained by normalizing the  $\rho_{nk}$  that are given by (10.46). We see that (this expression) involves expectations with respect to the variational distributions of the parameters, and these are easily evaluated to give

式(10.46)

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k} [(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Lambda}_k (\mathbf{x}_n - \boldsymbol{\mu}_k)] \\ = D \beta_k^{-1} + \nu_k (\mathbf{x}_n - \mathbf{m}_k)^T \mathbf{W}_k (\mathbf{x}_n - \mathbf{m}_k) \end{aligned} \quad (10.64)$$

$$\ln \tilde{\Lambda}_k \equiv \mathbb{E} [\ln |\boldsymbol{\Lambda}_k|] = \sum_{i=1}^D \psi \left( \frac{\nu_k + 1 - i}{2} \right) + D \ln 2 + \ln |\mathbf{W}_k| \quad (10.65)$$

$$\ln \tilde{\pi}_k \equiv \mathbb{E} [\ln \pi_k] = \psi(\alpha_k) - \psi(\hat{\alpha}) \quad (10.66)$$

where we have introduced definitions of  $\tilde{\Lambda}_k$  and  $\tilde{\pi}_k$ , and  $\psi(\cdot)$  is the digamma function defined by (B.25), with  $\hat{\alpha} = \sum_k \alpha_k$ . The results (10.65) and (10.66) follow from the standard properties of the Wishart and Dirichlet distributions.

If we substitute (10.64), (10.65), and (10.66) into (10.46) and make use of (10.49), we obtain the following result for the responsibilities

*Appendix B*  
知道  $p_{nk}$  如何  
计算而，进一步可推得  
 $r_{nk}$  的计算公式，即式(10.67)。  
现由变分推导得  
的  $r_{nk}$  与 EM 算法最优化  
的  $r_{nk}$  在计算公  
式是相似的，二者都有  
responsibility 的含义。

$$r_{nk} \propto \tilde{\pi}_k \tilde{\Lambda}_k^{1/2} \exp \left\{ -\frac{D}{2\beta_k} - \frac{\nu_k}{2} (\mathbf{x}_n - \mathbf{m}_k)^T \mathbf{W}_k (\mathbf{x}_n - \mathbf{m}_k) \right\}. \quad (10.67)$$

Notice the similarity to the corresponding result for the responsibilities in maximum likelihood EM, which from (9.13) can be written in the form

$$r_{nk} \propto \pi_k |\Lambda_k|^{1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Lambda_k (\mathbf{x}_n - \boldsymbol{\mu}_k) \right\} \quad (10.68)$$

where we have used the precision in place of the covariance to highlight the similarity to (10.67).

Thus the optimization of the variational posterior distribution involves cycling between two stages analogous to the E and M steps of the maximum likelihood EM algorithm. In the variational equivalent of the E step, we use the current distributions over the model parameters to evaluate the moments in (10.64), (10.65), and (10.66) and hence evaluate  $\mathbb{E}[z_{nk}] = r_{nk}$ . Then in the subsequent variational equivalent of the M step, we keep these responsibilities fixed and use them to re-compute the variational distribution over the parameters using (10.57) and (10.59). In each case, we see that the variational posterior distribution has the same functional form as the corresponding factor in the joint distribution (10.41). This is a general result and is a consequence of the choice of conjugate distributions.

Figure 10.6 shows the results of applying this approach to the rescaled Old Faithful data set for a Gaussian mixture model having  $K = 6$  components. We see that after convergence, there are only two components for which the expected values of the mixing coefficients are numerically distinguishable from their prior values. This effect can be understood qualitatively in terms of the automatic trade-off in a Bayesian model between fitting the data and the complexity of the model, in which the complexity penalty arises from components whose parameters are pushed away from their prior values. Components that take essentially no responsibility for explaining the data points have  $r_{nk} \approx 0$  and hence  $N_k \approx 0$ . From (10.58), we see that  $\alpha_k \approx \alpha_0$  and from (10.60)–(10.63) we see that the other parameters revert to their prior values. In principle such components are fitted slightly to the data points, but for broad priors this effect is too small to be seen numerically. For the variational Gaussian mixture model the expected values of the mixing coefficients in the posterior distribution are given by

$$\mathbb{E}[\pi_k] = \frac{\alpha_k + N_k}{K\alpha_0 + N}. \quad (10.69)$$

Consider a component for which  $N_k \approx 0$  and  $\alpha_k \approx \alpha_0$ . If the prior is broad so that  $\alpha_0 \rightarrow 0$ , then  $\mathbb{E}[\pi_k] \rightarrow 0$  and the component plays no role in the model, whereas if

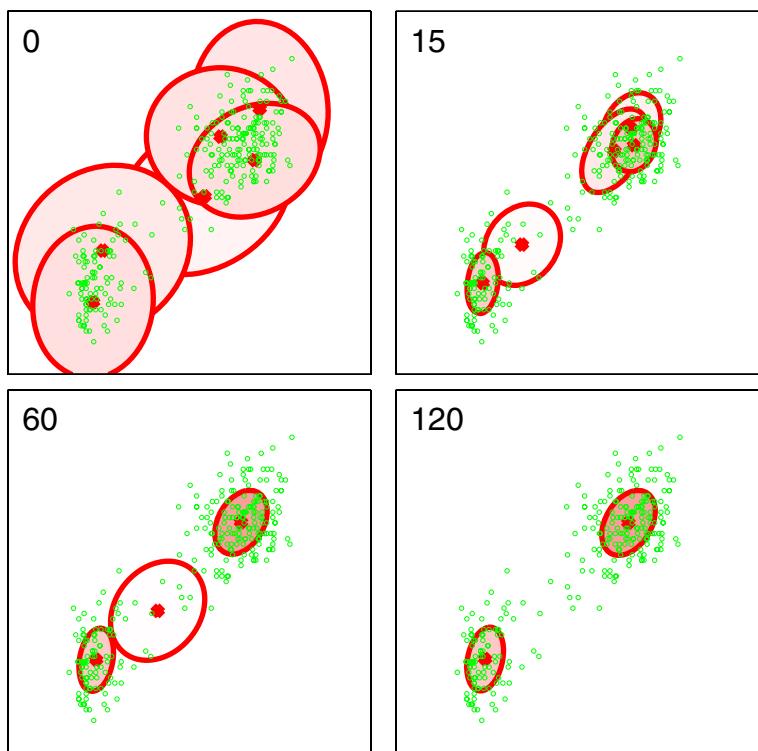
总结：从计算上看，上述是  
分后验分布的迭代优化  
EM算法相似，事实上，从  
逻辑上讲，EM算法的E步  
已退化了，因为这里对参数  
又加上了处理与隐变量  
已相同，所有的优化在选择  
上均属M步。 举例：

### Section 3.4

对没例子结果  
的理论分析

### Exercise 10.15

**Figure 10.6** Variational Bayesian mixture of  $K = 6$  Gaussians applied to the Old Faithful data set, in which the ellipses denote the one standard-deviation density contours for each of the components, and the density of red ink inside each ellipse corresponds to the mean value of the mixing coefficient for each component. The number in the top left of each diagram shows the number of iterations of variational inference. Components whose expected mixing coefficient are numerically indistinguishable from zero are not plotted.



the prior tightly constrains the mixing coefficients so that  $\alpha_0 \rightarrow \infty$ , then  $\mathbb{E}[\pi_k] \rightarrow 1/K$ .

In Figure 10.6, the prior over the mixing coefficients is a Dirichlet of the form (10.39). Recall from Figure 2.5 that for  $\alpha_0 < 1$  the prior favours solutions in which some of the mixing coefficients are zero. Figure 10.6 was obtained using  $\alpha_0 = 10^{-3}$ , and resulted in two components having nonzero mixing coefficients. If instead we choose  $\alpha_0 = 1$  we obtain three components with nonzero mixing coefficients, and for  $\alpha = 10$  all six components have nonzero mixing coefficients.]

As we have seen there is a close similarity between the variational solution for the Bayesian mixture of Gaussians and the EM algorithm for maximum likelihood. In fact if we consider the limit  $N \rightarrow \infty$  then the Bayesian treatment converges to the maximum likelihood EM algorithm. For anything other than very small data sets, the dominant computational cost of the variational algorithm for Gaussian mixtures arises from the evaluation of the responsibilities, together with the evaluation and inversion of the weighted data covariance matrices. These computations mirror precisely those that arise in the maximum likelihood EM algorithm, and so there is little computational overhead in using this Bayesian approach as compared to the traditional maximum likelihood one. There are, however, some substantial advantages. First of all, the singularities that arise in maximum likelihood when a Gaussian component ‘collapses’ onto a specific data point are absent in the Bayesian treatment.

★ 注意，为什么同样是EM的步骤第9章中就是MLE和EM算法，而这里说变分推断却是一种Bayesian方法呢？答案在于第9章进过EM

感觉就是  $N \rightarrow \infty$  时  
MAP 收敛于 MLE。  
证明见 Exercise 10.20.

总结：1. 当  $N \rightarrow \infty$  时，  
上述对高斯混合模型  
进行的变分推断（属于  
近似推断，也是一种贝叶  
斯方法）收敛于用EM算  
法进行极大似然估计；  
2. 变分推断的计算量比  
较而言并不高；  
3. 变分推断相比MLE还有①②③3个优势。

## 10.2. Illustration: Variational Mixture of Gaussians 481

### 见9.2.1

王相同, 参考了参数的近似后验分布, 因此是一种  
贝叶斯方法!

Indeed, these singularities are removed if we simply introduce a prior and then use a MAP estimate instead of maximum likelihood. Furthermore, there is no over-fitting if we choose a large number  $K$  of components in the mixture, as we saw in Figure 10.6. Finally, the variational treatment opens up the possibility of determining the optimal number of components in the mixture without resorting to techniques such as cross validation.

### 模型选择 Section 10.2.4

计算离分下界两个  
例外

这句话是指: 我们可以使用有限差分计算梯度, 通过检查梯度是否为0来判断当前的优化是否完成。

### 10.2.2 Variational lower bound

We can also straightforwardly evaluate the lower bound (10.3) for this model. In practice, it is useful to be able to monitor the bound during the re-estimation in order to test for convergence. It can also provide a valuable check on both the mathematical expressions for the solutions and their software implementation, because at each step of the iterative re-estimation procedure the value of this bound should not decrease. [We can take this a stage further to provide a deeper test of the correctness of both the mathematical derivation of the update equations and of their software implementation by using finite differences to check that each update does indeed give a (constrained) maximum of the bound (Svensén and Bishop, 2004).]

[For the variational mixture of Gaussians, the lower bound (10.3) is given by

$$\begin{aligned} \mathcal{L} &= \sum_{\mathbf{Z}} \iiint q(\mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})}{q(\mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})} \right\} d\boldsymbol{\pi} d\boldsymbol{\mu} d\boldsymbol{\Lambda} \\ &= \mathbb{E}[\ln p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] - \mathbb{E}[\ln q(\mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] \\ &= \mathbb{E}[\ln p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] + \mathbb{E}[\ln p(\mathbf{Z}|\boldsymbol{\pi})] + \mathbb{E}[\ln p(\boldsymbol{\pi})] + \mathbb{E}[\ln p(\boldsymbol{\mu}, \boldsymbol{\Lambda})] \\ &\quad - \mathbb{E}[\ln q(\mathbf{Z})] - \mathbb{E}[\ln q(\boldsymbol{\pi})] - \mathbb{E}[\ln q(\boldsymbol{\mu}, \boldsymbol{\Lambda})] \end{aligned} \quad (10.70)$$

注意, 除了  $\mathbf{X}$  是数据  
以外其他都是变量  
Exercise 10.16

where, to keep the notation uncluttered, we have omitted the  $\star$  superscript on the  $q$  distributions, along with the subscripts on the expectation operators because each expectation is taken with respect to all of the random variables in its argument. The various terms in the bound are easily evaluated to give the following results

$$\begin{aligned} \mathbb{E}[\ln p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] &= \frac{1}{2} \sum_{k=1}^K N_k \left\{ \ln \tilde{\Lambda}_k - D\beta_k^{-1} - \nu_k \text{Tr}(\mathbf{S}_k \mathbf{W}_k) \right. \\ &\quad \left. - \nu_k (\bar{\mathbf{x}}_k - \mathbf{m}_k)^T \mathbf{W}_k (\bar{\mathbf{x}}_k - \mathbf{m}_k) - D \ln(2\pi) \right\} \end{aligned} \quad (10.71)$$

$$\mathbb{E}[\ln p(\mathbf{Z}|\boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \ln \tilde{\pi}_k \quad (10.72)$$

$$\mathbb{E}[\ln p(\boldsymbol{\pi})] = \ln C(\boldsymbol{\alpha}_0) + (\alpha_0 - 1) \sum_{k=1}^K \ln \tilde{\pi}_k \quad (10.73)$$

$$\begin{aligned}\mathbb{E}[\ln p(\boldsymbol{\mu}, \boldsymbol{\Lambda})] &= \frac{1}{2} \sum_{k=1}^K \left\{ D \ln(\beta_0/2\pi) + \ln \tilde{\Lambda}_k - \frac{D\beta_0}{\beta_k} \right. \\ &\quad \left. - \beta_0 \nu_k (\mathbf{m}_k - \mathbf{m}_0)^T \mathbf{W}_k (\mathbf{m}_k - \mathbf{m}_0) \right\} + K \ln B(\mathbf{W}_0, \nu_0) \\ &\quad + \frac{(\nu_0 - D - 1)}{2} \sum_{k=1}^K \ln \tilde{\Lambda}_k - \frac{1}{2} \sum_{k=1}^K \nu_k \text{Tr}(\mathbf{W}_0^{-1} \mathbf{W}_k)\end{aligned}\quad (10.74)$$

$$\mathbb{E}[\ln q(\mathbf{Z})] = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \ln r_{nk} \quad (10.75)$$

$$\mathbb{E}[\ln q(\boldsymbol{\pi})] = \sum_{k=1}^K (\alpha_k - 1) \ln \tilde{\pi}_k + \ln C(\boldsymbol{\alpha}) \quad (10.76)$$

$$\mathbb{E}[\ln q(\boldsymbol{\mu}, \boldsymbol{\Lambda})] = \sum_{k=1}^K \left\{ \frac{1}{2} \ln \tilde{\Lambda}_k + \frac{D}{2} \ln \left( \frac{\beta_k}{2\pi} \right) - \frac{D}{2} - H[q(\boldsymbol{\Lambda}_k)] \right\} \quad (10.77)$$

where  $D$  is the dimensionality of  $\mathbf{x}$ ,  $H[q(\boldsymbol{\Lambda}_k)]$  is the entropy of the Wishart distribution given by (B.82), and the coefficients  $C(\boldsymbol{\alpha})$  and  $B(\mathbf{W}, \nu)$  are defined by (B.23) and (B.79), respectively. Note that the terms involving expectations of the logs of the  $q$  distributions simply represent the negative entropies of those distributions. Some simplifications and combination of terms can be performed when these expressions are summed to give the lower bound. However, we have kept the expressions separate for ease of understanding.

这是是说：10.2.1节是人  
造出来的角度是大化要分下界  
得到了相互的 re-estimation  
equations，另外还可以从参数之  
间的角度，是大化要分下界。

**Exercise 10.18**  
同样可以得到相同的  
迭代计算公式，但这就  
需要先对分布的函数形式  
进行假设。

Finally, it is worth noting that the lower bound provides an alternative approach for deriving the variational re-estimation equations obtained in Section 10.2.1. To do this we use the fact that, since the model has conjugate priors, the functional form of the factors in the variational posterior distribution is known, namely discrete for  $\mathbf{Z}$ , Dirichlet for  $\boldsymbol{\pi}$ , and Gaussian-Wishart for  $(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$ . By taking general parametric forms for these distributions we can derive the form of the lower bound as a function of the parameters of the distributions. Maximizing the bound with respect to these parameters then gives the required re-estimation equations.

### 10.2.3 Predictive density: 解决预测问题

In applications of the Bayesian mixture of Gaussians model we will often be interested in the predictive density for a new value  $\hat{\mathbf{x}}$  of the observed variable. Associated with this observation will be a corresponding latent variable  $\hat{\mathbf{z}}$ , and the predictive density is then given by

$$p(\hat{\mathbf{x}}|\mathbf{X}) = \sum_{\hat{\mathbf{z}}} \iiint p(\hat{\mathbf{x}}|\hat{\mathbf{z}}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\hat{\mathbf{z}}|\boldsymbol{\pi}) p(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}|\mathbf{X}) d\boldsymbol{\pi} d\boldsymbol{\mu} d\boldsymbol{\Lambda} \quad (10.78)$$

where  $p(\pi, \mu, \Lambda | \mathbf{X})$  is the (unknown) true posterior distribution of the parameters. Using (10.37) and (10.38) we can first perform the summation over  $\widehat{\mathbf{z}}$  to give

$$p(\widehat{\mathbf{x}} | \mathbf{X}) = \sum_{k=1}^K \iiint \pi_k \mathcal{N}(\widehat{\mathbf{x}} | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) p(\pi, \mu, \Lambda | \mathbf{X}) d\pi d\mu d\Lambda. \quad (10.79)$$

Because the remaining integrations are intractable, we approximate the predictive density by replacing the true posterior distribution  $p(\pi, \mu, \Lambda | \mathbf{X})$  with its variational approximation  $q(\pi)q(\mu, \Lambda)$  to give

可以看到，我们只用过前面  
计算参数的后验  $q(\pi, \mu, \Lambda)$  而  
没必要用后验  $q(\mu, \Lambda)$ ，因为  $q(\mu, \Lambda)$  是  
为了计算  $q(\pi, \mu, \Lambda)$  而存在的

**Exercise 10.19**  
中间量，即引入隐藏变量后，  
优化更简单。我们当然也可以不用式(10.78)  
而是：

$$\begin{aligned} p(\widehat{\mathbf{x}} | \mathbf{X}) &= \int p(\widehat{\mathbf{x}} | \mathbf{z}, \pi, \mu, \Lambda) p(\mathbf{z}, \pi, \mu, \Lambda | \mathbf{X}) d\mathbf{z} d\pi d\mu d\Lambda \\ &= \int p(\widehat{\mathbf{x}} | \mathbf{z}, \mu, \Lambda) p(\mathbf{z}, \pi, \mu, \Lambda | \mathbf{X}) d\mathbf{z} d\pi d\mu d\Lambda \end{aligned}$$

**Exercise 10.20**  
然后用  $q(\mathbf{z}, \pi, \mu, \Lambda)$   
 $= q(\mathbf{z})q(\pi, \mu, \Lambda)$  近似

**Section 10.1.4**  
 $p(\mathbf{z}, \pi, \mu, \Lambda | \mathbf{X})$  代入进  
级计算，但是没有必要  
这么做，直接像文中  
那样，只用到  $q(\mu, \Lambda)$   
简单又准确。

**Exercise 10.21**

$$p(\widehat{\mathbf{x}} | \mathbf{X}) \underset{\text{最终得预测值}}{\approx} \sum_{k=1}^K \iiint \pi_k \mathcal{N}(\widehat{\mathbf{x}} | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) q(\pi)q(\mu_k, \Lambda_k) d\pi d\mu_k d\Lambda_k \quad (10.80)$$

where we have made use of the factorization (10.55) and in each term we have implicitly integrated out all variables  $\{\boldsymbol{\mu}_j, \Lambda_j\}$  for  $j \neq k$ . The remaining integrations can now be evaluated analytically giving a mixture of Student's t-distributions

$$p(\widehat{\mathbf{x}} | \mathbf{X}) \underset{\text{预测分布}}{\approx} \frac{1}{\bar{\alpha}} \sum_{k=1}^K \alpha_k \text{St}(\widehat{\mathbf{x}} | \mathbf{m}_k, \mathbf{L}_k, \nu_k + 1 - D) \quad (10.81)$$

in which the  $k^{\text{th}}$  component has mean  $\mathbf{m}_k$ , and the precision is given by

$$\mathbf{L}_k = \frac{(\nu_k + 1 - D)\beta_k}{(1 + \beta_k)} \mathbf{W}_k \quad (10.82)$$

in which  $\nu_k$  is given by (10.63). When the size  $N$  of the data set is large the predictive distribution (10.81) reduces to a mixture of Gaussians.

#### 10.2.4 Determining the number of components 对模型选择的问题

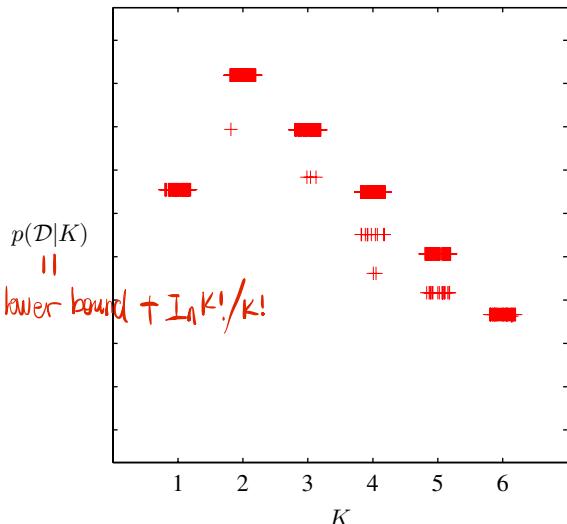
We have seen that the variational lower bound can be used to determine a posterior distribution over the number  $K$  of components in the mixture model. There is, however, one subtlety that needs to be addressed. For any given setting of the parameters in a Gaussian mixture model (except for specific degenerate settings), there will exist other parameter settings for which the density over the observed variables will be identical. These parameter values differ only through a re-labelling of the components. For instance, consider a mixture of two Gaussians and a single observed variable  $x$ , in which the parameters have the values  $\pi_1 = a, \pi_2 = b, \mu_1 = c, \mu_2 = d, \sigma_1 = e, \sigma_2 = f$ . Then the parameter values  $\pi_1 = b, \pi_2 = a, \mu_1 = d, \mu_2 = c, \sigma_1 = f, \sigma_2 = e$ , in which the two components have been exchanged, will by symmetry give rise to the same value of  $p(x)$ . If we have a mixture model comprising  $K$  components, then each parameter setting will be a member of a family of  $K!$  equivalent settings.

In the context of maximum likelihood, this redundancy is irrelevant because the parameter optimization algorithm (for example EM) will, depending on the initialization of the parameters, find one specific solution, and the other equivalent solutions play no role. In a Bayesian setting, however, we marginalize over all possible

这里说的“变分下界”是前文提到的变分下界 +  $\ln K! / K!$

Figure 10.7

Plot of the variational lower bound  $\mathcal{L}$  versus the number  $K$  of components in the Gaussian mixture model, for the Old Faithful data, showing a distinct peak at  $K = 2$  components. For each value of  $K$ , the model is trained from 100 different random starts, and the results shown as '+' symbols plotted with small random horizontal perturbations so that they can be distinguished. Note that some solutions find suboptimal local maxima, but that this happens infrequently.



10.3

是指参数数(从八到九)的经验分布

在括号里说这段叙述需  
结合 Exercise 10.22 来理解。  
详细见笔记。

实际上这已经给出了模型选  
择的一种方案，类似于 KIC, AIC  
Exercise 10.22  
model compare 的指标是：  
变分下界 +  $\ln K! / K!$

虽然，在变分下界上新增的项  $\ln K! / K!$   
就体现了模型的复杂程度

Section 3.4

下面给出了模型选择的另  
一种方案：混合模型的 K  
基于相关向量机的技术等

Exercise 10.23  
得到，即对数据解译性  
差的 components 会被自  
动剪枝。

parameter values. We have seen in Figure 10.2 that if the true posterior distribution is multimodal, variational inference based on the minimization of  $KL(q||p)$  will tend to approximate the distribution in the neighbourhood of one of the modes and ignore the others. Again, because equivalent modes have equivalent predictive densities, (究竟在 K 固定为某个值的情况下上面的推论是否成立。) this is of no concern provided we are considering a model having a specific number  $K$  of components. If, however, we wish to compare different values of  $K$ , then we need to take account of this multimodality. A simple approximate solution is to add a term  $\ln K! / K!$  onto the lower bound when used for model comparison and averaging<sup>33</sup>

Figure 10.7 shows a plot of the lower bound, including (the multimodality factor) versus the number  $K$  of components for the Old Faithful data set. It is worth emphasizing once again that maximum likelihood would lead to values of the likelihood function that increase monotonically with  $K$  (assuming the singular solutions have been avoided, and discounting the effects of local maxima) and so cannot be used to determine an appropriate model complexity. By contrast, Bayesian inference automatically makes the trade-off between model complexity and fitting the data.

This approach to the determination of  $K$  requires that a range of models having different  $K$  values be trained and compared. An alternative approach to determining a suitable value for  $K$  is to treat the mixing coefficients  $\pi$  as parameters and make point estimates of their values by maximizing the lower bound (Corduneanu and Bishop, 2001) with respect to  $\pi$  instead of maintaining a probability distribution over them as in the fully Bayesian approach. This leads to the re-estimation equation

$$\pi_k = \frac{1}{N} \sum_{n=1}^N r_{nk} \quad (10.83)$$

对  $\pi$  的迭代优化是嵌套到近似后验的迭代计算之中  
and (this maximization is interleaved with the variational updates for the  $q$  distribution over the remaining parameters.) Components that provide insufficient contribution

to explaining the data will have their mixing coefficients driven to zero during the optimization, and so they are effectively removed from the model through *automatic relevance determination*. This allows us to make a single training run in which we start with a relatively large initial value of  $K$ , and allow surplus components to be pruned out of the model. The origins of the sparsity when optimizing with respect to hyperparameters is discussed in detail in the context of the relevance vector machine.

### Section 7.2.2

- ① 介绍了什么是 induced factorizations，它们是怎样产生的？
- ② 介绍了 induced factorizations 在实际场景下的用处；
- ③ 给出了如何找 induced factorizations。

这是指 the true joint distribution 不是指 数据的真分布  $P_{\text{real}}$ ，而是  $P_{\text{model}}$ ，  
比如 GMM。高分后验分布  
是对应  $P_{\text{model}}$  后验分布  
的近似，而  $P_{\text{model}}$  则是  
对  $P_{\text{real}}$  的近似，这  
里存在两层近似。

#### 10.2.5 Induced factorizations

- ① [In deriving these variational update equations for the Gaussian mixture model, we assumed a particular factorization of the variational posterior distribution given by (10.42).] However, the optimal solutions for the various factors exhibit additional factorizations. In particular, the solution for  $q^*(\mu, \Lambda)$  is given by the product of an independent distribution  $q^*(\mu_k, \Lambda_k)$  over each of the components  $k$  of the mixture, whereas the variational posterior distribution  $q^*(\mathbf{Z})$  over the latent variables, given by (10.48), factorizes into an independent distribution  $q^*(\mathbf{z}_n)$  for each observation  $n$  (note that it does not further factorize with respect to  $k$  because, for each value of  $n$ , the  $z_{nk}$  are constrained to sum to one over  $k$ ).] (These additional factorizations) are a consequence of the interaction between the assumed factorization and the conditional independence properties of the true distribution, as characterized by the directed graph in Figure 10.5. *factorizations, 这样自然地导致出新的分解形式*

We shall refer to these additional factorizations as *induced factorizations* because they arise from an interaction between the factorization assumed in the variational posterior distribution and the conditional independence properties of the true joint distribution. [In a numerical implementation of the variational approach it is important to take account of such additional factorizations.] For instance, it would be very inefficient to maintain a full precision matrix for the Gaussian distribution over a set of variables if the optimal form for that distribution always had a diagonal precision matrix (corresponding to a factorization with respect to the individual variables described by that Gaussian).]

- ⑤ [Such induced factorizations can easily be detected using a simple graphical test based on d-separation as follows. We partition the latent variables into three disjoint groups A, B, C and then let us suppose that we are assuming a factorization between C and the remaining latent variables, so that

这是人为假设的 (10.84),  
而不管  $P(X, A, B, C)$  if  $q(A, B, C) = q(A, B)q(C)$ .  $P(A, B|C)$  是否成立 (10.84) 三.

$$\ln q^*(\mathbf{A}, \mathbf{B}) = \mathbb{E}_{\mathbf{C}}[\ln p(\mathbf{X}, \mathbf{A}, \mathbf{B}, \mathbf{C})] + \text{const}$$

Using the general result (10.9), together with the product rule for probabilities, we see that the optimal solution for  $q(\mathbf{A}, \mathbf{B})$  is given by

$$\begin{aligned} \ln q^*(\mathbf{A}, \mathbf{B}) &= \mathbb{E}_{\mathbf{C}}[\ln p(\mathbf{X}, \mathbf{A}, \mathbf{B}, \mathbf{C})] + \text{const} \\ &= \mathbb{E}_{\mathbf{C}}[\ln p(\mathbf{A}, \mathbf{B} | \mathbf{X}, \mathbf{C})] + \text{const}. \end{aligned} \quad (10.85)$$

We now ask whether this resulting solution will factorize between A and B, in other words whether  $q^*(\mathbf{A}, \mathbf{B}) = q^*(\mathbf{A})q^*(\mathbf{B})$ . This will happen if, and only if,  $\ln p(\mathbf{A}, \mathbf{B} | \mathbf{X}, \mathbf{C}) = \ln p(\mathbf{A} | \mathbf{X}, \mathbf{C}) + \ln p(\mathbf{B} | \mathbf{X}, \mathbf{C})$ , that is, if the conditional independence relation

也就是说是在  $P(X, A, B, C)$  中  $A \perp\!\!\!\perp B | X, C$  (10.86)

根据图模型中检查式 (10.86)

是否成立 (基于 d-separation)

is satisfied. We can test to see if this relation does hold, for any choice of  $\mathbf{A}$  and  $\mathbf{B}$  by making use of the d-separation criterion.

**笔记:** To illustrate this, consider again the Bayesian mixture of Gaussians represented by the directed graph in Figure 10.5, in which we are assuming a variational factorization given by (10.42). We can see immediately that the variational posterior distribution over the parameters must factorize between  $\boldsymbol{\pi}$  and the remaining parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\Lambda}$  because all paths connecting  $\boldsymbol{\pi}$  to either  $\boldsymbol{\mu}$  or  $\boldsymbol{\Lambda}$  must pass through one of the nodes  $\mathbf{z}_n$  all of which are in the conditioning set for our conditional independence test and all of which are head-to-tail with respect to such paths.

### 10.3. Variational Linear Regression

As a second illustration of variational inference, we return to the Bayesian linear regression model of Section 3.3. In the evidence framework, we approximated the integration over  $\alpha$  and  $\beta$  by making point estimates obtained by maximizing the log marginal likelihood. A fully Bayesian approach would integrate over the hyperparameters as well as over the parameters. Although exact integration is intractable, we can use variational methods to find a tractable approximation. In order to simplify the discussion, we shall suppose that the noise precision parameter  $\beta$  is known, and is fixed to its true value, although the framework is easily extended to include the distribution over  $\beta$ . For the linear regression model, the variational treatment will turn out to be equivalent to the evidence framework. Nevertheless, it provides a good exercise in the use of variational methods and will also lay the foundation for variational treatment of Bayesian logistic regression in Section 10.6.

Recall that the likelihood function for  $\mathbf{w}$ , and the prior over  $\mathbf{w}$ , are given by

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}_n, \beta^{-1}) \quad (10.87)$$

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \quad (10.88)$$

where  $\boldsymbol{\phi}_n = \boldsymbol{\phi}(\mathbf{x}_n)$ . We now introduce a prior distribution over  $\alpha$ . From our discussion in Section 2.3.6, we know that the conjugate prior for the precision of a Gaussian is given by a gamma distribution, and so we choose

$$p(\alpha) = \text{Gam}(\alpha|a_0, b_0) \quad (10.89)$$

where  $\text{Gam}(\cdot|\cdot, \cdot)$  is defined by (B.26). Thus the joint distribution of all the variables is given by

$$p(\mathbf{t}, \mathbf{w}, \alpha) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}|\alpha)p(\alpha). \quad (10.90)$$

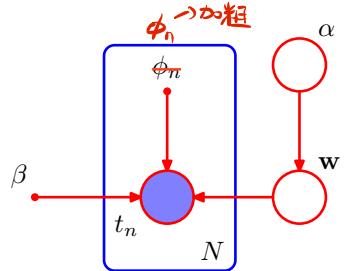
This can be represented as a directed graphical model as shown in Figure 10.8.

#### 10.3.1 Variational distribution

Our first goal is to find an approximation to the posterior distribution  $p(\mathbf{w}, \alpha|\mathbf{t})$ . To do this, we employ the variational framework of Section 10.1, with a variational

Exercise 10.26  
练习题 10.26 相当于在考虑  
分布的情况下将 10.3 节  
的结果重新推导一遍

**Figure 10.8** Probabilistic graphical model representing the joint distribution (10.90) for the Bayesian linear regression model.



posterior distribution given by the factorized expression

$$q(\mathbf{w}, \alpha) = q(\mathbf{w})q(\alpha). \quad (10.91)$$

We can find re-estimation equations for the factors in this distribution by making use of the general result (10.9). Recall that for each factor, we take the log of the joint distribution over all variables and then average with respect to those variables not in that factor. Consider first the distribution over  $\alpha$ . Keeping only terms that have a functional dependence on  $\alpha$ , we have

$$\begin{aligned} \ln q^*(\alpha) &= \ln p(\alpha) + \mathbb{E}_{\mathbf{w}} [\ln p(\mathbf{w}|\alpha)] + \text{const} \\ &= (a_0 - 1) \ln \alpha - b_0 \alpha + \frac{M}{2} \ln \alpha - \frac{\alpha}{2} \mathbb{E}[\mathbf{w}^T \mathbf{w}] + \text{const}. \end{aligned} \quad (10.92)$$

We recognize this as the log of a gamma distribution, and so identifying the coefficients of  $\alpha$  and  $\ln \alpha$  we obtain

$$q^*(\alpha) = \text{Gam}(\alpha|a_N, b_N) \quad (10.93)$$

where

$$a_N = a_0 + \frac{M}{2} \quad (10.94)$$

$$b_N = b_0 + \frac{1}{2} \mathbb{E}[\mathbf{w}^T \mathbf{w}]. \quad (10.95)$$

// Similarly, we can find the variational re-estimation equation for the posterior distribution over  $\mathbf{w}$ . Again, using the general result (10.9), and keeping only those terms that have a functional dependence on  $\mathbf{w}$ , we have

$$\ln q^*(\mathbf{w}) = \ln p(\mathbf{t}|\mathbf{w}) + \mathbb{E}_{\alpha} [\ln p(\mathbf{w}|\alpha)] + \text{const} \quad (10.96)$$

$$= -\frac{\beta}{2} \sum_{n=1}^N \{\mathbf{w}^T \boldsymbol{\phi}_n - t_n\}^2 - \frac{1}{2} \mathbb{E}[\alpha] \mathbf{w}^T \mathbf{w} + \text{const} \quad (10.97)$$

$$= -\frac{1}{2} \mathbf{w}^T (\mathbb{E}[\alpha] \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}) \mathbf{w} + \beta \mathbf{w}^T \boldsymbol{\Phi}^T \mathbf{t} + \text{const}. \quad (10.98)$$

Because this is a quadratic form, the distribution  $q^*(\mathbf{w})$  is Gaussian, and so we can complete the square in the usual way to identify the mean and covariance, giving

$$q^*(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \quad (10.99)$$

上面给出了q(w)的迭代  
计算式，下面给出q(w)的  
迭代计算公式

where

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t} \quad (10.100)$$

$$\mathbf{S}_N = (\mathbb{E}[\alpha] \mathbf{I} + \beta \Phi^T \Phi)^{-1}. \quad (10.101)$$

Note the close similarity to the posterior distribution (3.52) obtained when  $\alpha$  was treated as a fixed parameter. The difference is that here  $\alpha$  is replaced by its expectation  $\mathbb{E}[\alpha]$  under the variational distribution. Indeed, we have chosen to use the same notation for the covariance matrix  $\mathbf{S}_N$  in both cases.

Using the standard results (B.27), (B.38), and (B.39), we can obtain the required moments as follows

$$\mathbb{E}[\alpha] = a_N / b_N \quad (10.102)$$

$$\mathbb{E}[\mathbf{w}\mathbf{w}^T] = \mathbf{m}_N \mathbf{m}_N^T + \mathbf{S}_N. \quad (10.103)$$

**这个过程**

The evaluation of the variational posterior distribution begins by initializing the parameters of one of the distributions  $q(\mathbf{w})$  or  $q(\alpha)$ , and then alternately re-estimates these factors in turn until a suitable convergence criterion is satisfied (usually specified in terms of the lower bound to be discussed shortly). ]

[ It is instructive to relate the variational solution to that found using the evidence framework in Section 3.5. To do this consider the case  $a_0 = b_0 = 0$ , corresponding to the limit of an infinitely broad prior over  $\alpha$ . The mean of the variational posterior distribution  $q(\alpha)$  is then given by

$$\mathbb{E}[\alpha] = \frac{a_N}{b_N} = \frac{M/2}{\mathbb{E}[\mathbf{w}^T \mathbf{w}] / 2} = \frac{M}{\mathbf{m}_N^T \mathbf{m}_N + \text{Tr}(\mathbf{S}_N)}. \quad (10.104)$$

Comparison with (9.63) shows that in the case of this particularly simple model, the variational approach gives precisely the same expression as that obtained by maximizing the evidence function using EM except that the point estimate for  $\alpha$  is replaced by its expected value. Because the distribution  $q(\mathbf{w})$  depends on  $q(\alpha)$  only through the expectation  $\mathbb{E}[\alpha]$ , we see that the two approaches will give identical results for the case of an infinitely broad prior. ]

给出了变分推断先验取  $a_0=b_0=0$   
时，它与使用 EM max model evidence 结果等价的原因

### 10.3.2 Predictive distribution

The predictive distribution over  $t$ , given a new input  $\mathbf{x}$ , is easily evaluated for this model using the Gaussian variational posterior for the parameters

$$\begin{aligned} p(t|\mathbf{x}, \mathbf{t}) &= \int p(t|\mathbf{x}, \mathbf{w}) p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \\ &\simeq \int p(t|\mathbf{x}, \mathbf{w}) q(\mathbf{w}) d\mathbf{w} \\ &= \int \mathcal{N}(t|\mathbf{w}^T \phi(\mathbf{x}), \beta^{-1}) \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma^2(\mathbf{x})) \end{aligned} \quad (10.105)$$

where we have evaluated the integral by making use of the result (2.115) for the linear-Gaussian model. Here the input-dependent variance is given by

$$\sigma^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}). \quad (10.106)$$

Note that this takes the same form as the result (3.59) obtained with fixed  $\alpha$  except that now the expected value  $\mathbb{E}[\alpha]$  appears in the definition of  $\mathbf{S}_N$ .

### 10.3.3 Lower bound

给出高分段  
的计算公式

Another quantity of importance is the lower bound  $\mathcal{L}$  defined by

$$\begin{aligned} \mathcal{L}(q) &= \mathbb{E}[\ln p(\mathbf{w}, \alpha, \mathbf{t})] - \mathbb{E}[\ln q(\mathbf{w}, \alpha)] \\ &= \mathbb{E}_{\mathbf{w}}[\ln p(\mathbf{t}|\mathbf{w})] + \mathbb{E}_{\mathbf{w}, \alpha}[\ln p(\mathbf{w}|\alpha)] + \mathbb{E}_{\alpha}[\ln p(\alpha)] \\ &\quad - \mathbb{E}_{\mathbf{w}}[\ln q(\mathbf{w})] - \mathbb{E}[\ln q(\alpha)]. \end{aligned} \quad (10.107)$$

#### Exercise 10.27

Evaluation of the various terms is straightforward, making use of results obtained in previous chapters, and gives

$$\begin{aligned} \mathbb{E}[\ln p(\mathbf{t}|\mathbf{w})]_{\mathbf{w}} &= \frac{N}{2} \ln \left( \frac{\beta}{2\pi} \right) - \frac{\beta}{2} \mathbf{t}^T \mathbf{t} + \beta \mathbf{m}_N^T \Phi^T \mathbf{t} \\ &\quad - \frac{\beta}{2} \text{Tr} [\Phi^T \Phi (\mathbf{m}_N \mathbf{m}_N^T + \mathbf{S}_N)] \end{aligned} \quad (10.108)$$

$$\begin{aligned} \mathbb{E}[\ln p(\mathbf{w}|\alpha)]_{\mathbf{w}, \alpha} &= -\frac{M}{2} \ln(2\pi) + \frac{M}{2} (\psi(a_N) - \ln b_N) \\ &\quad - \frac{a_N}{2b_N} [\mathbf{m}_N^T \mathbf{m}_N + \text{Tr}(\mathbf{S}_N)] \end{aligned} \quad (10.109)$$

$$\begin{aligned} \mathbb{E}[\ln p(\alpha)]_{\alpha} &= a_0 \ln b_0 + (a_0 - 1) [\psi(a_N) - \ln b_N] \\ &\quad - b_0 \frac{a_N}{b_N} - \ln \Gamma(a_N) \end{aligned} \quad (10.110)$$

$$-\mathbb{E}[\ln q(\mathbf{w})]_{\mathbf{w}} = \frac{1}{2} \ln |\mathbf{S}_N| + \frac{M}{2} [1 + \ln(2\pi)] \quad (10.111)$$

$$-\mathbb{E}[\ln q(\alpha)]_{\alpha} = \ln \Gamma(a_N) - (a_N - 1)\psi(a_N) - \ln b_N + a_N. \quad (10.112) \quad ]$$

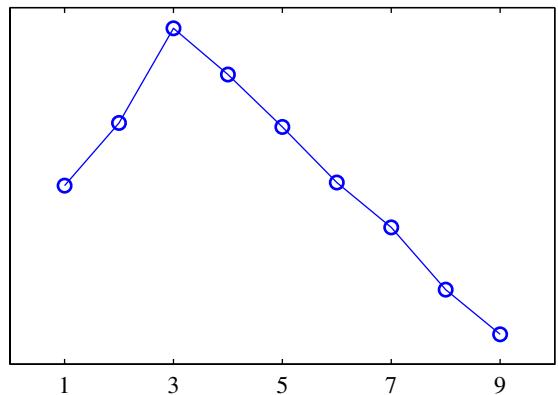
举例：

Figure 10.9 shows a plot of the lower bound  $\mathcal{L}(q)$  versus the degree of a polynomial model for a synthetic data set generated from a degree three polynomial. Here the prior parameters have been set to  $a_0 = b_0 = 0$ , corresponding to the noninformative prior  $p(\alpha) \propto 1/\alpha$ , which is uniform over  $\ln \alpha$  as discussed in Section 2.3.6. As we saw in Section 10.1, the quantity  $\mathcal{L}$  represents lower bound on the log marginal likelihood  $p(\mathbf{t}|M)$  for the model. If we assign equal prior probabilities  $p(M)$  to the different values of  $M$ , then we can interpret  $\mathcal{L}$  as an approximation to the posterior model probability  $p(M|\mathbf{t})$ . Thus the variational framework assigns the highest probability to the model with  $M = 3$ . This should be contrasted with the maximum likelihood result, which assigns ever smaller residual error to models of increasing complexity until the residual error is driven to zero, causing maximum likelihood to favour severely over-fitted models.

高分题断的高分段实际上 Condition on M,  
前面是为方便省略了未写即  $\ln p(\mathbf{t})$

$$\begin{aligned} p(x) &\propto 1/d, \quad y = 1/d \\ \text{根据概率密度函数} \\ d = e^y, \quad p(y) &\propto 1/d \cdot \frac{dy}{dy} \\ &= 1/e^y \cdot e^y \\ &= 1 \end{aligned}$$

**Figure 10.9** Plot of the lower bound  $\mathcal{L}$  versus the order  $M$  of the polynomial, for a polynomial model, in which a set of 10 data points is generated from a polynomial with  $M = 3$  sampled over the interval  $(-5, 5)$  with additive Gaussian noise of variance 0.09. The value of the bound gives the log probability of the model, and we see that the value of the bound peaks at  $M = 3$ , corresponding to the true model from which the data set was generated.



## 10.4. Exponential Family Distributions

In Chapter 2, we discussed the important role played by the exponential family of distributions and their conjugate priors. For many of the models discussed in this book, the complete-data likelihood is drawn from the exponential family. However, in general this will not be the case for the marginal likelihood function for the observed data. For example, in a mixture of Gaussians, the joint distribution of observations  $\mathbf{x}_n$  and corresponding hidden variables  $\mathbf{z}_n$  is a member of the exponential family, whereas the marginal distribution of  $\mathbf{x}_n$  is a mixture of Gaussians and hence is not.

隐藏量和参数的划分，但在数学上来说，二者没有本质上的区别。

Up to now we have grouped the variables in the model into observed variables and hidden variables. We now make a further distinction between latent variables, denoted  $\mathbf{Z}$ , and parameters, denoted  $\boldsymbol{\theta}$ , where parameters are *intensive* (fixed in number independent of the size of the data set), whereas latent variables are *extensive* (scale in number with the size of the data set). For example, in a Gaussian mixture model, the indicator variables  $z_{kn}$  (which specify which component  $k$  is responsible for generating data point  $\mathbf{x}_n$ ) represent the latent variables, whereas the means  $\mu_k$ , precisions  $\Lambda_k$  and mixing proportions  $\pi_k$  represent the parameters.

Consider the case of independent identically distributed data. We denote the data values by  $\mathbf{X} = \{\mathbf{x}_n\}$ , where  $n = 1, \dots, N$ , with corresponding latent variables  $\mathbf{Z} = \{\mathbf{z}_n\}$ . Now suppose that the joint distribution of observed and latent variables is a member of the exponential family, parameterized by natural parameters  $\boldsymbol{\eta}$  so that

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\eta}) = \prod_{n=1}^N h(\mathbf{x}_n, \mathbf{z}_n) g(\boldsymbol{\eta}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}_n, \mathbf{z}_n) \right\}. \quad (10.113)$$

We shall also use a conjugate prior for  $\boldsymbol{\eta}$ , which can be written as

$$p(\boldsymbol{\eta} | \nu_0, \chi_0) = f(\nu_0, \chi_0) g(\boldsymbol{\eta})^{\nu_0} \exp \left\{ \nu_0 \boldsymbol{\eta}^T \chi_0 \right\}. \quad (10.114)$$

Recall that the conjugate prior distribution can be interpreted as a prior number  $\nu_0$  of observations all having the value  $\chi_0$  for the  $\mathbf{u}$  vector. Now consider a variational

针对合隐藏量的后验分布及其  
拒绝检验，下面开始进行推断的推导

distribution that factorizes between the latent variables and the parameters, so that  $q(\mathbf{Z}, \boldsymbol{\eta}) = q(\mathbf{Z})q(\boldsymbol{\eta})$ . Using the general result (10.9), we can solve for the two factors as follows

$$\begin{aligned}\ln q^*(\mathbf{Z}) &= \mathbb{E}_{\boldsymbol{\eta}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\eta})] + \text{const} \\ &= \sum_{n=1}^N \{\ln h(\mathbf{x}_n, \mathbf{z}_n) + \mathbb{E}[\boldsymbol{\eta}^T] \mathbf{u}(\mathbf{x}_n, \mathbf{z}_n)\} + \text{const.}\end{aligned}\quad (10.115)$$

Thus we see that this decomposes into a sum of independent terms, one for each value of  $n$ , and hence the solution for  $q^*(\mathbf{Z})$  will factorize over  $n$  so that  $q^*(\mathbf{Z}) = \prod_n q^*(\mathbf{z}_n)$ . This is an example of an induced factorization. Taking the exponential of both sides, we have

$$q^*(\mathbf{z}_n) = h(\mathbf{x}_n, \mathbf{z}_n)g(\mathbb{E}[\boldsymbol{\eta}]) \exp\{\mathbb{E}[\boldsymbol{\eta}^T] \mathbf{u}(\mathbf{x}_n, \mathbf{z}_n)\} \quad (10.116)$$

where the normalization coefficient has been re-instated by comparison with the standard form for the exponential family.

// Similarly, for the variational distribution over the parameters, we have

$$\ln q^*(\boldsymbol{\eta}) = \ln p(\boldsymbol{\eta}|\boldsymbol{\nu}_0, \boldsymbol{\chi}_0) + \mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\eta})] + \text{const} \quad (10.117)$$

$$= \nu_0 \ln g(\boldsymbol{\eta}) + \boldsymbol{\nu}_{\boldsymbol{\eta}}^T \boldsymbol{\chi}_0 + \sum_{n=1}^N \{\ln g(\boldsymbol{\eta}) + \boldsymbol{\eta}^T \mathbb{E}_{\mathbf{z}_n}[\mathbf{u}(\mathbf{x}_n, \mathbf{z}_n)]\} + \text{const.} \quad (10.118)$$

Again, taking the exponential of both sides, and re-instating the normalization coefficient by inspection, we have

$$q^*(\boldsymbol{\eta}) = f(\nu_N, \boldsymbol{\chi}_N) g(\boldsymbol{\eta})^{\nu_N} \exp\left\{\boldsymbol{\nu}_{\boldsymbol{\eta}}^T \boldsymbol{\chi}_N\right\} \quad (10.119)$$

where we have defined

$$\nu_N = \nu_0 + N \quad (10.120)$$

$$\boldsymbol{\nu}_{\boldsymbol{\eta}} \boldsymbol{\chi}_N = \boldsymbol{\nu}_{\boldsymbol{\eta}} \boldsymbol{\chi}_0 + \sum_{n=1}^N \mathbb{E}_{\mathbf{z}_n}[\mathbf{u}(\mathbf{x}_n, \mathbf{z}_n)]. \quad (10.121)$$

Note that the solutions for  $q^*(\mathbf{z}_n)$  and  $q^*(\boldsymbol{\eta})$  are coupled, and so we solve them iteratively in a two-stage procedure. In the variational E step, we evaluate the expected sufficient statistics  $\mathbb{E}[\mathbf{u}(\mathbf{x}_n, \mathbf{z}_n)]$  using the current posterior distribution  $q(\mathbf{z}_n)$  over the latent variables and use this to compute a revised posterior distribution  $q(\boldsymbol{\eta})$  over the parameters. Then in the subsequent variational M step, we use this revised parameter posterior distribution to find the expected natural parameters  $\mathbb{E}[\boldsymbol{\eta}^T]$ , which gives rise to a revised variational distribution over the latent variables.

### 10.4.1 Variational message passing

We have illustrated the application of variational methods by considering a specific model, the Bayesian mixture of Gaussians, in some detail. This model can be

#### Section 10.2.5

最后分别给出了什么及  
如何通过迭代计算公式

described by the directed graph shown in Figure 10.5. Here we consider more generally the use of variational methods for models described by directed graphs and derive a number of widely applicable results.

The joint distribution corresponding to a directed graph can be written using the decomposition

$$p(\mathbf{x}) = \prod_i p(\mathbf{x}_i | \text{pa}_i) \quad (10.122)$$

where  $\mathbf{x}_i$  denotes the variable(s) associated with node  $i$ , and  $\text{pa}_i$  denotes the parent set corresponding to node  $i$ . Note that  $\mathbf{x}_i$  may be a latent variable or it may belong to the set of observed variables. Now consider a variational approximation in which the distribution  $q(\mathbf{x})$  is assumed to factorize with respect to the  $\mathbf{x}_i$  so that

$$q(\mathbf{x}) = \prod_i q_i(\mathbf{x}_i). \quad (10.123)$$

Note that for observed nodes, there is no factor  $q(\mathbf{x}_i)$  in the variational distribution. We now substitute (10.122) into our general result (10.9) to give

$$\ln q_j^*(\mathbf{x}_j) = \mathbb{E}_{i \neq j} \left[ \sum_i \ln p(\mathbf{x}_i | \text{pa}_i) \right] + \text{const.} \quad (10.124)$$

Any terms on the right-hand side that do not depend on  $\mathbf{x}_j$  can be absorbed into the additive constant. In fact, the only terms that do depend on  $\mathbf{x}_j$  are the conditional distribution for  $\mathbf{x}_j$  given by  $p(\mathbf{x}_j | \text{pa}_j)$  together with any other conditional distributions that have  $\mathbf{x}_j$  in the conditioning set. By definition, these conditional distributions correspond to the children of node  $j$ , and they therefore also depend on the *co-parents* of the child nodes, i.e., the other parents of the child nodes besides node  $\mathbf{x}_j$  itself. We see that the set of all nodes on which  $q_j^*(\mathbf{x}_j)$  depends corresponds to the Markov blanket of node  $\mathbf{x}_j$ , as illustrated in Figure 8.26. Thus the update of the factors in the variational posterior distribution represents a local calculation on the graph. This makes possible the construction of general purpose software for variational inference in which the form of the model does not need to be specified in advance (Bishop *et al.*, 2003).

结论：

**应用示例：** If we now specialize to the case of a model in which all of the conditional distributions have a conjugate-exponential structure, then the variational update procedure can be cast in terms of a local message passing algorithm (Winn and Bishop, 2005). In particular, the distribution associated with a particular node can be updated once that node has received messages from all of its parents and all of its children. This in turn requires that the children have already received messages from their co-parents. The evaluation of the lower bound can also be simplified because many of the required quantities are already evaluated as part of the message passing scheme. This distributed message passing formulation has good scaling properties and is well suited to large networks.

VI的特别 Variational  
message passing algorithm

对称 BP算子的特别

BP 10.7.2小节介绍该信息  
传递算法 Sum-product

algorithm / BP algorithm,

但这里并未详细地介绍 Variational  
message passing algorithm.