

## 10.5. Local Variational Methods

The variational framework discussed in Sections 10.1 and 10.2 can be considered a ‘global’ method in the sense that it directly seeks an approximation to the full posterior distribution over all random variables. An alternative ‘local’ approach involves finding bounds on functions over individual variables or groups of variables within a model. For instance, we might seek a bound on a conditional distribution  $p(y|x)$ , which is itself just one factor in a much larger probabilistic model specified by a directed graph. The purpose of introducing the bound of course is to simplify the resulting distribution. This local approximation can be applied to multiple variables in turn until a tractable approximation is obtained, and in Section 10.6.1 we shall give a practical example of this approach in the context of logistic regression. Here we focus on developing the bounds themselves.

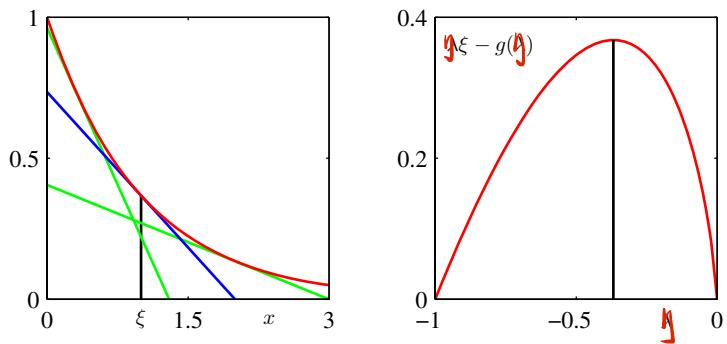
We have already seen in our discussion of the Kullback-Leibler divergence that the convexity of the logarithm function played a key role in developing the lower bound in the global variational approach. We have defined a (strictly) convex function as one for which every chord lies above the function. Convexity also plays a central role in the local variational framework. Note that our discussion will apply equally to concave functions with ‘min’ and ‘max’ interchanged and with lower bounds replaced by upper bounds.

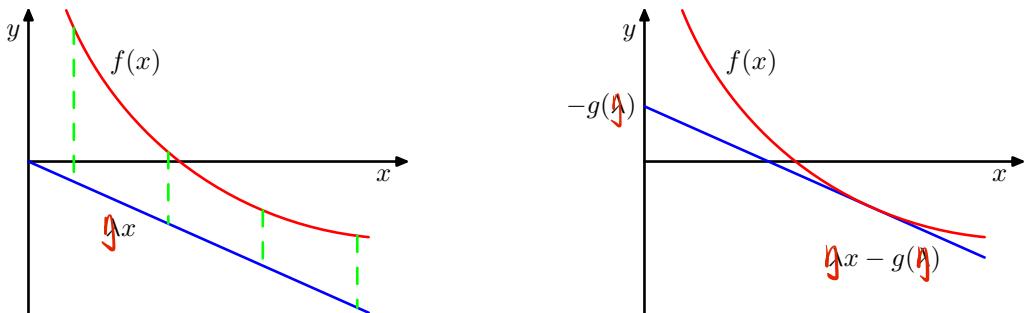
**举例子** Let us begin by considering a simple example, namely the function  $f(x) = \exp(-x)$ , which is a convex function of  $x$ , and which is shown in the left-hand plot of Figure 10.10. Our goal is to approximate  $f(x)$  by a simpler function, in particular a linear function of  $x$ . From Figure 10.10, we see that this linear function will be a lower bound on  $f(x)$  if it corresponds to a tangent. We can obtain the tangent line  $y(x)$  at a specific value of  $x$ , say  $x = \xi$ , by making a first order Taylor expansion

$$y(x) = f(\xi) + f'(\xi)(x - \xi) \quad (10.125)$$

so that  $y(x) \leq f(x)$  with equality when  $x = \xi$ . For our example function  $f(x) =$

**Figure 10.10** In the left-hand figure the red curve shows the function  $\exp(-x)$ , and the blue line shows the tangent at  $x = \xi$  defined by (10.125) with  $\xi = 1$ . This line has slope  $\frac{dy}{dx} = f'(\xi) = -\exp(-\xi)$ . Note that any other tangent line, for example the ones shown in green, will have a smaller value of  $y$  at  $x = \xi$ . The right-hand figure shows the corresponding plot of the function  $y = \xi - g(\frac{x-\xi}{\xi})$ , where  $g(\frac{x-\xi}{\xi})$  is given by (10.131), versus  $\frac{x-\xi}{\xi}$  for  $\xi = 1$ , in which the maximum corresponds to  $\frac{x-\xi}{\xi} = -\exp(-\xi) = -1/e$ .





**Figure 10.11** In the left-hand plot the red curve shows a convex function  $f(x)$ , and the blue line represents the linear function  $\lambda x$ , which is a lower bound on  $f(x)$  because  $f(x) > \lambda x$  for all  $x$ . For the given value of slope  $\lambda$  the contact point of the tangent line having the same slope is found by minimizing with respect to  $x$  the discrepancy (shown by the green dashed lines) given by  $f(x) - \lambda x$ . This defines the dual function  $g(\lambda)$ , which corresponds to the (negative of the) intercept of the tangent line having slope  $\lambda$ .

$\exp(-x)$ , we therefore obtain the tangent line in the form

$$y(x) = \exp(-\xi) - \exp(-\xi)(x - \xi) \quad (10.126)$$

which is a linear function parameterized by  $\xi$ . For consistency with subsequent discussion, let us define  $\lambda = -\exp(-\xi)$  so that

$$y(x, \lambda) = \lambda x - \lambda + \lambda \ln(-\lambda). \quad (10.127)$$

Different values of  $\lambda$  correspond to different tangent lines, and because all such lines are lower bounds on the function, we have  $f(x) \geq y(x, \lambda)$ . Thus we can write the function in the form

$$f(x) = \max_{\lambda} \{ \lambda x - \lambda + \lambda \ln(-\lambda) \}. \quad (10.128)$$

We have succeeded in approximating the convex function  $f(x)$  by a simpler, linear function  $y(x, \lambda)$ . The price we have paid is that we have introduced a variational parameter  $\lambda$ , and to obtain the tightest bound we must optimize with respect to  $\lambda$ . ]

【**拓展总结：** We can formulate this approach more generally using the framework of convex duality (Rockafellar, 1972; Jordan *et al.*, 1999). Consider the illustration of a convex function  $f(x)$  shown in the left-hand plot in Figure 10.11. In this example, the function  $\lambda x$  is a lower bound on  $f(x)$  but it is not the best lower bound that can be achieved by a linear function having slope  $\lambda$ , because the tightest bound is given by the tangent line. Let us write the equation of the tangent line, having slope  $\lambda$  as  $\lambda x - g(\lambda)$  where the (negative) intercept  $g(\lambda)$  clearly depends on the slope  $\lambda$  of the tangent. To determine the intercept, we note that the line must be moved vertically by an amount equal to the smallest vertical distance between the line and the function, as shown in Figure 10.11. Thus

$$\begin{aligned} g(\lambda) &= -\min_x \{f(x) - \lambda x\} \\ &= \max_x \{\lambda x - f(x)\}. \end{aligned} \quad (10.129)$$

$f(x) \leq g(\xi)$  in 2D //  
凸性

Now, instead of fixing  $\xi$  and varying  $x$ , we can consider a particular  $x$  and then adjust  $\xi$  until the tangent plane is tangent at that particular  $x$ . Because the  $y$  value of the tangent line at a particular  $x$  is maximized when that value coincides with its contact point, we have

$$f(x) = \max_{\xi} \{ \xi x - g(\xi) \}. \quad (10.130)$$

We see that the functions  $f(x)$  and  $g(\xi)$  play a dual role, and are related through (10.129) and (10.130).

Let us apply these duality relations to our simple example  $f(x) = \exp(-x)$ . From (10.129) we see that the maximizing value of  $x$  is given by  $\xi = -\ln(-\lambda)$ , and back-substituting we obtain the conjugate function  $g(\xi)$  in the form

$$g(\xi) = \lambda - \lambda \ln(-\lambda) \quad (10.131)$$

as obtained previously. The function  $\xi f(x) - g(\xi)$  is shown, for  $\xi = 1$  in the right-hand plot in Figure 10.10. As a check, we can substitute (10.131) into (10.130), which gives the maximizing value of  $\lambda = -\exp(-x)$ , and back-substituting then recovers the original function  $f(x) = \exp(-x)$ .

**Concave functions** For concave functions, we can follow a similar argument to obtain upper bounds, in which ‘max’ is replaced with ‘min’, so that

$$f(x) = \min_{\xi} \{ \xi x - g(\xi) \} \quad (10.132)$$

$$g(\xi) = \min_x \{ \xi x - f(x) \}. \quad (10.133)$$

**不具有凸性时的处理:** If the function of interest is not convex (or concave), then we cannot directly apply the method above to obtain a bound. However, we can first seek (invertible transformations) either of (the function) or of (its argument) which change it into a convex form. We then calculate the conjugate function and then transform back to the original variables.

[ An important example, which arises frequently in pattern recognition, is the logistic sigmoid function defined by

$$\sigma(x) = \frac{1}{1 + e^{-x}}. \quad (10.134)$$

As it stands this function is neither convex nor concave. However, if we take the logarithm we obtain a function which is concave, as is easily verified by finding the second derivative. From (10.133) the corresponding conjugate function then takes the form

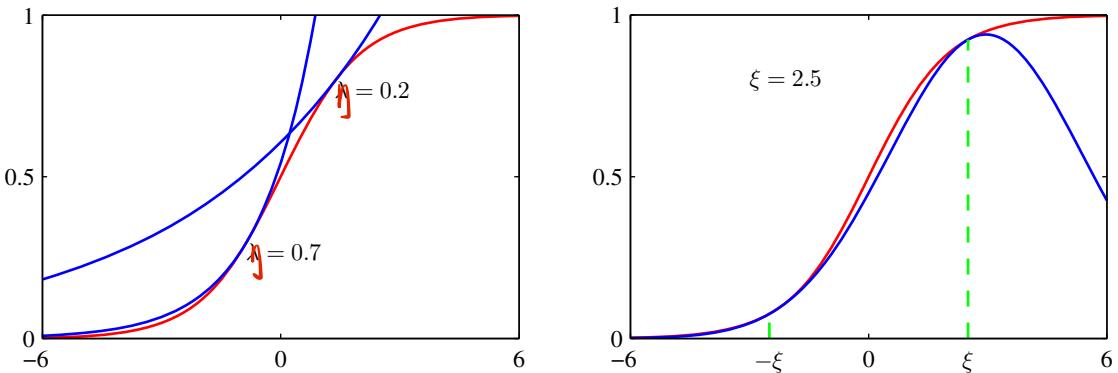
$$g(\lambda) = \min_x \{ \lambda x - f(x) \} = -\lambda \ln \lambda - (1 - \lambda) \ln(1 - \lambda) \quad (10.135)$$

### Exercise 10.30

### Appendix B

which we recognize as the binary entropy function for a variable whose probability of having the value 1 is  $\lambda$ . Using (10.132), we then obtain an upper bound on the log sigmoid

$$\ln \sigma(x) \leq \lambda x - g(\lambda) \quad (10.136)$$



**Figure 10.12** The left-hand plot shows the logistic sigmoid function  $\sigma(x)$  defined by (10.134) in red, together with two examples of the exponential upper bound (10.137) shown in blue. The right-hand plot shows the logistic sigmoid again in red together with the Gaussian lower bound (10.144) shown in blue. Here the parameter  $\xi = 2.5$ , and the bound is exact at  $x = \xi$  and  $x = -\xi$ , denoted by the dashed green lines.

and taking the exponential, we obtain an upper bound on the logistic sigmoid itself of the form

$$\sigma(x) \leq \exp(\eta x - g(\eta)) \quad (10.137)$$

which is plotted for two values of  $\eta$  on the left-hand plot in Figure 10.12.

// We can also obtain a lower bound on the sigmoid having the functional form of a Gaussian. To do this, we follow Jaakkola and Jordan (2000) and make transformations both of the input variable and of the function itself. First we take the log of the logistic function and then decompose it so that

$$\begin{aligned} \ln \sigma(x) &= -\ln(1 + e^{-x}) = -\ln \{e^{-x/2}(e^{x/2} + e^{-x/2})\} \\ &= x/2 - \ln(e^{x/2} + e^{-x/2}). \end{aligned} \quad (10.138)$$

### Exercise 10.31

We now note that the function  $f(x) = -\ln(e^{x/2} + e^{-x/2})$  is a convex function of the variable  $x^2$ , as can again be verified by finding the second derivative. This leads to a lower bound on  $f(x)$ , which is a linear function of  $x^2$  whose conjugate function is given by

$$g(\eta) = \max_{x^2} \left\{ \eta x^2 - f\left(\sqrt{x^2}\right) \right\}. \quad (10.139)$$

The stationarity condition leads to

$$0 = \eta - \frac{dx}{dx^2} \frac{d}{dx} f(x) = \eta + \frac{1}{4x} \tanh\left(\frac{x}{2}\right). \quad (10.140)$$

If we denote this value of  $x$ , corresponding to the contact point of the tangent line for this particular value of  $\eta$ , by  $\xi$ , then we have

$$\eta \lambda(\eta) = -\frac{1}{4\xi} \tanh\left(\frac{\xi}{2}\right) = -\frac{1}{2\xi} \left[ \sigma(\xi) - \frac{1}{2} \right] = -\lambda(\xi) \quad (10.141)$$

where we have defined  $\lambda = -\eta$  to maintain consistency with Jaakkola and Jordan (2000).

上面构造了  $\sigma(x)$   
in local 上界，下面开始  
构造其 local 下界

Instead of thinking of  $\lambda$  as the variational parameter, we can let  $\xi$  play this role as this leads to simpler expressions for the conjugate function, which is then given by

$$g(\xi) = -\lambda(\xi)\xi^2 - f(\xi) = -\lambda(\xi)\xi^2 + \ln(e^{\xi/2} + e^{-\xi/2}). \quad (10.142)$$

Hence the bound on  $f(x)$  can be written as

$$f(x) \geq g(\xi) = -\lambda(\xi)x^2 + \xi^2 - \ln(e^{\xi/2} + e^{-\xi/2}). \quad (10.143)$$

The bound on the sigmoid then becomes

$$\sigma(x) \geq \sigma(\xi) \exp\{(x - \xi)/2 - \lambda(\xi)(x^2 - \xi^2)\} \quad (10.144)$$

where  $\lambda(\xi)$  is defined by (10.141). This bound is illustrated in the right-hand plot of Figure 10.12. [We see that the bound has the form of the exponential of a quadratic function of  $x$ , which will prove useful when we seek Gaussian representations of posterior distributions defined through logistic sigmoid functions.]

### Section 4.5

~~上述推导并不直接推广到 softmax 函数中~~

### Section 4.3

The logistic sigmoid arises frequently in probabilistic models over binary variables because it is the function that transforms a log odds ratio into a posterior probability. The corresponding transformation for a multiclass distribution is given by the softmax function. Unfortunately, the lower bound derived here for the logistic sigmoid does not directly extend to the softmax. Gibbs (1997) proposes a method for constructing a Gaussian distribution that is conjectured to be a bound (although no rigorous proof is given), which may be used to apply local variational methods to multiclass problems.

[ We shall see an example of the use of local variational bounds in Sections 10.6.1. For the moment, however, it is instructive to consider in general terms how these bounds can be used. Suppose we wish to evaluate an integral of the form

$$I = \int \sigma(a)p(a) da \quad (10.145)$$

where  $\sigma(a)$  is the logistic sigmoid, and  $p(a)$  is a Gaussian probability density. Such integrals arise in Bayesian models when, for instance, we wish to evaluate the predictive distribution, in which case  $p(a)$  represents a posterior parameter distribution. Because the integral is intractable, we employ the variational bound (10.144), which we write in the form  $\sigma(a) \geq f(a, \xi)$  where  $\xi$  is a variational parameter. The integral now becomes the product of two exponential-quadratic functions and so can be integrated analytically to give a bound on  $I$

$$I \geq \int f(a, \xi)p(a) da = F(\xi). \quad (10.146)$$

We now have the freedom to choose the variational parameter  $\xi$ , which we do by finding the value  $\xi^*$  that maximizes the function  $F(\xi)$ . [The resulting value  $F(\xi^*)$  represents the tightest bound within this family of bounds and can be used as an approximation to  $I$ . This optimized bound, however, will in general not be exact.

【如何离散下界的应用  
举例】

指  $f(a)$  而不是  $F(\xi)$

Although the bound  $\sigma(a) \geq f(a, \xi)$  on the logistic sigmoid can be optimized exactly, the required choice for  $\xi$  depends on the value of  $a$ , so that the bound is exact for one value of  $a$  only. Because the quantity  $F(\xi)$  is obtained by integrating over all values of  $a$ , the value of  $\xi^*$  represents a compromise, weighted by the distribution  $p(a)$ .

## 10.6. Variational Logistic Regression

*对比 Laplace approximation*

We now illustrate the use of local variational methods by returning to the Bayesian logistic regression model studied in Section 4.5. There we focussed on the use of the Laplace approximation, while here we consider a variational treatment based on the approach of Jaakkola and Jordan (2000). Like the Laplace method, this also leads to a Gaussian approximation to the posterior distribution. However, the greater flexibility of the variational approximation leads to improved accuracy compared to the Laplace method. Furthermore (unlike the Laplace method), the variational approach is optimizing a well defined objective function given by a rigorous bound on the model evidence. Logistic regression has also been treated by Dydowski and Roberts (2005) from a Bayesian perspective using Monte Carlo sampling techniques.

*参数先验*  
 $P(w) = \mathcal{N}(w | m_0, S_0)$

### 10.6.1 Variational posterior distribution: 估算 $P(w|t)$ 的近似后验 $q(w)$

Here we shall make use of a variational approximation based on the local bounds introduced in Section 10.5. This allows the likelihood function for logistic regression, which is governed by the logistic sigmoid, to be approximated by the exponential of a quadratic form. It is therefore again convenient to choose a conjugate Gaussian prior of the form (4.140). For the moment, we shall treat the hyperparameters  $m_0$  and  $S_0$  as fixed constants. In Section 10.6.3, we shall demonstrate how the variational formalism can be extended to the case where there are unknown hyperparameters whose values are to be inferred from the data.

In the variational framework, we seek to maximize a lower bound on the marginal likelihood. For the Bayesian logistic regression model, the marginal likelihood takes the form

$$p(t) = \int p(t|w)p(w) dw = \int \left[ \prod_{n=1}^N p(t_n|w) \right] p(w) dw. \quad (10.147)$$

We first note that the conditional distribution for  $t$  can be written as

$$\begin{aligned} p(t|w) &= \sigma(a)^t \{1 - \sigma(a)\}^{1-t} \\ &= \left( \frac{1}{1 + e^{-a}} \right)^t \left( 1 - \frac{1}{1 + e^{-a}} \right)^{1-t} \\ &= e^{at} \frac{e^{-a}}{1 + e^{-a}} = e^{at} \sigma(-a) \end{aligned} \quad (10.148)$$

where  $a = w^T \phi$ . In order to obtain a lower bound on  $p(t)$ , we make use of the variational lower bound on the logistic sigmoid function given by (10.144), which

we reproduce here for convenience

$$\sigma(z) \geq \sigma(\xi) \exp \left\{ (z - \xi)/2 - \lambda(\xi)(z^2 - \xi^2) \right\} \quad (10.149)$$

where

$$\lambda(\xi) = \frac{1}{2\xi} \left[ \sigma(\xi) - \frac{1}{2} \right]. \quad (10.150)$$

We can therefore write

$$p(t|\mathbf{w}) = e^{at} \sigma(-a) \geq e^{at} \sigma(\xi) \exp \left\{ -(a + \xi)/2 - \lambda(\xi)(a^2 - \xi^2) \right\}. \quad (10.151)$$

{ Note that because this bound is applied to each of the terms in the likelihood function separately, there is a variational parameter  $\xi_n$  corresponding to each training set observation  $(\phi_n, t_n)$ . Using  $a = \mathbf{w}^\top \phi$ , and multiplying by the prior distribution, we obtain the following bound on the joint distribution of  $\mathbf{t}$  and  $\mathbf{w}$

$$p(\mathbf{t}, \mathbf{w}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}) \geq h(\mathbf{w}, \boldsymbol{\xi})p(\mathbf{w}) \quad (10.152)$$

where  $\boldsymbol{\xi}$  denotes the set  $\{\xi_n\}$  of variational parameters, and

$$\begin{aligned} h(\mathbf{w}, \boldsymbol{\xi}) &= \prod_{n=1}^N \sigma(\xi_n) \exp \left\{ \mathbf{w}^\top \phi_n t_n - (\mathbf{w}^\top \phi_n + \xi_n)/2 \right. \\ &\quad \left. - \lambda(\xi_n)([\mathbf{w}^\top \phi_n]^2 - \xi_n^2) \right\}. \end{aligned} \quad (10.153)$$

Evaluation of the exact posterior distribution would require normalization of the left-hand side of this inequality. Because this is intractable, we work instead with the right-hand side. Note that the function on the right-hand side cannot be interpreted as a probability density because it is not normalized. Once it is normalized to give a variational posterior distribution  $q(\mathbf{w})$ , however, it no longer represents a bound.

Because the logarithm function is monotonically increasing, the inequality  $A \geq B$  implies  $\ln A \geq \ln B$ . This gives a lower bound on the log of the joint distribution of  $\mathbf{t}$  and  $\mathbf{w}$  of the form

$$\begin{aligned} \ln \{p(\mathbf{t}|\mathbf{w})p(\mathbf{w})\} &\geq \ln p(\mathbf{w}) + \sum_{n=1}^N \left\{ \ln \sigma(\xi_n) + \mathbf{w}^\top \phi_n t_n \right. \\ &\quad \left. - (\mathbf{w}^\top \phi_n + \xi_n)/2 - \lambda(\xi_n)([\mathbf{w}^\top \phi_n]^2 - \xi_n^2) \right\}. \end{aligned} \quad (10.154)$$

Substituting for the prior  $p(\mathbf{w})$ , the right-hand side of this inequality becomes, as a function of  $\mathbf{w}$

$$\begin{aligned} &- \frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^\top \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) \\ &+ \sum_{n=1}^N \left\{ \mathbf{w}^\top \phi_n (t_n - 1/2) - \lambda(\xi_n) \mathbf{w}^\top (\phi_n \phi_n^\top) \mathbf{w} \right\} + \text{const.} \end{aligned} \quad (10.155)$$

這：式(10.152)右边是下界，但是一個概率分布，而一旦归一化为概率分布，则其不再为下界。

高次項得  
參數  $w$  的變動進行  
近似後驗  $q(w)$   
注意，這裡是五個而  
不是下界，只不過這  
一近似是由下界推  
導而來的

This is a quadratic function of  $w$ , and so we can obtain the corresponding variational approximation to the posterior distribution by identifying the linear and quadratic terms in  $w$ , giving a Gaussian variational posterior of the form

$$q(w) = \mathcal{N}(w | m_N, S_N) \quad (10.156)$$

where

$$m_N = S_N \left( S_0^{-1} m_0 + \sum_{n=1}^N (t_n - 1/2) \phi_n \right) \quad (10.157)$$

$$S_N^{-1} = S_0^{-1} + 2 \sum_{n=1}^N \lambda(\xi_n) \phi_n \phi_n^T. \quad (10.158)$$

$q(w)$  帶有高分參數  
 $\{\xi_n\}$  因此相比 Laplace  
approximation 能提供  
更灵活精準的近似

### Exercise 10.32

對應到 sequential learning (結合 10.6.2 \{\xi\_n\} 的學習)  
Note that the bound given by (10.149) applies only to the two-class problem and so this approach does not directly generalize to classification problems with  $K > 2$  classes. An alternative bound for the multiclass case has been explored by Gibbs (1997).

## 10.6.2 Optimizing the variational parameters

We now have a normalized Gaussian approximation to the posterior distribution, which we shall use shortly to evaluate the predictive distribution for new data points. First, however, we need to determine the variational parameters  $\{\xi_n\}$  by maximizing the lower bound on the marginal likelihood.

To do this, we substitute the inequality (10.152) back into the marginal likelihood to give

$$\ln p(t) = \ln \int p(t|w)p(w) dw \geq \ln \int h(w, \xi) p(w) dw = \mathcal{L}(\xi). \quad (10.159)$$

As with the optimization of the hyperparameter  $\alpha$  in the linear regression model of Section 3.5, there are two approaches to determining the  $\xi_n$ . In the first approach, we recognize that the function  $\mathcal{L}(\xi)$  is defined by an integration over  $w$  and so we can view  $w$  as a latent variable and invoke the EM algorithm. In the second approach, we integrate over  $w$  analytically and then perform a direct maximization over  $\xi$ . Let us begin by considering the EM approach.

① [The EM algorithm starts by choosing some initial values for the parameters  $\{\xi_n\}$ , which we denote collectively by  $\xi^{\text{old}}$ . In the E step of the EM algorithm,

這裡的 EM 算法並不要謬： $h(w, \xi)$  是  $p(t|w)$  的近似，但那並不是一個概率分布，因此严格来讲这里并不符合EM算法的使用场景，这里只是硬套EM的计算过程罢了。

we then use these parameter values to find the posterior distribution over  $\mathbf{w}$ , which is given by (10.156). In the M step, we then maximize the expected complete-data log likelihood which is given by

$$Q(\xi, \xi^{\text{old}}) = \mathbb{E} [\ln h(\mathbf{w}, \xi) p(\mathbf{w})] \quad (10.160)$$

where the expectation is taken with respect to the posterior distribution  $q(\mathbf{w})$  evaluated using  $\xi^{\text{old}}$ . Noting that  $p(\mathbf{w})$  does not depend on  $\xi$ , and substituting for  $h(\mathbf{w}, \xi)$  we obtain

$$Q(\xi, \xi^{\text{old}}) = \sum_{n=1}^N \left\{ \ln \sigma(\xi_n) - \xi_n/2 - \lambda(\xi_n) (\phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n - \xi_n^2) \right\} + \text{const} \quad (10.161)$$

where ‘const’ denotes terms that are independent of  $\xi$ . We now set the derivative with respect to  $\xi_n$  equal to zero. A few lines of algebra, making use of the definitions of  $\sigma(\xi)$  and  $\lambda(\xi)$ , then gives

$$0 = \lambda'(\xi_n) (\phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n - \xi_n^2). \quad (10.162)$$

We now note that  $\lambda'(\xi)$  is a monotonic function of  $\xi$  for  $\xi \geq 0$ , and that we can restrict attention to nonnegative values of  $\xi$  without loss of generality due to the symmetry of the bound around  $\xi = 0$ . Thus  $\lambda'(\xi) \neq 0$ , and hence we obtain the following re-estimation equations

$$(\xi_n^{\text{new}})^2 = \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n = \phi_n^T (\mathbf{S}_N + \mathbf{m}_N \mathbf{m}_N^T) \phi_n \quad (10.163)$$

where we have used (10.156).

**EM算法计算过程的总结** = Let us summarize the EM algorithm for finding the variational posterior distribution. We first initialize the variational parameters  $\xi^{\text{old}}$ . In the E step, we evaluate the posterior distribution over  $\mathbf{w}$  given by (10.156), in which the mean and covariance are defined by (10.157) and (10.158). In the M step, we then use this variational posterior to compute a new value for  $\xi$  given by (10.163). The E and M steps are repeated until a suitable convergence criterion is satisfied, which in practice typically requires only a few iterations.

[An alternative approach to obtaining re-estimation equations for  $\xi$  is to note that in the integral over  $\mathbf{w}$  in the definition (10.159) of the lower bound  $\mathcal{L}(\xi)$ , the integrand has a Gaussian-like form and so the integral can be evaluated analytically. Having evaluated the integral, we can then differentiate with respect to  $\xi_n$ . It turns out that this gives rise to exactly the same re-estimation equations as does the EM approach given by (10.163).]

As we have emphasized already, in the application of variational methods it is useful to be able to evaluate the lower bound  $\mathcal{L}(\xi)$  given by (10.159). The integration over  $\mathbf{w}$  can be performed analytically by noting that  $p(\mathbf{w})$  is Gaussian and  $h(\mathbf{w}, \xi)$  is the exponential of a quadratic function of  $\mathbf{w}$ . Thus, by completing the square and making use of the standard result for the normalization coefficient of a Gaussian distribution, we can obtain a closed form solution which takes the form

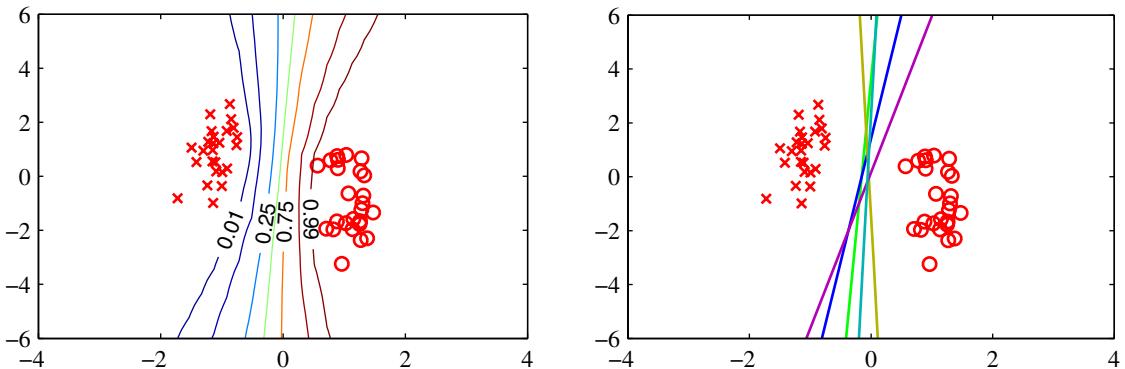
如何解析计算公式

Exercise 10.35

Exercise 10.34

如何解析计算公式

Exercise 10.35



**Figure 10.13** Illustration of the Bayesian approach to logistic regression for a simple linearly separable data set. The plot on the left shows the predictive distribution obtained using variational inference. We see that the decision boundary lies roughly mid way between the clusters of data points, and that the contours of the predictive distribution splay out away from the data reflecting the greater uncertainty in the classification of such regions. The plot on the right shows the decision boundaries corresponding to five samples of the parameter vector  $w$  drawn from the posterior distribution  $p(w|t)$ .

$$\begin{aligned} \mathcal{L}(\xi) = & \frac{1}{2} \ln \frac{|\mathbf{S}_N|}{|\mathbf{S}_0|} + \frac{1}{2} \mathbf{m}_N^T \mathbf{S}_N^{-1} \mathbf{m}_N - \frac{1}{2} \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0 \\ & + \sum_{n=1}^N \left\{ \ln \sigma(\xi_n) - \frac{1}{2} \xi_n + \lambda(\xi_n) \xi_n^2 \right\}. \end{aligned} \quad (10.164) \boxed{1}$$

如何进行 sequential learning, This variational framework can also be applied to situations in which the data is arriving sequentially (Jaakkola and Jordan, 2000). In this case we maintain a Gaussian posterior distribution over  $w$ , which is initialized using the prior  $p(w)$ . As each data point arrives, the posterior is updated by making use of the bound (10.151) and then normalized to give an updated posterior distribution.

**预测分布** The predictive distribution is obtained by marginalizing over the posterior distribution, and takes the same form as for the Laplace approximation discussed in Section 4.5.2. Figure 10.13 shows the variational predictive distributions for a synthetic data set. This example provides interesting insights into the concept of ‘large margin’, which was discussed in Section 7.1 and which has qualitatively similar behaviour to the Bayesian solution.

### 10.6.3 Inference of hyperparameters

So far, we have treated the hyperparameter  $\alpha$  in the prior distribution as a known constant. We now extend the Bayesian logistic regression model to allow the value of this parameter to be inferred from the data set. This can be achieved by combining the global and local variational approximations into a single framework, so as to maintain a lower bound on the marginal likelihood at each stage. Such a combined approach was adopted by Bishop and Svensén (2003) in the context of a Bayesian treatment of the hierarchical mixture of experts model.

注意，10.6.1和10.6.2是一脉  
相承的： $m_0, S_0$ 是超参数，  
直接给定， $\{\xi_n\}$ 是差分  
超参数，在10.6.2小节中进行  
了优化，参数由后验给定  
实际上 Condition on  $m_0, S_0, \{\xi_n\}$

在10.6.1小节给定，而10.6.3小节则自成体系：前文中的超参数  $m_0, S_0$  对应这里的参数  $\alpha$ ，并引入  
先验分布  $P(\alpha) = \text{Gam}(\alpha | a_0, b_0)$ ，参数由后验给定此时为  $(\alpha | \mathbf{x}, \mathbf{y})$ ，多了个参数  $\alpha$ ；  
超参数  $a_0, b_0$  直接给定剩下的部分超参数  $\{\xi_n\}$  则仍然通过迭代得到。

Specifically, we consider once again a simple isotropic Gaussian prior distribution of the form

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}). \quad (10.165)$$

Our analysis is readily extended to more general Gaussian priors, for instance if we wish to associate a different hyperparameter with different subsets of the parameters  $w_j$ . As usual, we consider a conjugate hyperprior over  $\alpha$  given by a gamma distribution

$$p(\alpha) = \text{Gam}(\alpha|a_0, b_0) \quad (10.166)$$

governed by the constants  $a_0$  and  $b_0$ .

The marginal likelihood for this model now takes the form

$$p(\mathbf{t}) = \iint p(\mathbf{w}, \alpha, \mathbf{t}) d\mathbf{w} d\alpha \quad (10.167)$$

where the joint distribution is given by

$$p(\mathbf{w}, \alpha, \mathbf{t}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}|\alpha)p(\alpha). \quad (10.168)$$

We are now faced with an analytically intractable integration over  $\mathbf{w}$  and  $\alpha$ , which we shall tackle by using both the (local) and (global) variational approaches in the same model.

上面给出了背景条件，下面开始推导 // To begin with, we introduce a variational distribution  $q(\mathbf{w}, \alpha)$ , and then apply the decomposition (10.2), which in this instance takes the form

$$\ln p(\mathbf{t}) = \mathcal{L}(q) + \text{KL}(q\|p) \quad (10.169)$$

where the lower bound  $\mathcal{L}(q)$  and the Kullback-Leibler divergence  $\text{KL}(q\|p)$  are defined by

$$\mathcal{L}(q) = \iint q(\mathbf{w}, \alpha) \ln \left\{ \frac{p(\mathbf{w}, \alpha, \mathbf{t})}{q(\mathbf{w}, \alpha)} \right\} d\mathbf{w} d\alpha \quad (10.170)$$

$$\text{KL}(q\|p) = - \iint q(\mathbf{w}, \alpha) \ln \left\{ \frac{p(\mathbf{w}, \alpha|\mathbf{t})}{q(\mathbf{w}, \alpha)} \right\} d\mathbf{w} d\alpha. \quad (10.171)$$

At this point, the lower bound  $\mathcal{L}(q)$  is still intractable due to the form of the likelihood factor  $p(\mathbf{t}|\mathbf{w})$ . We therefore apply the local variational bound to each of the logistic sigmoid factors as before. This allows us to use the inequality (10.152) and place a lower bound on  $\mathcal{L}(q)$ , which will therefore also be a lower bound on the log marginal likelihood

$$\begin{aligned} \ln p(\mathbf{t}) &\stackrel{\text{①}}{\geq} \mathcal{L}(q) \stackrel{\text{②}}{\geq} \tilde{\mathcal{L}}(q, \xi) \\ &= \iint q(\mathbf{w}, \alpha) \ln \left\{ \frac{h(\mathbf{w}, \xi)p(\mathbf{w}|\alpha)p(\alpha)}{q(\mathbf{w}, \alpha)} \right\} d\mathbf{w} d\alpha. \end{aligned} \quad (10.172)$$

Next we assume that the variational distribution factorizes between parameters and hyperparameters so that

$$q(\mathbf{w}, \alpha) = q(\mathbf{w})q(\alpha). \quad (10.173)$$

VI框架下， $q(\mathbf{w})$ 的迭代公式

With this factorization we can appeal to the general result (10.9) to find expressions for the optimal factors // Consider first the distribution  $q(\mathbf{w})$ . Discarding terms that are independent of  $\mathbf{w}$ , we have

$$\begin{aligned}\ln q(\mathbf{w}) &= \mathbb{E}_\alpha [\ln \{h(\mathbf{w}, \xi)p(\mathbf{w}|\alpha)p(\alpha)\}] + \text{const} \\ &= \ln h(\mathbf{w}, \xi) + \mathbb{E}_\alpha [\ln p(\mathbf{w}|\alpha)] + \text{const}.\end{aligned}$$

We now substitute for  $\ln h(\mathbf{w}, \xi)$  using (10.153), and for  $\ln p(\mathbf{w}|\alpha)$  using (10.165), giving

$$\ln q(\mathbf{w}) = -\frac{\mathbb{E}[\alpha]}{2}\mathbf{w}^T\mathbf{w} + \sum_{n=1}^N \{(t_n - 1/2)\mathbf{w}^T\boldsymbol{\phi}_n - \lambda(\xi_n)\mathbf{w}^T\boldsymbol{\phi}_n\boldsymbol{\phi}_n^T\mathbf{w}\} + \text{const.}$$

We see that this is a quadratic function of  $\mathbf{w}$  and so the solution for  $q(\mathbf{w})$  will be Gaussian. Completing the square in the usual way, we obtain

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N) \quad (10.174)$$

where we have defined

$$\boldsymbol{\Sigma}_N^{-1} \boldsymbol{\mu}_N = \sum_{n=1}^N (t_n - 1/2) \boldsymbol{\phi}_n \quad (10.175)$$

$$\boldsymbol{\Sigma}_N^{-1} = \mathbb{E}[\alpha]\mathbf{I} + 2 \sum_{n=1}^N \lambda(\xi_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T. \quad (10.176)$$

VI框架下， $q(\alpha)$ 的迭代公式 //

Similarly, the optimal solution for the factor  $q(\alpha)$  is obtained from

$$\ln q(\alpha) = \mathbb{E}_{\mathbf{w}} [\ln p(\mathbf{w}|\alpha)] + \ln p(\alpha) + \text{const.}$$

Substituting for  $\ln p(\mathbf{w}|\alpha)$  using (10.165), and for  $\ln p(\alpha)$  using (10.166), we obtain

$$\ln q(\alpha) = \frac{M}{2} \ln \alpha - \frac{\alpha}{2} \mathbb{E} [\mathbf{w}^T \mathbf{w}] + (a_0 - 1) \ln \alpha - b_0 \alpha + \text{const.}$$

We recognize this as the log of a gamma distribution, and so we obtain

$$q(\alpha) = \text{Gam}(\alpha|a_N, b_N) = \frac{1}{\Gamma(a_N)} a_N^{b_N} \alpha^{a_N - 1} e^{-b_N \alpha} \quad (10.177)$$

where

$$a_N = a_0 + \frac{M}{2} \quad (10.178)$$

$$b_N = b_0 + \frac{1}{2} \mathbb{E}_{\mathbf{w}} [\mathbf{w}^T \mathbf{w}]. \quad (10.179)$$

**local variational 框架下 //** We also need to optimize the variational parameters  $\xi_n$ , and this is also done by maximizing the lower bound  $\tilde{\mathcal{L}}(q, \xi)$ . Omitting terms that are independent of  $\xi$ , and integrating over  $\alpha$ , we have

VI (global Variational) 中  
的，因此这里得到的也是  
迭代公式。

$$\tilde{\mathcal{L}}(q, \xi) = \int q(\mathbf{w}) \ln h(\mathbf{w}, \xi) d\mathbf{w} + \text{const.} \quad (10.180)$$

10.160

Note that this has precisely the same form as (10.159), and so we can again appeal to our earlier result (10.163), which can be obtained by direct optimization of the marginal likelihood function, leading to re-estimation equations of the form

$$(\xi_n^{\text{new}})^2 = \phi_n^T (\Sigma_N + \mu_N \mu_N^T) \phi_n. \quad (10.181)$$

总结前文并给出会  
用到的一些矩阵计算公式  
Appendix B

// We have obtained re-estimation equations for the three quantities  $q(\mathbf{w})$ ,  $q(\alpha)$ , and  $\xi$ , and so after making suitable initializations, we can cycle through these quantities, updating each in turn. The required moments are given by

$$\mathbb{E}[\alpha] = \frac{a_N}{b_N} \quad (10.182)$$

$$\mathbb{E}[\mathbf{w}^T \mathbf{w}] = \Sigma_N + \mu_N \mu_N^T. \quad (10.183)$$

## 10.7. Expectation Propagation

We conclude this chapter by discussing an alternative form of deterministic approximate inference, known as *expectation propagation or EP* (Minka, 2001a; Minka, 2001b). As with the variational Bayes methods discussed so far, this too is based on the minimization of a Kullback-Leibler divergence but now of the (reverse) form, which gives the approximation rather different properties.

Consider for a moment the problem of minimizing  $\text{KL}(p||q)$  with respect to  $q(\mathbf{z})$  when  $p(\mathbf{z})$  is a fixed distribution and  $q(\mathbf{z})$  is a member of the exponential family and so, from (2.194), can be written in the form

$$q(\mathbf{z}) = h(\mathbf{z})g(\boldsymbol{\eta}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{z}) \}. \quad (10.184)$$

As a function of  $\boldsymbol{\eta}$ , the Kullback-Leibler divergence then becomes

$$\text{KL}(p||q) = -\ln g(\boldsymbol{\eta}) - \boldsymbol{\eta}^T \mathbb{E}_{p(\mathbf{z})}[\mathbf{u}(\mathbf{z})] + \text{const} \quad (10.185)$$

where the constant terms are independent of the natural parameters  $\boldsymbol{\eta}$ . We can minimize  $\text{KL}(p||q)$  within this family of distributions by setting the gradient with respect to  $\boldsymbol{\eta}$  to zero, giving

$$-\nabla \ln g(\boldsymbol{\eta}) = \mathbb{E}_{p(\mathbf{z})}[\mathbf{u}(\mathbf{z})]. \quad (10.186)$$

However, we have already seen in (2.226) that the negative gradient of  $\ln g(\boldsymbol{\eta})$  is given by the expectation of  $\mathbf{u}(\mathbf{z})$  under the distribution  $q(\mathbf{z})$ . Equating these two results, we obtain

$$\mathbb{E}_{q(\mathbf{z})}[\mathbf{u}(\mathbf{z})] = \mathbb{E}_{p(\mathbf{z})}[\mathbf{u}(\mathbf{z})]. \quad (10.187)$$

We see that the optimum solution simply corresponds to matching the expected sufficient statistics. So, for instance, if  $q(\mathbf{z})$  is a Gaussian  $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  then we minimize the Kullback-Leibler divergence by setting the mean  $\boldsymbol{\mu}$  of  $q(\mathbf{z})$  equal to the mean of the distribution  $p(\mathbf{z})$  and the covariance  $\boldsymbol{\Sigma}$  equal to the covariance of  $p(\mathbf{z})$ . This is sometimes called *moment matching*. An example of this was seen in Figure 10.3(a). [

Now let us exploit this result to obtain a practical algorithm for approximate inference. For many probabilistic models, the joint distribution of (data  $\mathcal{D}$ ) and (hidden variables (including parameters)  $\boldsymbol{\theta}$ ) comprises a product of factors in the form

注意，这里f\_i(\theta)包括隐变量和参数。而前文  
介绍VI时的含义相同。同时注意以  $p(\mathcal{D}, \boldsymbol{\theta}) = \prod_i f_i(\boldsymbol{\theta})$  所以在式(10.188)右侧未  
 $\boldsymbol{\theta}$  对应上一自然段的 $f_i(\boldsymbol{\theta})$ 是错误。  
少写观测数据，是已知的。  
呈水平状态。

$$p(\mathcal{D}, \boldsymbol{\theta}) = \prod_i f_i(\boldsymbol{\theta}). \quad (10.188)$$

举列见例(10.188)  
形式化场景

This would arise, for example, in a model for independent, identically distributed data in which there is one factor  $f_n(\boldsymbol{\theta}) = p(\mathbf{x}_n|\boldsymbol{\theta})$  for each data point  $\mathbf{x}_n$ , along with a factor  $f_0(\boldsymbol{\theta}) = p(\boldsymbol{\theta})$  corresponding to the prior. More generally, it would also apply to any model defined by a directed probabilistic graph in which each factor is a conditional distribution corresponding to one of the nodes, or an undirected graph in which each factor is a clique potential. We are interested in evaluating the posterior distribution  $p(\boldsymbol{\theta}|\mathcal{D})$  for the purpose of making predictions, as well as the model evidence  $p(\mathcal{D})$  for the purpose of model comparison. From (10.188) the posterior is given by

任务1  $p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{p(\mathcal{D})} \prod_i f_i(\boldsymbol{\theta}) \quad (10.189)$

and the model evidence is given by

任务2  $p(\mathcal{D}) = \int \prod_i f_i(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (10.190)$

Here we are considering continuous variables, but the following discussion applies equally to discrete variables with integrals replaced by summations. We shall suppose that the marginalization over  $\boldsymbol{\theta}$ , along with the marginalizations with respect to the posterior distribution required to make predictions, are intractable so that some form of approximation is required. [

Expectation propagation is based on an approximation to the posterior distribution which is also given by a product of factors

EP算法则推导

$$q(\boldsymbol{\theta}) = \frac{1}{Z} \prod_i \tilde{f}_i(\boldsymbol{\theta}) \quad (10.191)$$

in which each factor  $\tilde{f}_i(\boldsymbol{\theta})$  in the approximation corresponds to one of the factors  $f_i(\boldsymbol{\theta})$  in the true posterior (10.189), and the factor  $1/Z$  is the normalizing constant needed to ensure that the left-hand side of (10.191) integrates to unity. In order to obtain a practical algorithm, we need to constrain the factors  $\tilde{f}_i(\boldsymbol{\theta})$  in some way, and in particular we shall assume that they come from the exponential family. The product of the factors will therefore also be from the exponential family and so can

$p(\boldsymbol{\theta}|\mathcal{D})$ 的近似，注意  
高f\_i(\theta)与f\_{i'}(\theta)是一对  
应的，此外，f\_i(\theta)的函数  
表达式是已知的，不过相后  
的计算intractable；我们  
限制于旧源自指数族  
分布，其作用于f\_i(\theta)相差  
一个常数因子的近似，  
即  $f_i(\theta) \approx C f_i(\theta)$ 。

be described by a finite set of sufficient statistics. For example, if each of the  $\tilde{f}_i(\theta)$  is a Gaussian, then the overall approximation  $q(\theta)$  will also be Gaussian.

// Ideally we would like to determine the  $\tilde{f}_i(\theta)$  by minimizing the Kullback-Leibler divergence between the true posterior and the approximation given by

$$\text{优化目标} \quad \text{KL}(p\|q) = \text{KL}\left(\frac{1}{p(\mathcal{D})} \prod_i f_i(\theta) \middle\| \frac{1}{Z} \prod_i \tilde{f}_i(\theta)\right). \quad (10.192)$$

**优化困难的原因及一种解决方案的问题**

Note that this is the reverse form of KL divergence compared with that used in variational inference. In general, this minimization will be intractable because the KL divergence involves averaging with respect to the true distribution. As a rough approximation, we could instead minimize the KL divergences between the corresponding pairs  $f_i(\theta)$  and  $\tilde{f}_i(\theta)$  of factors. This represents a much simpler problem to solve, and has the advantage that the algorithm is noniterative. However, because each factor is individually approximated, the product of the factors could well give a poor approximation.

下面开始介绍  
印算法

// Expectation propagation makes a much better approximation by optimizing each factor in turn in the context of all of the remaining factors. It starts by initializing the factors  $\tilde{f}_i(\theta)$ , and then cycles through the factors refining them one at a time. This is similar in spirit to the update of factors in the variational Bayes framework considered earlier. Suppose we wish to refine factor  $\tilde{f}_j(\theta)$ . We first remove this factor from the product to give  $\prod_{i \neq j} \tilde{f}_i(\theta)$ . Conceptually, we will now determine a revised form of the factor  $\tilde{f}_j(\theta)$  by ensuring that the product

$$q^{\text{new}}(\theta) \propto \tilde{f}_j(\theta) \prod_{i \neq j} \tilde{f}_i(\theta) \quad (10.193)$$

是新得更新的量  
并修正书中相应的叙述。 } 已知保持不变的是

$$\text{is as close as possible to } \tilde{f}_j(\theta) \prod_{i \neq j} \tilde{f}_i(\theta) \quad (10.194)$$

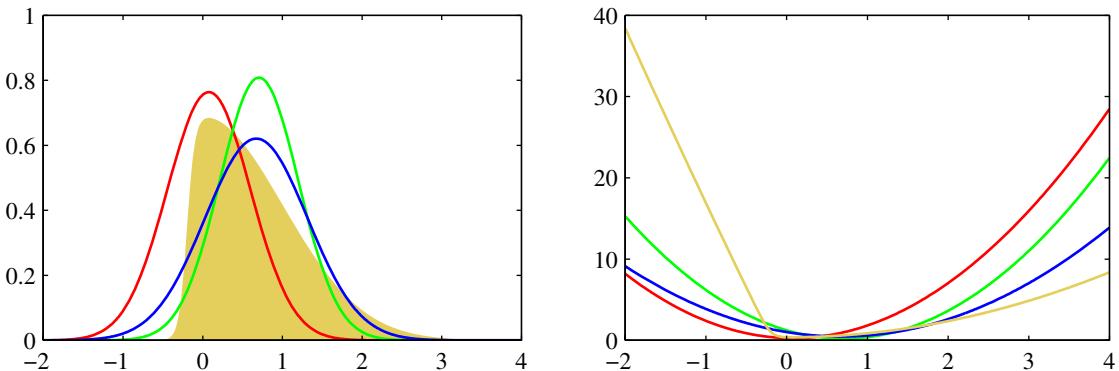
in which we keep fixed all of the factors  $\tilde{f}_i(\theta)$  for  $i \neq j$ . This ensures that the approximation is most accurate in the regions of high posterior probability as defined by the remaining factors. We shall see an example of this effect when we apply EP to the ‘clutter problem’. To achieve this, we first remove the factor  $\tilde{f}_j(\theta)$  from the current approximation to the posterior by defining the unnormalized distribution

$$q^{\setminus j}(\theta) = \frac{q(\theta)}{\tilde{f}_j(\theta)}. \quad (10.195)$$

Note that we could instead find  $q^{\setminus j}(\theta)$  from the product of factors  $i \neq j$ , although in practice division is usually easier. This is now combined with the factor  $f_j(\theta)$  to give a distribution

$$\frac{1}{Z_j} f_j(\theta) q^{\setminus j}(\theta) \quad (10.196)$$

### Section 10.7.1



**Figure 10.14** Illustration of the expectation propagation approximation using a Gaussian distribution for the example considered earlier in Figures 4.14 and 10.1. The left-hand plot shows the original distribution (yellow) along with the Laplace (red), global variational (green), and EP (blue) approximations, and the right-hand plot shows the corresponding negative logarithms of the distributions. Note that the EP distribution is broader than that of variational inference, as a consequence of the different form of KL divergence.

obtained by

where  $Z_j$  is the normalization constant given by

$$Z_j = \int f_j(\theta) q^{\backslash j}(\theta) d\theta. \quad (10.197)$$

We now determine a revised factor  $\tilde{f}_j(\theta)$  by minimizing the Kullback-Leibler divergence

$$\text{KL}\left(\frac{f_j(\theta)q^{\backslash j}(\theta)}{Z_j} \parallel \frac{\tilde{f}_j(\theta)q^{\backslash j}(\theta)}{K_j}\right) \quad \text{KL}\left(\frac{f_j(\theta)q^{\backslash j}(\theta)}{Z_j} \parallel q^{\text{new}}(\theta)\right). \quad (10.198)$$

This is easily solved because the approximating distribution  $q^{\text{new}}(\theta)$  is from the exponential family, and so we can appeal to the result (10.187), which tells us that the parameters of  $q^{\text{new}}(\theta)$  are obtained by matching its expected sufficient statistics to the corresponding moments of (10.196). We shall assume that this is a tractable operation. For example, if we choose  $q(\theta)$  to be a Gaussian distribution  $\mathcal{N}(\theta|\mu, \Sigma)$ , then  $\mu$  is set equal to the mean of the (unnormalized) distribution  $f_j(\theta)q^{\backslash j}(\theta)$ , and  $\Sigma$  is set to its covariance. More generally, it is straightforward to obtain the required expectations for any member of the exponential family, provided it can be normalized, because the expected statistics can be related to the derivatives of the normalization coefficient, as given by (2.226). The EP approximation is illustrated in Figure 10.14.

$\tilde{f}_j(\theta) / K_j$ ,  $K_j$  为未知而通过  $\min \text{KL}$ , 我们能得 到  $\tilde{f}_j(\theta) / K_j$ , 二者分配实际上 是 住意的, 为此, 我们令  $K_j = Z_j$ , 而非而文提

到 “matching zeroth-order” From (10.193), we see that the revised factor  $\tilde{f}_j(\theta)$  can be found by taking  $q^{\text{new}}(\theta)$  and dividing out the remaining factors so that

$$\tilde{f}_j(\theta) = K_j \frac{q^{\text{new}}(\theta)}{q^{\backslash j}(\theta)} \quad (10.199)$$

where we have used (10.195). The coefficient  $K_j$  is determined by multiplying both

$\int \frac{f_j(\theta)q^{\backslash j}(\theta)}{Z_j} d\theta = 1 = \int \frac{\tilde{f}_j(\theta)q^{\backslash j}(\theta)}{K_j} d\theta$ , 我们方面强行取  $K_j = Z_j$ , 这使得  $\tilde{f}_j(\theta) \approx c f_j(\theta)$  的  $c = 1$ .

sides of (10.199) by  $q^{\backslash j}(\theta)$  and integrating to give

$$K_j = \int \tilde{f}_j'(\theta) q^{\backslash j}(\theta) d\theta \quad (10.200)$$

where we have used the fact that  $q^{\text{new}}(\theta)$  is normalized. The value of  $K$  can therefore be found by matching zeroth-order moments X, 原因并非如此, 而是强行使其相等。

$$\int \tilde{f}_j'(\theta) q^{\backslash j}(\theta) d\theta = \int f_j(\theta) q^{\backslash j}(\theta) d\theta. \quad (10.201)$$

Combining this with (10.197), we then see that  $K_j = Z_j$  and so can be found by evaluating the integral in (10.197). }

In practice, several passes are made through the set of factors, revising each

迭代更新各因子直到满足终止条件。最后回归到任务 1.2 对  $p(\theta|D)$ ,  $p(D)$  的计算。 factor in turn. The posterior distribution  $p(\theta|D)$  is then approximated using (10.191), and the model evidence  $p(D)$  can be approximated by using (10.190) with the factors  $f_i(\theta)$  replaced by their approximations  $\tilde{f}_i(\theta)$ .

### Expectation Propagation

We are given a joint distribution over observed data  $D$  and stochastic variables  $\theta$  in the form of a product of factors D是已知数据，θ是随机变量，包括隐变量及参数，未知。

$$p(D, \theta) = \prod_i f_i(\theta) \quad (10.202)$$

and we wish to approximate the posterior distribution  $p(\theta|D)$  by a distribution of the form

$$\underline{q(\theta)} = \frac{1}{Z} \prod_i \tilde{f}_i(\theta). \quad (10.203)$$

E步一化 We also wish to approximate the model evidence  $p(D)$ .

1. Initialize all of the approximating factors  $\tilde{f}_i(\theta)$ .
2. Initialize the posterior approximation by setting

$$q(\theta) \propto \prod_i \tilde{f}_i(\theta). \quad (10.204)$$

3. Until convergence:

- (a) Choose a factor  $\tilde{f}_j(\theta)$  to refine.
- (b) Remove  $\tilde{f}_j(\theta)$  from the posterior by division 于 j 因子迭代更新为 f~j 因子

$$\text{未归一化 } \underline{q^{\backslash j}(\theta)} = \frac{q(\theta)}{\tilde{f}_j(\theta)}. \quad \text{即式(10.207)} \quad (10.205)$$

给定结果后，下一次再对其进行更新时，计算时除以的并不一定要是 f~j 因子，而可以相差任意一个常数因子 C，

当 $q$ 为指数族分布时，通过 matching 最小化式(10.198)，  
就算没不仅仅局限于指数族分布，  
只不过非指数族分布时  
通过 moment matching 可以(10.198)了。比如，10.7.2 节就入  
要求 $q$ 为指数族分布。

moment

但显然

指数族分布，

没有理由

是这样

但为什么

能这样

呢？

因为就入

- (c) Evaluate the new posterior by setting the sufficient statistics (moments) of  $q^{\text{new}}(\theta)$  equal to those of  $q^j(\theta)f_j(\theta)$ , including evaluation of the normalization constant

已归一化

1) 是的。但因为我们取  $k_j = z_j$ ，因此优化得到的系数总是对  $f_j(\theta)$  的直接正比，而不会

相差常数因子，即  $f_j(\theta) \propto \tilde{f}_j(\theta)$ ，这也是

- (d) Evaluate and store the new factor

$$\tilde{f}_j'(\theta) = Z_j \frac{q^{\text{new}}(\theta)}{q^j(\theta)}. \quad (10.207)$$

$$\tilde{f}_j(\theta) \leftarrow \tilde{f}_j'(\theta)$$

4. Evaluate the approximation to the model evidence

$$p(\mathcal{D}) \simeq \int \prod_i \tilde{f}_i(\theta) d\theta. \quad (10.208)$$

EP算法的特例：ADF

只更新一轮

A special case of EP, known as *assumed density filtering* (ADF) or *moment matching* (Maybeck, 1982; Lauritzen, 1992; Boyen and Koller, 1998; Opper and Winther, 1999), is obtained by initializing all of the approximating factors except the first to unity and then making one pass through the factors updating each of them once. Assumed density filtering can be appropriate for on-line learning in which data points are arriving in a sequence and we need to learn from each data point and then discard it before considering the next point. However, in a batch setting we have the opportunity to re-use the data points many times in order to achieve improved accuracy, and it is this idea that is exploited in expectation propagation. Furthermore, if we apply ADF to batch data, the results will have an undesirable dependence on the (arbitrary) order in which the data points are considered, which again EP can overcome.

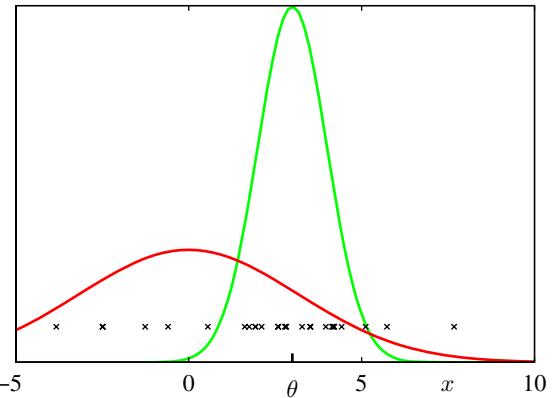
EP不保证一定收敛

One disadvantage of expectation propagation is that there is no guarantee that the iterations will converge. However, for approximations  $q(\theta)$  in the exponential family, if the iterations do converge, the resulting solution will be a stationary point of a particular energy function (Minka, 2001a), although each iteration of EP does not necessarily decrease the value of this energy function. This is in contrast to variational Bayes, which iteratively maximizes a lower bound on the log marginal likelihood, in which each iteration is guaranteed not to decrease the bound. It is possible to optimize the EP cost function directly, in which case it is guaranteed to converge, although the resulting algorithms can be slower and more complex to implement.

Another difference between variational Bayes and EP arises from the form of KL divergence that is minimized by the two algorithms, because the former minimizes  $\text{KL}(q||p)$  whereas the latter minimizes  $\text{KL}(p||q)$ . As we saw in Figure 10.3, for distributions  $p(\theta)$  which are multimodal, minimizing  $\text{KL}(p||q)$  can lead to poor approximations. In particular, if EP is applied to mixtures the results are not sensible because the approximation tries to capture all of the modes of the posterior distribution. Conversely, in logistic-type models, EP often out-performs both local variational methods and the Laplace approximation (Kuss and Rasmussen, 2006).

即再次更新第3个factor时，3(a)与  
3选择第3个factor  $\rightarrow c f_j(\theta)$ 。 $P_{S12}$  给  
若在 clutter problem 中使用公式，实际上暗  
含有先迭代前通过选择适当的 C 使  $q(\theta)$  已满足

**Figure 10.15** Illustration of the clutter problem for a data space dimensionality of  $D = 1$ . Training data points, denoted by the crosses, are drawn from a mixture of two Gaussians with components shown in red and green. The goal is to infer the mean of the green Gaussian from the observed data.



### 10.7.1 Example: The clutter problem

Following Minka (2001b), we illustrate the EP algorithm using a simple example in which the goal is to infer the mean  $\theta$  of a multivariate Gaussian distribution over a variable  $x$  given a set of observations drawn from that distribution. To make the problem more interesting, the observations are embedded in background clutter, which itself is also Gaussian distributed, as illustrated in Figure 10.15. The distribution of observed values  $x$  is therefore a mixture of Gaussians, which we take to be of the form

$$p(\mathbf{x}|\boldsymbol{\theta}) = (1-w)\mathcal{N}(\mathbf{x}|\boldsymbol{\theta}, \mathbf{I}) + w\mathcal{N}(\mathbf{x}|\mathbf{0}, a\mathbf{I}) \quad (10.209)$$

where  $w$  is the proportion of background clutter and is assumed to be known. The prior over  $\boldsymbol{\theta}$  is taken to be Gaussian

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, b\mathbf{I}) \quad (10.210)$$

注意,  $a, b, w$  都已知 and Minka (2001a) chooses the parameter values  $a = 10$ ,  $b = 100$  and  $w = 0.5$ . The joint distribution of  $N$  observations  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and  $\boldsymbol{\theta}$  is given by

$$p(\mathcal{D}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\theta}) \quad (10.211)$$

and so the posterior distribution comprises a mixture of  $2^N$  Gaussians. { Thus the computational cost of solving this problem exactly would grow exponentially with the size of the data set, and so an exact solution is intractable for moderately large  $N$ . }

To apply EP to the clutter problem, we first identify the factors  $f_0(\boldsymbol{\theta}) = p(\boldsymbol{\theta})$  and  $f_n(\boldsymbol{\theta}) = p(\mathbf{x}_n | \boldsymbol{\theta})$ . Next we select an approximating distribution from the exponential family, and for this example it is convenient to choose a spherical Gaussian

$$\begin{aligned} q(\boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{\theta} | \tilde{\mathbf{m}}, v\mathbf{I}). \\ &= \frac{1}{2} \tilde{\mathbf{m}}^\top \tilde{\mathbf{m}} / v \end{aligned} \quad (10.212)$$

The factor approximations will therefore take the form of exponential-quadratic functions of the form

$$\tilde{f}_n(\boldsymbol{\theta}) = s_n \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_n, v_n \mathbf{I}) \quad (10.213)$$

where  $n = 1, \dots, N$ , and we set  $\tilde{f}_0(\boldsymbol{\theta})$  equal to the prior  $p(\boldsymbol{\theta})$ . Note that the use of  $\mathcal{N}(\boldsymbol{\theta} | \cdot, \cdot)$  does not imply that the right-hand side is a well-defined Gaussian density (in fact, as we shall see, the variance parameter  $v_n$  can be negative) but is simply a convenient shorthand notation. The approximations  $\tilde{f}_n(\boldsymbol{\theta})$ , for  $n = 1, \dots, N$ , can be initialized to unity, corresponding to  $s_n = (2\pi v_n)^{D/2}$ ,  $v_n \rightarrow \infty$  and  $\mathbf{m}_n = \mathbf{0}$ , where  $D$  is the dimensionality of  $\mathbf{x}$  and hence of  $\boldsymbol{\theta}$ . The initial  $q(\boldsymbol{\theta})$ , defined by (10.191), is therefore equal to the prior.

We then iteratively refine the factors by taking one factor  $f_n(\boldsymbol{\theta})$  at a time and applying (10.205), (10.206), and (10.207). Note that we do not need to revise the term  $f_0(\boldsymbol{\theta})$  because an EP update will leave this term unchanged. Here we state the results and leave the reader to fill in the details.

**Exercise 10.37** 没问题，但证明不清楚，已经知道修正了。 **Exercise 10.38**

证明式(10.214)-(10.216) 注意，这里通过将  $\tilde{f}_n(\boldsymbol{\theta})$  代入  $\mathbf{m}^n = \mathbf{m} + v^n v_n^{-1} (\mathbf{m} - \mathbf{m}_n) P(\boldsymbol{\theta})$  并随着计算不断更新，而  $\mathbf{m}^n$  又是迭代值，使得  $q^{(n)}(\boldsymbol{\theta})_{(v^n)^{-1}} = \{v^{-1} - v_n^{-1}\}$  已归一化，且与迭代步数无关。 上面这段对应前文印算法3(b)

Next we evaluate the normalization constant  $Z_n$  using (10.206) to give

$$Z_n = (1 - w) \mathcal{N}(\mathbf{x}_n | \mathbf{m}^n, (v^n + 1) \mathbf{I}) + w \mathcal{N}(\mathbf{x}_n | \mathbf{0}, a \mathbf{I}). \quad (10.216)$$

Similarly, we compute the mean and variance of  $q^{new}(\boldsymbol{\theta})$  by finding the mean and variance of  $q^{(n)}(\boldsymbol{\theta}) f_n(\boldsymbol{\theta})$  to give

$$\begin{aligned} \mathbf{m}^{new} &= \mathbf{m}^n + \rho_n \frac{v^n}{v^n + 1} (\mathbf{x}_n - \mathbf{m}^n) \\ v^{new} &= v^n - \rho_n \frac{(v^n)^2}{v^n + 1} + \rho_n (1 - \rho_n) \frac{(v^n)^2 \|\mathbf{x}_n - \mathbf{m}^n\|^2}{D(v^n + 1)^2} \end{aligned} \quad (10.217) \quad (10.218)$$

where the quantity

$$\rho_n = 1 - \frac{w}{Z_n} \mathcal{N}(\mathbf{x}_n | \mathbf{0}, a \mathbf{I}) \quad (10.219)$$

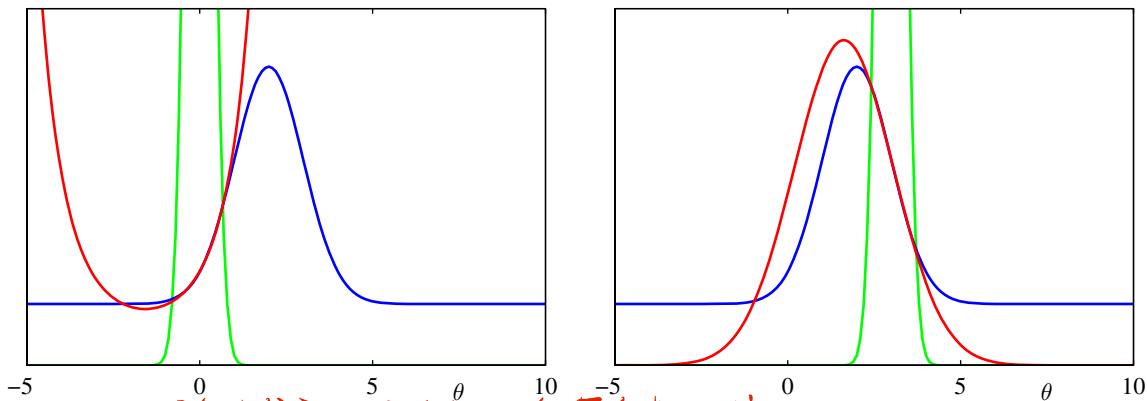
has a simple interpretation as the probability of the point  $\mathbf{x}_n$  not being clutter. Then we use (10.207) to compute the refined factor  $\tilde{f}_n(\boldsymbol{\theta})$  whose parameters are given by

$$\text{注意，这是 } \tilde{f}_n(\boldsymbol{\theta}), \mathbf{m}_n, v_n^{-1} = (v^{new})^{-1} - (v^n)^{-1} \quad (10.220)$$

$$\text{注意，这是更新后的值，而式 } \mathbf{m}'_n = \mathbf{m}^n + (v'_n + v^n)(v^n)^{-1}(\mathbf{m}^{new} - \mathbf{m}^n) \quad (10.221)$$

$$(10.214), (10.215) 中的 \mathbf{m}'_n, v'_n 是 \text{注意，这是更新前的值，且更新前的 } s'_n = \frac{Z_n}{(2\pi v'_n)^{D/2} \mathcal{N}(\mathbf{m}'_n | \mathbf{m}^n, (v'_n + v^n) \mathbf{I})}. \quad (10.222)$$

This refinement process is repeated until a suitable termination criterion is satisfied, for instance that the maximum change in parameter values resulting from a complete iteration is less than a specified tolerance. 即印算法第3步是一个迭代计算步，需迭代直至收敛。为了表示区别，我们加上上标“'”。



绿色曲线高的地方，红色曲线和蓝色曲线比较一致

**Figure 10.16** Examples of the approximation of specific factors for a one-dimensional version of the clutter problem, showing  $f_n(\theta)$  in blue,  $\tilde{f}_n(\theta)$  in red, and  $q^{<n}(\theta)$  in green. Notice that the current form for  $q^{<n}(\theta)$  controls the range of  $\theta$  over which  $\tilde{f}_n(\theta)$  will be a good approximation to  $f_n(\theta)$ .

上面这阶段对应印算法 3(d)

不仅仅是针  
对最后一次，  
或 (10.223),  $\mathbf{m}^{\text{new}}, \mathbf{v}^{\text{new}} =$  最后一次更新时，式 (10.217), (10.218) 的结果  $\mathbf{v}_N$ , 即  $\mathbf{v}^{\text{new}}$  的均值与方差  
(10.224) 实际上  $\mathbf{v} =$  最后一次更新前最近似后验  $q(\theta)$  的均值

pass through all factors is less than some threshold. Finally, we use (10.208) to evaluate the approximation to the model evidence, given by

$$p(D) \approx (2\pi v^{\text{new}})^{D/2} \exp(B/2) \prod_{n=1}^N \{s_n(2\pi v_n)^{-D/2}\} \quad (10.223)$$

实际上给出  $s_n, v_n, m_n$ : where 由最新的一次给定，注意这里指的是传播而不是本次更新的 factor。  
3 次更新完

factor  $\tilde{f}_n(\theta)$  后，回正关系更新近似后验

$q(\theta)$  的计算公式，即式 (10.223), 或

其中  $s_n(\theta)$  均为局部值，得到的  $v_n$

可再嵌入到 factor  $\tilde{f}_{n+1}(\theta)$  的更新中。 Examples of factor approximations for the clutter problem with a one-dimensional parameter space  $\theta$  are shown in Figure 10.16. Note that the factor approximations can have infinite or even negative values for the ‘variance’ parameter  $v_n$ . This simply corresponds to approximations that curve upwards instead of downwards and are not necessarily problematic provided the overall approximate posterior  $q(\theta)$  has positive variance. Figure 10.17 compares the performance of EP with variational Bayes (mean field theory) and the Laplace approximation on the clutter problem.

## 10.7.2 Expectation propagation on graphs

对 10.19 与进一步分解的情况，  
进而引用概率图模型描述。

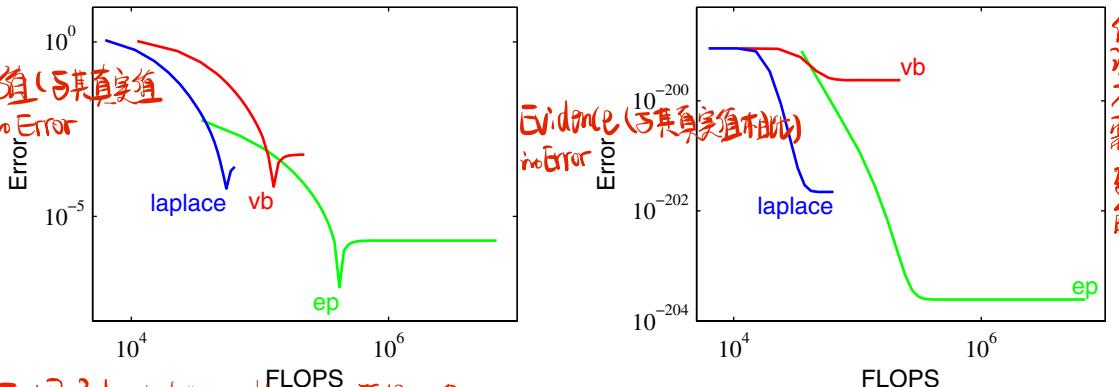
So far in our general discussion of EP, we have allowed the factors  $f_i(\theta)$  in the distribution  $p(\theta)$  to be functions of all of the components of  $\theta$ , and similarly for the approximating factors  $\tilde{f}(\theta)$  in the approximating distribution  $q(\theta)$ . We now consider situations in which the factors depend only on subsets of the variables. Such restrictions can be conveniently expressed using the framework of probabilistic graphical models, as discussed in Chapter 8. Here we use a factor graph representation because this encompasses both directed and undirected graphs.

这里对两个概念进行梳理和说明：belief propagation (BP) 算法全书并未详细叙述，只是在 P403 中提到 BP 是 sum-product 算法的特例。我们还以为本书提到 BP 算法时就是在谈 sum-product 算法。此外，8.4.7 小节介绍了 loopy BP 算法，它是一种近似推断算法，需要不断迭代更新，但 sum-product 算法

是一种精确推断算法，只需要传递信息传递而无需多轮迭代更新。另外需要说明的是，BP在求积分，而EP则是在求近似分布，这两个主题似乎并不相关，但注意，之所以要求近似分布，是因为我们虽然知道分布的函数形式

## 514 10. APPROXIMATE INFERENCE

却并不知道归一化常量  $Z$ ，而它的计算就是一个积分问题，而从另一角度，BP求积分的应用之一正是计算归一化常量  $Z$  以获得完整的概率分布。Posterior mean 分布，也就是说在某种程度上EP和BP是相同的，否则就证明两者特例这一说了。



可以看到，对 clutter problem，EP算得更准。

Figure 10.17 Comparison of expectation propagation, variational inference, and the Laplace approximation on the clutter problem. The left-hand plot shows the error in the predicted posterior mean versus the number of floating point operations, and the right-hand plot shows the corresponding results for the model evidence.

对近似分布的限制：fully factorized

我们提到，当被近似分布为 tree-structure 时，reduce to 精确推断算法。举例 We shall focus on the case in which the approximating distribution is fully factorized, and we shall show that in this case expectation propagation reduces to loopy belief propagation (Minka, 2001a). To start with, we show this in the context of a simple example, and then we shall explore the general case.

First of all, recall from (10.17) that if we minimize the Kullback-Leibler divergence  $KL(p\|q)$  with respect to a factorized distribution  $q$  then the optimal solution for each factor is simply the corresponding marginal of  $p$ . (10.17) 给出了解析解

Now consider the factor graph shown on the left in Figure 10.18, which was introduced earlier in the context of the sum-product algorithm. The joint distribution is given by

$$p(\mathbf{x}) = f_a(x_1, x_2)f_b(x_2, x_3)f_c(x_2, x_4). \quad (10.225)$$

Section 8.4.4  
这节的含义与前文不同，要  
示随机变量，对应前文介  
绍EP算三时的。此外，

前面在介绍EP算三时讨论  
的是参数后验PCID 的近似

似，事实上，对任意分布均  
可用EP进行近似，无论就  
是优化  $KL(p\|q)$ ，更侧重  
对任意分布也均可进行VI，  
关键是优化  $KL(p\|q)$ ，  
只不过在被近似分布为后  
验PCID 时，存在ELBO  
那一套方法。

$$q(\mathbf{x}) \propto \tilde{f}_a(x_1, x_2)\tilde{f}_b(x_2, x_3)\tilde{f}_c(x_2, x_4). \quad (10.226)$$

Note that normalization constants have been omitted, and these can be re-instated at the end by local normalization, as is generally done in belief propagation. Now suppose we restrict attention to approximations in which the factors themselves factorize with respect to the individual variables so that 假设近似分布 fully factorized

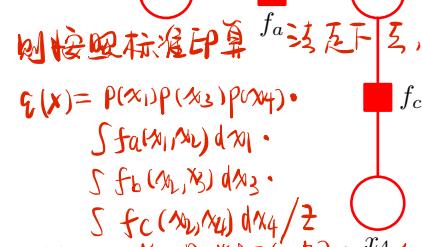
$$q(\mathbf{x}) \propto \tilde{f}_{a1}(x_1)\tilde{f}_{a2}(x_2)\tilde{f}_{b2}(x_2)\tilde{f}_{b3}(x_3)\tilde{f}_{c2}(x_2)\tilde{f}_{c4}(x_4) \quad (10.227)$$

which corresponds to the factor graph shown on the right in Figure 10.18. Because the individual factors are factorized, the overall distribution  $q(\mathbf{x})$  is itself fully factorized.

上面介绍完问题 // Now we apply the EP algorithm using the fully factorized approximation. Suppose that we have initialized all of the factors and that we choose to refine factor 与假设后，下面开始推导计算式

按BP算法进行两级信息传递后计算各子 factor:

$$\begin{aligned} \tilde{f}_{a1}(x_1) &\propto \int p(x) dx_{x_1} \stackrel{\text{def}}{=} p(x_1) & \tilde{f}_{a2}(x_2) &\propto \int f_a(x_1, x_2) dx_{x_1} \\ \tilde{f}_{b3}(x_3) &\propto \int p(x) dx_{x_3} \stackrel{\text{def}}{=} p(x_3) & \tilde{f}_{b2}(x_2) &\propto \int f_b(x_2, x_3) dx_{x_3} \\ \tilde{f}_{c4}(x_4) &\propto \int_x \int p(x) dx_{x_4} \stackrel{\text{def}}{=} p(x_4) & \tilde{f}_{c2}(x_2) &\propto \int f_c(x_2, x_4) dx_{x_4} \end{aligned}$$



$$q(x) = p(x_1)p(x_2)p(x_3)p(x_4) \cdot$$

$$\int f_a(x_1, x_2) dx_{x_1} \cdot$$

$$\int f_b(x_2, x_3) dx_{x_3} \cdot$$

$$\int f_c(x_3, x_4) dx_{x_4} / 2$$

注意  $f_b(x)$  并不是精确解  $p(x)$ , 但 BP 算法最终得到的是精确解。

## 10.7. Expectation Propagation

515

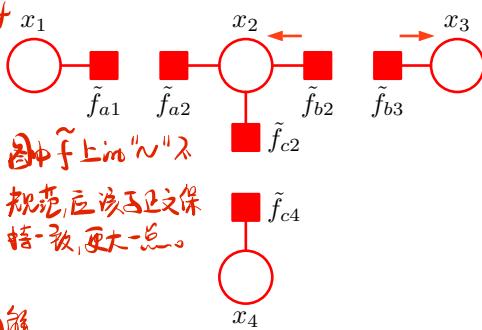


Figure 10.18 On the left is a simple factor graph from Figure 8.51 and reproduced here for convenience. On the right is the corresponding factorized approximation.

$\tilde{f}_b(x_2, x_3) = \tilde{f}_{b2}(x_2)\tilde{f}_{b3}(x_3)$ . We first remove this factor from the approximating distribution to give

$$q^b(\mathbf{x}) \propto \tilde{f}_{a1}(x_1)\tilde{f}_{a2}(x_2)\tilde{f}_{c2}(x_2)\tilde{f}_{c4}(x_4) \quad (10.228)$$

and we then multiply this by the exact factor  $f_b(x_2, x_3)$  to give

$$\hat{p}(\mathbf{x}) = q^b(\mathbf{x})f_b(x_2, x_3) = \tilde{f}_{a1}(x_1)\tilde{f}_{a2}(x_2)\tilde{f}_{c2}(x_2)\tilde{f}_{c4}(x_4)f_b(x_2, x_3). \quad (10.229)$$

注意, 这里是基本式 (10.17)  
求解地  $\min_{q^{\text{new}}} \text{KL}(\hat{p} \parallel q^{\text{new}})$ ,  
得到式 (10.232) - (10.233) 后,  
再从中“识别”到  $P_{\text{sq}}$   
在印算法中想要更新的项,  
即式 (10.24), (10.25). 这里  
不能像  $P_{\text{sq}}$  3(c) 中说的那样  
用 moment matching 求解, 因  
为这里  $q$  并不一定为后  
验分布。  
  
We now find  $q^{\text{new}}(\mathbf{x})$  by minimizing the Kullback-Leibler divergence  $\text{KL}(\hat{p} \parallel q^{\text{new}})$ . The result, as noted above, is that  $q^{\text{new}}(\mathbf{z})$  comprises the product of factors, one for each variable  $x_i$ , in which each factor is given by the corresponding marginal of  $\hat{p}(\mathbf{x})$ . These four marginals are given by 這個  $q^{\text{new}}$  的含義與前文不同。前文在介紹印算法時  
 $\hat{p}(x_1) \propto \tilde{f}_{a1}(x_1)$   $q^{\text{new}}$  對應這個  $\hat{p}$ , 而這裡的  $q^{\text{new}}$  則對應前面等更  
 $\hat{p}(x_2) \propto \tilde{f}_{a2}(x_2)\tilde{f}_{c2}(x_2) \sum_{x_3} f_b(x_2, x_3)$  新加近似後驗, 也就是我所說的  $q$ 。這裏的  
 $\hat{p}(x_3) \propto \sum_{x_2} \{f_b(x_2, x_3)\tilde{f}_{a2}(x_2)\tilde{f}_{c2}(x_2)\}$  不像前面那樣為要更新而直接加上  
 $\hat{p}(x_4) \propto \tilde{f}_{c4}(x_4)$  上標“!”了。  
  
(10.230)  
(10.231)  
(10.232)  
(10.233)

and  $q^{\text{new}}(\mathbf{x})$  is obtained by multiplying these marginals together. We see that the only factors in  $q(\mathbf{x})$  that change when we update  $\tilde{f}_b(x_2, x_3)$  are those that involve the variables in  $f_b$  namely  $x_2$  and  $x_3$ . To obtain the refined factor  $\tilde{f}_b(x_2, x_3) = \tilde{f}_{b2}(x_2)\tilde{f}_{b3}(x_3)$  we simply divide  $q^{\text{new}}(\mathbf{x})$  by  $q^b(\mathbf{x})$ , which gives

$$\tilde{f}_{b2}(x_2) \propto \sum_{x_3} f_b(x_2, x_3) \quad (10.234)$$

$$\tilde{f}_{b3}(x_3) \propto \sum_{x_2} \{f_b(x_2, x_3)\tilde{f}_{a2}(x_2)\tilde{f}_{c2}(x_2)\}. \quad (10.235)$$

关于EP与BP关系的详细论述请参见笔记。这里将BP的两轮  
佳量计算得因子factor元，若仍按标准BP的步聚计算 $f(x), q(x)$

仍然只是PIW的近似解，而不含像BP那

样得到精确解。Figure 10.18 处备注了本例的详细结果。

## Section 8.4.4

These are precisely the messages obtained using (belief propagation) in which messages from variable nodes to factor nodes have been folded into the messages from factor nodes to variable nodes. In particular,  $\tilde{f}_{b2}(x_2)$  corresponds to the message  $\mu_{f_b \rightarrow x_2}(x_2)$  sent by factor node  $f_b$  to variable node  $x_2$  and is given by (8.81). Similarly, if we substitute (8.78) into (8.79), we obtain (10.235) in which  $\tilde{f}_{a2}(x_2)$  corresponds to  $\mu_{f_a \rightarrow x_2}(x_2)$  and  $\tilde{f}_{c2}(x_2)$  corresponds to  $\mu_{f_c \rightarrow x_2}(x_2)$ , giving the message  $\tilde{f}_{b3}(x_3)$  which corresponds to  $\mu_{f_b \rightarrow x_3}(x_3)$ .

这是说，若按既标准的EP  
算法进行更新，会同时更新  
 $\tilde{f}_b(x)$ ,  $\tilde{f}_b(x)$ 两个量，对应从同  
一个因子结点发出的两个不同  
方向的信息，见Figure 10.18中标  
注的箭头。注意，这两个信息  
均是 factor node to variable  
node message，是不同边上的  
信息，这跟计算不是本地  
涉及 variable node to factor  
node message.

This result differs slightly from standard belief propagation in that (messages are passed in both directions at the same time.) We can easily modify the EP procedure to give the standard form of the sum-product algorithm by updating just one of the factors at a time, for instance if we refine only  $\tilde{f}_{b3}(x_3)$ , then  $\tilde{f}_{b2}(x_2)$  is unchanged by definition, while the refined version of  $\tilde{f}_{b3}(x_3)$  is again given by (10.235). If we are refining only one term at a time, then we can choose the order in which the refinements are done as we wish. In particular, for a tree-structured graph we can follow a two-pass update scheme, corresponding to the standard belief propagation schedule, which will result in exact inference of the variable and factor marginals.

The initialization of the approximation factors in this case is unimportant. ] tree-structure

Now let us consider a general factor graph corresponding to the distribution  $p(\theta) = \prod f_i(\theta_i)$  [ 条易例子，从这里开始随机变量 又从 X 变回了。  
注意，θ之间可以包含相同的变量  $\theta_k$  ]  
(10.236) 沿用原有保  
持一致的情况 F, 等 factor 以  
初简化对计  
算并没有影响，  
按实际情况一通  
即而泡明。

where  $\theta_i$  represents the subset of variables associated with factor  $f_i$ . We approximate this using a fully factorized distribution of the form

$$q(\theta) \propto \prod_i \prod_k \tilde{f}_{ik}(\theta_k) \quad (10.237)$$

remove 整个 factor; 但是加粗 得同 message passing 形式上 相同的计算公式，即式(10.240)  
得同 message passing 形式上 相同的计算公式，即式(10.240)  
同一时刻只能更新一个子 factor。比如，BP 的 第 2 轮 佳量时，对某一个因子节点，而同时更新其  
和 then multiply by the exact factor  $f_j(\theta_j)$ . To determine the refined term  $\tilde{f}_{jl}(\theta_l)$ , we need only consider the functional dependence on  $\theta_l$ , and so we simply find the factor, 它们计算相互之间不  
受影响。

$$q^{(j)}(\theta) \propto \prod_{i \neq j} \prod_k \tilde{f}_{ik}(\theta_k) \quad (10.238)$$

再次强调其他 factor 也可能含有变量  $\theta_l$   
and then multiply by the exact factor  $f_j(\theta_j)$ . To determine the refined term  $\tilde{f}_{jl}(\theta_l)$ , we need only consider the functional dependence on  $\theta_l$ , and so we simply find the factor, 它们计算相互之间不  
受影响。

$$q^{(j)}(\theta) f_j(\theta_j). \quad (10.239)$$

Up to a multiplicative constant, this involves taking the marginal of  $f_j(\theta_j)$  multiplied by any terms from  $q^{(j)}(\theta)$  that are functions of any of the variables in  $\theta_j$ . Terms that correspond to other factors  $\tilde{f}_i(\theta_i)$  for  $i \neq j$  will cancel between numerator and denominator when we subsequently divide by  $q^{(j)}(\theta)$ . We therefore obtain

$$\tilde{f}_{jl}(\theta_l) \propto \sum_{\theta_m \neq l \in \theta_j} f_j(\theta_j) \prod_k \prod_{m \neq l} \tilde{f}_{km}(\theta_m). \quad (10.240)$$

这正是对如图懂得式  
(10.240)叙述得很简洁，  
对式(10.240)本身已作详  
细的说明，具体推导过程  
同上面举的例子。在笔记中  
也作了详细的推导，对式  
(10.240)也作了更清晰的表达。

We recognize this as the sum-product rule in the form in which messages from variable nodes to factor nodes have been eliminated, as illustrated by the example shown in Figure 8.50. The quantity  $\tilde{f}_{jm}(\theta_m)$  corresponds to the message  $\mu_{f_j \rightarrow \theta_m}(\theta_m)$ , which factor node  $j$  sends to variable node  $m$ , and the product over  $k$  in (10.24) is over all factors that depend on the variables  $\theta_m$  that have variables (other than variable  $\theta_l$ ) in common with factor  $f_j(\theta_j)$ . In other words, to compute the outgoing message from a factor node, we take the product of all the incoming messages from other factor nodes, multiply by the local factor, and then marginalize.

(Thus, the sum-product algorithm arises as a special case of expectation propagation if we use an approximating distribution that is fully factorized.) This suggests that more flexible approximating distributions, corresponding to partially disconnected graphs, could be used to achieve higher accuracy. Another generalization is to group factors  $f_i(\theta_i)$  together into sets and to refine all the factors in a set together at each iteration. Both of these approaches can lead to improvements in accuracy (Minka, 2001b). In general, the problem of choosing the best combination of grouping and disconnection is an open research issue.] 10.4 小节但没详述 本书

总结：① VI/VB:  $\min_{q \in \mathcal{Q}} KL(q||p)$ , We have seen that variational message passing and expectation propagation optimize two different forms of the Kullback-Leibler divergence. Minka (2005) has shown that a broad range of message passing algorithms can be derived from a common framework involving minimization of members of the alpha family of divergences, given by (10.19). These include variational message passing, loopy belief propagation, and expectation propagation, as well as a range of other algorithms, which we do not have space to discuss here, such as tree-reweighted message passing (Wainwright et al., 2005), fractional belief propagation (Wiegerinck and Heskes, 2003), and power EP (Minka, 2004).

② EP:  $\min_{q \in \mathcal{Q}} KL(p||q)$ , 应用在图模型上时，即等价于本小节所介绍的 message passing 算法；③ 更一般地，有  $\min_{q \in \mathcal{Q}} D_\alpha(p||q)$ ，当应用在图模型上时，也存在相应的 message passing 算法。

## Exercises

- 10.1 (\*) **www** Verify that the log marginal distribution of the observed data  $\ln p(\mathbf{X})$  can be decomposed into two terms in the form (10.2) where  $\mathcal{L}(q)$  is given by (10.3) and  $KL(q||p)$  is given by (10.4).
- 10.2 (\*) Use the properties  $\mathbb{E}[z_1] = m_1$  and  $\mathbb{E}[z_2] = m_2$  to solve the simultaneous equations (10.13) and (10.15), and hence show that, provided the original distribution  $p(\mathbf{z})$  is nonsingular, the unique solution for the means of the factors in the approximation distribution is given by  $\mathbb{E}[z_1] = \mu_1$  and  $\mathbb{E}[z_2] = \mu_2$ .
- 10.3 (\*\*) **www** Consider a factorized variational distribution  $q(\mathbf{Z})$  of the form (10.5). By using the technique of Lagrange multipliers, verify that minimization of the Kullback-Leibler divergence  $KL(p||q)$  with respect to one of the factors  $q_i(\mathbf{Z}_i)$ , keeping all other factors fixed, leads to the solution (10.17).
- 10.4 (\*\*) Suppose that  $p(\mathbf{x})$  is some fixed distribution and that we wish to approximate it using a Gaussian distribution  $q(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . By writing down the form of the KL divergence  $KL(p||q)$  for a Gaussian  $q(\mathbf{x})$  and then differentiating, show that

minimization of  $\text{KL}(p\|q)$  with respect to  $\mu$  and  $\Sigma$  leads to the result that  $\mu$  is given by the expectation of  $\mathbf{x}$  under  $p(\mathbf{x})$  and that  $\Sigma$  is given by the covariance.

- 10.5** (\*\*) **www** Consider a model in which the set of all hidden stochastic variables, denoted collectively by  $\mathbf{Z}$ , comprises some latent variables  $\mathbf{z}$  together with some model parameters  $\theta$ . Suppose we use a variational distribution that factorizes between latent variables and parameters so that  $q(\mathbf{z}, \theta) = q_{\mathbf{z}}(\mathbf{z})q_{\theta}(\theta)$ , in which the distribution  $q_{\theta}(\theta)$  is approximated by a point estimate of the form  $q_{\theta}(\theta) = \delta(\theta - \theta_0)$  where  $\theta_0$  is a vector of free parameters. Show that variational optimization of this factorized distribution is equivalent to an EM algorithm, in which the E step optimizes  $q_{\mathbf{z}}(\mathbf{z})$ , and the M step maximizes the expected complete-data log posterior distribution of  $\theta$  with respect to  $\theta_0$ .
- 10.6** (\*\*) The alpha family of divergences is defined by (10.19). Show that the Kullback-Leibler divergence  $\text{KL}(p\|q)$  corresponds to  $\alpha \rightarrow 1$ . This can be done by writing  $p^{\epsilon} = \exp(\epsilon \ln p) = 1 + \epsilon \ln p + O(\epsilon^2)$  and then taking  $\epsilon \rightarrow 0$ . Similarly show that  $\text{KL}(q\|p)$  corresponds to  $\alpha \rightarrow -1$ .
- 10.7** (\*\*) Consider the problem of inferring the mean and precision of a univariate Gaussian using a factorized variational approximation, as considered in Section 10.1.3. Show that the factor  $q_{\mu}(\mu)$  is a Gaussian of the form  $\mathcal{N}(\mu|\mu_N, \lambda_N^{-1})$  with mean and precision given by (10.26) and (10.27), respectively. Similarly show that the factor  $q_{\tau}(\tau)$  is a gamma distribution of the form  $\text{Gam}(\tau|a_N, b_N)$  with parameters given by (10.29) and (10.30).
- 10.8** (\*) Consider the variational posterior distribution for the precision of a univariate Gaussian whose parameters are given by (10.29) and (10.30). By using the standard results for the mean and variance of the gamma distribution given by (B.27) and (B.28), show that if we let  $N \rightarrow \infty$ , this variational posterior distribution has a mean given by the inverse of the maximum likelihood estimator for the variance of the data, and a variance that goes to zero.
- 10.9** (\*\*) By making use of the standard result  $\mathbb{E}[\tau] = a_N/b_N$  for the mean of a gamma distribution, together with (10.26), (10.27), (10.29), and (10.30), derive the result (10.33) for the reciprocal of the expected precision in the factorized variational treatment of a univariate Gaussian.
- 10.10** (\*) **www** Derive the decomposition given by (10.34) that is used to find approximate posterior distributions over models using variational inference.
- 10.11** (\*\*) **www** By using a Lagrange multiplier to enforce the normalization constraint on the distribution  $q(m)$ , show that the maximum of the lower bound (10.35) is given by (10.36).
- 10.12** (\*\*) Starting from the joint distribution (10.41), and applying the general result (10.9), show that the optimal variational distribution  $q^*(\mathbf{Z})$  over the latent variables for the Bayesian mixture of Gaussians is given by (10.48) by verifying the steps given in the text.

- 10.13** (\*\*) **www** Starting from (10.54), derive the result (10.59) for the optimum variational posterior distribution over  $\mu_k$  and  $\Lambda_k$  in the Bayesian mixture of Gaussians, and hence verify the expressions for the parameters of this distribution given by (10.60)–(10.63).
- 10.14** (\*\*) Using the distribution (10.59), verify the result (10.64).
- 10.15** (\*) Using the result (B.17), show that the expected value of the mixing coefficients in the variational mixture of Gaussians is given by (10.69).
- 10.16** (\*\*) **www** Verify the results (10.71) and (10.72) for the first two terms in the lower bound for the variational Gaussian mixture model given by (10.70).
- 10.17** (\*\*\*) Verify the results (10.73)–(10.77) for the remaining terms in the lower bound for the variational Gaussian mixture model given by (10.70).
- 10.18** (\*\*\*\*) In this exercise, we shall derive the variational re-estimation equations for the Gaussian mixture model by direct differentiation of the lower bound. To do this we assume that the variational distribution has the factorization defined by (10.42) and (10.55) with factors given by (10.48), (10.57), and (10.59). Substitute these into (10.70) and hence obtain the lower bound as a function of the parameters of the variational distribution. Then, by maximizing the bound with respect to these parameters, derive the re-estimation equations for the factors in the variational distribution, and show that these are the same as those obtained in Section 10.2.1.
- 10.19** (\*\*) Derive the result (10.81) for the predictive distribution in the variational treatment of the Bayesian mixture of Gaussians model.
- 10.20** (\*\*) **www** This exercise explores the variational Bayes solution for the mixture of Gaussians model when the size  $N$  of the data set is large and shows that it reduces (as we would expect) to the maximum likelihood solution based on EM derived in Chapter 9. Note that results from Appendix B may be used to help answer this exercise. First show that the posterior distribution  $q^*(\Lambda_k)$  of the precisions becomes sharply peaked around the maximum likelihood solution. Do the same for the posterior distribution of the means  $q^*(\mu_k|\Lambda_k)$ . Next consider the posterior distribution  $q^*(\pi)$  for the mixing coefficients and show that this too becomes sharply peaked around the maximum likelihood solution. Similarly, show that the responsibilities become equal to the corresponding maximum likelihood values for large  $N$ , by making use of the following asymptotic result for the digamma function for large  $x$

$$\psi(x) = \ln x + O(1/x). \quad (10.241)$$

Finally, by making use of (10.80), show that for large  $N$ , the predictive distribution becomes a mixture of Gaussians.

- 10.21** (\*) Show that the number of equivalent parameter settings due to interchange symmetries in a mixture model with  $K$  components is  $K!$ .

## 10. APPROXIMATE INFERENCE

解答表述得不严谨，结果也有误，请见别处。参数又从人问及验分布

- 10.22** (\*\*) We have seen that each mode of the posterior distribution in a Gaussian mixture model is a member of a family of  $K!$  equivalent modes. Suppose that the result of running the variational inference algorithm is an approximate posterior distribution  $q$  that is localized in the neighbourhood of one of the modes. We can then approximate the full posterior distribution as a mixture of  $K!$  such  $q$  distributions, once centred on each mode and having equal mixing coefficients. Show that if we assume negligible overlap between the components of the  $q$  mixture, the resulting lower bound differs from that for a single component  $q$  distribution through the addition of an extra term  $\ln K!/K!$

- 10.23** (\*\*) **www** Consider a variational Gaussian mixture model in which there is no prior distribution over mixing coefficients  $\{\pi_k\}$ . Instead, the mixing coefficients are treated as parameters, whose values are to be found by maximizing the variational lower bound on the log marginal likelihood. Show that maximizing this lower bound with respect to the mixing coefficients, using a Lagrange multiplier to enforce the constraint that the mixing coefficients sum to one, leads to the re-estimation result (10.83). Note that there is no need to consider all of the terms in the lower bound but only the dependence of the bound on the  $\{\pi_k\}$ .

- 10.24** (\*\*) **www** We have seen in Section 10.2 that the singularities arising in the maximum likelihood treatment of Gaussian mixture models do not arise in a Bayesian treatment. Discuss whether such singularities would arise if the Bayesian model were solved using maximum posterior (MAP) estimation.

- 10.25** (\*\*) The variational treatment of the Bayesian mixture of Gaussians, discussed in Section 10.2, made use of a factorized approximation (10.5) to the posterior distribution. As we saw in Figure 10.2, the factorized assumption causes the variance of the posterior distribution to be under-estimated for certain directions in parameter space. Discuss qualitatively the effect this will have on the variational approximation to the model evidence, and how this effect will vary with the number of components in the mixture. Hence explain whether the variational Gaussian mixture will tend to under-estimate or over-estimate the optimal number of components.

- 10.26** (\*\*\*) Extend the variational treatment of Bayesian linear regression to include a gamma hyperprior  $\text{Gam}(\beta|c_0, d_0)$  over  $\beta$  and solve variationally, by assuming a factorized variational distribution of the form  $q(\mathbf{w})q(\alpha)q(\beta)$ . Derive the variational update equations for the three factors in the variational distribution and also obtain an expression for the lower bound and for the predictive distribution.

- 10.27** (\*\*) By making use of the formulae given in Appendix B show that the variational lower bound for the linear basis function regression model, ~~defined by (10.107)~~, can be written in the form (10.107) with the various terms defined by (10.108)–(10.112).

- 10.28** (\*\*\*) Rewrite the model for the Bayesian mixture of Gaussians, introduced in Section 10.2, as a conjugate model from the exponential family, as discussed in Section 10.4. Hence use the general results (10.115) and (10.119) to derive the specific results (10.48), (10.57), and (10.59).

不是指元为  $1/K$   
而是指  $K!$  个  $q$  distributions  
混合成的平均互信息后验  
的混合系数为  $1/K!$ 。

- 10.29** (★) **www** Show that the function  $f(x) = \ln(x)$  is concave for  $0 < x < \infty$  by computing its second derivative. Determine the form of the dual function  $g(\lambda)$  defined by (10.133), and verify that minimization of  $\lambda x - g(\lambda)$  with respect to  $\lambda$  according to (10.132) indeed recovers the function  $\ln(x)$ .
- 10.30** (★) By evaluating the second derivative, show that the log logistic function  $f(x) = -\ln(1 + e^{-x})$  is concave. Derive the variational upper bound (10.137) directly by making a ~~second~~ order Taylor expansion of the log logistic function around a point  $x = \xi$ . ~~first~~
- 10.31** (★★) By finding the second derivative with respect to  $x$ , show that the function  $f(x) = -\ln(e^{x/2} + e^{-x/2})$  is a concave function of  $x$ . Now consider the second derivatives with respect to the variable  $x^2$  and hence show that it is a convex function of  $x^2$ . Plot graphs of  $f(x)$  against  $x$  and against  $x^2$ . Derive the lower bound (10.144) on the logistic sigmoid function directly by making a first order Taylor series expansion of the function  $f(x)$  in the variable  $x^2$  centred on the value  $\xi^2$ .
- 10.32** (★★) **www** Consider the variational treatment of logistic regression with sequential learning in which data points are arriving one at a time and each must be processed and discarded before the next data point arrives. Show that a Gaussian approximation to the posterior distribution can be maintained through the use of the lower bound (10.151), in which the distribution is initialized using the prior, and as each data point is absorbed its corresponding variational parameter  $\xi_n$  is optimized.
- 10.33** (★) By differentiating the quantity  $Q(\xi, \xi^{\text{old}})$  defined by (10.161) with respect to the variational parameter  $\xi_n$  show that the update equation for  $\xi_n$  for the Bayesian logistic regression model is given by (10.163).
- 10.34** (★★) In this exercise we derive re-estimation equations for the variational parameters  $\xi$  in the Bayesian logistic regression model of Section 4.5 by direct maximization of the lower bound given by (10.164). To do this set the derivative of  $\mathcal{L}(\xi)$  with respect to  $\xi_n$  equal to zero, making use of the result (3.117) for the derivative of the log of a determinant, together with the expressions (10.157) and (10.158) which define the mean and covariance of the variational posterior distribution  $q(\mathbf{w})$ .
- 10.35** (★★) Derive the result (10.164) for the lower bound  $\mathcal{L}(\xi)$  in the variational logistic regression model. This is most easily done by substituting the expressions for the Gaussian prior  $q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$ , together with the lower bound  $h(\mathbf{w}, \xi)$  on the likelihood function, into the integral (10.159) which defines  $\mathcal{L}(\xi)$ . Next gather together the terms which depend on  $\mathbf{w}$  in the exponential and complete the square to give a Gaussian integral, which can then be evaluated by invoking the standard result for the normalization coefficient of a multivariate Gaussian. Finally take the logarithm to obtain (10.164).
- 10.36** (★★) Consider the ADF approximation scheme discussed in Section 10.7, and show that inclusion of the factor  $f_j(\theta)$  leads to an update of the model evidence of the form

$$p_j(\mathcal{D}) \simeq p_{j-1}(\mathcal{D}) Z_j \quad (10.242)$$

where  $Z_j$  is the normalization constant defined by (10.197). By applying this result recursively, and initializing with  $p_0(\mathcal{D}) = 1$ , derive the result

$$p(\mathcal{D}) \simeq \prod_j Z_j. \quad (10.243)$$

- 10.37** (\*) www Consider the expectation propagation algorithm from Section 10.7, and suppose that one of the factors  $f_0(\boldsymbol{\theta})$  in the definition (10.188) has the same exponential family functional form as the approximating distribution  $q(\boldsymbol{\theta})$ . Show that if the factor  $\tilde{f}_0(\boldsymbol{\theta})$  is initialized to be  $f_0(\boldsymbol{\theta})$ , then an EP update to refine  $\tilde{f}_0(\boldsymbol{\theta})$  leaves  $\tilde{f}_0(\boldsymbol{\theta})$  unchanged. This situation typically arises when one of the factors is the prior  $p(\boldsymbol{\theta})$ , and so we see that the prior factor can be incorporated once exactly and does not need to be refined.
- 10.38** (\*\*\*) In this exercise and the next, we shall verify the results (10.214)–(10.224) for the expectation propagation algorithm applied to the clutter problem. Begin by using the division formula (10.205) to derive the expressions (10.214) and (10.215) by completing the square inside the exponential to identify the mean and variance. Also, show that the normalization constant  $Z_n$ , defined by (10.206), is given for the clutter problem by (10.216). This can be done by making use of the general result (2.115).
- 10.39** (\*\*\*) Show that the mean and variance of  $q^{\text{new}}(\boldsymbol{\theta})$  for EP applied to the clutter problem are given by (10.217) and (10.218). To do this, first prove the following results for the expectations of  $\boldsymbol{\theta}$  and  $\boldsymbol{\theta}\boldsymbol{\theta}^T$  under  $q^{\text{new}}(\boldsymbol{\theta})$

$$\mathbb{E}[\boldsymbol{\theta}] = \mathbf{m}^{\backslash n} + v^{\backslash n} \nabla_{\mathbf{m}^{\backslash n}} \ln Z_n \quad (10.244)$$

$$\mathbb{E}[\boldsymbol{\theta}^T \boldsymbol{\theta}] = 2(v^{\backslash n})^2 \nabla_{v^{\backslash n}} \ln Z_n + 2\mathbb{E}[\boldsymbol{\theta}]^T \mathbf{m}^{\backslash n} - \|\mathbf{m}^{\backslash n}\|^2 \quad (10.245)$$

+ v^{\backslash n} D

and then make use of the result (10.216) for  $Z_n$ . Next, prove the results (10.220)–(10.222) by using (10.207) and completing the square in the exponential. Finally, use (10.208) to derive the result (10.223).