

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

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

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

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

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

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

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

Appendix B

## 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

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

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

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

and the model evidence is given by

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

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

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

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

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

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

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

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  $f_i(\boldsymbol{\theta})$  and  $\tilde{f}_i(\boldsymbol{\theta})$  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.

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  $\tilde{f}_i(\boldsymbol{\theta})$ , 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  $\tilde{f}_j(\boldsymbol{\theta})$ . We first remove this factor from the product to give  $\prod_{i \neq j} \tilde{f}_i(\boldsymbol{\theta})$ . Conceptually, we will now determine a revised form of the factor  $\tilde{f}_j(\boldsymbol{\theta})$  by ensuring that the product

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

is as close as possible to

$$f_j(\boldsymbol{\theta}) \prod_{i \neq j} \tilde{f}_i(\boldsymbol{\theta}) \quad (10.194)$$

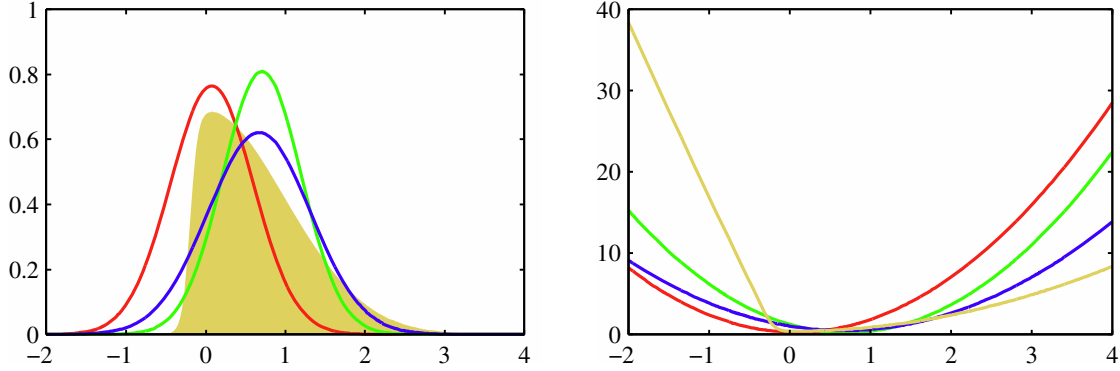
in which we keep fixed all of the factors  $\tilde{f}_i(\boldsymbol{\theta})$  for  $i \neq j$ . This ensures that the approximation is most accurate in the regions of high posterior probability as defined by the remaining factors. We shall see an example of this effect when we apply EP to the ‘clutter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

Section 10.7.1

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

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

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



**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 variational inference, as a consequence of the different form of KL divergence.

where  $Z_j$  is the normalization constant given by

$$Z_j = \int f_j(\boldsymbol{\theta})q^{\setminus j}(\boldsymbol{\theta}) d\boldsymbol{\theta}. \tag{10.197}$$

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

$$\text{KL} \left( \frac{f_j(\boldsymbol{\theta})q^{\setminus j}(\boldsymbol{\theta})}{Z_j} \parallel q^{\text{new}}(\boldsymbol{\theta}) \right). \tag{10.198}$$

This is easily solved because the approximating distribution  $q^{\text{new}}(\boldsymbol{\theta})$  is from the exponential family, and so we can appeal to the result (10.187), which tells us that the parameters of  $q^{\text{new}}(\boldsymbol{\theta})$  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  $q(\boldsymbol{\theta})$  to be a Gaussian distribution  $\mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $\boldsymbol{\mu}$  is set equal to the mean of the (unnormalized) distribution  $f_j(\boldsymbol{\theta})q^{\setminus j}(\boldsymbol{\theta})$ , and  $\boldsymbol{\Sigma}$  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.

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

$$\tilde{f}_j(\boldsymbol{\theta}) = K \frac{q^{\text{new}}(\boldsymbol{\theta})}{q^{\setminus j}(\boldsymbol{\theta})} \tag{10.199}$$

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

sides of (10.199) by  $q^{\setminus i}(\boldsymbol{\theta})$  and integrating to give

$$K = \int \tilde{f}_j(\boldsymbol{\theta}) q^{\setminus j}(\boldsymbol{\theta}) d\boldsymbol{\theta} \tag{10.200}$$

where we have used the fact that  $q^{\text{new}}(\boldsymbol{\theta})$  is normalized. The value of  $K$  can therefore be found by matching zeroth-order moments

$$\int \tilde{f}_j(\boldsymbol{\theta}) q^{\setminus j}(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int f_j(\boldsymbol{\theta}) q^{\setminus j}(\boldsymbol{\theta}) d\boldsymbol{\theta}. \tag{10.201}$$

Combining this with (10.197), we then see that  $K = Z_j$  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 factor in turn. The posterior distribution  $p(\boldsymbol{\theta}|\mathcal{D})$  is then approximated using (10.191), and the model evidence  $p(\mathcal{D})$  can be approximated by using (10.190) with the factors  $f_i(\boldsymbol{\theta})$  replaced by their approximations  $\tilde{f}_i(\boldsymbol{\theta})$ .

### Expectation Propagation

We are given a joint distribution over observed data  $\mathcal{D}$  and stochastic variables  $\boldsymbol{\theta}$  in the form of a product of factors

$$p(\mathcal{D}, \boldsymbol{\theta}) = \prod_i f_i(\boldsymbol{\theta}) \tag{10.202}$$

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

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

We also wish to approximate the model evidence  $p(\mathcal{D})$ .

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

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

3. Until convergence:
  - (a) Choose a factor  $\tilde{f}_j(\boldsymbol{\theta})$  to refine.
  - (b) Remove  $\tilde{f}_j(\boldsymbol{\theta})$  from the posterior by division

$$q^{\setminus j}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_j(\boldsymbol{\theta})}. \tag{10.205}$$

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

$$Z_j = \int q^{\setminus j}(\boldsymbol{\theta})f_j(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (10.206)$$

- (d) Evaluate and store the new factor

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

4. Evaluate the approximation to the model evidence

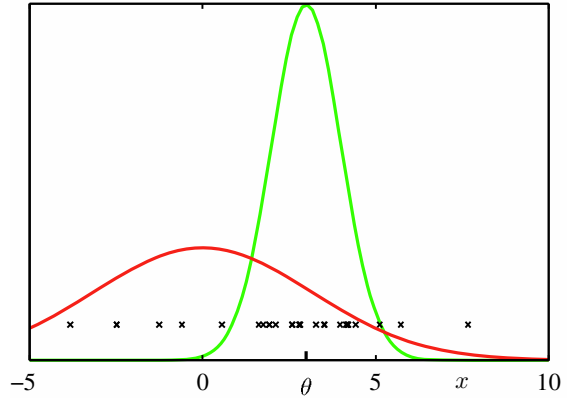
$$p(\mathcal{D}) \simeq \int \prod_i \tilde{f}_i(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (10.208)$$

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.

One disadvantage of expectation propagation is that there is no guarantee that the iterations will converge. However, for approximations  $q(\boldsymbol{\theta})$  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  $p(\boldsymbol{\theta})$  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).

**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  $\mathbf{x}$  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  $\mathbf{x}$  is therefore a mixture of Gaussians, which we take to be of the form

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

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

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

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  $\theta$  is given by

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

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

$$q(\theta) = \mathcal{N}(\theta|\mathbf{m}, v\mathbf{I}). \tag{10.212}$$

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*

First we remove the current estimate  $f_n(\boldsymbol{\theta})$  from  $q(\boldsymbol{\theta})$  by division using (10.205) to give  $q^{\setminus n}(\boldsymbol{\theta})$ , which has mean and inverse variance given by

*Exercise 10.38*

$$\mathbf{m}^{\setminus n} = \mathbf{m} + v^{\setminus n} v_n^{-1} (\mathbf{m} - \mathbf{m}_n) \quad (10.214)$$

$$(v^{\setminus n})^{-1} = v^{-1} - v_n^{-1}. \quad (10.215)$$

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

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

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

*Exercise 10.39*

$$\mathbf{m} = \mathbf{m}^{\setminus n} + \rho_n \frac{v^{\setminus n}}{v^{\setminus n} + 1} (\mathbf{x}_n - \mathbf{m}^{\setminus n}) \quad (10.217)$$

$$v = v^{\setminus n} - \rho_n \frac{(v^{\setminus n})^2}{v^{\setminus n} + 1} + \rho_n (1 - \rho_n) \frac{(v^{\setminus n})^2 \|\mathbf{x}_n - \mathbf{m}^{\setminus n}\|^2}{D(v^{\setminus n} + 1)^2} \quad (10.218)$$

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

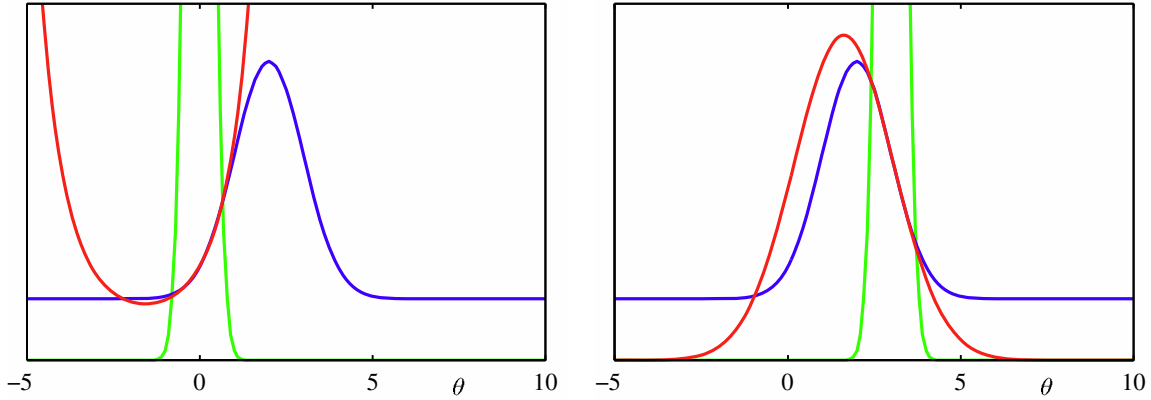
$$v_n^{-1} = (v^{\text{new}})^{-1} - (v^{\setminus n})^{-1} \quad (10.220)$$

$$\mathbf{m}_n = \mathbf{m}^{\setminus n} + (v_n + v^{\setminus n})(v^{\setminus n})^{-1} (\mathbf{m}^{\text{new}} - \mathbf{m}^{\setminus n}) \quad (10.221)$$

$$s_n = \frac{Z_n}{(2\pi v_n)^{D/2} \mathcal{N}(\mathbf{m}_n | \mathbf{m}^{\setminus n}, (v_n + v^{\setminus n}) \mathbf{I})}. \quad (10.222)$$

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





**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  $f_n(\theta)$ .

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

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

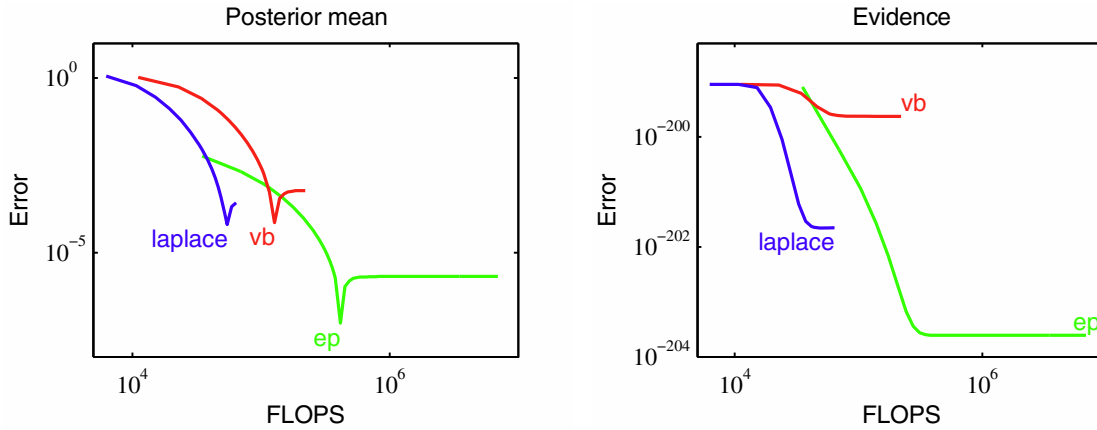
where

$$B = \frac{(\mathbf{m}^{\text{new}})^T \mathbf{m}^{\text{new}}}{v} - \sum_{n=1}^N \frac{\mathbf{m}_n^T \mathbf{m}_n}{v_n}. \tag{10.224}$$

Examples 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  $q(\boldsymbol{\theta})$  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

So far in our general discussion of EP, we have allowed the factors  $f_i(\boldsymbol{\theta})$  in the distribution  $p(\boldsymbol{\theta})$  to be functions of all of the components of  $\boldsymbol{\theta}$ , and similarly for the approximating factors  $\tilde{f}_i(\boldsymbol{\theta})$  in the approximating distribution  $q(\boldsymbol{\theta})$ . 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.



**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.

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  $p$ .

*Section 8.4.4*

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

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

We seek an approximation  $q(\mathbf{x})$  that has the same factorization, so that

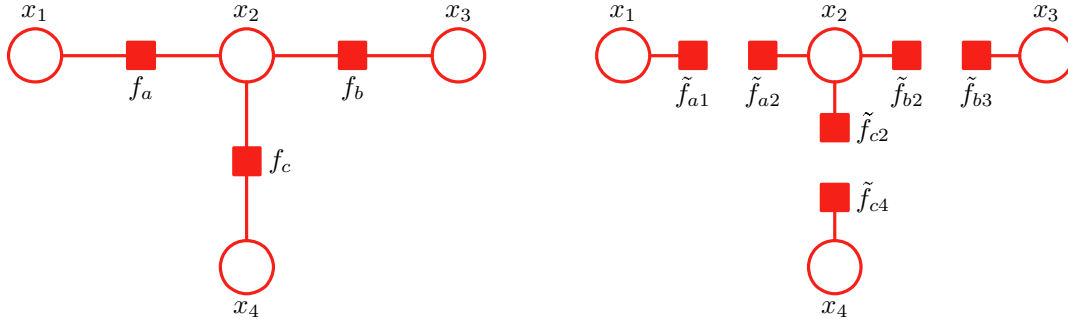
$$q(\mathbf{x}) \propto \tilde{f}_a(x_1, x_2) \tilde{f}_b(x_2, x_3) \tilde{f}_c(x_2, x_4). \tag{10.226}$$

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

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

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.

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



**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

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

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

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

We now find  $q^{\text{new}}(\mathbf{x})$  by minimizing the Kullback-Leibler divergence  $\text{KL}(\hat{p}||q^{\text{new}})$ . 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

$$\hat{p}(x_1) \propto \tilde{f}_{a1}(x_1) \quad (10.230)$$

$$\hat{p}(x_2) \propto \tilde{f}_{a2}(x_2)\tilde{f}_{c2}(x_2) \sum_{x_3} f_b(x_2, x_3) \quad (10.231)$$

$$\hat{p}(x_3) \propto \sum_{x_2} \left\{ f_b(x_2, x_3)\tilde{f}_{a2}(x_2)\tilde{f}_{c2}(x_2) \right\} \quad (10.232)$$

$$\hat{p}(x_4) \propto \tilde{f}_{c4}(x_4) \quad (10.233)$$

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  $f_b(x_2, x_3)$  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  $q^{\setminus b}(\mathbf{x})$ , which gives

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

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

## 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  $\mu_{f_b \rightarrow x_3}(x_3)$ .

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.

Now let us consider a general factor graph corresponding to the distribution

$$p(\boldsymbol{\theta}) = \prod_i f_i(\boldsymbol{\theta}_i) \quad (10.236)$$

where  $\boldsymbol{\theta}_i$  represents the subset of variables associated with factor  $f_i$ . We approximate this using a fully factorized distribution of the form

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

where  $\theta_k$  corresponds to an individual variable node. Suppose that we wish to refine the particular term  $\tilde{f}_{jl}(\theta_l)$  keeping all other terms fixed. We first remove the term  $\tilde{f}_j(\boldsymbol{\theta}_j)$  from  $q(\boldsymbol{\theta})$  to give

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

and then multiply by the exact factor  $f_j(\boldsymbol{\theta}_j)$ . 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 corresponding marginal of

$$q^{\setminus j}(\boldsymbol{\theta}) f_j(\boldsymbol{\theta}_j). \quad (10.239)$$

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

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

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 (10.240) 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.

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).

---

## 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  $\text{KL}(q||p)$  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  $\text{KL}(p||q)$  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  $\text{KL}(p||q)$  for a Gaussian  $q(\mathbf{x})$  and then differentiating, show that