

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.

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 $f(x) =$

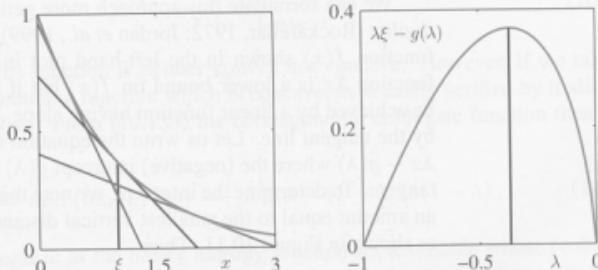


Figure 10.10 In the left-hand figure the red curve shows the function $\exp(-x)$, and the blue line shows the tangent at $x = \xi$ defined by (10.125) with $\xi = 1$. This line has slope $\lambda = f'(\xi) = -\exp(-\xi)$. Note that any other tangent line, for example the ones shown in green, will have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

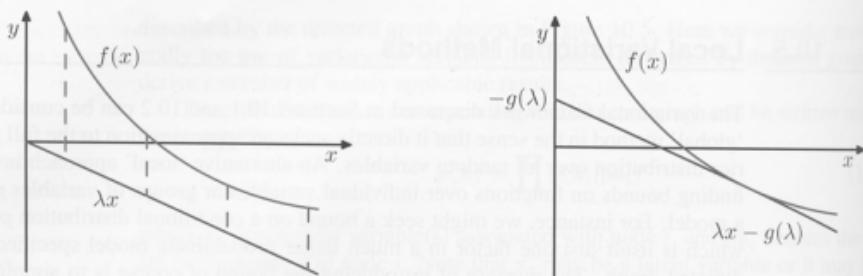


Figure 10.11 In the left-hand plot the red curve shows a convex function $f(x)$, and the blue line represents the linear function λx , which is a lower bound on $f(x)$ because $f(x) > \lambda x$ for all x . For the given value of slope λ the contact point of the tangent line having the same slope is 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 λ .

$\exp(-x)$, 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 ξ . 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 λ 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 λ , and to obtain the tightest bound we must optimize with respect to λ .

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 λx 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 λ , because the tightest bound is given by the tangent line. Let us write the equation of the tangent line, having slope λ as $\lambda x - g(\lambda)$ where the (negative) intercept $g(\lambda)$ clearly depends on the slope λ 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

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

Now, instead of fixing λ and varying x , we can consider a particular x and then adjust λ 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

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

We see that the functions $f(x)$ and $g(\lambda)$ 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 $g(\lambda)$ in the form

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

as obtained previously. The function $\lambda\xi - g(\lambda)$ 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 $f(x) = \exp(-x)$.

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

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

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

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

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

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

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

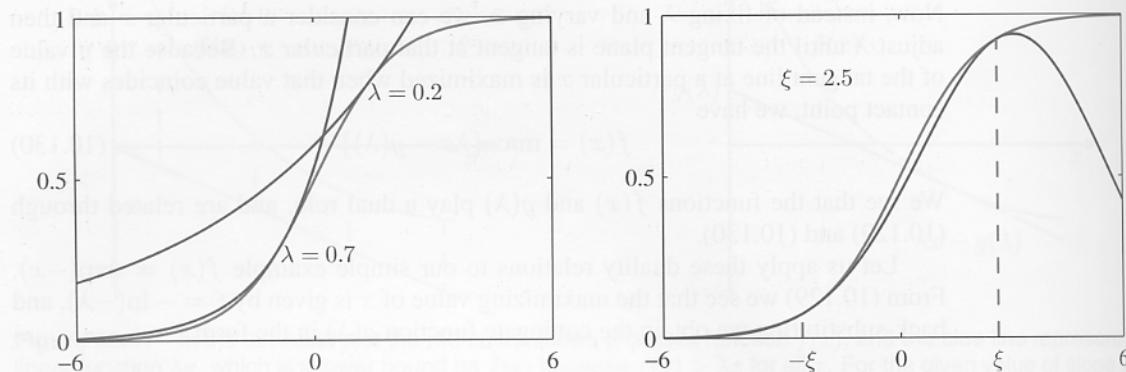


Figure 10.12 The left-hand plot shows the logistic sigmoid function $\sigma(x)$ defined by (10.134) in red, together with two examples of the exponential upper bound (10.137) shown in blue. The right-hand plot shows the logistic sigmoid again in red together with the Gaussian lower bound (10.144) shown 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

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

which is plotted for two values of λ 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

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

The stationarity condition leads to

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

If we denote this value of x , corresponding to the contact point of the tangent line for this particular value of λ , by ξ , then we have

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

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

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

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

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

The bound on the sigmoid then becomes

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

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.

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 ξ 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 ξ , which we do by finding the value ξ^* 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.

Although the bound $\sigma(a) \geq f(a, \xi)$ on the logistic sigmoid can be optimized exactly, the required choice for ξ 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 ξ^* represents a compromise, weighted by the distribution $p(a)$.

10.6. Variational Logistic Regression

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 Dybowski and Roberts (2005) from a Bayesian perspective using Monte Carlo sampling techniques.

10.6.1 Variational posterior distribution

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

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

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

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

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

we reproduce here for convenience

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

where

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

We can therefore write

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

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

$$p(t, w) = p(t|w)p(w) \geq h(w, \xi)p(w) \quad (10.152)$$

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

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

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 $q(w)$, 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 t and w of the form

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

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

$$\begin{aligned} &-\frac{1}{2}(w - m_0)^T S_0^{-1}(w - m_0) \\ &+ \sum_{n=1}^N \{ w^T \phi_n (t_n - 1/2) - \lambda(\xi_n) w^T (\phi_n \phi_n^T) w \} + \text{const.} \end{aligned} \quad (10.155)$$

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

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) \quad (10.156)$$

where

$$\mathbf{m}_N = \mathbf{S}_N \left(\mathbf{S}_0^{-1} \mathbf{m}_0 + \sum_{n=1}^N (t_n - 1/2) \boldsymbol{\phi}_n \right) \quad (10.157)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + 2 \sum_{n=1}^N \lambda(\xi_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T. \quad (10.158)$$

As with the Laplace framework, we have again obtained a Gaussian approximation to the posterior distribution. However, the additional flexibility provided by the variational parameters $\{\xi_n\}$ leads to improved accuracy in the approximation (Jaakkola and Jordan, 2000).

Here we have considered a batch learning context in which all of the training data is available at once. However, Bayesian methods are intrinsically well suited to sequential learning in which the data points are processed one at a time and then discarded. The formulation of this variational approach for the sequential case is straightforward.

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

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

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

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

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 ξ^{old} . Noting that $p(\mathbf{w})$ does not depend on ξ , 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 ξ . We now set the derivative with respect to ξ_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 ξ for $\xi \geq 0$, and that we can restrict attention to nonnegative values of ξ 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).

Let us summarize the EM algorithm for finding the variational posterior distribution. We first initialize the variational parameters ξ^{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 ξ 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 ξ 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 ξ_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

Exercise 10.33

Exercise 10.34

Exercise 10.35

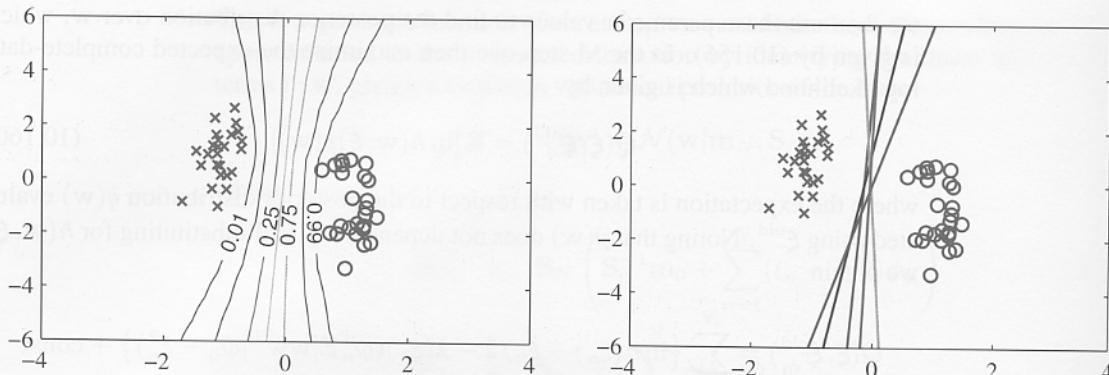


Figure 10.13 Illustration of the Bayesian approach to logistic regression for a simple linearly separable data set. The plot on the left shows the predictive distribution obtained using variational inference. We see that the decision boundary lies roughly mid way between the clusters of data points, 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 $p(w|t)$.

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

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.

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 α 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.