

# Pattern Recognition and Machine Learning Solutions to the Exercises: Tutors' Edition

### Markus Svensén and Christopher M. Bishop

Copyright © 2002–2009

This is the solutions manual (Tutors' Edition) for the book *Pattern Recognition and Machine Learning* (PRML; published by Springer in 2006). This release was created September 8, 2009. Any future releases (e.g. with corrections to errors) will be announced on the PRML web-site (see below) and published via Springer.

#### PLEASE DO NOT DISTRIBUTE

Most of the solutions in this manual are intended as a resource for tutors teaching courses based on PRML and the value of this resource would be greatly diminished if was to become generally available. All tutors who want a copy should contact Springer directly.

The authors would like to express their gratitude to the various people who have provided feedback on earlier releases of this document.

The authors welcome all comments, questions and suggestions about the solutions as well as reports on (potential) errors in text or formulae in this document; please send any such feedback to

prml-fb@microsoft.com

Further information about PRML is available from

http://research.microsoft.com/~cmbishop/PRML

# **Contents**

Contents	5
Chapter 1: Introduction	7
Chapter 2: Probability Distributions	28
Chapter 3: Linear Models for Regression	62
Chapter 4: Linear Models for Classification	78
Chapter 5: Neural Networks	93
Chapter 6: Kernel Methods	114
Chapter 7: Sparse Kernel Machines	128
Chapter 8: Graphical Models	136
Chapter 9: Mixture Models and EM	150
Chapter 10: Approximate Inference	163
Chapter 11: Sampling Methods	198
Chapter 12: Continuous Latent Variables	207
Chapter 13: Sequential Data	223
Chapter 14: Combining Models	

# 6 CONTENTS

# **Chapter 1** Introduction

**1.1** Substituting (1.1) into (1.2) and then differentiating with respect to  $w_i$  we obtain

$$\sum_{n=1}^{N} \left( \sum_{j=0}^{M} w_j x_n^j - t_n \right) x_n^i = 0.$$
 (1)

Re-arranging terms then gives the required result.

**1.2** For the regularized sum-of-squares error function given by (1.4) the corresponding linear equations are again obtained by differentiation, and take the same form as (1.122), but with  $A_{ij}$  replaced by  $\widetilde{A}_{ij}$ , given by

$$\widetilde{A}_{ij} = A_{ij} + \lambda I_{ij}. (2)$$

**1.3** Let us denote apples, oranges and limes by a, o and l respectively. The marginal probability of selecting an apple is given by

$$p(a) = p(a|r)p(r) + p(a|b)p(b) + p(a|g)p(g)$$

$$= \frac{3}{10} \times 0.2 + \frac{1}{2} \times 0.2 + \frac{3}{10} \times 0.6 = 0.34$$
(3)

where the conditional probabilities are obtained from the proportions of apples in each box.

To find the probability that the box was green, given that the fruit we selected was an orange, we can use Bayes' theorem

$$p(g|o) = \frac{p(o|g)p(g)}{p(o)}. (4)$$

The denominator in (4) is given by

$$p(o) = p(o|r)p(r) + p(o|b)p(b) + p(o|g)p(g)$$

$$= \frac{4}{10} \times 0.2 + \frac{1}{2} \times 0.2 + \frac{3}{10} \times 0.6 = 0.36$$
(5)

from which we obtain

$$p(g|o) = \frac{3}{10} \times \frac{0.6}{0.36} = \frac{1}{2}.$$
 (6)

1.4 We are often interested in finding the most probable value for some quantity. In the case of probability distributions over discrete variables this poses little problem. However, for continuous variables there is a subtlety arising from the nature of probability densities and the way they transform under non-linear changes of variable.

#### 8 Solution 1.4

Consider first the way a function f(x) behaves when we change to a new variable y where the two variables are related by x = g(y). This defines a new function of y given by

$$\widetilde{f}(y) = f(g(y)). \tag{7}$$

Suppose f(x) has a mode (i.e. a maximum) at  $\widehat{x}$  so that  $f'(\widehat{x}) = 0$ . The corresponding mode of  $\widetilde{f}(y)$  will occur for a value  $\widehat{y}$  obtained by differentiating both sides of (7) with respect to y

$$\widetilde{f}'(\widehat{y}) = f'(g(\widehat{y}))g'(\widehat{y}) = 0.$$
(8)

Assuming  $g'(\widehat{y}) \neq 0$  at the mode, then  $f'(g(\widehat{y})) = 0$ . However, we know that  $f'(\widehat{x}) = 0$ , and so we see that the locations of the mode expressed in terms of each of the variables x and y are related by  $\widehat{x} = g(\widehat{y})$ , as one would expect. Thus, finding a mode with respect to the variable x is completely equivalent to first transforming to the variable y, then finding a mode with respect to y, and then transforming back to x.

Now consider the behaviour of a probability density  $p_x(x)$  under the change of variables x=g(y), where the density with respect to the new variable is  $p_y(y)$  and is given by ((1.27)). Let us write g'(y)=s|g'(y)| where  $s\in\{-1,+1\}$ . Then ((1.27)) can be written

$$p_y(y) = p_x(g(y))sg'(y).$$

Differentiating both sides with respect to y then gives

$$p'_{y}(y) = sp'_{x}(g(y))\{g'(y)\}^{2} + sp_{x}(g(y))g''(y).$$
(9)

Due to the presence of the second term on the right hand side of (9) the relationship  $\widehat{x}=g(\widehat{y})$  no longer holds. Thus the value of x obtained by maximizing  $p_x(x)$  will not be the value obtained by transforming to  $p_y(y)$  then maximizing with respect to y and then transforming back to x. This causes modes of densities to be dependent on the choice of variables. In the case of linear transformation, the second term on the right hand side of (9) vanishes, and so the location of the maximum transforms according to  $\widehat{x}=g(\widehat{y})$ .

This effect can be illustrated with a simple example, as shown in Figure 1. We begin by considering a Gaussian distribution  $p_x(x)$  over x with mean  $\mu=6$  and standard deviation  $\sigma=1$ , shown by the red curve in Figure 1. Next we draw a sample of N=50,000 points from this distribution and plot a histogram of their values, which as expected agrees with the distribution  $p_x(x)$ .

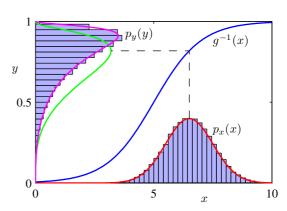
Now consider a non-linear change of variables from x to y given by

$$x = g(y) = \ln(y) - \ln(1 - y) + 5. \tag{10}$$

The inverse of this function is given by

$$y = g^{-1}(x) = \frac{1}{1 + \exp(-x + 5)} \tag{11}$$

Figure 1 Example of the transformation of the mode of a density under a non-linear change of variables, illustrating the different behaviour compared to a simple function. See the text for details.



which is a logistic sigmoid function, and is shown in Figure 1 by the blue curve.

If we simply transform  $p_x(x)$  as a function of x we obtain the green curve  $p_x(g(y))$  shown in Figure 1, and we see that the mode of the density  $p_x(x)$  is transformed via the sigmoid function to the mode of this curve. However, the density over y transforms instead according to (1.27) and is shown by the magenta curve on the left side of the diagram. Note that this has its mode shifted relative to the mode of the green curve.

To confirm this result we take our sample of 50,000 values of x, evaluate the corresponding values of y using (11), and then plot a histogram of their values. We see that this histogram matches the magenta curve in Figure 1 and not the green curve!

**1.5** Expanding the square we have

$$\mathbb{E}[(f(x) - \mathbb{E}[f(x)])^{2}] = \mathbb{E}[f(x)^{2} - 2f(x)\mathbb{E}[f(x)] + \mathbb{E}[f(x)]^{2}] 
= \mathbb{E}[f(x)^{2}] - 2\mathbb{E}[f(x)]\mathbb{E}[f(x)] + \mathbb{E}[f(x)]^{2} 
= \mathbb{E}[f(x)^{2}] - \mathbb{E}[f(x)]^{2}$$

as required.

**1.6** The definition of covariance is given by (1.41) as

$$cov[x, y] = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y].$$

Using (1.33) and the fact that p(x,y)=p(x)p(y) when x and y are independent, we obtain

$$\begin{split} \mathbb{E}[xy] &= \sum_{x} \sum_{y} p(x,y) xy \\ &= \sum_{x} p(x) x \sum_{y} p(y) y \\ &= \mathbb{E}[x] \mathbb{E}[y] \end{split}$$

#### **10** Solutions 1.7–1.8

and hence cov[x, y] = 0. The case where x and y are continuous variables is analogous, with (1.33) replaced by (1.34) and the sums replaced by integrals.

1.7 The transformation from Cartesian to polar coordinates is defined by

$$x = r\cos\theta \tag{12}$$

$$y = r\sin\theta \tag{13}$$

and hence we have  $x^2+y^2=r^2$  where we have used the well-known trigonometric result (2.177). Also the Jacobian of the change of variables is easily seen to be

$$\frac{\partial(x,y)}{\partial(r,\theta)} = \begin{vmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{vmatrix} \\
= \begin{vmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{vmatrix} = r$$

where again we have used (2.177). Thus the double integral in (1.125) becomes

$$I^{2} = \int_{0}^{2\pi} \int_{0}^{\infty} \exp\left(-\frac{r^{2}}{2\sigma^{2}}\right) r \, \mathrm{d}r \, \mathrm{d}\theta \tag{14}$$

$$= 2\pi \int_0^\infty \exp\left(-\frac{u}{2\sigma^2}\right) \frac{1}{2} du \tag{15}$$

$$= \pi \left[ \exp\left(-\frac{u}{2\sigma^2}\right) \left(-2\sigma^2\right) \right]_0^{\infty} \tag{16}$$

$$= 2\pi\sigma^2 \tag{17}$$

where we have used the change of variables  $r^2 = u$ . Thus

$$I = \left(2\pi\sigma^2\right)^{1/2}.$$

Finally, using the transformation  $y=x-\mu$ , the integral of the Gaussian distribution becomes

$$\int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) dx = \frac{1}{(2\pi\sigma^2)^{1/2}} \int_{-\infty}^{\infty} \exp\left(-\frac{y^2}{2\sigma^2}\right) dy$$
$$= \frac{I}{(2\pi\sigma^2)^{1/2}} = 1$$

as required.

**1.8** From the definition (1.46) of the univariate Gaussian distribution, we have

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \left(\frac{1}{2\pi\sigma^2}\right)^{1/2} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} x \,\mathrm{d}x. \tag{18}$$

Now change variables using  $y = x - \mu$  to give

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \left(\frac{1}{2\pi\sigma^2}\right)^{1/2} \exp\left\{-\frac{1}{2\sigma^2}y^2\right\} (y+\mu) \,\mathrm{d}y. \tag{19}$$

We now note that in the factor  $(y + \mu)$  the first term in y corresponds to an odd integrand and so this integral must vanish (to show this explicitly, write the integral as the sum of two integrals, one from  $-\infty$  to 0 and the other from 0 to  $\infty$  and then show that these two integrals cancel). In the second term,  $\mu$  is a constant and pulls outside the integral, leaving a normalized Gaussian distribution which integrates to 1, and so we obtain (1.49).

To derive (1.50) we first substitute the expression (1.46) for the normal distribution into the normalization result (1.48) and re-arrange to obtain

$$\int_{-\infty}^{\infty} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} dx = \left(2\pi\sigma^2\right)^{1/2}.$$
 (20)

We now differentiate both sides of (20) with respect to  $\sigma^2$  and then re-arrange to obtain

$$\left(\frac{1}{2\pi\sigma^2}\right)^{1/2} \int_{-\infty}^{\infty} \exp\left\{-\frac{1}{2\sigma^2} (x-\mu)^2\right\} (x-\mu)^2 \, \mathrm{d}x = \sigma^2$$
 (21)

which directly shows that

$$\mathbb{E}[(x-\mu)^2] = \text{var}[x] = \sigma^2. \tag{22}$$

Now we expand the square on the left-hand side giving

$$\mathbb{E}[x^2] - 2\mu \mathbb{E}[x] + \mu^2 = \sigma^2.$$

Making use of (1.49) then gives (1.50) as required.

Finally, (1.51) follows directly from (1.49) and (1.50)

$$\mathbb{E}[x^2] - \mathbb{E}[x]^2 = \left(\mu^2 + \sigma^2\right) - \mu^2 = \sigma^2.$$

**1.9** For the univariate case, we simply differentiate (1.46) with respect to x to obtain

$$\frac{\mathrm{d}}{\mathrm{d}x} \mathcal{N}\left(x|\mu, \sigma^2\right) = -\mathcal{N}\left(x|\mu, \sigma^2\right) \frac{x-\mu}{\sigma^2}.$$

Setting this to zero we obtain  $x = \mu$ .

Similarly, for the multivariate case we differentiate (1.52) with respect to x to obtain

$$\begin{split} \frac{\partial}{\partial \mathbf{x}} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= -\frac{1}{2} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \nabla_{\mathbf{x}} \left\{ (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= -\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}), \end{split}$$

where we have used (C.19), (C.20)<sup>1</sup> and the fact that  $\Sigma^{-1}$  is symmetric. Setting this

derivative equal to 0, and left-multiplying by  $\Sigma$ , leads to the solution  $x = \mu$ .

NOTE: In the 1<sup>st</sup> printing of PRML, there are mistakes in (C.20), all instances of x (vector) denominators should be x (scalar).

#### **Solutions 1.10–1.11**

**1.10** Since x and z are independent, their joint distribution factorizes p(x,z) = p(x)p(z), and so

$$\mathbb{E}[x+z] = \iint (x+z)p(x)p(z) \,dx \,dz \tag{23}$$

$$= \int xp(x) dx + \int zp(z) dz$$
 (24)

$$= \mathbb{E}[x] + \mathbb{E}[z]. \tag{25}$$

Similarly for the variances, we first note that

$$(x+z-\mathbb{E}[x+z])^2 = (x-\mathbb{E}[x])^2 + (z-\mathbb{E}[z])^2 + 2(x-\mathbb{E}[x])(z-\mathbb{E}[z])$$
 (26)

where the final term will integrate to zero with respect to the factorized distribution p(x)p(z). Hence

$$\operatorname{var}[x+z] = \iint (x+z-\mathbb{E}[x+z])^2 p(x) p(z) \, \mathrm{d}x \, \mathrm{d}z$$

$$= \int (x-\mathbb{E}[x])^2 p(x) \, \mathrm{d}x + \int (z-\mathbb{E}[z])^2 p(z) \, \mathrm{d}z$$

$$= \operatorname{var}(x) + \operatorname{var}(z). \tag{27}$$

For discrete variables the integrals are replaced by summations, and the same results are again obtained.

**1.11** We use  $\ell$  to denote  $\ln p(\mathbf{X}|\mu, \sigma^2)$  from (1.54). By standard rules of differentiation we obtain

$$\frac{\partial \ell}{\partial \mu} = \frac{1}{\sigma^2} \sum_{n=1}^{N} (x_n - \mu).$$

Setting this equal to zero and moving the terms involving  $\mu$  to the other side of the equation we get

$$\frac{1}{\sigma^2} \sum_{n=1}^{N} x_n = \frac{1}{\sigma^2} N \mu$$

and by multiplying ing both sides by  $\sigma^2/N$  we get (1.55).

Similarly we have

$$\frac{\partial \ell}{\partial \sigma^2} = \frac{1}{2(\sigma^2)^2} \sum_{n=1}^{N} (x_n - \mu)^2 - \frac{N}{2} \frac{1}{\sigma^2}$$

and setting this to zero we obtain

$$\frac{N}{2} \frac{1}{\sigma^2} = \frac{1}{2(\sigma^2)^2} \sum_{n=1}^{N} (x_n - \mu)^2.$$

Multiplying both sides by  $2(\sigma^2)^2/N$  and substituting  $\mu_{\rm ML}$  for  $\mu$  we get (1.56).

**1.12** If m=n then  $x_nx_m=x_n^2$  and using (1.50) we obtain  $\mathbb{E}[x_n^2]=\mu^2+\sigma^2$ , whereas if  $n\neq m$  then the two data points  $x_n$  and  $x_m$  are independent and hence  $\mathbb{E}[x_nx_m]=\mathbb{E}[x_n]\mathbb{E}[x_m]=\mu^2$  where we have used (1.49). Combining these two results we obtain (1.130).

Next we have

$$\mathbb{E}[\mu_{\mathrm{ML}}] = \frac{1}{N} \sum_{n=1}^{N} \mathbb{E}[x_n] = \mu \tag{28}$$

using (1.49).

Finally, consider  $\mathbb{E}[\sigma_{\mathrm{ML}}^2]$ . From (1.55) and (1.56), and making use of (1.130), we have

$$\mathbb{E}[\sigma_{\text{ML}}^{2}] = \mathbb{E}\left[\frac{1}{N}\sum_{n=1}^{N}\left(x_{n} - \frac{1}{N}\sum_{m=1}^{N}x_{m}\right)^{2}\right]$$

$$= \frac{1}{N}\sum_{n=1}^{N}\mathbb{E}\left[x_{n}^{2} - \frac{2}{N}x_{n}\sum_{m=1}^{N}x_{m} + \frac{1}{N^{2}}\sum_{m=1}^{N}\sum_{l=1}^{N}x_{m}x_{l}\right]$$

$$= \left\{\mu^{2} + \sigma^{2} - 2\left(\mu^{2} + \frac{1}{N}\sigma^{2}\right) + \mu^{2} + \frac{1}{N}\sigma^{2}\right\}$$

$$= \left(\frac{N-1}{N}\right)\sigma^{2}$$
(29)

as required.

**1.13** In a similar fashion to solution 1.12, substituting  $\mu$  for  $\mu_{\rm ML}$  in (1.56) and using (1.49) and (1.50) we have

$$\mathbb{E}_{\{\mathbf{x}_n\}} \left[ \frac{1}{N} \sum_{n=1}^{N} (x_n - \mu)^2 \right] = \frac{1}{N} \sum_{n=1}^{N} \mathbb{E}_{\mathbf{x}_n} \left[ x_n^2 - 2x_n \mu + \mu^2 \right]$$
$$= \frac{1}{N} \sum_{n=1}^{N} \left( \mu^2 + \sigma^2 - 2\mu \mu + \mu^2 \right)$$
$$= \sigma^2$$

**1.14** Define

$$w_{ij}^{S} = \frac{1}{2}(w_{ij} + w_{ji})$$
  $w_{ij}^{A} = \frac{1}{2}(w_{ij} - w_{ji}).$  (30)

from which the (anti)symmetry properties follow directly, as does the relation  $w_{ij}=w_{ij}^{\rm S}+w_{ij}^{\rm A}$ . We now note that

$$\sum_{i=1}^{D} \sum_{j=1}^{D} w_{ij}^{A} x_i x_j = \frac{1}{2} \sum_{i=1}^{D} \sum_{j=1}^{D} w_{ij} x_i x_j - \frac{1}{2} \sum_{i=1}^{D} \sum_{j=1}^{D} w_{ji} x_i x_j = 0$$
 (31)

from which we obtain (1.132). The number of independent components in  $w_{ij}^{\rm S}$  can be found by noting that there are  $D^2$  parameters in total in this matrix, and that entries off the leading diagonal occur in constrained pairs  $w_{ij}=w_{ji}$  for  $j\neq i$ . Thus we start with  $D^2$  parameters in the matrix  $w_{ij}^{\rm S}$ , subtract D for the number of parameters on the leading diagonal, divide by two, and then add back D for the leading diagonal and we obtain  $(D^2-D)/2+D=D(D+1)/2$ .

**1.15** The redundancy in the coefficients in (1.133) arises from interchange symmetries between the indices  $i_k$ . Such symmetries can therefore be removed by enforcing an ordering on the indices, as in (1.134), so that only one member in each group of equivalent configurations occurs in the summation.

To derive (1.135) we note that the number of independent parameters n(D, M) which appear at order M can be written as

$$n(D,M) = \sum_{i_1=1}^{D} \sum_{i_2=1}^{i_1} \cdots \sum_{i_M=1}^{i_{M-1}} 1$$
 (32)

which has M terms. This can clearly also be written as

$$n(D,M) = \sum_{i_1=1}^{D} \left\{ \sum_{i_2=1}^{i_1} \cdots \sum_{i_M=1}^{i_{M-1}} 1 \right\}$$
 (33)

where the term in braces has M-1 terms which, from (32), must equal  $n(i_1, M-1)$ . Thus we can write

$$n(D,M) = \sum_{i_1=1}^{D} n(i_1, M-1)$$
(34)

which is equivalent to (1.135).

To prove (1.136) we first set D=1 on both sides of the equation, and make use of 0!=1, which gives the value 1 on both sides, thus showing the equation is valid for D=1. Now we assume that it is true for a specific value of dimensionality D and then show that it must be true for dimensionality D+1. Thus consider the left-hand side of (1.136) evaluated for D+1 which gives

$$\sum_{i=1}^{D+1} \frac{(i+M-2)!}{(i-1)!(M-1)!} = \frac{(D+M-1)!}{(D-1)!M!} + \frac{(D+M-1)!}{D!(M-1)!}$$

$$= \frac{(D+M-1)!D + (D+M-1)!M}{D!M!}$$

$$= \frac{(D+M)!}{D!M!}$$
(35)

which equals the right hand side of (1.136) for dimensionality D+1. Thus, by induction, (1.136) must hold true for all values of D.

Finally we use induction to prove (1.137). For M=2 we find obtain the standard result  $n(D,2)=\frac{1}{2}D(D+1)$ , which is also proved in Exercise 1.14. Now assume that (1.137) is correct for a specific order M-1 so that

$$n(D, M-1) = \frac{(D+M-2)!}{(D-1)!(M-1)!}.$$
(36)

Substituting this into the right hand side of (1.135) we obtain

$$n(D,M) = \sum_{i=1}^{D} \frac{(i+M-2)!}{(i-1)!(M-1)!}$$
(37)

which, making use of (1.136), gives

$$n(D,M) = \frac{(D+M-1)!}{(D-1)!M!}$$
(38)

and hence shows that (1.137) is true for polynomials of order M. Thus by induction (1.137) must be true for all values of M.

**1.16 NOTE**: In the 1<sup>st</sup> printing of PRML, this exercise contains two typographical errors. On line 4, M6th should be  $M^{th}$  and on the l.h.s. of (1.139), N(d, M) should be N(D, M).

The result (1.138) follows simply from summing up the coefficients at all order up to and including order M. To prove (1.139), we first note that when M=0 the right hand side of (1.139) equals 1, which we know to be correct since this is the number of parameters at zeroth order which is just the constant offset in the polynomial. Assuming that (1.139) is correct at order M, we obtain the following result at order M+1

$$N(D, M+1) = \sum_{m=0}^{M+1} n(D, m)$$

$$= \sum_{m=0}^{M} n(D, m) + n(D, M+1)$$

$$= \frac{(D+M)!}{D!M!} + \frac{(D+M)!}{(D-1)!(M+1)!}$$

$$= \frac{(D+M)!(M+1) + (D+M)!D}{D!(M+1)!}$$

$$= \frac{(D+M+1)!}{D!(M+1)!}$$

which is the required result at order M+1.

Now assume  $M \gg D$ . Using Stirling's formula we have

$$n(D,M) \simeq \frac{(D+M)^{D+M}e^{-D-M}}{D! M^M e^{-M}}$$

$$= \frac{M^{D+M}e^{-D}}{D! M^M} \left(1 + \frac{D}{M}\right)^{D+M}$$

$$\simeq \frac{M^D e^{-D}}{D!} \left(1 + \frac{D(D+M)}{M}\right)$$

$$\simeq \frac{(1+D)e^{-D}}{D!} M^D$$

which grows like  $M^D$  with M. The case where  $D\gg M$  is identical, with the roles of D and M exchanged. By numerical evaluation we obtain N(10,3)=286 and N(100,3)=176,851.

**1.17** Using integration by parts we have

$$\Gamma(x+1) = \int_0^\infty u^x e^{-u} du$$

$$= \left[ -e^{-u} u^x \right]_0^\infty + \int_0^\infty x u^{x-1} e^{-u} du = 0 + x \Gamma(x).$$
 (39)

For x = 1 we have

$$\Gamma(1) = \int_0^\infty e^{-u} du = \left[ -e^{-u} \right]_0^\infty = 1.$$
 (40)

If x is an integer we can apply proof by induction to relate the gamma function to the factorial function. Suppose that  $\Gamma(x+1)=x!$  holds. Then from the result (39) we have  $\Gamma(x+2)=(x+1)\Gamma(x+1)=(x+1)!$ . Finally,  $\Gamma(1)=1=0!$ , which completes the proof by induction.

**1.18** On the right-hand side of (1.142) we make the change of variables  $u = r^2$  to give

$$\frac{1}{2}S_D \int_0^\infty e^{-u} u^{D/2-1} \, \mathrm{d}u = \frac{1}{2}S_D \Gamma(D/2) \tag{41}$$

where we have used the definition (1.141) of the Gamma function. On the left hand side of (1.142) we can use (1.126) to obtain  $\pi^{D/2}$ . Equating these we obtain the desired result (1.143).

The volume of a sphere of radius 1 in D-dimensions is obtained by integration

$$V_D = S_D \int_0^1 r^{D-1} \, \mathrm{d}r = \frac{S_D}{D}.$$
 (42)

For D=2 and D=3 we obtain the following results

$$S_2 = 2\pi,$$
  $S_3 = 4\pi,$   $V_2 = \pi a^2,$   $V_3 = \frac{4}{3}\pi a^3.$  (43)

**1.19** The volume of the cube is  $(2a)^D$ . Combining this with (1.143) and (1.144) we obtain (1.145). Using Stirling's formula (1.146) in (1.145) the ratio becomes, for large D,

$$\frac{\text{volume of sphere}}{\text{volume of cube}} = \left(\frac{\pi e}{2D}\right)^{D/2} \frac{1}{D}$$
 (44)

which goes to 0 as  $D\to\infty$ . The distance from the center of the cube to the mid point of one of the sides is a, since this is where it makes contact with the sphere. Similarly the distance to one of the corners is  $a\sqrt{D}$  from Pythagoras' theorem. Thus the ratio is  $\sqrt{D}$ .

**1.20** Since  $p(\mathbf{x})$  is radially symmetric it will be roughly constant over the shell of radius r and thickness  $\epsilon$ . This shell has volume  $S_D r^{D-1} \epsilon$  and since  $\|\mathbf{x}\|^2 = r^2$  we have

$$\int_{\text{shell}} p(\mathbf{x}) \, d\mathbf{x} \simeq p(r) S_D r^{D-1} \epsilon \tag{45}$$

from which we obtain (1.148). We can find the stationary points of p(r) by differentiation

$$\frac{\mathrm{d}}{\mathrm{d}r}p(r) \propto \left[ (D-1)r^{D-2} + r^{D-1} \left( -\frac{r}{\sigma^2} \right) \right] \exp\left( -\frac{r^2}{2\sigma^2} \right) = 0. \tag{46}$$

Solving for r, and using  $D \gg 1$ , we obtain  $\hat{r} \simeq \sqrt{D}\sigma$ .

Next we note that

$$p(\widehat{r} + \epsilon) \propto (\widehat{r} + \epsilon)^{D-1} \exp\left[-\frac{(\widehat{r} + \epsilon)^2}{2\sigma^2}\right]$$

$$= \exp\left[-\frac{(\widehat{r} + \epsilon)^2}{2\sigma^2} + (D - 1)\ln(\widehat{r} + \epsilon)\right]. \tag{47}$$

We now expand p(r) around the point  $\hat{r}$ . Since this is a stationary point of p(r) we must keep terms up to second order. Making use of the expansion  $\ln(1+x) = x - x^2/2 + O(x^3)$ , together with  $D \gg 1$ , we obtain (1.149).

Finally, from (1.147) we see that the probability density at the origin is given by

$$p(\mathbf{x} = \mathbf{0}) = \frac{1}{(2\pi\sigma^2)^{1/2}}$$

while the density at  $\|\mathbf{x}\| = \hat{r}$  is given from (1.147) by

$$p(\|\mathbf{x}\| = \hat{r}) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{\hat{r}^2}{2\sigma^2}\right) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{D}{2}\right)$$

where we have used  $\hat{r} \simeq \sqrt{D}\sigma$ . Thus the ratio of densities is given by  $\exp(D/2)$ .

**1.21** Since the square root function is monotonic for non-negative numbers, we can take the square root of the relation  $a \le b$  to obtain  $a^{1/2} \le b^{1/2}$ . Then we multiply both sides by the non-negative quantity  $a^{1/2}$  to obtain  $a \le (ab)^{1/2}$ .

The probability of a misclassification is given, from (1.78), by

$$p(\text{mistake}) = \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) \, d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) \, d\mathbf{x}$$
$$= \int_{\mathcal{R}_1} p(\mathcal{C}_2 | \mathbf{x}) p(\mathbf{x}) \, d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathcal{C}_1 | \mathbf{x}) p(\mathbf{x}) \, d\mathbf{x}. \tag{48}$$

Since we have chosen the decision regions to minimize the probability of misclassification we must have  $p(C_2|\mathbf{x}) \leq p(C_1|\mathbf{x})$  in region  $\mathcal{R}_1$ , and  $p(C_1|\mathbf{x}) \leq p(C_2|\mathbf{x})$  in region  $\mathcal{R}_2$ . We now apply the result  $a \leq b \Rightarrow a^{1/2} \leq b^{1/2}$  to give

$$p(\text{mistake}) \leqslant \int_{\mathcal{R}_1} \{ p(\mathcal{C}_1 | \mathbf{x}) p(\mathcal{C}_2 | \mathbf{x}) \}^{1/2} p(\mathbf{x}) \, d\mathbf{x}$$

$$+ \int_{\mathcal{R}_2} \{ p(\mathcal{C}_1 | \mathbf{x}) p(\mathcal{C}_2 | \mathbf{x}) \}^{1/2} p(\mathbf{x}) \, d\mathbf{x}$$

$$= \int \{ p(\mathcal{C}_1 | \mathbf{x}) p(\mathbf{x}) p(\mathcal{C}_2 | \mathbf{x}) p(\mathbf{x}) \}^{1/2} \, d\mathbf{x}$$

$$(49)$$

since the two integrals have the same integrand. The final integral is taken over the whole of the domain of  $\mathbf{x}$ .

- 1.22 Substituting  $L_{kj} = 1 \delta_{kj}$  into (1.81), and using the fact that the posterior probabilities sum to one, we find that, for each  $\mathbf{x}$  we should choose the class j for which  $1 p(\mathcal{C}_j | \mathbf{x})$  is a minimum, which is equivalent to choosing the j for which the posterior probability  $p(\mathcal{C}_j | \mathbf{x})$  is a maximum. This loss matrix assigns a loss of one if the example is misclassified, and a loss of zero if it is correctly classified, and hence minimizing the expected loss will minimize the misclassification rate.
- **1.23** From (1.81) we see that for a general loss matrix and arbitrary class priors, the expected loss is minimized by assigning an input x to class the j which minimizes

$$\sum_{k} L_{kj} p(\mathcal{C}_k | \mathbf{x}) = \frac{1}{p(\mathbf{x})} \sum_{k} L_{kj} p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)$$

and so there is a direct trade-off between the priors  $p(C_k)$  and the loss matrix  $L_{kj}$ .

**1.24** A vector  $\mathbf{x}$  belongs to class  $\mathcal{C}_k$  with probability  $p(\mathcal{C}_k|\mathbf{x})$ . If we decide to assign  $\mathbf{x}$  to class  $\mathcal{C}_j$  we will incur an expected loss of  $\sum_k L_{kj} p(\mathcal{C}_k|\mathbf{x})$ , whereas if we select the reject option we will incur a loss of  $\lambda$ . Thus, if

$$j = \arg\min_{l} \sum_{k} L_{kl} p(\mathcal{C}_k | \mathbf{x})$$
 (50)

then we minimize the expected loss if we take the following action

choose 
$$\begin{cases} \text{class } j, & \text{if } \min_{l} \sum_{k} L_{kl} p(\mathcal{C}_{k} | \mathbf{x}) < \lambda; \\ \text{reject,} & \text{otherwise.} \end{cases}$$
 (51)

For a loss matrix  $L_{kj} = 1 - I_{kj}$  we have  $\sum_k L_{kl} p(\mathcal{C}_k | \mathbf{x}) = 1 - p(\mathcal{C}_l | \mathbf{x})$  and so we reject unless the smallest value of  $1 - p(\mathcal{C}_l | \mathbf{x})$  is less than  $\lambda$ , or equivalently if the largest value of  $p(\mathcal{C}_l | \mathbf{x})$  is less than  $1 - \lambda$ . In the standard reject criterion we reject if the largest posterior probability is less than  $\theta$ . Thus these two criteria for rejection are equivalent provided  $\theta = 1 - \lambda$ .

**1.25** The expected squared loss for a vectorial target variable is given by

$$\mathbb{E}[L] = \iint \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{t}, \mathbf{x}) \, d\mathbf{x} \, d\mathbf{t}.$$

Our goal is to choose y(x) so as to minimize  $\mathbb{E}[L]$ . We can do this formally using the calculus of variations to give

$$\frac{\delta \mathbb{E}[L]}{\delta \mathbf{y}(\mathbf{x})} = \int 2(\mathbf{y}(\mathbf{x}) - \mathbf{t})p(\mathbf{t}, \mathbf{x}) \, d\mathbf{t} = 0.$$

Solving for y(x), and using the sum and product rules of probability, we obtain

$$\mathbf{y}(\mathbf{x}) = \frac{\int \mathbf{t} p(\mathbf{t}, \mathbf{x}) \, d\mathbf{t}}{\int p(\mathbf{t}, \mathbf{x}) \, d\mathbf{t}} = \int \mathbf{t} p(\mathbf{t} | \mathbf{x}) \, d\mathbf{t}$$

which is the conditional average of t conditioned on x. For the case of a scalar target variable we have

$$y(\mathbf{x}) = \int t p(t|\mathbf{x}) \, \mathrm{d}t$$

which is equivalent to (1.89).

**1.26** NOTE: In the 1<sup>st</sup> printing of PRML, there is an error in equation (1.90); the integrand of the second integral should be replaced by  $var[t|\mathbf{x}]p(\mathbf{x})$ .

We start by expanding the square in (1.151), in a similar fashion to the univariate case in the equation preceding (1.90),

$$\begin{split} \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 &= \|\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}\|^2 \\ &= \|\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}]\|^2 + (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^{\mathrm{T}} (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}) \\ &+ (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^{\mathrm{T}} (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}]) + \|\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}\|^2. \end{split}$$

Following the treatment of the univariate case, we now substitute this into (1.151) and perform the integral over t. Again the cross-term vanishes and we are left with

$$\mathbb{E}[L] = \int \|\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}]\|^2 p(\mathbf{x}) \, d\mathbf{x} + \int \text{var}[\mathbf{t}|\mathbf{x}] p(\mathbf{x}) \, d\mathbf{x}$$

from which we see directly that the function  $\mathbf{y}(\mathbf{x})$  that minimizes  $\mathbb{E}[L]$  is given by  $\mathbb{E}[\mathbf{t}|\mathbf{x}]$ .

**1.27** Since we can choose  $y(\mathbf{x})$  independently for each value of  $\mathbf{x}$ , the minimum of the expected  $L_q$  loss can be found by minimizing the integrand given by

$$\int |y(\mathbf{x}) - t|^q p(t|\mathbf{x}) \, \mathrm{d}t \tag{52}$$

for each value of x. Setting the derivative of (52) with respect to y(x) to zero gives the stationarity condition

$$\int q|y(\mathbf{x}) - t|^{q-1} \operatorname{sign}(y(\mathbf{x}) - t)p(t|\mathbf{x}) dt$$

$$= q \int_{-\infty}^{y(\mathbf{x})} |y(\mathbf{x}) - t|^{q-1} p(t|\mathbf{x}) dt - q \int_{y(\mathbf{x})}^{\infty} |y(\mathbf{x}) - t|^{q-1} p(t|\mathbf{x}) dt = 0$$

which can also be obtained directly by setting the functional derivative of (1.91) with respect to  $y(\mathbf{x})$  equal to zero. It follows that  $y(\mathbf{x})$  must satisfy

$$\int_{-\infty}^{y(\mathbf{x})} |y(\mathbf{x}) - t|^{q-1} p(t|\mathbf{x}) \, \mathrm{d}t = \int_{y(\mathbf{x})}^{\infty} |y(\mathbf{x}) - t|^{q-1} p(t|\mathbf{x}) \, \mathrm{d}t. \tag{53}$$

For the case of q = 1 this reduces to

$$\int_{-\infty}^{y(\mathbf{x})} p(t|\mathbf{x}) \, \mathrm{d}t = \int_{y(\mathbf{x})}^{\infty} p(t|\mathbf{x}) \, \mathrm{d}t.$$
 (54)

which says that  $y(\mathbf{x})$  must be the conditional median of t.

For  $q \to 0$  we note that, as a function of t, the quantity  $|y(\mathbf{x}) - t|^q$  is close to 1 everywhere except in a small neighbourhood around  $t = y(\mathbf{x})$  where it falls to zero. The value of (52) will therefore be close to 1, since the density p(t) is normalized, but reduced slightly by the 'notch' close to  $t = y(\mathbf{x})$ . We obtain the biggest reduction in (52) by choosing the location of the notch to coincide with the largest value of p(t), i.e. with the (conditional) mode.

**1.28** From the discussion of the introduction of Section 1.6, we have

$$h(p^2) = h(p) + h(p) = 2 h(p).$$

We then assume that for all  $k \leq K$ ,  $h(p^k) = k h(p)$ . For k = K + 1 we have

$$h(p^{K+1}) = h(p^K p) = h(p^K) + h(p) = K h(p) + h(p) = (K+1) h(p).$$

Moreover,

$$h(p^{n/m}) = n h(p^{1/m}) = \frac{n}{m} m h(p^{1/m}) = \frac{n}{m} h(p^{m/m}) = \frac{n}{m} h(p)$$

and so, by continuity, we have that  $h(p^x) = x h(p)$  for any real number x.

Now consider the positive real numbers p and q and the real number x such that  $p=q^x$ . From the above discussion, we see that

$$\frac{h(p)}{\ln(p)} = \frac{h(q^x)}{\ln(q^x)} = \frac{x h(q)}{x \ln(q)} = \frac{h(q)}{\ln(q)}$$

and hence  $h(p) \propto \ln(p)$ .

**1.29** The entropy of an M-state discrete variable x can be written in the form

$$H(x) = -\sum_{i=1}^{M} p(x_i) \ln p(x_i) = \sum_{i=1}^{M} p(x_i) \ln \frac{1}{p(x_i)}.$$
 (55)

The function ln(x) is concave and so we can apply Jensen's inequality in the form (1.115) but with the inequality reversed, so that

$$H(x) \leqslant \ln\left(\sum_{i=1}^{M} p(x_i) \frac{1}{p(x_i)}\right) = \ln M.$$
 (56)

**1.30 NOTE**: In PRML, there is a minus sign ('-') missing on the l.h.s. of (1.103). From (1.113) we have

$$KL(p||q) = -\int p(x) \ln q(x) dx + \int p(x) \ln p(x) dx.$$
 (57)

Using (1.46) and (1.48)–(1.50), we can rewrite the first integral on the r.h.s. of (57)

$$-\int p(x) \ln q(x) dx = \int \mathcal{N}(x|\mu, \sigma^2) \frac{1}{2} \left( \ln(2\pi s^2) + \frac{(x-m)^2}{s^2} \right) dx$$

$$= \frac{1}{2} \left( \ln(2\pi s^2) + \frac{1}{s^2} \int \mathcal{N}(x|\mu, \sigma^2) (x^2 - 2xm + m^2) dx \right)$$

$$= \frac{1}{2} \left( \ln(2\pi s^2) + \frac{\sigma^2 + \mu^2 - 2\mu m + m^2}{s^2} \right). \tag{58}$$

The second integral on the r.h.s. of (57) we recognize from (1.103) as the negative differential entropy of a Gaussian. Thus, from (57), (58) and (1.110), we have

$$KL(p||q) = \frac{1}{2} \left( \ln(2\pi s^2) + \frac{\sigma^2 + \mu^2 - 2\mu m + m^2}{s^2} - 1 - \ln(2\pi\sigma^2) \right)$$
 次迎关注公众号  $\phi^2$  方  $\phi^2$  方  $\phi^2$  方  $\phi^2$  方  $\phi^2$  方  $\phi^2$  方  $\phi^2$ 

**1.31** We first make use of the relation  $I(\mathbf{x}; \mathbf{y}) = H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x})$  which we obtained in (1.121), and note that the mutual information satisfies  $I(\mathbf{x}; \mathbf{y}) \ge 0$  since it is a form of Kullback-Leibler divergence. Finally we make use of the relation (1.112) to obtain the desired result (1.152).

To show that statistical independence is a sufficient condition for the equality to be satisfied, we substitute  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$  into the definition of the entropy, giving

$$H(\mathbf{x}, \mathbf{y}) = \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}, \mathbf{y}) \, d\mathbf{x} \, d\mathbf{y}$$

$$= \iint p(\mathbf{x}) p(\mathbf{y}) \left\{ \ln p(\mathbf{x}) + \ln p(\mathbf{y}) \right\} \, d\mathbf{x} \, d\mathbf{y}$$

$$= \int p(\mathbf{x}) \ln p(\mathbf{x}) \, d\mathbf{x} + \int p(\mathbf{y}) \ln p(\mathbf{y}) \, d\mathbf{y}$$

$$= H(\mathbf{x}) + H(\mathbf{y}).$$

To show that statistical independence is a necessary condition, we combine the equality condition

$$H(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) + H(\mathbf{y})$$

with the result (1.112) to give

$$H(\mathbf{y}|\mathbf{x}) = H(\mathbf{y}).$$

We now note that the right-hand side is independent of  $\mathbf{x}$  and hence the left-hand side must also be constant with respect to  $\mathbf{x}$ . Using (1.121) it then follows that the mutual information  $I[\mathbf{x}, \mathbf{y}] = 0$ . Finally, using (1.120) we see that the mutual information is a form of KL divergence, and this vanishes only if the two distributions are equal, so that  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$  as required.

**1.32** When we make a change of variables, the probability density is transformed by the Jacobian of the change of variables. Thus we have

$$p(\mathbf{x}) = p(\mathbf{y}) \left| \frac{\partial y_i}{\partial x_j} \right| = p(\mathbf{y}) |\mathbf{A}|$$
 (59)

where  $|\cdot|$  denotes the determinant. Then the entropy of y can be written

$$H(\mathbf{y}) = -\int p(\mathbf{y}) \ln p(\mathbf{y}) \, d\mathbf{y} = -\int p(\mathbf{x}) \ln \left\{ p(\mathbf{x}) |\mathbf{A}|^{-1} \right\} \, d\mathbf{x} = H(\mathbf{x}) + \ln |\mathbf{A}|$$
(60)

as required.

**1.33** The conditional entropy H(y|x) can be written

$$H(y|x) = -\sum_{i} \sum_{j} p(y_i|x_j) p(x_j) \ln p(y_i|x_j)$$
 (61)

which equals 0 by assumption. Since the quantity  $-p(y_i|x_j) \ln p(y_i|x_j)$  is nonnegative each of these terms must vanish for any value  $x_j$  such that  $p(x_j) \neq 0$ . However, the quantity  $p \ln p$  only vanishes for p=0 or p=1. Thus the quantities  $p(y_i|x_j)$  are all either 0 or 1. However, they must also sum to 1, since this is a normalized probability distribution, and so precisely one of the  $p(y_i|x_j)$  is 1, and the rest are 0. Thus, for each value  $x_j$  there is a unique value  $y_i$  with non-zero probability.

**1.34** Obtaining the required functional derivative can be done simply by inspection. However, if a more formal approach is required we can proceed as follows using the techniques set out in Appendix D. Consider first the functional

$$I[p(x)] = \int p(x)f(x) dx.$$

Under a small variation  $p(x) \rightarrow p(x) + \epsilon \eta(x)$  we have

$$I[p(x) + \epsilon \eta(x)] = \int p(x)f(x) dx + \epsilon \int \eta(x)f(x) dx$$

and hence from (D.3) we deduce that the functional derivative is given by

$$\frac{\delta I}{\delta p(x)} = f(x).$$

Similarly, if we define

$$J[p(x)] = \int p(x) \ln p(x) \, \mathrm{d}x$$

then under a small variation  $p(x) \rightarrow p(x) + \epsilon \eta(x)$  we have

$$J[p(x) + \epsilon \eta(x)] = \int p(x) \ln p(x) dx$$
$$+\epsilon \left\{ \int \eta(x) \ln p(x) dx + \int p(x) \frac{1}{p(x)} \eta(x) dx \right\} + O(\epsilon^2)$$

and hence

$$\frac{\delta J}{\delta p(x)} = p(x) + 1.$$

Using these two results we obtain the following result for the functional derivative

$$-\ln p(x) - 1 + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2$$
.

Re-arranging then gives (1.108).

To eliminate the Lagrange multipliers we substitute (1.108) into each of the three constraints (1.105), (1.106) and (1.107) in turn. The solution is most easily obtained

by comparison with the standard form of the Gaussian, and noting that the results

$$\lambda_1 = 1 - \frac{1}{2} \ln \left( 2\pi \sigma^2 \right) \tag{62}$$

$$\lambda_2 = 0 \tag{63}$$

$$\lambda_3 = \frac{1}{2\sigma^2} \tag{64}$$

do indeed satisfy the three constraints.

Note that there is a typographical error in the question, which should read "Use calculus of variations to show that the stationary point of the functional shown just before (1.108) is given by (1.108)".

For the multivariate version of this derivation, see Exercise 2.14.

**1.35** NOTE: In PRML, there is a minus sign ('-') missing on the l.h.s. of (1.103).

Substituting the right hand side of (1.109) in the argument of the logarithm on the right hand side of (1.103), we obtain

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

$$= -\int p(x) \left( -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x-\mu)^2}{2\sigma^2} \right) dx$$

$$= \frac{1}{2} \left( \ln(2\pi\sigma^2) + \frac{1}{\sigma^2} \int p(x) (x-\mu)^2 dx \right)$$

$$= \frac{1}{2} \left( \ln(2\pi\sigma^2) + 1 \right),$$

where in the last step we used (1.107).

**1.36** Consider (1.114) with  $\lambda = 0.5$  and  $b = a + 2\epsilon$  (and hence  $a = b - 2\epsilon$ ),

$$0.5f(a) + 0.5f(b) > f(0.5a + 0.5b)$$

$$= 0.5f(0.5a + 0.5(a + 2\epsilon)) + 0.5f(0.5(b - 2\epsilon) + 0.5b)$$

$$= 0.5f(a + \epsilon) + 0.5f(b - \epsilon)$$

We can rewrite this as

$$f(b) - f(b - \epsilon) > f(a + \epsilon) - f(a)$$

We then divide both sides by  $\epsilon$  and let  $\epsilon \to 0$ , giving

$$f'(b) > f'(a).$$

Since this holds at all points, it follows that  $f''(x) \ge 0$  everywhere.

To show the implication in the other direction, we make use of Taylor's theorem (with the remainder in Lagrange form), according to which there exist an  $x^*$  such that

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x^*)(x - x_0)^2.$$

Since we assume that f''(x) > 0 everywhere, the third term on the r.h.s. will always be positive and therefore

$$f(x) > f(x_0) + f'(x_0)(x - x_0)$$

Now let  $x_0 = \lambda a + (1 - \lambda)b$  and consider setting x = a, which gives

$$f(a) > f(x_0) + f'(x_0)(a - x_0)$$
  
=  $f(x_0) + f'(x_0)((1 - \lambda)(a - b)).$  (65)

Similarly, setting x = b gives

$$f(b) > f(x_0) + f'(x_0)(\lambda(b-a)).$$
 (66)

Multiplying (65) by  $\lambda$  and (66) by  $1-\lambda$  and adding up the results on both sides, we obtain

$$\lambda f(a) + (1 - \lambda)f(b) > f(x_0) = f(\lambda a + (1 - \lambda)b)$$

as required.

**1.37** From (1.104), making use of (1.111), we have

$$H[\mathbf{x}, \mathbf{y}] = -\iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}, \mathbf{y}) \, d\mathbf{x} \, d\mathbf{y}$$

$$= -\iint p(\mathbf{x}, \mathbf{y}) \ln (p(\mathbf{y}|\mathbf{x})p(\mathbf{x})) \, d\mathbf{x} \, d\mathbf{y}$$

$$= -\iint p(\mathbf{x}, \mathbf{y}) (\ln p(\mathbf{y}|\mathbf{x}) + \ln p(\mathbf{x})) \, d\mathbf{x} \, d\mathbf{y}$$

$$= -\iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) \, d\mathbf{x} \, d\mathbf{y} - \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}) \, d\mathbf{x} \, d\mathbf{y}$$

$$= -\iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) \, d\mathbf{x} \, d\mathbf{y} - \int p(\mathbf{x}) \ln p(\mathbf{x}) \, d\mathbf{x} \, d\mathbf{y}$$

$$= -\iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) \, d\mathbf{x} \, d\mathbf{y} - \int p(\mathbf{x}) \ln p(\mathbf{x}) \, d\mathbf{x}$$

$$= H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}].$$

**1.38** From (1.114) we know that the result (1.115) holds for M=1. We now suppose that it holds for some general value M and show that it must therefore hold for M+1. Consider the left hand side of (1.115)

$$f\left(\sum_{i=1}^{M+1} \lambda_i x_i\right) = f\left(\lambda_{M+1} x_{M+1} + \sum_{i=1}^{M} \lambda_i x_i\right)$$

$$= f\left(\lambda_{M+1} x_{M+1} + (1 - \lambda_{M+1}) \sum_{i=1}^{M} \eta_i x_i\right)$$
(67)

where we have defined

$$\eta_i = \frac{\lambda_i}{1 - \lambda_{M+1}}. (69)$$

We now apply (1.114) to give

$$f\left(\sum_{i=1}^{M+1} \lambda_i x_i\right) \leqslant \lambda_{M+1} f(x_{M+1}) + (1 - \lambda_{M+1}) f\left(\sum_{i=1}^{M} \eta_i x_i\right). \tag{70}$$

We now note that the quantities  $\lambda_i$  by definition satisfy

$$\sum_{i=1}^{M+1} \lambda_i = 1 \tag{71}$$

and hence we have

$$\sum_{i=1}^{M} \lambda_i = 1 - \lambda_{M+1} \tag{72}$$

Then using (69) we see that the quantities  $\eta_i$  satisfy the property

$$\sum_{i=1}^{M} \eta_i = \frac{1}{1 - \lambda_{M+1}} \sum_{i=1}^{M} \lambda_i = 1.$$
 (73)

Thus we can apply the result (1.115) at order M and so (70) becomes

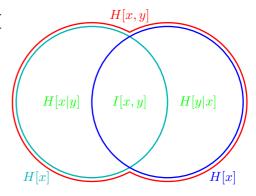
$$f\left(\sum_{i=1}^{M+1} \lambda_i x_i\right) \leqslant \lambda_{M+1} f(x_{M+1}) + (1 - \lambda_{M+1}) \sum_{i=1}^{M} \eta_i f(x_i) = \sum_{i=1}^{M+1} \lambda_i f(x_i)$$
 (74)

where we have made use of (69).

**1.39** From Table 1.3 we obtain the marginal probabilities by summation and the conditional probabilities by normalization, to give

	y			y	
$x \mid 0 \mid 2/3$	0 1			0	1
1 1/3	1/3 2/3	$\boldsymbol{x}$	0	1	1/2
			1	0	1/2
p(x)	p(y)		1	p(x	y)
	y				
	0 1				
	$x \mid 0 \mid 1/2 \mid 1/2 \mid$				
	p(y x)				

Figure 2 Diagram showing the relationship between marginal, conditional and joint entropies and the mutual information.



From these tables, together with the definitions

$$H(x) = -\sum_{i} p(x_i) \ln p(x_i)$$
 (75)

$$H(x|y) = -\sum_{i} \sum_{j} p(x_i, y_j) \ln p(x_i|y_j)$$
 (76)

and similar definitions for H(y) and H(y|x), we obtain the following results

(a) 
$$H(x) = \ln 3 - \frac{2}{3} \ln 2$$

**(b)** 
$$H(y) = \ln 3 - \frac{2}{3} \ln 2$$

(c) 
$$H(y|x) = \frac{2}{3} \ln 2$$

(d) 
$$H(x|y) = \frac{2}{3} \ln 2$$

(e) 
$$H(x,y) = \ln 3$$

(f) 
$$I(x;y) = \ln 3 - \frac{4}{3} \ln 2$$

where we have used (1.121) to evaluate the mutual information. The corresponding diagram is shown in Figure 2.

#### **1.40** The arithmetic and geometric means are defined as

$$ar{x}_{\mathrm{A}} = rac{1}{K} \sum_{k}^{K} x_{k} \quad ext{and} \quad ar{x}_{\mathrm{G}} = \left(\prod_{k}^{K} x_{k}\right)^{1/K},$$

respectively. Taking the logarithm of  $\bar{x}_A$  and  $\bar{x}_G$ , we see that

$$\ln \bar{x}_{\mathrm{A}} = \ln \left( \frac{1}{K} \sum_{k}^{K} x_{k} \right) \quad \text{and} \quad \ln \bar{x}_{\mathrm{G}} = \frac{1}{K} \sum_{k}^{K} \ln x_{k}.$$

By matching f with  $\ln$  and  $\lambda_i$  with 1/K in (1.115), taking into account that the logarithm is concave rather than convex and the inequality therefore goes the other way, we obtain the desired result.

**1.41** From the product rule we have  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$ , and so (1.120) can be written as

$$I(\mathbf{x}; \mathbf{y}) = -\iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} + \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) \, d\mathbf{x} \, d\mathbf{y}$$

$$= -\int p(\mathbf{y}) \ln p(\mathbf{y}) \, d\mathbf{y} + \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) \, d\mathbf{x} \, d\mathbf{y}$$

$$= H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x}). \tag{77}$$

## **Chapter 2** Probability Distributions

**2.1** From the definition (2.2) of the Bernoulli distribution we have

$$\sum_{x \in \{0,1\}} p(x|\mu) = p(x=0|\mu) + p(x=1|\mu)$$

$$= (1-\mu) + \mu = 1$$

$$\sum_{x \in \{0,1\}} xp(x|\mu) = 0.p(x=0|\mu) + 1.p(x=1|\mu) = \mu$$

$$\sum_{x \in \{0,1\}} (x-\mu)^2 p(x|\mu) = \mu^2 p(x=0|\mu) + (1-\mu)^2 p(x=1|\mu)$$

$$= \mu^2 (1-\mu) + (1-\mu)^2 \mu = \mu (1-\mu).$$

The entropy is given by

$$\begin{split} \mathrm{H}[x] &= -\sum_{x \in \{0,1\}} p(x|\mu) \ln p(x|\mu) \\ &= -\sum_{x \in \{0,1\}} \mu^x (1-\mu)^{1-x} \left\{ x \ln \mu + (1-x) \ln (1-\mu) \right\} \\ &= -(1-\mu) \ln (1-\mu) - \mu \ln \mu. \end{split}$$

**2.2** The normalization of (2.261) follows from

$$p(x = +1|\mu) + p(x = -1|\mu) = \left(\frac{1+\mu}{2}\right) + \left(\frac{1-\mu}{2}\right) = 1.$$

The mean is given by

$$\mathbb{E}[x] = \left(\frac{1+\mu}{2}\right) - \left(\frac{1-\mu}{2}\right) = \mu.$$

To evaluate the variance we use

$$\mathbb{E}[x^2] = \left(\frac{1-\mu}{2}\right) + \left(\frac{1+\mu}{2}\right) = 1$$

from which we have

$$var[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = 1 - \mu^2.$$

Finally the entropy is given by

$$H[x] = -\sum_{x=-1}^{x=+1} p(x|\mu) \ln p(x|\mu)$$
$$= -\left(\frac{1-\mu}{2}\right) \ln \left(\frac{1-\mu}{2}\right) - \left(\frac{1+\mu}{2}\right) \ln \left(\frac{1+\mu}{2}\right).$$

**2.3** Using the definition (2.10) we have

$$\binom{N}{n} + \binom{N}{n-1} = \frac{N!}{n!(N-n)!} + \frac{N!}{(n-1)!(N+1-n)!}$$

$$= \frac{(N+1-n)N! + nN!}{n!(N+1-n)!} = \frac{(N+1)!}{n!(N+1-n)!}$$

$$= \binom{N+1}{n}.$$
(78)

To prove the binomial theorem (2.263) we note that the theorem is trivially true for N=0. We now assume that it holds for some general value N and prove its correctness for N+1, which can be done as follows

$$(1+x)^{N+1} = (1+x) \sum_{n=0}^{N} {N \choose n} x^{n}$$

$$= \sum_{n=0}^{N} {N \choose n} x^{n} + \sum_{n=1}^{N+1} {N \choose n-1} x^{n}$$

$$= {N \choose 0} x^{0} + \sum_{n=1}^{N} \left\{ {N \choose n} + {N \choose n-1} \right\} x^{n} + {N \choose N} x^{N+1}$$

$$= {N+1 \choose 0} x^{0} + \sum_{n=1}^{N} {N+1 \choose n} x^{n} + {N+1 \choose N+1} x^{N+1}$$

$$= \sum_{n=0}^{N+1} {N+1 \choose n} x^{n}$$

$$(79)$$

which completes the inductive proof. Finally, using the binomial theorem, the normalization condition (2.264) for the binomial distribution gives

$$\sum_{n=0}^{N} {N \choose n} \mu^n (1-\mu)^{N-n} = (1-\mu)^N \sum_{n=0}^{N} {N \choose n} \left(\frac{\mu}{1-\mu}\right)^n$$
$$= (1-\mu)^N \left(1+\frac{\mu}{1-\mu}\right)^N = 1$$
 (80)

as required.

**2.4** Differentiating (2.264) with respect to  $\mu$  we obtain

$$\sum_{n=1}^{N} \binom{N}{n} \mu^{n} (1-\mu)^{N-n} \left[ \frac{n}{\mu} - \frac{(N-n)}{(1-\mu)} \right] = 0.$$

Multiplying through by  $\mu(1-\mu)$  and re-arranging we obtain (2.11).

If we differentiate (2.264) twice with respect to  $\mu$  we obtain

$$\sum_{n=1}^{N} \binom{N}{n} \mu^n (1-\mu)^{N-n} \left\{ \left[ \frac{n}{\mu} - \frac{(N-n)}{(1-\mu)} \right]^2 - \frac{n}{\mu^2} - \frac{(N-n)}{(1-\mu)^2} \right\} = 0.$$

We now multiply through by  $\mu^2(1-\mu)^2$  and re-arrange, making use of the result (2.11) for the mean of the binomial distribution, to obtain

$$\mathbb{E}[n^2] = N\mu(1-\mu) + N^2\mu^2.$$

Finally, we use (1.40) to obtain the result (2.12) for the variance.

**2.5** Making the change of variable t = y + x in (2.266) we obtain

$$\Gamma(a)\Gamma(b) = \int_0^\infty x^{a-1} \left\{ \int_x^\infty \exp(-t)(t-x)^{b-1} dt \right\} dx.$$
 (81)

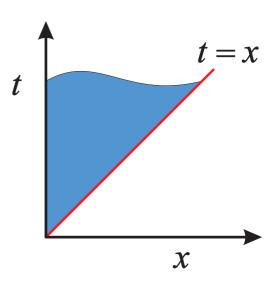
We now exchange the order of integration, taking care over the limits of integration

$$\Gamma(a)\Gamma(b) = \int_0^\infty \int_0^t x^{a-1} \exp(-t)(t-x)^{b-1} dx dt.$$
 (82)

The change in the limits of integration in going from (81) to (82) can be understood by reference to Figure 3. Finally we change variables in the x integral using  $x=t\mu$  to give

$$\Gamma(a)\Gamma(b) = \int_0^\infty \exp(-t)t^{a-1}t^{b-1}t \,dt \int_0^1 \mu^{a-1}(1-\mu)^{b-1} \,d\mu$$
$$= \Gamma(a+b)\int_0^1 \mu^{a-1}(1-\mu)^{b-1} \,d\mu. \tag{83}$$

Figure 3 Plot of the region of integration of (81) in (x, t) space.



**2.6** From (2.13) the mean of the beta distribution is given by

$$\mathbb{E}[\mu] = \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{(a+1)-1} (1-\mu)^{b-1} d\mu.$$

Using the result (2.265), which follows directly from the normalization condition for the Beta distribution, we have

$$\mathbb{E}[\mu] = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+1+b)}{\Gamma(a+1)\Gamma(b)} = \frac{a}{a+b}$$

where we have used the property  $\Gamma(x+1)=x\Gamma(x)$ . We can find the variance in the same way, by first showing that

$$\mathbb{E}[\mu^{2}] = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_{0}^{1} \frac{\Gamma(a+2+b)}{\Gamma(a+2)\Gamma(b)} \mu^{(a+2)-1} (1-\mu)^{b-1} d\mu$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+2+b)}{\Gamma(a+2)\Gamma(b)} = \frac{a}{(a+b)} \frac{a+1}{(a+1+b)}.$$
(84)

Now we use the result (1.40), together with the result (2.15) to derive the result (2.16) for  $var[\mu]$ . Finally, we obtain the result (2.269) for the mode of the beta distribution simply by setting the derivative of the right hand side of (2.13) with respect to  $\mu$  to zero and re-arranging.

**2.7 NOTE**: In PRML, the exercise text contains a typographical error. On the third line, "mean value of x" should be "mean value of  $\mu$ ".

Using the result (2.15) for the mean of a Beta distribution we see that the prior mean is a/(a+b) while the posterior mean is (a+n)/(b+n+n). The maximum likelihood estimate for  $\mu$  is given by the relative frequency n/(n+m) of observations

of x=1. Thus the posterior mean will lie between the prior mean and the maximum likelihood solution provided the following equation is satisfied for  $\lambda$  in the interval (0,1)

$$\lambda \frac{a}{a+b} + (1-\lambda) \frac{n}{n+m} = \frac{a+n}{a+b+n+m}.$$

which represents a convex combination of the prior mean and the maximum likelihood estimator. This is a linear equation for  $\lambda$  which is easily solved by re-arranging terms to give

$$\lambda = \frac{1}{1 + (n+m)/(a+b)}.$$

Since a > 0, b > 0, n > 0, and m > 0, it follows that the term (n + m)/(a + b) lies in the range  $(0, \infty)$  and hence  $\lambda$  must lie in the range (0, 1).

**2.8** To prove the result (2.270) we use the product rule of probability

$$\mathbb{E}_{y} \left[ \mathbb{E}_{x}[x|y] \right] = \int \left\{ \int x p(x|y) \, \mathrm{d}x \right\} p(y) \, \mathrm{d}y$$
$$= \iint x p(x,y) \, \mathrm{d}x \, \mathrm{d}y = \int x p(x) \, \mathrm{d}x = \mathbb{E}_{x}[x]. \tag{85}$$

For the result (2.271) for the conditional variance we make use of the result (1.40), as well as the relation (85), to give

$$\mathbb{E}_{y} \left[ \operatorname{var}_{x}[x|y] \right] + \operatorname{var}_{y} \left[ \mathbb{E}_{x}[x|y] \right] = \mathbb{E}_{y} \left[ \mathbb{E}_{x}[x^{2}|y] - \mathbb{E}_{x}[x|y]^{2} \right]$$

$$+ \mathbb{E}_{y} \left[ \mathbb{E}_{x}[x|y]^{2} \right] - \mathbb{E}_{y} \left[ \mathbb{E}_{x}[x|y] \right]^{2}$$

$$= \mathbb{E}_{x}[x^{2}] - \mathbb{E}_{x}[x]^{2} = \operatorname{var}_{x}[x]$$

where we have made use of  $\mathbb{E}_y\left[\mathbb{E}_x[x^2|y]\right] = \mathbb{E}_x[x^2]$  which can be proved by analogy with (85).

**2.9** When we integrate over  $\mu_{M-1}$  the lower limit of integration is 0, while the upper limit is  $1 - \sum_{j=1}^{M-2} \mu_j$  since the remaining probabilities must sum to one (see Figure 2.4). Thus we have

$$p_{M-1}(\mu_1, \dots, \mu_{M-2}) = \int_0^{1 - \sum_{j=1}^{M-2} \mu_j} p_M(\mu_1, \dots, \mu_{M-1}) \, \mathrm{d}\mu_{M-1}$$

$$= C_M \left[ \prod_{k=1}^{M-2} \mu_k^{\alpha_k - 1} \right] \int_0^{1 - \sum_{j=1}^{M-2} \mu_j} \mu_{M-1}^{\alpha_{M-1} - 1} \left( 1 - \sum_{j=1}^{M-1} \mu_j \right)^{\alpha_M - 1} \, \mathrm{d}\mu_{M-1}.$$

In order to make the limits of integration equal to 0 and 1 we change integration variable from  $\mu_{M-1}$  to t using

$$\mu_{M-1} = t \left( 1 - \sum_{j=1}^{M-2} \mu_j \right)$$

which gives

$$p_{M-1}(\mu_{1}, \dots, \mu_{M-2})$$

$$= C_{M} \left[ \prod_{k=1}^{M-2} \mu_{k}^{\alpha_{k}-1} \right] \left( 1 - \sum_{j=1}^{M-2} \mu_{j} \right)^{\alpha_{M-1} + \alpha_{M} - 1} \int_{0}^{1} t^{\alpha_{M-1} - 1} (1 - t)^{\alpha_{M} - 1} dt$$

$$= C_{M} \left[ \prod_{k=1}^{M-2} \mu_{k}^{\alpha_{k} - 1} \right] \left( 1 - \sum_{j=1}^{M-2} \mu_{j} \right)^{\alpha_{M-1} + \alpha_{M} - 1} \frac{\Gamma(\alpha_{M-1}) \Gamma(\alpha_{M})}{\Gamma(\alpha_{M-1} + \alpha_{M})}$$
(86)

where we have used (2.265). The right hand side of (86) is seen to be a normalized Dirichlet distribution over M-1 variables, with coefficients  $\alpha_1, \ldots, \alpha_{M-2}, \alpha_{M-1} + \alpha_M$ , (note that we have effectively combined the final two categories) and we can identify its normalization coefficient using (2.38). Thus

$$C_{M} = \frac{\Gamma(\alpha_{1} + \ldots + \alpha_{M})}{\Gamma(\alpha_{1}) \ldots \Gamma(\alpha_{M-2}) \Gamma(\alpha_{M-1} + \alpha_{M})} \cdot \frac{\Gamma(\alpha_{M-1} + \alpha_{M})}{\Gamma(\alpha_{M-1}) \Gamma(\alpha_{M})}$$

$$= \frac{\Gamma(\alpha_{1} + \ldots + \alpha_{M})}{\Gamma(\alpha_{1}) \ldots \Gamma(\alpha_{M})}$$
(87)

as required.

**2.10** Using the fact that the Dirichlet distribution (2.38) is normalized we have

$$\int \prod_{k=1}^{M} \mu_k^{\alpha_k - 1} \, \mathrm{d}\boldsymbol{\mu} = \frac{\Gamma(\alpha_1) \cdots \Gamma(\alpha_M)}{\Gamma(\alpha_0)}$$
 (88)

where  $\int d\boldsymbol{\mu}$  denotes the integral over the (M-1)-dimensional simplex defined by  $0 \leqslant \mu_k \leqslant 1$  and  $\sum_k \mu_k = 1$ . Now consider the expectation of  $\mu_j$  which can be written

$$\mathbb{E}[\mu_j] = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_M)} \int \mu_j \prod_{k=1}^M \mu_k^{\alpha_k - 1} \, \mathrm{d}\boldsymbol{\mu}$$
$$= \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_M)} \cdot \frac{\Gamma(\alpha_1) \cdots \Gamma(\alpha_j + 1) \cdots \Gamma(\alpha_M)}{\Gamma(\alpha_0 + 1)} = \frac{\alpha_j}{\alpha_0}$$

where we have made use of (88), noting that the effect of the extra factor of  $\mu_j$  is to increase the coefficient  $\alpha_j$  by 1, and then made use of  $\Gamma(x+1)=x\Gamma(x)$ . By similar reasoning we have

$$\operatorname{var}[\mu_j] = \mathbb{E}[\mu_j^2] - \mathbb{E}[\mu_j]^2 = \frac{\alpha_j(\alpha_j + 1)}{\alpha_0(\alpha_0 + 1)} - \frac{\alpha_j^2}{\alpha_0^2}$$
$$= \frac{\alpha_j(\alpha_0 - \alpha_j)}{\alpha_0^2(\alpha_0 + 1)}.$$

Likewise, for  $j \neq l$  we have

$$cov[\mu_j \mu_l] = \mathbb{E}[\mu_j \mu_l] - \mathbb{E}[\mu_j] \mathbb{E}[\mu_l] = \frac{\alpha_j \alpha_l}{\alpha_0 (\alpha_0 + 1)} - \frac{\alpha_j}{\alpha_0} \frac{\alpha_l}{\alpha_0} 
= -\frac{\alpha_j \alpha_l}{\alpha_0^2 (\alpha_0 + 1)}.$$

**2.11** We first of all write the Dirichlet distribution (2.38) in the form

$$Dir(\boldsymbol{\mu}|\boldsymbol{\alpha}) = K(\boldsymbol{\alpha}) \prod_{k=1}^{M} \mu_k^{\alpha_k - 1}$$

where

$$K(\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_M)}.$$

Next we note the following relation

$$\frac{\partial}{\partial \alpha_j} \prod_{k=1}^M \mu_k^{\alpha_k - 1} = \frac{\partial}{\partial \alpha_j} \prod_{k=1}^M \exp\left((\alpha_k - 1) \ln \mu_k\right)$$

$$= \prod_{k=1}^M \ln \mu_j \exp\left\{(\alpha_k - 1) \ln \mu_k\right\}$$

$$= \ln \mu_j \prod_{k=1}^M \mu_k^{\alpha_k - 1}$$

from which we obtain

$$\mathbb{E}[\ln \mu_j] = K(\boldsymbol{\alpha}) \int_0^1 \cdots \int_0^1 \ln \mu_j \prod_{k=1}^M \mu_k^{\alpha_k - 1} \, \mathrm{d}\mu_1 \dots \, \mathrm{d}\mu_M$$

$$= K(\boldsymbol{\alpha}) \frac{\partial}{\partial \alpha_j} \int_0^1 \cdots \int_0^1 \prod_{k=1}^M \mu_k^{\alpha_k - 1} \, \mathrm{d}\mu_1 \dots \, \mathrm{d}\mu_M$$

$$= K(\boldsymbol{\alpha}) \frac{\partial}{\partial \mu_j} \frac{1}{K(\boldsymbol{\alpha})}$$

$$= -\frac{\partial}{\partial \mu_j} \ln K(\boldsymbol{\alpha}).$$

Finally, using the expression for  $K(\alpha)$ , together with the definition of the digamma function  $\psi(\cdot)$ , we have

$$\mathbb{E}[\ln \mu_i] = \psi(\alpha_i) - \psi(\alpha_0).$$

**2.12** The normalization of the uniform distribution is proved trivially

$$\int_{a}^{b} \frac{1}{b-a} \, \mathrm{d}x = \frac{b-a}{b-a} = 1.$$

For the mean of the distribution we have

$$\mathbb{E}[x] = \int_a^b \frac{1}{b-a} x \, \mathrm{d}x = \left[\frac{x^2}{2(b-a)}\right]_a^b = \frac{b^2 - a^2}{2(b-a)} = \frac{a+b}{2}.$$

The variance can be found by first evaluating

$$\mathbb{E}[x^2] = \int_a^b \frac{1}{b-a} x^2 \, \mathrm{d}x = \left[\frac{x^3}{3(b-a)}\right]_a^b = \frac{b^3 - a^3}{3(b-a)} = \frac{a^2 + ab + b^2}{3}$$

and then using (1.40) to give

$$var[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \frac{a^2 + ab + b^2}{3} - \frac{(a+b)^2}{4} = \frac{(b-a)^2}{12}.$$

**2.13** Note that this solution is the multivariate version of Solution 1.30.

From (1.113) we have

$$\mathrm{KL}(p||q) = -\int p(\mathbf{x}) \ln q(\mathbf{x}) \, \mathrm{d}\mathbf{x} + \int p(\mathbf{x}) \ln p(\mathbf{x}) \, \mathrm{d}\mathbf{x}.$$

Using (2.43), (2.57), (2.59) and (2.62), we can rewrite the first integral on the r.h.s. of () as

$$-\int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x}$$

$$= \int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}^{2}) \frac{1}{2} \left( D \ln(2\pi) + \ln |\mathbf{L}| + (\mathbf{x} - \mathbf{m})^{\mathrm{T}} \mathbf{L}^{-1} (\mathbf{x} - \mathbf{m}) \right) d\mathbf{x}$$

$$= \frac{1}{2} \left( D \ln(2\pi) + \ln |\mathbf{L}| + \mathrm{Tr}[\mathbf{L}^{-1} (\boldsymbol{\mu} \boldsymbol{\mu}^{\mathrm{T}} + \boldsymbol{\Sigma})] - \boldsymbol{\mu} \mathbf{L}^{-1} \mathbf{m} - \mathbf{m}^{\mathrm{T}} \mathbf{L}^{-1} \boldsymbol{\mu} + \mathbf{m}^{\mathrm{T}} \mathbf{L}^{-1} \mathbf{m} \right). \tag{89}$$

The second integral on the r.h.s. of () we recognize from (1.104) as the negative differential entropy of a multivariate Gaussian. Thus, from (), (89) and (B.41), we have

$$KL(p||q) = \frac{1}{2} \left( \ln \frac{|\mathbf{L}|}{|\mathbf{\Sigma}|} + Tr[\mathbf{L}^{-1}(\boldsymbol{\mu}\boldsymbol{\mu}^{T} + \mathbf{\Sigma})] - \boldsymbol{\mu}^{T}\mathbf{L}^{-1}\mathbf{m} - \mathbf{m}^{T}\mathbf{L}^{-1}\boldsymbol{\mu} + \mathbf{m}^{T}\mathbf{L}^{-1}\mathbf{m} - D \right)$$

**2.14** As for the univariate Gaussian considered in Section 1.6, we can make use of Lagrange multipliers to enforce the constraints on the maximum entropy solution. Note that we need a single Lagrange multiplier for the normalization constraint (2.280), a D-dimensional vector  $\mathbf{m}$  of Lagrange multipliers for the D constraints given by (2.281), and a  $D \times D$  matrix  $\mathbf{L}$  of Lagrange multipliers to enforce the  $D^2$  constraints represented by (2.282). Thus we maximize

$$\widetilde{H}[p] = -\int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} + \lambda \left( \int p(\mathbf{x}) d\mathbf{x} - 1 \right)$$

$$+ \mathbf{m}^{\mathrm{T}} \left( \int p(\mathbf{x}) \mathbf{x} d\mathbf{x} - \boldsymbol{\mu} \right)$$

$$+ \mathrm{Tr} \left\{ \mathbf{L} \left( \int p(\mathbf{x}) (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} d\mathbf{x} - \boldsymbol{\Sigma} \right) \right\}.$$
 (90)

By functional differentiation (Appendix D) the maximum of this functional with respect to  $p(\mathbf{x})$  occurs when

$$0 = -1 - \ln p(\mathbf{x}) + \lambda + \mathbf{m}^{\mathrm{T}} \mathbf{x} + \mathrm{Tr} \{ \mathbf{L} (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \}.$$

Solving for  $p(\mathbf{x})$  we obtain

$$p(\mathbf{x}) = \exp\left\{\lambda - 1 + \mathbf{m}^{\mathrm{T}}\mathbf{x} + (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}}\mathbf{L}(\mathbf{x} - \boldsymbol{\mu})\right\}. \tag{91}$$

We now find the values of the Lagrange multipliers by applying the constraints. First we complete the square inside the exponential, which becomes

$$\lambda - 1 + \left(\mathbf{x} - \boldsymbol{\mu} + \frac{1}{2}\mathbf{L}^{-1}\mathbf{m}\right)^{\mathrm{T}}\mathbf{L}\left(\mathbf{x} - \boldsymbol{\mu} + \frac{1}{2}\mathbf{L}^{-1}\mathbf{m}\right) + \boldsymbol{\mu}^{\mathrm{T}}\mathbf{m} - \frac{1}{4}\mathbf{m}^{\mathrm{T}}\mathbf{L}^{-1}\mathbf{m}.$$

We now make the change of variable

$$\mathbf{y} = \mathbf{x} - \boldsymbol{\mu} + \frac{1}{2} \mathbf{L}^{-1} \mathbf{m}.$$

The constraint (2.281) then becomes

$$\int \exp\left\{\lambda - 1 + \mathbf{y}^{\mathrm{T}}\mathbf{L}\mathbf{y} + \boldsymbol{\mu}^{\mathrm{T}}\mathbf{m} - \frac{1}{4}\mathbf{m}^{\mathrm{T}}\mathbf{L}^{-1}\mathbf{m}\right\} \left(\mathbf{y} + \boldsymbol{\mu} - \frac{1}{2}\mathbf{L}^{-1}\mathbf{m}\right) d\mathbf{y} = \boldsymbol{\mu}.$$

In the final parentheses, the term in y vanishes by symmetry, while the term in  $\mu$  simply integrates to  $\mu$  by virtue of the normalization constraint (2.280) which now takes the form

$$\int \exp\left\{\lambda - 1 + \mathbf{y}^{\mathrm{T}}\mathbf{L}\mathbf{y} + \boldsymbol{\mu}^{\mathrm{T}}\mathbf{m} - \frac{1}{4}\mathbf{m}^{\mathrm{T}}\mathbf{L}^{-1}\mathbf{m}\right\} \,\mathrm{d}\mathbf{y} = 1.$$

and hence we have

$$-\frac{1}{2}\mathbf{L}^{-1}\mathbf{m} = \mathbf{0}$$

where again we have made use of the constraint (2.280). Thus  $\mathbf{m} = \mathbf{0}$  and so the density becomes

$$p(\mathbf{x}) = \exp\left\{\lambda - 1 + (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \mathbf{L} (\mathbf{x} - \boldsymbol{\mu})\right\}.$$

Substituting this into the final constraint (2.282), and making the change of variable  $\mathbf{x} - \boldsymbol{\mu} = \mathbf{z}$  we obtain

$$\int \exp\left\{\lambda - 1 + \mathbf{z}^{\mathrm{T}} \mathbf{L} \mathbf{z}\right\} \mathbf{z} \mathbf{z}^{\mathrm{T}} d\mathbf{x} = \mathbf{\Sigma}.$$

Applying an analogous argument to that used to derive (2.64) we obtain  $\mathbf{L} = -\frac{1}{2}\Sigma$ . Finally, the value of  $\lambda$  is simply that value needed to ensure that the Gaussian distribution is correctly normalized, as derived in Section 2.3, and hence is given by

$$\lambda - 1 = \ln \left\{ \frac{1}{(2\pi)^{D/2}} \frac{1}{|\mathbf{\Sigma}|^{1/2}} \right\}.$$

**2.15** From the definitions of the multivariate differential entropy (1.104) and the multivariate Gaussian distribution (2.43), we get

$$H[\mathbf{x}] = -\int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}$$

$$= \int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \frac{1}{2} \left( D \ln(2\pi) + \ln |\boldsymbol{\Sigma}| + (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right) d\mathbf{x}$$

$$= \frac{1}{2} \left( D \ln(2\pi) + \ln |\boldsymbol{\Sigma}| + \mathrm{Tr} \left[ \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \right] \right)$$

$$= \frac{1}{2} \left( D \ln(2\pi) + \ln |\boldsymbol{\Sigma}| + D \right)$$

**2.16** We have  $p(x_1) = \mathcal{N}(x_1|\mu_1, \tau_1^{-1})$  and  $p(x_2) = \mathcal{N}(x_2|\mu_2, \tau_2^{-1})$ . Since  $x = x_1 + x_2$  we also have  $p(x|x_2) = \mathcal{N}(x|\mu_1 + x_2, \tau_1^{-1})$ . We now evaluate the convolution integral given by (2.284) which takes the form

$$p(x) = \left(\frac{\tau_1}{2\pi}\right)^{1/2} \left(\frac{\tau_2}{2\pi}\right)^{1/2} \int_{-\infty}^{\infty} \exp\left\{-\frac{\tau_1}{2}(x - \mu_1 - x_2)^2 - \frac{\tau_2}{2}(x_2 - \mu_2)^2\right\} dx_2.$$
(92)

Since the final result will be a Gaussian distribution for p(x) we need only evaluate its precision, since, from (1.110), the entropy is determined by the variance or equivalently the precision, and is independent of the mean. This allows us to simplify the calculation by ignoring such things as normalization constants.

We begin by considering the terms in the exponent of (92) which depend on  $x_2$  which are given by

$$-\frac{1}{2}x_2^2(\tau_1+\tau_2)+x_2\left\{\tau_1(x-\mu_1)+\tau_2\mu_2\right\}$$

$$= -\frac{1}{2}(\tau_1+\tau_2)\left\{x_2-\frac{\tau_1(x-\mu_1)+\tau_2\mu_2}{\tau_1+\tau_2}\right\}^2+\frac{\left\{\tau_1(x-\mu_1)+\tau_2\mu_2\right\}^2}{2(\tau_1+\tau_2)}$$

where we have completed the square over  $x_2$ . When we integrate out  $x_2$ , the first term on the right hand side will simply give rise to a constant factor independent of x. The second term, when expanded out, will involve a term in  $x^2$ . Since the precision of x is given directly in terms of the coefficient of  $x^2$  in the exponent, it is only such terms that we need to consider. There is one other term in  $x^2$  arising from the original exponent in (92). Combining these we have

$$-\frac{\tau_1}{2}x^2 + \frac{\tau_1^2}{2(\tau_1 + \tau_2)}x^2 = -\frac{1}{2}\frac{\tau_1\tau_2}{\tau_1 + \tau_2}x^2$$

from which we see that x has precision  $\tau_1 \tau_2/(\tau_1 + \tau_2)$ .

We can also obtain this result for the precision directly by appealing to the general result (2.115) for the convolution of two linear-Gaussian distributions.

The entropy of x is then given, from (1.110), by

$$H[x] = \frac{1}{2} \ln \left\{ \frac{2\pi(\tau_1 + \tau_2)}{\tau_1 \tau_2} \right\}.$$

**2.17** We can use an analogous argument to that used in the solution of Exercise 1.14. Consider a general square matrix  $\Lambda$  with elements  $\Lambda_{ij}$ . Then we can always write  $\Lambda = \Lambda^{A} + \Lambda^{S}$  where

$$\Lambda_{ij}^{S} = \frac{\Lambda_{ij} + \Lambda_{ji}}{2}, \qquad \Lambda_{ij}^{A} = \frac{\Lambda_{ij} - \Lambda_{ji}}{2}$$
(93)

and it is easily verified that  $\Lambda^{\rm S}$  is symmetric so that  $\Lambda^{\rm S}_{ij}=\Lambda^{\rm S}_{ji}$ , and  $\Lambda^{\rm A}$  is antisymmetric so that  $\Lambda^{\rm A}_{ij}=-\Lambda^{\rm S}_{ji}$ . The quadratic form in the exponent of a D-dimensional multivariate Gaussian distribution can be written

$$\frac{1}{2} \sum_{i=1}^{D} \sum_{j=1}^{D} (x_i - \mu_i) \Lambda_{ij} (x_j - \mu_j)$$
(94)

where  $\Lambda = \Sigma^{-1}$  is the precision matrix. When we substitute  $\Lambda = \Lambda^{A} + \Lambda^{S}$  into (94) we see that the term involving  $\Lambda^{A}$  vanishes since for every positive term there is an equal and opposite negative term. Thus we can always take  $\Lambda$  to be symmetric.

**2.18** We start by pre-multiplying both sides of (2.45) by  $\mathbf{u}_i^{\dagger}$ , the conjugate transpose of  $\mathbf{u}_i$ . This gives us

$$\mathbf{u}_{i}^{\dagger} \mathbf{\Sigma} \mathbf{u}_{i} = \lambda_{i} \mathbf{u}_{i}^{\dagger} \mathbf{u}_{i}. \tag{95}$$

Next consider the conjugate transpose of (2.45) and post-multiply it by  $\mathbf{u}_i$ , which gives us

$$\mathbf{u}_i^{\dagger} \mathbf{\Sigma}^{\dagger} \mathbf{u}_i = \lambda_i^* \mathbf{u}_i^{\dagger} \mathbf{u}_i. \tag{96}$$

where  $\lambda_i^*$  is the complex conjugate of  $\lambda_i$ . We now subtract (95) from (96) and use the fact the  $\Sigma$  is real and symmetric and hence  $\Sigma = \Sigma^{\dagger}$ , to get

$$0 = (\lambda_i^* - \lambda_i) \mathbf{u}_i^{\dagger} \mathbf{u}_i.$$

Hence  $\lambda_i^* = \lambda_i$  and so  $\lambda_i$  must be real.

Now consider

$$\mathbf{u}_i^{\mathrm{T}} \mathbf{u}_j \lambda_j = \mathbf{u}_i^{\mathrm{T}} \mathbf{\Sigma} \mathbf{u}_j$$

$$= \mathbf{u}_i^{\mathrm{T}} \mathbf{\Sigma}^{\mathrm{T}} \mathbf{u}_j$$

$$= (\mathbf{\Sigma} \mathbf{u}_i)^{\mathrm{T}} \mathbf{u}_j$$

$$= \lambda_i \mathbf{u}_i^{\mathrm{T}} \mathbf{u}_j,$$

where we have used (2.45) and the fact that  $\Sigma$  is symmetric. If we assume that  $0 \neq \lambda_i \neq \lambda_j \neq 0$ , the only solution to this equation is that  $\mathbf{u}_i^{\mathrm{T}} \mathbf{u}_j = 0$ , i.e., that  $\mathbf{u}_i$  and  $\mathbf{u}_j$  are orthogonal.

If  $0 \neq \lambda_i = \lambda_j \neq 0$ , any linear combination of  $\mathbf{u}_i$  and  $\mathbf{u}_j$  will be an eigenvector with eigenvalue  $\lambda = \lambda_i = \lambda_j$ , since, from (2.45),

$$\Sigma(a\mathbf{u}_i + b\mathbf{u}_j) = a\lambda_i\mathbf{u}_i + b\lambda_j\mathbf{u}_j$$
$$= \lambda(a\mathbf{u}_i + b\mathbf{u}_j).$$

Assuming that  $\mathbf{u}_i \neq \mathbf{u}_j$ , we can construct

$$\mathbf{u}_{\alpha} = a\mathbf{u}_i + b\mathbf{u}_j$$
$$\mathbf{u}_{\beta} = c\mathbf{u}_i + d\mathbf{u}_j$$

such that  $\mathbf{u}_{\alpha}$  and  $\mathbf{u}_{\beta}$  are mutually orthogonal and of unit length. Since  $\mathbf{u}_{i}$  and  $\mathbf{u}_{j}$  are orthogonal to  $\mathbf{u}_{k}$  ( $k \neq i, k \neq j$ ), so are  $\mathbf{u}_{\alpha}$  and  $\mathbf{u}_{\beta}$ . Thus,  $\mathbf{u}_{\alpha}$  and  $\mathbf{u}_{\beta}$  satisfy (2.46).

Finally, if  $\lambda_i = 0$ ,  $\Sigma$  must be singular, with  $\mathbf{u}_i$  lying in the nullspace of  $\Sigma$ . In this case,  $\mathbf{u}_i$  will be orthogonal to the eigenvectors projecting onto the rowspace of  $\Sigma$  and we can chose  $\|\mathbf{u}_i\| = 1$ , so that (2.46) is satisfied. If more than one eigenvalue equals zero, we can chose the corresponding eigenvectors arbitrily, as long as they remain in the nullspace of  $\Sigma$ , and so we can chose them to satisfy (2.46).

**2.19** We can write the r.h.s. of (2.48) in matrix form as

$$\sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^{\mathrm{T}} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^{\mathrm{T}} = \mathbf{M},$$

where U is a  $D \times D$  matrix with the eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_D$  as its columns and  $\Lambda$  is a diagonal matrix with the eigenvalues  $\lambda_1, \dots, \lambda_D$  along its diagonal.

Thus we have

$$\mathbf{U}^{\mathrm{T}}\mathbf{M}\mathbf{U} = \mathbf{U}^{\mathrm{T}}\mathbf{U}\mathbf{\Lambda}\mathbf{U}^{\mathrm{T}}\mathbf{U} = \mathbf{\Lambda}.$$

However, from (2.45)–(2.47), we also have that

$$\mathbf{U}^{\mathrm{T}} \mathbf{\Sigma} \mathbf{U} = \mathbf{U}^{\mathrm{T}} \mathbf{\Lambda} \mathbf{U} = \mathbf{U}^{\mathrm{T}} \mathbf{U} \mathbf{\Lambda} = \mathbf{\Lambda},$$

and so  $\mathbf{M} = \mathbf{\Sigma}$  and (2.48) holds.

Moreover, since  $\mathbf{U}$  is orthonormal,  $\mathbf{U}^{-1} = \mathbf{U}^{\mathrm{T}}$  and so

$$\boldsymbol{\Sigma}^{-1} = \left(\mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^{\mathrm{T}}\right)^{-1} = \left(\mathbf{U}^{\mathrm{T}}\right)^{-1}\boldsymbol{\Lambda}^{-1}\mathbf{U}^{-1} = \mathbf{U}\boldsymbol{\Lambda}^{-1}\mathbf{U}^{\mathrm{T}} = \sum_{i=1}^{D}\lambda_{i}\mathbf{u}_{i}\mathbf{u}_{i}^{\mathrm{T}}.$$

**2.20** Since  $\mathbf{u}_1, \dots, \mathbf{u}_D$  constitute a basis for  $\mathbb{R}^D$ , we can write

$$\mathbf{a} = \hat{a}_1 \mathbf{u}_1 + \hat{a}_2 \mathbf{u}_2 + \ldots + \hat{a}_D \mathbf{u}_D,$$

where  $\hat{a}_1, \dots, \hat{a}_D$  are coefficients obtained by projecting  $\mathbf{a}$  on  $\mathbf{u}_1, \dots, \mathbf{u}_D$ . Note that they typically do *not* equal the elements of  $\mathbf{a}$ .

Using this we can write

$$\mathbf{a}^{\mathrm{T}} \mathbf{\Sigma} \mathbf{a} = (\hat{a}_{1} \mathbf{u}_{1}^{\mathrm{T}} + \ldots + \hat{a}_{D} \mathbf{u}_{D}^{\mathrm{T}}) \mathbf{\Sigma} (\hat{a}_{1} \mathbf{u}_{1} + \ldots + \hat{a}_{D} \mathbf{u}_{D})$$

and combining this result with (2.45) we get

$$(\hat{a}_1\mathbf{u}_1^{\mathrm{T}} + \ldots + \hat{a}_D\mathbf{u}_D^{\mathrm{T}})(\hat{a}_1\lambda_1\mathbf{u}_1 + \ldots + \hat{a}_D\lambda_D\mathbf{u}_D).$$

Now, since  $\mathbf{u}_i^{\mathrm{T}} \mathbf{u}_j = 1$  only if i = j, and 0 otherwise, this becomes

$$\hat{a}_1^2\lambda_1+\ldots+\hat{a}_D^2\lambda_D$$

and since  ${\bf a}$  is real, we see that this expression will be strictly positive for any non-zero  ${\bf a}$ , if all eigenvalues are strictly positive. It is also clear that if an eigenvalue,  $\lambda_i$ , is zero or negative, there exist a vector  ${\bf a}$  (e.g.  ${\bf a}={\bf u}_i$ ), for which this expression will be less than or equal to zero. Thus, that a matrix has eigenvectors which are all strictly positive is a sufficient and necessary condition for the matrix to be positive definite.

**2.21** A  $D \times D$  matrix has  $D^2$  elements. If it is symmetric then the elements not on the leading diagonal form pairs of equal value. There are D elements on the diagonal so the number of elements not on the diagonal is  $D^2 - D$  and only half of these are independent giving

$$\frac{D^2 - D}{2}.$$

If we now add back the D elements on the diagonal we get

$$\frac{D^2 - D}{2} + D = \frac{D(D+1)}{2}.$$

**2.22** Consider a matrix M which is symmetric, so that  $\mathbf{M}^{\mathrm{T}}=\mathbf{M}$ . The inverse matrix  $\mathbf{M}^{-1}$  satisfies

$$\mathbf{M}\mathbf{M}^{-1} = \mathbf{I}.$$

Taking the transpose of both sides of this equation, and using the relation (C.1), we obtain

$$\left(\mathbf{M}^{-1}\right)^{\mathrm{T}}\mathbf{M}^{\mathrm{T}}=\mathbf{I}^{\mathrm{T}}=\mathbf{I}$$

since the identity matrix is symmetric. Making use of the symmetry condition for M we then have

$$\left(\mathbf{M}^{-1}
ight)^{\mathrm{T}}\mathbf{M}=\mathbf{I}$$

and hence, from the definition of the matrix inverse,

$$\left(\mathbf{M}^{-1}\right)^{\mathrm{T}} = \mathbf{M}^{-1}$$

and so  $M^{-1}$  is also a symmetric matrix.

**2.23** Recall that the transformation (2.51) diagonalizes the coordinate system and that the quadratic form (2.44), corresponding to the square of the Mahalanobis distance, is then given by (2.50). This corresponds to a shift in the origin of the coordinate system and a rotation so that the hyper-ellipsoidal contours along which the Mahalanobis distance is constant become axis aligned. The volume contained within any one such contour is unchanged by shifts and rotations. We now make the further transformation  $z_i = \lambda_i^{1/2} y_i$  for  $i = 1, \ldots, D$ . The volume within the hyper-ellipsoid then becomes

$$\int \prod_{i=1}^{D} dy_i = \prod_{i=1}^{D} \lambda_i^{1/2} \int \prod_{i=1}^{D} dz_i = |\mathbf{\Sigma}|^{1/2} V_D \Delta^D$$

where we have used the property that the determinant of  $\Sigma$  is given by the product of its eigenvalues, together with the fact that in the z coordinates the volume has become a sphere of radius  $\Delta$  whose volume is  $V_D \Delta^D$ .

**2.24** Multiplying the left hand side of (2.76) by the matrix (2.287) trivially gives the identity matrix. On the right hand side consider the four blocks of the resulting partitioned matrix:

upper left

$$\mathbf{A}\mathbf{M} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} = \mathbf{I}$$

upper right

$$\begin{aligned} -\mathbf{A}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} + \mathbf{B}\mathbf{D}^{-1} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ &= -(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} + \mathbf{B}\mathbf{D}^{-1} \\ &= -\mathbf{B}\mathbf{D}^{-1} + \mathbf{B}\mathbf{D}^{-1} = \mathbf{0} \end{aligned}$$

## 欢迎关途。公员。机器学习与算法之道 CM-DD-1CM-CM-CM-0

lower right

$$-CMBD^{-1} + DD^{-1} + DD^{-1}CMBD^{-1} = DD^{-1} = I.$$

Thus the right hand side also equals the identity matrix.

**2.25** We first of all take the joint distribution  $p(\mathbf{x}_a, \mathbf{x}_b, \mathbf{x}_c)$  and marginalize to obtain the distribution  $p(\mathbf{x}_a, \mathbf{x}_b)$ . Using the results of Section 2.3.2 this is again a Gaussian distribution with mean and covariance given by

$$oldsymbol{\mu} = egin{pmatrix} oldsymbol{\mu}_a \ oldsymbol{\mu}_b \end{pmatrix}, \qquad oldsymbol{\Sigma} = egin{pmatrix} oldsymbol{\Sigma}_{aa} & oldsymbol{\Sigma}_{ab} \ oldsymbol{\Sigma}_{ba} & oldsymbol{\Sigma}_{bb} \end{pmatrix}.$$

From Section 2.3.1 the distribution  $p(\mathbf{x}_a, \mathbf{x}_b)$  is then Gaussian with mean and covariance given by (2.81) and (2.82) respectively.

**2.26** Multiplying the left hand side of (2.289) by (A + BCD) trivially gives the identity matrix I. On the right hand side we obtain

$$\begin{split} &(\mathbf{A} + \mathbf{B}\mathbf{C}\mathbf{D})(\mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D}\mathbf{A}^{-1}) \\ &= \ \mathbf{I} + \mathbf{B}\mathbf{C}\mathbf{D}\mathbf{A}^{-1} - \mathbf{B}(\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D}\mathbf{A}^{-1} \\ &- \mathbf{B}\mathbf{C}\mathbf{D}\mathbf{A}^{-1}\mathbf{B}(\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D}\mathbf{A}^{-1} \\ &= \ \mathbf{I} + \mathbf{B}\mathbf{C}\mathbf{D}\mathbf{A}^{-1} - \mathbf{B}\mathbf{C}(\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B})(\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D}\mathbf{A}^{-1} \\ &= \ \mathbf{I} + \mathbf{B}\mathbf{C}\mathbf{D}\mathbf{A}^{-1} - \mathbf{B}\mathbf{C}\mathbf{D}\mathbf{A}^{-1} = \mathbf{I} \end{split}$$

**2.27** From  $\mathbf{y} = \mathbf{x} + \mathbf{z}$  we have trivially that  $\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{z}]$ . For the covariance we have

$$cov[\mathbf{y}] = \mathbb{E}\left[ (\mathbf{x} - \mathbb{E}[\mathbf{x}] + \mathbf{y} - \mathbb{E}[\mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x}] + \mathbf{y} - \mathbb{E}[\mathbf{y}])^{\mathrm{T}} \right]$$

$$= \mathbb{E}\left[ (\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^{\mathrm{T}} \right] + \mathbb{E}\left[ (\mathbf{y} - \mathbb{E}[\mathbf{y}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])^{\mathrm{T}} \right]$$

$$+ \mathbb{E}\left[ (\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])^{\mathrm{T}} \right] + \mathbb{E}\left[ (\mathbf{y} - \mathbb{E}[\mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^{\mathrm{T}} \right]$$

$$= cov[\mathbf{x}] + cov[\mathbf{z}]$$

where we have used the independence of  $\mathbf{x}$  and  $\mathbf{z}$ , together with  $\mathbb{E}\left[(\mathbf{x} - \mathbb{E}[\mathbf{x}])\right] = \mathbb{E}\left[(\mathbf{z} - \mathbb{E}[\mathbf{z}])\right] = 0$ , to set the third and fourth terms in the expansion to zero. For 1-dimensional variables the covariances become variances and we obtain the result of Exercise 1.10 as a special case.

**2.28** For the marginal distribution  $p(\mathbf{x})$  we see from (2.92) that the mean is given by the upper partition of (2.108) which is simply  $\mu$ . Similarly from (2.93) we see that the covariance is given by the top left partition of (2.105) and is therefore given by  $\Lambda^{-1}$ . Now consider the conditional distribution  $p(\mathbf{y}|\mathbf{x})$ . Applying the result (2.81) for the conditional mean we obtain

$$\boldsymbol{\mu}_{\mathbf{y}|\mathbf{x}} = \mathbf{A}\boldsymbol{\mu} + \mathbf{b} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) = \mathbf{A}\mathbf{x} + \mathbf{b}.$$

Similarly applying the result (2.82) for the covariance of the conditional distribution we have

$$\mathrm{cov}[\mathbf{y}|\mathbf{x}] = \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}} - \mathbf{A}\boldsymbol{\Lambda}^{-1}\boldsymbol{\Lambda}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}} = \mathbf{L}^{-1}$$

as required.

### 2.29 We first define

$$\mathbf{X} = \mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A} \tag{97}$$

and

$$\mathbf{W} = -\mathbf{L}\mathbf{A}$$
, and thus  $\mathbf{W}^{\mathrm{T}} = -\mathbf{A}^{\mathrm{T}}\mathbf{L}^{\mathrm{T}} = -\mathbf{A}^{\mathrm{T}}\mathbf{L}$ , (98)

since L is symmetric. We can use (97) and (98) to re-write (2.104) as

$$\mathbf{R} = \left( egin{array}{cc} \mathbf{X} & \mathbf{W}^{\mathrm{T}} \ \mathbf{W} & \mathbf{L} \end{array} 
ight)$$

and using (2.76) we get

$$\begin{pmatrix} \mathbf{X} & \mathbf{W}^{\mathrm{T}} \\ \mathbf{W} & \mathbf{L} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{W}^{\mathrm{T}}\mathbf{L}^{-1} \\ -\mathbf{L}^{-1}\mathbf{W}\mathbf{M} & \mathbf{L}^{-1} + \mathbf{L}^{-1}\mathbf{W}\mathbf{M}\mathbf{W}^{\mathrm{T}}\mathbf{L}^{-1} \end{pmatrix}$$

where now

$$\mathbf{M} = \left(\mathbf{X} - \mathbf{W}^{\mathrm{T}} \mathbf{L}^{-1} \mathbf{W}\right)^{-1}.$$

Substituting X and W using (97) and (98), respectively, we get

$$\mathbf{M} = \left(\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A} - \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{L}^{-1}\mathbf{L}\mathbf{A}\right)^{-1} = \mathbf{\Lambda}^{-1},$$
$$-\mathbf{M}\mathbf{W}^{\mathrm{T}}\mathbf{L}^{-1} = \mathbf{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{L}^{-1} = \mathbf{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}$$

and

$$\mathbf{L}^{-1} + \mathbf{L}^{-1}\mathbf{W}\mathbf{M}\mathbf{W}^{\mathrm{T}}\mathbf{L}^{-1} = \mathbf{L}^{-1} + \mathbf{L}^{-1}\mathbf{L}\mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{L}^{-1}$$
$$= \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}},$$

as required.

**2.30** Substituting the leftmost expression of (2.105) for  $\mathbb{R}^{-1}$  in (2.107), we get

$$\begin{pmatrix} \mathbf{\Lambda}^{-1} & \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathrm{T}} \\ \mathbf{A}\mathbf{\Lambda}^{-1} & \mathbf{S}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1} \mathbf{A}^{\mathrm{T}} \end{pmatrix} \begin{pmatrix} \mathbf{\Lambda}\boldsymbol{\mu} - \mathbf{A}^{\mathrm{T}} \mathbf{S} \mathbf{b} \\ \mathbf{S} \mathbf{b} \end{pmatrix}$$

$$= \begin{pmatrix} \mathbf{\Lambda}^{-1} \left( \mathbf{\Lambda}\boldsymbol{\mu} - \mathbf{A}^{\mathrm{T}} \mathbf{S} \mathbf{b} \right) + \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathrm{T}} \mathbf{S} \mathbf{b} \\ \mathbf{A}\mathbf{\Lambda}^{-1} \left( \mathbf{\Lambda}\boldsymbol{\mu} - \mathbf{A}^{\mathrm{T}} \mathbf{S} \mathbf{b} \right) + \left( \mathbf{S}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1} \mathbf{A}^{\mathrm{T}} \mathbf{S} \mathbf{b} \right) \end{pmatrix}$$

$$= \begin{pmatrix} \boldsymbol{\mu} - \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathrm{T}} \mathbf{S} \mathbf{b} + \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathrm{T}} \mathbf{S} \mathbf{b} \\ \mathbf{A}\boldsymbol{\mu} - \mathbf{A}\mathbf{\Lambda}^{-1} \mathbf{A}^{\mathrm{T}} \mathbf{S} \mathbf{b} + \mathbf{b} + \mathbf{A}\mathbf{\Lambda}^{-1} \mathbf{A}^{\mathrm{T}} \mathbf{S} \mathbf{b} \end{pmatrix}$$

$$= \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} - \mathbf{b} \end{pmatrix}$$

- 2.31 Since  $\mathbf{y} = \mathbf{x} + \mathbf{z}$  we can write the conditional distribution of  $\mathbf{y}$  given  $\mathbf{x}$  in the form  $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}_{\mathbf{z}} + \mathbf{x}, \boldsymbol{\Sigma}_{\mathbf{z}})$ . This gives a decomposition of the joint distribution of  $\mathbf{x}$  and  $\mathbf{y}$  in the form  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$  where  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{x}})$ . This therefore takes the form of (2.99) and (2.100) in which we can identify  $\boldsymbol{\mu} \to \boldsymbol{\mu}_{\mathbf{x}}$ ,  $\boldsymbol{\Lambda}^{-1} \to \boldsymbol{\Sigma}_{\mathbf{x}}$ ,  $\boldsymbol{A} \to \mathbf{I}$ ,  $\mathbf{b} \to \boldsymbol{\mu}_{\mathbf{z}}$  and  $\mathbf{L}^{-1} \to \boldsymbol{\Sigma}_{\mathbf{z}}$ . We can now obtain the marginal distribution  $p(\mathbf{y})$  by making use of the result (2.115) from which we obtain  $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}_{\mathbf{x}} + \boldsymbol{\mu}_{\mathbf{z}}, \boldsymbol{\Sigma}_{\mathbf{z}} + \boldsymbol{\Sigma}_{\mathbf{x}})$ . Thus both the means and the covariances are additive, in agreement with the results of Exercise 2.27.
- **2.32** The quadratic form in the exponential of the joint distribution is given by

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})^{\mathrm{T}} \mathbf{L}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}). \tag{99}$$

We now extract all of those terms involving x and assemble them into a standard Gaussian quadratic form by completing the square

$$= -\frac{1}{2}\mathbf{x}^{\mathrm{T}}(\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})\mathbf{x} + \mathbf{x}^{\mathrm{T}}\left[\mathbf{\Lambda}\boldsymbol{\mu} + \mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y} - \mathbf{b})\right] + \text{const}$$

$$= -\frac{1}{2}(\mathbf{x} - \mathbf{m})^{\mathrm{T}}(\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})(\mathbf{x} - \mathbf{m})$$

$$+\frac{1}{2}\mathbf{m}^{\mathrm{T}}(\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})\mathbf{m} + \text{const}$$
(100)

where

$$\mathbf{m} = (\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1} \left[ \mathbf{\Lambda} \boldsymbol{\mu} + \mathbf{A}^{\mathrm{T}} \mathbf{L} (\mathbf{y} - \mathbf{b}) \right].$$

We can now perform the integration over x which eliminates the first term in (100). Then we extract the terms in y from the final term in (100) and combine these with the remaining terms from the quadratic form (99) which depend on y to give

$$= -\frac{1}{2}\mathbf{y}^{\mathrm{T}}\left\{\mathbf{L} - \mathbf{L}\mathbf{A}(\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{L}\right\}\mathbf{y}$$

$$+\mathbf{y}^{\mathrm{T}}\left[\left\{\mathbf{L} - \mathbf{L}\mathbf{A}(\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{L}\right\}\mathbf{b}$$

$$+\mathbf{L}\mathbf{A}(\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}\mathbf{\Lambda}\boldsymbol{\mu}\right]. \tag{101}$$

We can identify the precision of the marginal distribution p(y) from the second order term in y. To find the corresponding covariance, we take the inverse of the precision and apply the Woodbury inversion formula (2.289) to give

$$\left\{\mathbf{L} - \mathbf{L}\mathbf{A}(\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{L}\right\}^{-1} = \mathbf{L}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}$$
 (102)

which corresponds to (2.110).

Next we identify the mean  $\nu$  of the marginal distribution. To do this we make use of (102) in (101) and then complete the square to give

$$-\frac{1}{2}(\mathbf{y} - \boldsymbol{\nu})^{\mathrm{T}} \left(\mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}\right)^{-1} (\mathbf{y} - \boldsymbol{\nu}) + \mathrm{const}$$

where

$$\boldsymbol{\nu} = \left(\mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}\right) \left[ (\mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})^{-1}\mathbf{b} + \mathbf{L}\mathbf{A}(\boldsymbol{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}\boldsymbol{\Lambda}\boldsymbol{\mu} \right].$$

Now consider the two terms in the square brackets, the first one involving  ${\bf b}$  and the second involving  ${\boldsymbol \mu}$ . The first of these contribution simply gives  ${\bf b}$ , while the term in  ${\boldsymbol \mu}$  can be written

$$= (\mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})\mathbf{L}\mathbf{A}(\boldsymbol{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}\boldsymbol{\Lambda}\boldsymbol{\mu}$$

$$= \mathbf{A}(\mathbf{I} + \boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})(\mathbf{I} + \boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}\boldsymbol{\Lambda}^{-1}\boldsymbol{\Lambda}\boldsymbol{\mu} = \mathbf{A}\boldsymbol{\mu}$$

where we have used the general result  $(\mathbf{BC})^{-1} = \mathbf{C}^{-1}\mathbf{B}^{-1}$ . Hence we obtain (2.109).

**2.33** To find the conditional distribution  $p(\mathbf{x}|\mathbf{y})$  we start from the quadratic form (99) corresponding to the joint distribution  $p(\mathbf{x}, \mathbf{y})$ . Now, however, we treat  $\mathbf{y}$  as a constant and simply complete the square over  $\mathbf{x}$  to give

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})^{\mathrm{T}} \mathbf{L}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})$$

$$= -\frac{1}{2} \mathbf{x}^{\mathrm{T}} (\boldsymbol{\Lambda} + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A}) \mathbf{x} + \mathbf{x}^{\mathrm{T}} \left\{ \boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{A} \mathbf{L}(\mathbf{y} - \mathbf{b}) \right\} + \text{const}$$

$$= -\frac{1}{2} (\mathbf{x} - \mathbf{m})^{\mathrm{T}} (\boldsymbol{\Lambda} + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A}) (\mathbf{x} - \mathbf{m})$$

where, as in the solution to Exercise 2.32, we have defined

$$\mathbf{m} = (\mathbf{\Lambda} + \mathbf{A}^{\mathrm{T}}\mathbf{L}\mathbf{A})^{-1}\left\{\mathbf{\Lambda}\boldsymbol{\mu} + \mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y} - \mathbf{b})\right\}$$

from which we obtain directly the mean and covariance of the conditional distribution in the form (2.111) and (2.112).

**2.34** Differentiating (2.118) with respect to  $\Sigma$  we obtain two terms:

$$-\frac{N}{2}\frac{\partial}{\partial \mathbf{\Sigma}}\ln|\mathbf{\Sigma}| - \frac{1}{2}\frac{\partial}{\partial \mathbf{\Sigma}}\sum_{n=1}^{N}(\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}}\mathbf{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}).$$

For the first term, we can apply (C.28) directly to get

$$-\frac{N}{2}\frac{\partial}{\partial \mathbf{\Sigma}}\ln|\mathbf{\Sigma}| = -\frac{N}{2}\left(\mathbf{\Sigma}^{-1}\right)^{\mathrm{T}} = -\frac{N}{2}\mathbf{\Sigma}^{-1}.$$

For the second term, we first re-write the sum

$$\sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = N \mathrm{Tr} \left[ \boldsymbol{\Sigma}^{-1} \mathbf{S} \right],$$

where

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}}.$$

Using this together with (C.21), in which  $x = \Sigma_{ij}$  (element (i, j) in  $\Sigma$ ), and properties of the trace we get

$$\frac{\partial}{\partial \Sigma_{ij}} \sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = N \frac{\partial}{\partial \Sigma_{ij}} \mathrm{Tr} \left[ \boldsymbol{\Sigma}^{-1} \mathbf{S} \right] 
= N \mathrm{Tr} \left[ \frac{\partial}{\partial \Sigma_{ij}} \boldsymbol{\Sigma}^{-1} \mathbf{S} \right] 
= -N \mathrm{Tr} \left[ \boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \Sigma_{ij}} \boldsymbol{\Sigma}^{-1} \mathbf{S} \right] 
= -N \mathrm{Tr} \left[ \frac{\partial \boldsymbol{\Sigma}}{\partial \Sigma_{ij}} \boldsymbol{\Sigma}^{-1} \mathbf{S} \boldsymbol{\Sigma}^{-1} \right] 
= -N \left( \boldsymbol{\Sigma}^{-1} \mathbf{S} \boldsymbol{\Sigma}^{-1} \right)_{ij}$$

where we have used (C.26). Note that in the last step we have ignored the fact that  $\Sigma_{ij} = \Sigma_{ji}$ , so that  $\partial \mathbf{\Sigma}/\partial \Sigma_{ij}$  has a 1 in position (i,j) only and 0 everywhere else. Treating this result as valid nevertheless, we get

$$-\frac{1}{2}\frac{\partial}{\partial \mathbf{\Sigma}} \sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = \frac{N}{2} \mathbf{\Sigma}^{-1} \mathbf{S} \mathbf{\Sigma}^{-1}.$$

Combining the derivatives of the two terms and setting the result to zero, we obtain

$$\frac{N}{2} \mathbf{\Sigma}^{-1} = \frac{N}{2} \mathbf{\Sigma}^{-1} \mathbf{S} \mathbf{\Sigma}^{-1}.$$

Re-arrangement then yields

$$\Sigma = S$$

as required.

**2.35** NOTE: In PRML, this exercise contains a typographical error;  $\mathbb{E}\left[\mathbf{x}_{n}\mathbf{x}_{m}\right]$  should be  $\mathbb{E}\left[\mathbf{x}_{n}\mathbf{x}_{m}^{\mathrm{T}}\right]$  on the l.h.s. of (2.291).

The derivation of (2.62) is detailed in the text between (2.59) (page 82) and (2.62) (page 83).

If m=n then, using (2.62) we have  $\mathbb{E}[\mathbf{x}_n\mathbf{x}_n^{\mathrm{T}}]=\boldsymbol{\mu}\boldsymbol{\mu}^{\mathrm{T}}+\boldsymbol{\Sigma}$ , whereas if  $n\neq m$  then the two data points  $\mathbf{x}_n$  and  $\mathbf{x}_m$  are independent and hence  $\mathbb{E}[\mathbf{x}_n\mathbf{x}_m]=\boldsymbol{\mu}\boldsymbol{\mu}^{\mathrm{T}}$  where we have used (2.59). Combining these results we obtain (2.291). From (2.59) and

(2.62) we then have

$$\mathbb{E}\left[\mathbf{\Sigma}_{\mathrm{ML}}\right] = \frac{1}{N} \sum_{n=1}^{N} \mathbb{E}\left[\left(\mathbf{x}_{n} - \frac{1}{N} \sum_{m=1}^{N} \mathbf{x}_{m}\right) \left(\mathbf{x}_{n}^{\mathrm{T}} - \frac{1}{N} \sum_{l=1}^{N} \mathbf{x}_{l}^{\mathrm{T}}\right)\right]$$

$$= \frac{1}{N} \sum_{n=1}^{N} \mathbb{E}\left[\mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - \frac{2}{N} \mathbf{x}_{n} \sum_{m=1}^{N} \mathbf{x}_{m}^{\mathrm{T}} + \frac{1}{N^{2}} \sum_{m=1}^{N} \sum_{l=1}^{N} \mathbf{x}_{m} \mathbf{x}_{l}^{\mathrm{T}}\right]$$

$$= \left\{\mu \boldsymbol{\mu}^{\mathrm{T}} + \boldsymbol{\Sigma} - 2\left(\boldsymbol{\mu} \boldsymbol{\mu}^{\mathrm{T}} + \frac{1}{N} \boldsymbol{\Sigma}\right) + \boldsymbol{\mu} \boldsymbol{\mu}^{\mathrm{T}} + \frac{1}{N} \boldsymbol{\Sigma}\right\}$$

$$= \left(\frac{N-1}{N}\right) \boldsymbol{\Sigma}$$
(103)

as required.

**2.36 NOTE**: In the 1<sup>st</sup> printing of PRML, there are mistakes that affect this solution. The sign in (2.129) is incorrect, and this equation should read

$$\theta^{(N)} = \theta^{(N-1)} - a_{N-1} z(\theta^{(N-1)}).$$

Then, in order to be consistent with the assumption that  $f(\theta) > 0$  for  $\theta > \theta^*$  and  $f(\theta) < 0$  for  $\theta < \theta^*$  in Figure 2.10, we should find the root of the expected *negative* log likelihood. This lead to sign changes in (2.133) and (2.134), but in (2.135), these are cancelled against the change of sign in (2.129), so in effect, (2.135) remains unchanged. Also,  $\mathbf{x}_n$  should be  $x_n$  on the l.h.s. of (2.133). Finally, the labels  $\mu$  and  $\mu_{\mathrm{ML}}$  in Figure 2.11 should be interchanged and there are corresponding changes to the caption (see errata on the PRML web site for details).

Consider the expression for  $\sigma_{(N)}^2$  and separate out the contribution from observation  $x_N$  to give

$$\sigma_{(N)}^{2} = \frac{1}{N} \sum_{n=1}^{N} (x_{n} - \mu)^{2}$$

$$= \frac{1}{N} \sum_{n=1}^{N-1} (x_{n} - \mu)^{2} + \frac{(x_{N} - \mu)^{2}}{N}$$

$$= \frac{N-1}{N} \sigma_{(N-1)}^{2} + \frac{(x_{N} - \mu)^{2}}{N}$$

$$= \sigma_{(N-1)}^{2} - \frac{1}{N} \sigma_{(N-1)}^{2} + \frac{(x_{N} - \mu)^{2}}{N}$$

$$= \sigma_{(N-1)}^{2} + \frac{1}{N} \left\{ (x_{N} - \mu)^{2} - \sigma_{(N-1)}^{2} \right\}. \tag{104}$$

If we substitute the expression for a Gaussian distribution into the result (2.135) for

the Robbins-Monro procedure applied to maximizing likelihood, we obtain

$$\sigma_{(N)}^{2} = \sigma_{(N-1)}^{2} + a_{N-1} \frac{\partial}{\partial \sigma_{(N-1)}^{2}} \left\{ -\frac{1}{2} \ln \sigma_{(N-1)}^{2} - \frac{(x_{N} - \mu)^{2}}{2\sigma_{(N-1)}^{2}} \right\}$$

$$= \sigma_{(N-1)}^{2} + a_{N-1} \left\{ -\frac{1}{2\sigma_{(N-1)}^{2}} + \frac{(x_{N} - \mu)^{2}}{2\sigma_{(N-1)}^{4}} \right\}$$

$$= \sigma_{(N-1)}^{2} + \frac{a_{N-1}}{2\sigma_{(N-1)}^{4}} \left\{ (x_{N} - \mu)^{2} - \sigma_{(N-1)}^{2} \right\}. \tag{105}$$

Comparison of (105) with (104) allows us to identify

$$a_{N-1} = \frac{2\sigma_{(N-1)}^4}{N}.$$

**2.37 NOTE**: In PRML, this exercise requires the additional assumption that we can use the known true mean,  $\mu$ , in (2.122). Furthermore, for the derivation of the Robbins-Monro sequential estimation formula, we assume that the covariance matrix is restricted to be diagonal. Starting from (2.122), we have

$$\Sigma_{\mathrm{ML}}^{(N)} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}}$$

$$= \frac{1}{N} \sum_{n=1}^{N-1} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}}$$

$$+ \frac{1}{N} (\mathbf{x}_{N} - \boldsymbol{\mu}) (\mathbf{x}_{N} - \boldsymbol{\mu})^{\mathrm{T}}$$

$$= \frac{N-1}{N} \Sigma_{\mathrm{ML}}^{(N-1)} + \frac{1}{N} (\mathbf{x}_{N} - \boldsymbol{\mu}) (\mathbf{x}_{N} - \boldsymbol{\mu})^{\mathrm{T}}$$

$$= \Sigma_{\mathrm{ML}}^{(N-1)} + \frac{1}{N} \left( (\mathbf{x}_{N} - \boldsymbol{\mu}) (\mathbf{x}_{N} - \boldsymbol{\mu})^{\mathrm{T}} - \Sigma_{\mathrm{ML}}^{(N-1)} \right). \quad (106)$$

From Solution 2.34, we know that

$$\begin{split} & \frac{\partial}{\partial \boldsymbol{\Sigma}_{\mathrm{ML}}^{(N-1)}} \ln p(\mathbf{x}_{N} | \boldsymbol{\mu}, \boldsymbol{\Sigma}_{\mathrm{ML}}^{(N-1)}) \\ & = \frac{1}{2} \left( \boldsymbol{\Sigma}_{\mathrm{ML}}^{(N-1)} \right)^{-1} \left( (\mathbf{x}_{N} - \boldsymbol{\mu}) (\mathbf{x}_{N} - \boldsymbol{\mu})^{\mathrm{T}} - \boldsymbol{\Sigma}_{\mathrm{ML}}^{(N-1)} \right) \left( \boldsymbol{\Sigma}_{\mathrm{ML}}^{(N-1)} \right)^{-1} \\ & = \frac{1}{2} \left( \boldsymbol{\Sigma}_{\mathrm{ML}}^{(N-1)} \right)^{-2} \left( (\mathbf{x}_{N} - \boldsymbol{\mu}) (\mathbf{x}_{N} - \boldsymbol{\mu})^{\mathrm{T}} - \boldsymbol{\Sigma}_{\mathrm{ML}}^{(N-1)} \right) \end{split}$$

where we have used the assumption that  $\mathbf{\Sigma}_{\mathrm{ML}}^{(N-1)}$ , and hence  $\left(\mathbf{\Sigma}_{\mathrm{ML}}^{(N-1)}\right)^{-1}$ , is diag-

onal. If we substitute this into the multivariate form of (2.135), we get

$$\Sigma_{\mathrm{ML}}^{(N)} = \Sigma_{\mathrm{ML}}^{(N-1)} + \mathbf{A}_{N-1} \frac{1}{2} \left( \Sigma_{\mathrm{ML}}^{(N-1)} \right)^{-2} \left( (\mathbf{x}_N - \boldsymbol{\mu}) \left( \mathbf{x}_N - \boldsymbol{\mu} \right)^{\mathrm{T}} - \Sigma_{\mathrm{ML}}^{(N-1)} \right)$$
(107)

where  $A_{N-1}$  is a matrix of coefficients corresponding to  $a_{N-1}$  in (2.135). By comparing (106) with (107), we see that if we choose

$$\mathbf{A}_{N-1} = \frac{2}{N} \left( \mathbf{\Sigma}_{\mathrm{ML}}^{(N-1)} \right)^{2}.$$

we recover (106). Note that if the covariance matrix was restricted further, to the form  $\sigma^2 \mathbf{I}$ , i.e. a spherical Gaussian, the coefficient in (107) would again become a scalar.

**2.38** The exponent in the posterior distribution of (2.140) takes the form

$$-\frac{1}{2\sigma_0^2}(\mu - \mu_0)^2 - \frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2$$
$$= -\frac{\mu^2}{2} \left(\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}\right) + \mu \left(\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} \sum_{n=1}^N x_n\right) + \text{const.}$$

where 'const.' denotes terms independent of  $\mu$ . Following the discussion of (2.71) we see that the variance of the posterior distribution is given by

$$\frac{1}{\sigma_N^2} = \frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}.$$

Similarly the mean is given by

$$\mu_{N} = \left(\frac{N}{\sigma^{2}} + \frac{1}{\sigma_{0}^{2}}\right)^{-1} \left(\frac{\mu_{0}}{\sigma_{0}^{2}} + \frac{1}{\sigma^{2}} \sum_{n=1}^{N} x_{n}\right)$$

$$= \frac{\sigma^{2}}{N\sigma_{0}^{2} + \sigma^{2}} \mu_{0} + \frac{N\sigma_{0}^{2}}{N\sigma_{0}^{2} + \sigma^{2}} \mu_{ML}.$$
(108)
$$(109)$$

**2.39** From (2.142), we see directly that

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} = \frac{1}{\sigma_0^2} + \frac{N-1}{\sigma^2} + \frac{1}{\sigma^2} = \frac{1}{\sigma_{N-1}^2} + \frac{1}{\sigma^2}.$$
 (110)

We also note for later use, that

$$\frac{1}{\sigma_N^2} = \frac{\sigma^2 + N\sigma_0^2}{\sigma_0^2 \sigma^2} = \frac{\sigma^2 + \sigma_{N-1}^2}{\sigma_{N-1}^2 \sigma^2}$$
(111)

and similarly

$$\frac{1}{\sigma_{N-1}^2} = \frac{\sigma^2 + (N-1)\sigma_0^2}{\sigma_0^2 \sigma^2}.$$
 (112)

Using (2.143), we can rewrite (2.141) as

$$\begin{array}{rcl} \mu_N & = & \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{\sigma_0^2 \sum_{n=1}^N x_n}{N\sigma_0^2 + \sigma^2} \\ & = & \frac{\sigma^2 \mu_0 + \sigma_0^2 \sum_{n=1}^{N-1} x_n}{N\sigma_0^2 + \sigma^2} + \frac{\sigma_0^2 x_N}{N\sigma_0^2 + \sigma^2}. \end{array}$$

Using (2.141), (111) and (112), we can rewrite the first term of this expression as

$$\frac{\sigma_N^2}{\sigma_{N-1}^2} \frac{\sigma^2 \mu_0 + \sigma_0^2 \sum_{n=1}^{N-1} x_n}{(N-1)\sigma_0^2 + \sigma^2} = \frac{\sigma_N^2}{\sigma_{N-1}^2} \mu_{N-1}.$$

Similarly, using (111), the second term can be rewritten as

$$\frac{\sigma_N^2}{\sigma^2} x_N$$

and so

$$\mu_N = \frac{\sigma_N^2}{\sigma_{N-1}^2} \mu_{N-1} + \frac{\sigma_N^2}{\sigma^2} x_N.$$
 (113)

Now consider

$$p(\mu|\mu_{N}, \sigma_{N}^{2}) = p(\mu|\mu_{N-1}, \sigma_{N-1}^{2}) p(x_{N}|\mu, \sigma^{2})$$

$$= \mathcal{N}(\mu|\mu_{N-1}, \sigma_{N-1}^{2}) \mathcal{N}(x_{N}|\mu, \sigma^{2})$$

$$\propto \exp\left\{-\frac{1}{2} \left(\frac{\mu_{N-1}^{2} - 2\mu\mu_{N-1} + \mu^{2}}{\sigma_{N-1}^{2}} + \frac{x_{N}^{2} - 2x_{N}\mu + \mu^{2}}{\sigma^{2}}\right)\right\}$$

$$= \exp\left\{-\frac{1}{2} \left(\frac{\sigma^{2}(\mu_{N-1}^{2} - 2\mu\mu_{N-1} + \mu^{2})}{\sigma_{N-1}^{2}\sigma^{2}} + \frac{\sigma_{N-1}^{2}(x_{N}^{2} - 2x_{N}\mu + \mu^{2})}{\sigma_{N-1}^{2}\sigma^{2}}\right)\right\}$$

$$= \exp\left\{-\frac{1}{2} \frac{(\sigma_{N-1}^{2} + \sigma^{2})\mu^{2} - 2(\sigma^{2}\mu_{N-1} + \sigma_{N-1}^{2}x_{N})\mu}{\sigma_{N-1}^{2}\sigma^{2}}\right\} + C,$$

where C accounts for all the remaining terms that are independent of  $\mu$ . From this, we can directly read off

$$\frac{1}{\sigma_N^2} = \frac{\sigma^2 + \sigma_{N-1}^2}{\sigma_{N-1}^2 \sigma^2} = \frac{1}{\sigma_{N-1}^2} + \frac{1}{\sigma^2}$$

and

$$\mu_{N} = \frac{\sigma^{2}\mu_{N-1} + \sigma_{N-1}^{2}x_{N}}{\sigma_{N-1}^{2} + \sigma^{2}}$$

$$= \frac{\sigma^{2}}{\sigma_{N-1}^{2} + \sigma^{2}}\mu_{N-1} + \frac{\sigma_{N-1}^{2}}{\sigma_{N-1}^{2} + \sigma^{2}}x_{N}$$

$$= \frac{\sigma_{N}^{2}}{\sigma_{N-1}^{2}}\mu_{N-1} + \frac{\sigma_{N}^{2}}{\sigma^{2}}x_{N}$$

and so we have recovered (110) and (113).

The posterior distribution is proportional to the product of the prior and the likelihood function

$$p(\boldsymbol{\mu}|\mathbf{X}) \propto p(\boldsymbol{\mu}) \prod_{n=1}^{N} p(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

Thus the posterior is proportional to an exponential of a quadratic form in  $\mu$  given by

$$-\frac{1}{2}(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^{\mathrm{T}} \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0) - \frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

$$= -\frac{1}{2} \boldsymbol{\mu}^{\mathrm{T}} \left( \boldsymbol{\Sigma}_0^{-1} + N \boldsymbol{\Sigma}^{-1} \right) \boldsymbol{\mu} + \boldsymbol{\mu}^{\mathrm{T}} \left( \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\Sigma}^{-1} \sum_{n=1}^{N} \mathbf{x}_n \right) + \text{const}$$

where 'const.' denotes terms independent of  $\mu$ . Using the discussion following (2.71) we see that the mean and covariance of the posterior distribution are given by

$$\mu_{N} = \left(\Sigma_{0}^{-1} + N\Sigma^{-1}\right)^{-1} \left(\Sigma_{0}^{-1}\mu_{0} + \Sigma^{-1}N\mu_{\mathrm{ML}}\right)$$
(114)  
$$\Sigma_{N}^{-1} = \Sigma_{0}^{-1} + N\Sigma^{-1}$$
(115)

$$\boldsymbol{\Sigma}_{N}^{-1} = \boldsymbol{\Sigma}_{0}^{-1} + N \boldsymbol{\Sigma}^{-1} \tag{115}$$

where  $\mu_{
m ML}$  is the maximum likelihood solution for the mean given by

$$\boldsymbol{\mu}_{\mathrm{ML}} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_{n}.$$

If we consider the integral of the Gamma distribution over  $\tau$  and make the change of variable  $b\tau = u$  we have

$$\int_0^\infty \operatorname{Gam}(\tau|a,b) \, d\tau = \frac{1}{\Gamma(a)} \int_0^\infty b^a \tau^{a-1} \exp(-b\tau) \, d\tau$$
$$= \frac{1}{\Gamma(a)} \int_0^\infty b^a u^{a-1} \exp(-u) b^{1-a} b^{-1} \, du$$

# 欢迎关注公众号@机器学习与算法之道

where we have used the definition (1.141) of the Gamma function

**2.42** We can use the same change of variable as in the previous exercise to evaluate the mean of the Gamma distribution

$$\mathbb{E}[\tau] = \frac{1}{\Gamma(a)} \int_0^\infty b^a \tau^{a-1} \tau \exp(-b\tau) d\tau$$
$$= \frac{1}{\Gamma(a)} \int_0^\infty b^a u^a \exp(-u) b^{-a} b^{-1} du$$
$$= \frac{\Gamma(a+1)}{b\Gamma(a)} = \frac{a}{b}$$

where we have used the recurrence relation  $\Gamma(a+1)=a\Gamma(a)$  for the Gamma function. Similarly we can find the variance by first evaluating

$$\mathbb{E}[\tau^{2}] = \frac{1}{\Gamma(a)} \int_{0}^{\infty} b^{a} \tau^{a-1} \tau^{2} \exp(-b\tau) d\tau$$

$$= \frac{1}{\Gamma(a)} \int_{0}^{\infty} b^{a} u^{a+1} \exp(-u) b^{-a-1} b^{-1} du$$

$$= \frac{\Gamma(a+2)}{b^{2} \Gamma(a)} = \frac{(a+1)\Gamma(a+1)}{b^{2} \Gamma(a)} = \frac{a(a+1)}{b^{2}}$$

and then using

$$\operatorname{var}[\tau] = \mathbb{E}[\tau^2] - \mathbb{E}[\tau]^2 = \frac{a(a+1)}{b^2} - \frac{a^2}{b^2} = \frac{a}{b^2}.$$

Finally, the mode of the Gamma distribution is obtained simply by differentiation

$$\frac{\mathrm{d}}{\mathrm{d}\tau} \left\{ \tau^{a-1} \exp(-b\tau) \right\} = \left[ \frac{a-1}{\tau} - b \right] \tau^{a-1} \exp(-b\tau) = 0$$

from which we obtain

$$mode[\tau] = \frac{a-1}{b}.$$

Notice that the mode only exists if  $a \ge 1$ , since  $\tau$  must be a non-negative quantity. This is also apparent in the plot of Figure 2.13.

**2.43** To prove the normalization of the distribution (2.293) consider the integral

$$I = \int_{-\infty}^{\infty} \exp\left(-\frac{|x|^q}{2\sigma^2}\right) dx = 2\int_{0}^{\infty} \exp\left(-\frac{x^q}{2\sigma^2}\right) dx$$

and make the change of variable

$$u = \frac{x^q}{2\sigma^2}$$

Using the definition (1.141) of the Gamma function, this gives

$$I = 2 \int_0^\infty \frac{2\sigma^2}{q} (2\sigma^2 u)^{(1-q)/q} \exp(-u) du = \frac{2(2\sigma^2)^{1/q} \Gamma(1/q)}{q}$$

from which the normalization of (2.293) follows.

For the given noise distribution, the conditional distribution of the target variable given the input variable is

$$p(t|\mathbf{x}, \mathbf{w}, \sigma^2) = \frac{q}{2(2\sigma^2)^{1/q}\Gamma(1/q)} \exp\left(-\frac{|t - y(\mathbf{x}, \mathbf{w})|^q}{2\sigma^2}\right).$$

The likelihood function is obtained by taking products of factors of this form, over all pairs  $\{x_n, t_n\}$ . Taking the logarithm, and discarding additive constants, we obtain the desired result.

#### **2.44** From Bayes' theorem we have

$$p(\mu, \lambda | \mathbf{X}) \propto p(\mathbf{X} | \mu, \lambda) p(\mu, \lambda),$$

where the factors on the r.h.s. are given by (2.152) and (2.154), respectively. Writing this out in full, we get

$$p(\mu, \lambda) \propto \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^N \exp\left\{\lambda\mu \sum_{n=1}^N x_n - \frac{\lambda}{2} \sum_{n=1}^N x_n^2\right\}$$
$$(\beta\lambda)^{1/2} \exp\left[-\frac{\beta\lambda}{2} \left(\mu^2 - 2\mu\mu_0 + \mu_0^2\right)\right] \lambda^{a-1} \exp\left(-b\lambda\right),$$

where we have used the defintions of the Gaussian and Gamma distributions and we have ommitted terms independent of  $\mu$  and  $\lambda$ . We can rearrange this to obtain

$$\lambda^{N/2} \lambda^{a-1} \exp \left\{ -\left( b + \frac{1}{2} \sum_{n=1}^{N} x_n^2 + \frac{\beta}{2} \mu_0^2 \right) \lambda \right\}$$
$$(\lambda(N+\beta))^{1/2} \exp \left[ -\frac{\lambda(N+\beta)}{2} \left( \mu^2 - \frac{2}{N+\beta} \left\{ \beta \mu_0 + \sum_{n=1}^{N} x_n \right\} \mu \right) \right]$$

and by completing the square in the argument of the second exponential,

$$\lambda^{N/2} \lambda^{a-1} \exp \left\{ -\left( b + \frac{1}{2} \sum_{n=1}^{N} x_n^2 + \frac{\beta}{2} \mu_0^2 - \frac{\left(\beta \mu_0 + \sum_{n=1}^{N} x_n\right)^2}{2(N+\beta)} \right) \lambda \right\}$$

$$\left(\lambda (N+\beta)\right)^{1/2} \exp \left[ -\frac{\lambda (N+\beta)}{2} \left( \mu - \frac{\beta \mu_0 + \sum_{n=1}^{N} x_n}{N+\beta} \right) \right]$$

we arrive at an (unnormalised) Gaussian-Gamma distribution,

$$\mathcal{N}\left(\mu|\mu_N, ((N+\beta)\lambda)^{-1}\right) \operatorname{Gam}\left(\lambda|a_N, b_N\right),$$

with parameters

$$\mu_{N} = \frac{\beta \mu_{0} + \sum_{n=1}^{N} x_{n}}{N + \beta}$$

$$a_{N} = a + \frac{N}{2}$$

$$b_{N} = b + \frac{1}{2} \sum_{n=1}^{N} x_{n}^{2} + \frac{\beta}{2} \mu_{0}^{2} - \frac{N + \beta}{2} \mu_{N}^{2}.$$

**2.45** We do this, as in the univariate case, by considering the likelihood function of  $\Lambda$  for a given data set  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ :

$$\prod_{n=1}^{N} \mathcal{N}(\mathbf{x}_{n} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \propto |\boldsymbol{\Lambda}|^{N/2} \exp \left( -\frac{1}{2} \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Lambda} (\mathbf{x}_{n} - \boldsymbol{\mu}) \right) 
= |\boldsymbol{\Lambda}|^{N/2} \exp \left( -\frac{1}{2} \mathrm{Tr} \left[ \boldsymbol{\Lambda} \mathbf{S} \right] \right),$$

where  $S = \sum_{n} (\mathbf{x}_{n} - \boldsymbol{\mu})(\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}}$ . By simply comparing with (2.155), we see that the functional dependence on  $\boldsymbol{\Lambda}$  is indeed the same and thus a product of this likelihood and a Wishart prior will result in a Wishart posterior.

**2.46** From (2.158), we have

$$\begin{split} \int_0^\infty \frac{b^a e^{(-b\tau)} \tau^{a-1}}{\Gamma(a)} \left(\frac{\tau}{2\pi}\right)^{1/2} \exp\left\{-\frac{\tau}{2} (x-\mu)^2\right\} \, \mathrm{d}\tau \\ &= \frac{b^a}{\Gamma(a)} \left(\frac{1}{2\pi}\right)^{1/2} \int_0^\infty \tau^{a-1/2} \exp\left\{-\tau \left(b + \frac{(x-\mu)^2}{2}\right)\right\} \, \mathrm{d}\tau. \end{split}$$

We now make the proposed change of variable  $z=\tau\Delta$ , where  $\Delta=b+(x-\mu)^2/2$ , yielding

$$\frac{b^{a}}{\Gamma(a)} \left(\frac{1}{2\pi}\right)^{1/2} \Delta^{-a-1/2} \int_{0}^{\infty} z^{a-1/2} \exp(-z) dz$$

$$= \frac{b^{a}}{\Gamma(a)} \left(\frac{1}{2\pi}\right)^{1/2} \Delta^{-a-1/2} \Gamma(a+1/2)$$

where we have used the definition of the Gamma function (1.141). Finally, we substitute  $b + (x - \mu)^2/2$  for  $\Delta$ ,  $\nu/2$  for a and  $\nu/2\lambda$  for b:

$$\begin{split} &\frac{\Gamma(-a+1/2)}{\Gamma(a)} \, b^a \left(\frac{1}{2\pi}\right)^{1/2} \Delta^{a-1/2} \\ &= \frac{\Gamma\left((\nu+1)/2\right)}{\Gamma(\nu/2)} \left(\frac{\nu}{2\lambda}\right)^{\nu/2} \left(\frac{1}{2\pi}\right)^{1/2} \left(\frac{\nu}{2\lambda} + \frac{(x-\mu)^2}{2}\right)^{-(\nu+1)/2} \\ &= \frac{\Gamma\left((\nu+1)/2\right)}{\Gamma(\nu/2)} \left(\frac{\nu}{2\lambda}\right)^{\nu/2} \left(\frac{1}{2\pi}\right)^{1/2} \left(\frac{\nu}{2\lambda}\right)^{-(\nu+1)/2} \left(1 + \frac{\lambda(x-\mu)^2}{\nu}\right)^{-(\nu+1)/2} \\ &= \frac{\Gamma\left((\nu+1)/2\right)}{\Gamma(\nu/2)} \left(\frac{\lambda}{\nu\pi}\right)^{1/2} \left(1 + \frac{\lambda(x-\mu)^2}{\nu}\right)^{-(\nu+1)/2} \end{split}$$

**2.47** Ignoring the normalization constant, we write (2.159) as

$$\operatorname{St}(x|\mu,\lambda,\nu) \propto \left[1 + \frac{\lambda(x-\mu)^2}{\nu}\right]^{-(\nu-1)/2}$$

$$= \exp\left(-\frac{\nu-1}{2}\ln\left[1 + \frac{\lambda(x-\mu)^2}{\nu}\right]\right). \tag{116}$$

For large  $\nu$ , we make use of the Taylor expansion for the logarithm in the form

$$ln(1+\epsilon) = \epsilon + O(\epsilon^2)$$
(117)

to re-write (116) as

$$\exp\left(-\frac{\nu-1}{2}\ln\left[1+\frac{\lambda(x-\mu)^2}{\nu}\right]\right)$$

$$= \exp\left(-\frac{\nu-1}{2}\left[\frac{\lambda(x-\mu)^2}{\nu}+O(\nu^{-2})\right]\right)$$

$$= \exp\left(-\frac{\lambda(x-\mu)^2}{2}+O(\nu^{-1})\right).$$

We see that in the limit  $\nu \to \infty$  this becomes, up to an overall constant, the same as a Gaussian distribution with mean  $\mu$  and precision  $\lambda$ . Since the Student distribution is normalized to unity for all values of  $\nu$  it follows that it must remain normalized in this limit. The normalization coefficient is given by the standard expression (2.42) for a univariate Gaussian.

**2.48** Substituting expressions for the Gaussian and Gamma distributions into (2.161), we have

$$\begin{aligned} \operatorname{St}(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Lambda},\nu) &=& \int_0^\infty \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu},(\eta\boldsymbol{\Lambda})^{-1}\right) \operatorname{Gam}(\eta|\nu/2,\nu/2) \,\mathrm{d}\eta \\ &=& \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} \frac{|\boldsymbol{\Lambda}|^{1/2}}{(2\pi)^{D/2}} \int_0^\infty \eta^{D/2} \eta^{\nu/2-1} e^{-\nu\eta/2} e^{-\eta\Delta^2/2} \,\mathrm{d}\eta. \end{aligned}$$

Now we make the change of variable

$$\tau = \eta \left[ \frac{\nu}{2} + \frac{1}{2} \Delta^2 \right]^{-1}$$

which gives

$$\operatorname{St}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) = \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} \frac{|\boldsymbol{\Lambda}|^{1/2}}{(2\pi)^{D/2}} \left[ \frac{\nu}{2} + \frac{1}{2} \Delta^2 \right]^{-D/2 - \nu/2}$$
$$\int_0^\infty \tau^{D/2 + \nu/2 - 1} e^{-\tau} d\tau$$
$$= \frac{\Gamma(\nu/2 + d/2)}{\Gamma(\nu/2)} \frac{|\boldsymbol{\Lambda}|^{1/2}}{(\pi \nu)^{D/2}} \left[ 1 + \frac{\Delta^2}{\nu} \right]^{-D/2 - \nu/2}$$

as required.

The correct normalization of the multivariate Student's t-distribution follows directly from the fact that the Gaussian and Gamma distributions are normalized. From (2.161) we have

$$\begin{split} \int \mathrm{St} \left( \mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\nu} \right) \, \mathrm{d} \mathbf{x} &= \int \!\! \int \mathcal{N} \left( \mathbf{x} | \boldsymbol{\mu}, (\eta \boldsymbol{\Lambda})^{-1} \right) \mathrm{Gam} \left( \eta | \boldsymbol{\nu} / 2, \boldsymbol{\nu} / 2 \right) \, \mathrm{d} \boldsymbol{\eta} \, \mathrm{d} \mathbf{x} \\ &= \int \!\! \int \mathcal{N} \left( \mathbf{x} | \boldsymbol{\mu}, (\eta \boldsymbol{\Lambda})^{-1} \right) \, \mathrm{d} \mathbf{x} \, \mathrm{Gam} \left( \eta | \boldsymbol{\nu} / 2, \boldsymbol{\nu} / 2 \right) \, \mathrm{d} \boldsymbol{\eta} \\ &= \int \mathrm{Gam} \left( \eta | \boldsymbol{\nu} / 2, \boldsymbol{\nu} / 2 \right) \, \mathrm{d} \boldsymbol{\eta} = 1. \end{split}$$

**2.49** If we make the change of variable  $z = x - \mu$ , we can write

$$\mathbb{E}[\mathbf{x}] = \int \mathrm{St}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\nu}) \mathbf{x} \, \mathrm{d}\mathbf{x} = \int \mathrm{St}(\mathbf{z}|\mathbf{0}, \boldsymbol{\Lambda}, \boldsymbol{\nu}) (\mathbf{z} + \boldsymbol{\mu}) \, \mathrm{d}\mathbf{z}.$$

In the factor  $(\mathbf{z} + \boldsymbol{\mu})$  the first term vanishes as a consequence of the fact that the zero-mean Student distribution is an even function of  $\mathbf{z}$  that is  $\mathrm{St}(-\mathbf{z}|\mathbf{0}, \boldsymbol{\Lambda}, \nu) = \mathrm{St}(-\mathbf{z}|\mathbf{0}, \boldsymbol{\Lambda}, \nu)$ . This leaves the second term, which equals  $\boldsymbol{\mu}$  since the Student distribution is normalized.

The covariance of the multivariate Student can be re-expressed by using the expression for the multivariate Student distribution as a convolution of a Gaussian with a Gamma distribution given by (2.161) which gives

$$\operatorname{cov}[\mathbf{x}] = \int \operatorname{St}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu)(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} d\mathbf{x}$$

$$= \int_{0}^{\infty} \int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \eta \boldsymbol{\Lambda})(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} d\mathbf{x} \operatorname{Gam}(\eta|\nu/2, \nu/2) d\eta$$

$$= \int_{0}^{\infty} \eta^{-1} \boldsymbol{\Lambda}^{-1} \operatorname{Gam}(\eta|\nu/2, \nu/2) d\eta$$

where we have used the standard result for the covariance of a multivariate Gaussian. We now substitute for the Gamma distribution using (2.146) to give

$$\begin{aligned} \operatorname{cov}[\mathbf{x}] &= \frac{1}{\Gamma(\nu/2)} \left(\frac{\nu}{2}\right)^{\nu/2} \int_0^\infty e^{-\nu\eta/2} \eta^{\nu/2-2} \, \mathrm{d}\eta \mathbf{\Lambda}^{-1} \\ &= \frac{\nu}{2} \frac{\Gamma(\nu/2-2)}{\Gamma(\nu/2)} \mathbf{\Lambda}^{-1} \\ &= \frac{\nu}{\nu-2} \mathbf{\Lambda}^{-1} \end{aligned}$$

where we have used the integral representation for the Gamma function, together with the standard result  $\Gamma(1+x)=x\Gamma(x)$ .

The mode of the Student distribution is obtained by differentiation

$$\nabla_{\mathbf{x}} \mathrm{St}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) = \frac{\Gamma(\nu/2 + D/2)}{\Gamma(\nu/2)} \frac{|\boldsymbol{\Lambda}|^{1/2}}{(\pi\nu)^{D/2}} \left[ 1 + \frac{\Delta^2}{\nu} \right]^{-D/2 - \nu/2 - 1} \frac{1}{\nu} \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}).$$

Provided  $\Lambda$  is non-singular we therefore obtain

$$mode[\mathbf{x}] = \boldsymbol{\mu}.$$

**2.50** Just like in univariate case (Exercise 2.47), we ignore the normalization coefficient, which leaves us with

$$\left[1 + \frac{\Delta^2}{\nu}\right]^{-\nu/2 - D/2} = \exp\left\{-\left(\frac{\nu}{2} + \frac{D}{2}\right) \ln\left[1 + \frac{\Delta^2}{\nu}\right]\right\}$$

where  $\Delta^2$  is the squared Mahalanobis distance given by

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}).$$

Again we make use of (117) to give

$$\exp\left\{-\left(\frac{\nu}{2} + \frac{D}{2}\right)\ln\left[1 + \frac{\Delta^2}{\nu}\right]\right\} = \exp\left\{-\frac{\Delta^2}{2} + O(1/\nu)\right\}.$$

As in the univariate case, in the limit  $\nu \to \infty$  this becomes, up to an overall constant, the same as a Gaussian distribution, here with mean  $\mu$  and precision  $\Lambda$ ; the univariate normalization argument also applies in the multivariate case.

**2.51** Using the relation (2.296) we have

$$1 = \exp(iA) \exp(-iA) = (\cos A + i\sin A)(\cos A - i\sin A) = \cos^2 A + \sin^2 A.$$

Similarly, we have

$$\cos(A - B) = \Re \exp\{i(A - B)\}$$

$$= \Re \exp(iA) \exp(-iB)$$

$$= \Re(\cos A + i \sin A)(\cos B - i \sin B)$$

$$= \cos A \cos B + \sin A \sin B.$$

Finally

$$\sin(A - B) = \Im \exp\{i(A - B)\}$$

$$= \Im \exp(iA) \exp(-iB)$$

$$= \Im(\cos A + i \sin A)(\cos B - i \sin B)$$

$$= \sin A \cos B - \cos A \sin B.$$

**2.52** Expressed in terms of  $\xi$  the von Mises distribution becomes

$$p(\xi) \propto \exp\left\{m\cos(m^{-1/2}\xi)\right\}.$$

For large m we have  $\cos(m^{-1/2}\xi) = 1 - m^{-1}\xi^2/2 + O(m^{-2})$  and so

$$p(\xi) \propto \exp\left\{-\xi^2/2\right\}$$

and hence  $p(\theta) \propto \exp\{-m(\theta - \theta_0)^2/2\}$ .

**2.53** Using (2.183), we can write (2.182) as

$$\sum_{n=1}^{N} (\cos \theta_0 \sin \theta_n - \cos \theta_n \sin \theta_0) = \cos \theta_0 \sum_{n=1}^{N} \sin \theta_n - \sin \sum_{n=1}^{N} \cos \theta_n = 0.$$

Rearranging this, we get

$$\frac{\sum_{n} \sin \theta_{n}}{\sum_{n} \cos \theta_{n}} = \frac{\sin \theta_{0}}{\cos \theta_{0}} = \tan \theta_{0},$$

which we can solve w.r.t.  $\theta_0$  to obtain (2.184).

**2.54** Differentiating the von Mises distribution (2.179) we have

$$p'(\theta) = -\frac{1}{2\pi I_0(m)} \exp\left\{m\cos(\theta - \theta_0)\right\} \sin(\theta - \theta_0)$$

which vanishes when  $\theta=\theta_0$  or when  $\theta=\theta_0+\pi\,(\mathrm{mod}2\pi)$ . Differentiating again we have

$$p''(\theta) = -\frac{1}{2\pi I_0(m)} \exp\{m\cos(\theta - \theta_0)\} \left[\sin^2(\theta - \theta_0) + \cos(\theta - \theta_0)\right].$$

Since  $I_0(m) > 0$  we see that  $p''(\theta) < 0$  when  $\theta = \theta_0$ , which therefore represents a maximum of the density, while  $p''(\theta) > 0$  when  $\theta = \theta_0 + \pi \pmod{2\pi}$ , which is therefore a minimum.

**2.55 NOTE**: In the 1<sup>st</sup> printing of PRML, equation (2.187), which will be the starting point for this solution, contains a typo. The "–" on the r.h.s. should be a "+", as is easily seen from (2.178) and (2.185).

From (2.169) and (2.184), we see that  $\bar{\theta} = \theta_0^{\rm ML}$ . Using this together with (2.168) and (2.177), we can rewrite (2.187) as follows:

$$A(m_{\text{ML}}) = \left(\frac{1}{N} \sum_{n=1}^{N} \cos \theta_{n}\right) \cos \theta_{0}^{\text{ML}} + \left(\frac{1}{N} \sum_{n=1}^{N} \sin \theta_{n}\right) \sin \theta_{0}^{\text{ML}}$$

$$= \bar{r} \cos \bar{\theta} \cos \theta_{0}^{\text{ML}} + \bar{r} \sin \bar{\theta} \sin \theta_{0}^{\text{ML}}$$

$$= \bar{r} \left(\cos^{2} \theta_{0}^{\text{ML}} + \sin^{2} \theta_{0}^{\text{ML}}\right)$$

$$= \bar{r}.$$

**2.56** We can most conveniently cast distributions into standard exponential family form by taking the exponential of the logarithm of the distribution. For the Beta distribution (2.13) we have

$$\operatorname{Beta}(\mu|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \exp\left\{ (a-1)\ln\mu + (b-1)\ln(1-\mu) \right\}$$

which we can identify as being in standard exponential form (2.194) with

$$h(\mu) = 1 \tag{118}$$

$$g(a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$
(119)

$$\mathbf{u}(\mu) = \begin{pmatrix} \ln \mu \\ \ln(1-\mu) \end{pmatrix} \tag{120}$$

$$\eta(a,b) = \begin{pmatrix} a-1\\b-1 \end{pmatrix}. \tag{121}$$

Applying the same approach to the gamma distribution (2.146) we obtain

$$\operatorname{Gam}(\lambda|a,b) = \frac{b^a}{\Gamma(a)} \exp\left\{ (a-1) \ln \lambda - b\lambda \right\}.$$

from which it follows that

$$h(\lambda) = 1 \tag{122}$$

$$g(a,b) = \frac{b^a}{\Gamma(a)} \tag{123}$$

$$\mathbf{u}(\lambda) = \begin{pmatrix} \lambda \\ \ln \lambda \end{pmatrix} \tag{124}$$

$$\eta(a,b) = \begin{pmatrix} -b \\ a-1 \end{pmatrix}.$$
(125)

Finally, for the von Mises distribution (2.179) we make use of the identity (2.178) to give

$$p(\theta|\theta_0, m) = \frac{1}{2\pi I_0(m)} \exp\{m\cos\theta\cos\theta_0 + m\sin\theta\sin\theta_0\}$$

from which we find

$$h(\theta) = 1 \tag{126}$$

$$g(\theta_0, m) = \frac{1}{2\pi I_0(m)} \tag{127}$$

$$\mathbf{u}(\theta) = \begin{pmatrix} \cos \theta \\ \sin \theta \end{pmatrix} \tag{128}$$

$$n(\theta) = 1$$

$$g(\theta_0, m) = \frac{1}{2\pi I_0(m)}$$

$$\mathbf{u}(\theta) = \begin{pmatrix} \cos \theta \\ \sin \theta \end{pmatrix}$$

$$\eta(\theta_0, m) = \begin{pmatrix} m \cos \theta_0 \\ m \sin \theta_0 \end{pmatrix}.$$
(128)

Starting from (2.43), we can rewrite the argument of the exponential as

$$-\frac{1}{2}\mathrm{Tr}\left[\boldsymbol{\Sigma}^{-1}\mathbf{x}\mathbf{x}^{\mathrm{T}}\right] + \boldsymbol{\mu}^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}.$$

The last term is indepedent of x but depends on  $\mu$  and  $\Sigma$  and so should go into  $g(\eta)$ . The second term is already an inner product and can be kept as is. To deal with the first term, we define the  $D^2$ -dimensional vectors  $\mathbf{z}$  and  $\boldsymbol{\lambda}$ , which consist of the columns of  $\mathbf{x}\mathbf{x}^{\mathrm{T}}$  and  $\mathbf{\Sigma}^{-1}$ , respectively, stacked on top of each other. Now we can write the multivariate Gaussian distribution on the form (2.194), with

$$\eta = \begin{bmatrix} \mathbf{\Sigma}^{-1} \boldsymbol{\mu} \\ -\frac{1}{2} \boldsymbol{\lambda} \end{bmatrix} \\
\mathbf{u}(\mathbf{x}) = \begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} \\
h(\mathbf{x}) = (2\pi)^{-D/2} \\
g(\boldsymbol{\eta}) = |\mathbf{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2} \boldsymbol{\mu}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \boldsymbol{\mu}\right).$$

**2.58** Taking the first derivative of (2.226) we obtain, as in the text,

$$-\nabla \ln g(\boldsymbol{\eta}) = g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp\left\{\boldsymbol{\eta}^{\mathrm{T}} \mathbf{u}(\mathbf{x})\right\} \mathbf{u}(\mathbf{x}) \, \mathrm{d}\mathbf{x}$$

Taking the gradient again gives

$$-\nabla\nabla \ln g(\boldsymbol{\eta}) = g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp\left\{\boldsymbol{\eta}^{\mathrm{T}} \mathbf{u}(\mathbf{x})\right\} \mathbf{u}(\mathbf{x}) \mathbf{u}(\mathbf{x})^{\mathrm{T}} d\mathbf{x}$$
$$+\nabla g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp\left\{\boldsymbol{\eta}^{\mathrm{T}} \mathbf{u}(\mathbf{x})\right\} \mathbf{u}(\mathbf{x}) d\mathbf{x}$$
$$= \mathbb{E}[\mathbf{u}(\mathbf{x}) \mathbf{u}(\mathbf{x})^{\mathrm{T}}] - \mathbb{E}[\mathbf{u}(\mathbf{x})] \mathbb{E}[\mathbf{u}(\mathbf{x})^{\mathrm{T}}]$$
$$= \cos[\mathbf{u}(\mathbf{x})]$$

where we have used the result (2.226).

2.59

$$\int \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right) dx = \frac{1}{\sigma} \int f(y) \frac{dx}{dy} dy$$
$$= \frac{1}{\sigma} \int f(y) \sigma dy$$
$$= \frac{\sigma}{\sigma} \int f(y) dy = 1,$$

since f(x) integrates to 1.

**2.60** The value of the density  $p(\mathbf{x})$  at a point  $\mathbf{x}_n$  is given by  $h_{j(n)}$ , where the notation j(n)denotes that data point  $x_n$  falls within region j. Thus the log likelihood function takes the form

$$\sum_{n=1}^{N} \ln p(\mathbf{x}_n) = \sum_{n=1}^{N} \ln h_{j(n)}.$$

We now need to take account of the constraint that  $p(\mathbf{x})$  must integrate to unity. Since  $p(\mathbf{x})$  has the constant value  $h_i$  over region i, which has volume  $\Delta_i$ , the normalization constraint becomes  $\sum_i h_i \Delta_i = 1$ . Introducing a Lagrange multiplier  $\lambda$  we then minimize the function

$$\sum_{n=1}^{N} \ln h_{j(n)} + \lambda \left( \sum_{i} h_{i} \Delta_{i} - 1 \right)$$

with respect to  $h_k$  to give

$$0 = \frac{n_k}{h} + \lambda \Delta_k$$

 $0 = \frac{n_k}{h_k} + \lambda \Delta_k$   $\text{The points falling within region is falling within region is falling within region is falling within region in the property of the points falling within region in the property of the points falling within region in the property of the points falling within region in the property of the p$ both sides by  $h_k$ , summing over k and making use of the normalization constraint,

#### **62** Solutions 2.61–3.1

we obtain  $\lambda = -N$ . Eliminating  $\lambda$  then gives our final result for the maximum likelihood solution for  $h_k$  in the form

$$h_k = \frac{n_k}{N} \frac{1}{\Delta_k}.$$

Note that, for equal sized bins  $\Delta_k = \Delta$  we obtain a bin height  $h_k$  which is proportional to the fraction of points falling within that bin, as expected.

**2.61** From (2.246) we have

$$p(\mathbf{x}) = \frac{K}{NV(\rho)}$$

where  $V(\rho)$  is the volume of a D-dimensional hypersphere with radius  $\rho$ , where in turn  $\rho$  is the distance from  ${\bf x}$  to its  $K^{\rm th}$  nearest neighbour in the data set. Thus, in polar coordinates, if we consider sufficiently large values for the radial coordinate r, we have

$$p(\mathbf{x}) \propto r^{-D}$$
.

If we consider the integral of  $p(\mathbf{x})$  and note that the volume element  $d\mathbf{x}$  can be written as  $r^{D-1} dr$ , we get

$$\int p(\mathbf{x}) \, d\mathbf{x} \propto \int r^{-D} r^{D-1} \, dr = \int r^{-1} \, dr$$

which diverges logarithmically.

## **Chapter 3** Linear Models for Regression

**3.1 NOTE**: In the 1<sup>st</sup> printing of PRML, there is a 2 missing in the denominator of the argument to the 'tanh' function in equation (3.102).

Using (3.6), we have

$$2\sigma(2a) - 1 = \frac{2}{1 + e^{-2a}} - 1$$

$$= \frac{2}{1 + e^{-2a}} - \frac{1 + e^{-2a}}{1 + e^{-2a}}$$

$$= \frac{1 - e^{-2a}}{1 + e^{-2a}}$$

$$= \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

$$= \tanh(a)$$

If we now take  $a_j = (x - \mu_j)/2s$ , we can rewrite (3.101) as

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M} w_j \sigma(2a_j)$$

$$= w_0 + \sum_{j=1}^{M} \frac{w_j}{2} (2\sigma(2a_j) - 1 + 1)$$

$$= u_0 + \sum_{j=1}^{M} u_j \tanh(a_j),$$

where  $u_j = w_j/2$ , for j = 1, ..., M, and  $u_0 = w_0 + \sum_{j=1}^{M} w_j/2$ .

**3.2** We first write

$$\Phi(\Phi^{T}\Phi)^{-1}\Phi^{T}\mathbf{v} = \Phi\widetilde{\mathbf{v}}$$

$$= \varphi_{1}\widetilde{v}^{(1)} + \varphi_{2}\widetilde{v}^{(2)} + \ldots + \varphi_{M}\widetilde{v}^{(M)}$$

where  $\varphi_m$  is the *m*-th column of  $\Phi$  and  $\widetilde{\mathbf{v}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{v}$ . By comparing this with the least squares solution in (3.15), we see that

$$\mathbf{y} = \mathbf{\Phi} \mathbf{w}_{\mathrm{ML}} = \mathbf{\Phi} (\mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi})^{-1} \mathbf{\Phi}^{\mathrm{T}} \mathbf{t}$$

corresponds to a projection of **t** onto the space spanned by the columns of  $\Phi$ . To see that this is indeed an orthogonal projection, we first note that for any column of  $\Phi$ ,  $\varphi_i$ ,

$$\boldsymbol{\Phi}(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi})^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\varphi}_{j} = \left[\boldsymbol{\Phi}(\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi})^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\right]_{i} = \boldsymbol{\varphi}_{j}$$

and therefore

$$(\mathbf{y} - \mathbf{t})^{\mathrm{T}} \boldsymbol{\varphi}_{i} = (\mathbf{\Phi} \mathbf{w}_{\mathrm{ML}} - \mathbf{t})^{\mathrm{T}} \boldsymbol{\varphi}_{i} = \mathbf{t}^{\mathrm{T}} (\mathbf{\Phi} (\mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi})^{-1} \mathbf{\Phi}^{\mathrm{T}} - \mathbf{I})^{\mathrm{T}} \boldsymbol{\varphi}_{i} = 0$$

and thus  $(\mathbf{y} - \mathbf{t})$  is ortogonal to every column of  $\Phi$  and hence is orthogonal to S.

**3.3** If we define  $\mathbf{R} = \operatorname{diag}(r_1, \dots, r_N)$  to be a diagonal matrix containing the weighting coefficients, then we can write the weighted sum-of-squares cost function in the form

$$E_D(\mathbf{w}) = \frac{1}{2} (\mathbf{t} - \mathbf{\Phi} \mathbf{w})^{\mathrm{T}} \mathbf{R} (\mathbf{t} - \mathbf{\Phi} \mathbf{w}).$$

Setting the derivative with respect to w to zero, and re-arranging, then gives

$$\mathbf{w}^{\star} = \left(\mathbf{\Phi}^{\mathrm{T}} \mathbf{R} \mathbf{\Phi}\right)^{-1} \mathbf{\Phi}^{\mathrm{T}} \mathbf{R} \mathbf{t}$$

which reduces to the standard solution (3.15) for the case  $\mathbf{R} = \mathbf{I}$ .

If we compare (3.104) with (3.10)–(3.12), we see that  $r_n$  can be regarded as a precision (inverse variance) parameter, particular to the data point  $(\mathbf{x}_n, t_n)$ , that either replaces or scales  $\beta$ .

#### **64 Solution 3.4**

Alternatively,  $r_n$  can be regarded as an *effective* number of replicated observations of data point  $(\mathbf{x}_n, t_n)$ ; this becomes particularly clear if we consider (3.104) with  $r_n$  taking positive integer values, although it is valid for any  $r_n > 0$ .

#### **3.4** Let

$$\widetilde{y}_n = w_0 + \sum_{i=1}^D w_i (x_{ni} + \epsilon_{ni})$$

$$= y_n + \sum_{i=1}^D w_i \epsilon_{ni}$$

where  $y_n = y(x_n, \mathbf{w})$  and  $\epsilon_{ni} \sim \mathcal{N}(0, \sigma^2)$  and we have used (3.105). From (3.106) we then define

$$\widetilde{E} = \frac{1}{2} \sum_{n=1}^{N} {\{\widetilde{y}_n - t_n\}}^2$$

$$= \frac{1}{2} \sum_{n=1}^{N} {\{\widetilde{y}_n^2 - 2\widetilde{y}_n t_n + t_n^2\}}$$

$$= \frac{1}{2} \sum_{n=1}^{N} {\{y_n^2 + 2y_n \sum_{i=1}^{D} w_i \epsilon_{ni} + \left(\sum_{i=1}^{D} w_i \epsilon_{ni}\right)^2}$$

$$-2t_n y_n - 2t_n \sum_{i=1}^{D} w_i \epsilon_{ni} + t_n^2 }.$$

If we take the expectation of  $\widetilde{E}$  under the distribution of  $\epsilon_{ni}$ , we see that the second and fifth terms disappear, since  $\mathbb{E}[\epsilon_{ni}] = 0$ , while for the third term we get

$$\mathbb{E}\left[\left(\sum_{i=1}^{D} w_i \epsilon_{ni}\right)^2\right] = \sum_{i=1}^{D} w_i^2 \sigma^2$$

since the  $\epsilon_{ni}$  are all independent with variance  $\sigma^2$ .

From this and (3.106) we see that

$$\mathbb{E}\left[\widetilde{E}\right] = E_D + \frac{1}{2} \sum_{i=1}^{D} w_i^2 \sigma^2,$$

as required.

**3.5** We can rewrite (3.30) as

$$\frac{1}{2} \left( \sum_{j=1}^{M} |w_j|^q - \eta \right) \leqslant 0$$

where we have incorporated the 1/2 scaling factor for convenience. Clearly this does not affect the constraint.

Employing the technique described in Appendix E, we can combine this with (3.12) to obtain the Lagrangian function

$$L(\mathbf{w}, \lambda) = \frac{1}{2} \sum_{n=1}^{N} \{t_n - \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \left( \sum_{j=1}^{M} |w_j|^q - \eta \right)$$

and by comparing this with (3.29) we see immediately that they are identical in their dependence on  $\mathbf{w}$ .

Now suppose we choose a specific value of  $\lambda > 0$  and minimize (3.29). Denoting the resulting value of  $\mathbf{w}$  by  $\mathbf{w}^*(\lambda)$ , and using the KKT condition (E.11), we see that the value of  $\eta$  is given by

$$\eta = \sum_{j=1}^{M} |w_j^{\star}(\lambda)|^q.$$

**3.6** We first write down the log likelihood function which is given by

$$\ln L(\mathbf{W}, \boldsymbol{\Sigma}) = -\frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^{N} (\mathbf{t}_n - \mathbf{W}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n))^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{W}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n)).$$

First of all we set the derivative with respect to W equal to zero, giving

$$0 = -\sum_{n=1}^{N} \mathbf{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{W}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n)) \boldsymbol{\phi}(\mathbf{x}_n)^{\mathrm{T}}.$$

Multiplying through by  $\Sigma$  and introducing the design matrix  $\Phi$  and the target data matrix T we have

$$\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi}\mathbf{W} = \mathbf{\Phi}^{\mathrm{T}}\mathbf{T}$$

Solving for W then gives (3.15) as required.

The maximum likelihood solution for  $\Sigma$  is easily found by appealing to the standard result from Chapter 2 giving

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{t}_n - \mathbf{W}_{\mathrm{ML}}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{\mathrm{ML}}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n))^{\mathrm{T}}.$$

as required. Since we are finding a joint maximum with respect to both W and  $\Sigma$  we see that it is  $W_{\rm ML}$  which appears in this expression, as in the standard result for an unconditional Gaussian distribution.

#### **3.7** From Bayes' theorem we have

$$p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{t}|\mathbf{w})p(\mathbf{w}),$$

where the factors on the r.h.s. are given by (3.10) and (3.48), respectively. Writing this out in full, we get

$$p(\mathbf{w}|\mathbf{t}) \propto \left[ \prod_{n=1}^{N} \mathcal{N} \left( t_{n} | \mathbf{w}^{T} \boldsymbol{\phi}(\mathbf{x}_{n}), \beta^{-1} \right) \right] \mathcal{N} \left( \mathbf{w} | \mathbf{m}_{0}, \mathbf{S}_{0} \right)$$

$$\propto \exp \left( -\frac{\beta}{2} (\mathbf{t} - \boldsymbol{\Phi} \mathbf{w})^{T} (\mathbf{t} - \boldsymbol{\Phi} \mathbf{w}) \right)$$

$$= \exp \left( -\frac{1}{2} (\mathbf{w} - \mathbf{m}_{0})^{T} \mathbf{S}_{0}^{-1} (\mathbf{w} - \mathbf{m}_{0}) \right)$$

$$= \exp \left( -\frac{1}{2} (\mathbf{w}^{T} \left( \mathbf{S}_{0}^{-1} + \beta \boldsymbol{\Phi}^{T} \boldsymbol{\Phi} \right) \mathbf{w} - \beta \mathbf{t}^{T} \boldsymbol{\Phi} \mathbf{w} - \beta \mathbf{w}^{T} \boldsymbol{\Phi}^{T} \mathbf{t} + \beta \mathbf{t}^{T} \mathbf{t} \right)$$

$$= \exp \left( -\frac{1}{2} (\mathbf{w}^{T} \left( \mathbf{S}_{0}^{-1} + \beta \boldsymbol{\Phi}^{T} \boldsymbol{\Phi} \right) \mathbf{w} - \left( \mathbf{S}_{0}^{-1} \mathbf{m}_{0} + \beta \boldsymbol{\Phi}^{T} \mathbf{t} \right)^{T} \mathbf{w} \right)$$

$$= \exp \left( -\frac{1}{2} (\mathbf{w}^{T} \left( \mathbf{S}_{0}^{-1} + \beta \boldsymbol{\Phi}^{T} \boldsymbol{\Phi} \right) \mathbf{w} - \left( \mathbf{S}_{0}^{-1} \mathbf{m}_{0} + \beta \boldsymbol{\Phi}^{T} \mathbf{t} \right)^{T} \mathbf{w} \right)$$

$$= \exp \left( -\frac{1}{2} (\mathbf{w} - \mathbf{m}_{N})^{T} \mathbf{S}_{N}^{-1} (\mathbf{w} - \mathbf{m}_{N}) \right)$$

$$= \exp \left( -\frac{1}{2} \left( \beta \mathbf{t}^{T} \mathbf{t} + \mathbf{m}_{0}^{T} \mathbf{S}_{0}^{-1} \mathbf{m}_{0} - \mathbf{m}_{N}^{T} \mathbf{S}_{N}^{-1} \mathbf{m}_{N} \right) \right)$$

where we have used (3.50) and (3.51) when completing the square in the last step. The first exponential corrsponds to the posterior, unnormalized Gaussian distribution over w, while the second exponential is independent of w and hence can be absorbed into the normalization factor.

#### **3.8** Combining the prior

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

and the likelihood

$$p(t_{N+1}|\mathbf{x}_{N+1},\mathbf{w}) = \left(\frac{\beta}{2\pi}\right)^{1/2} \exp\left(-\frac{\beta}{2}(t_{N+1} - \mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}_{N+1})^{2}\right)$$
(130)

where  $\phi_{N+1} = \phi(\mathbf{x}_{N+1})$ , we obtain a posterior of the form

$$p(\mathbf{w}|t_{N+1}, \mathbf{x}_{N+1}, \mathbf{m}_N, \mathbf{S}_N)$$

$$\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^{\mathrm{T}} \mathbf{S}_N^{-1}(\mathbf{w} - \mathbf{m}_N) - \frac{1}{2}\beta(t_{N+1} - \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}_{N+1})^2\right).$$

We can expand the argument of the exponential, omitting the -1/2 factors, as follows

$$(\mathbf{w} - \mathbf{m}_{N})^{\mathrm{T}} \mathbf{S}_{N}^{-1} (\mathbf{w} - \mathbf{m}_{N}) + \beta (t_{N+1} - \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}_{N+1})^{2}$$

$$= \mathbf{w}^{\mathrm{T}} \mathbf{S}_{N}^{-1} \mathbf{w} - 2 \mathbf{w}^{\mathrm{T}} \mathbf{S}_{N}^{-1} \mathbf{m}_{N}$$

$$+ \beta \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}_{N+1}^{\mathrm{T}} \boldsymbol{\phi}_{N+1} \mathbf{w} - 2\beta \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}_{N+1} t_{N+1} + \text{const}$$

$$= \mathbf{w}^{\mathrm{T}} (\mathbf{S}_{N}^{-1} + \beta \boldsymbol{\phi}_{N+1} \boldsymbol{\phi}_{N+1}^{\mathrm{T}}) \mathbf{w} - 2 \mathbf{w}^{\mathrm{T}} (\mathbf{S}_{N}^{-1} \mathbf{m}_{N} + \beta \boldsymbol{\phi}_{N+1} t_{N+1}) + \text{const},$$

where const denotes remaining terms independent of w. From this we can read off the desired result directly,

$$p(\mathbf{w}|t_{N+1}, \mathbf{x}_{N+1}, \mathbf{m}_N, \mathbf{S}_N) = \mathcal{N}(\mathbf{w}|\mathbf{m}_{N+1}, \mathbf{S}_{N+1}),$$

with

$$\mathbf{S}_{N+1}^{-1} = \mathbf{S}_{N}^{-1} + \beta \phi_{N+1} \phi_{N+1}^{\mathrm{T}}. \tag{131}$$

and

$$\mathbf{m}_{N+1} = \mathbf{S}_{N+1}(\mathbf{S}_N^{-1}\mathbf{m}_N + \beta \phi_{N+1}t_{N+1}). \tag{132}$$

**3.9** Identifying (2.113) with (3.49) and (2.114) with (130), such that

$$\mathbf{x} \Rightarrow \mathbf{w} \quad \boldsymbol{\mu} \Rightarrow \mathbf{m}_N \quad \boldsymbol{\Lambda}^{-1} \Rightarrow \mathbf{S}_N$$
 $\mathbf{y} \Rightarrow t_{N+1} \quad \mathbf{A} \Rightarrow \boldsymbol{\phi}(\mathbf{x}_{N+1})^{\mathrm{T}} = \boldsymbol{\phi}_{N+1}^{\mathrm{T}} \quad \mathbf{b} \Rightarrow \mathbf{0} \quad \mathbf{L}^{-1} \Rightarrow \beta \mathbf{I},$ 

(2.116) and (2.117) directly give

$$p(\mathbf{w}|t_{N+1}, \mathbf{x}_{N+1}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_{N+1}, \mathbf{S}_{N+1})$$

where  $S_{N+1}$  and  $m_{N+1}$  are given by (131) and (132), respectively.

**3.10** Using (3.3), (3.8) and (3.49), we can re-write (3.57) as

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \int \mathcal{N}(t|\boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}}\mathbf{w}, \beta^{-1})\mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) d\mathbf{w}.$$

By matching the first factor of the integrand with (2.114) and the second factor with (2.113), we obtain the desired result directly from (2.115).

**3.11** From (3.59) we have

$$\sigma_{N+1}^{2}(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^{\mathrm{T}} \mathbf{S}_{N+1} \phi(\mathbf{x})$$
 (133)

where  $S_{N+1}$  is given by (131). From (131) and (3.110) we get

$$\mathbf{S}_{N+1} = \left(\mathbf{S}_{N}^{-1} + \beta \phi_{N+1} \phi_{N+1}^{\mathrm{T}}\right)^{-1}$$

$$= \mathbf{S}_{N} - \frac{\left(\mathbf{S}_{N} \phi_{N+1} \beta^{1/2}\right) \left(\beta^{1/2} \phi_{N+1}^{\mathrm{T}} \mathbf{S}_{N}\right)}{1 + \beta \phi_{N+1}^{\mathrm{T}} \mathbf{S}_{N} \phi_{N+1}}$$

$$= \mathbf{S}_{N} - \frac{\beta \mathbf{S}_{N} \phi_{N+1} \phi_{N+1}^{\mathrm{T}} \mathbf{S}_{N}}{1 + \beta \phi_{N+1}^{\mathrm{T}} \mathbf{S}_{N} \phi_{N+1}}.$$

Using this and (3.59), we can rewrite (133) as

$$\sigma_{N+1}^{2}(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^{\mathrm{T}} \left( \mathbf{S}_{N} - \frac{\beta \mathbf{S}_{N} \phi_{N+1} \phi_{N+1}^{\mathrm{T}} \mathbf{S}_{N}}{1 + \beta \phi_{N+1}^{\mathrm{T}} \mathbf{S}_{N} \phi_{N+1}} \right) \phi(\mathbf{x})$$

$$= \sigma_{N}^{2}(\mathbf{x}) - \frac{\beta \phi(\mathbf{x})^{\mathrm{T}} \mathbf{S}_{N} \phi_{N+1} \phi_{N+1}^{\mathrm{T}} \mathbf{S}_{N} \phi(\mathbf{x})}{1 + \beta \phi_{N+1}^{\mathrm{T}} \mathbf{S}_{N} \phi_{N+1}}.$$
(134)

Since  $\mathbf{S}_N$  is positive definite, the numerator and denominator of the second term in (134) will be non-negative and positive, respectively, and hence  $\sigma_{N+1}^2(\mathbf{x}) \leqslant \sigma_N^2(\mathbf{x})$ .

**3.12** It is easiest to work in log space. The log of the posterior distribution is given by

$$\ln p(\mathbf{w}, \beta | \mathbf{t}) = \ln p(\mathbf{w}, \beta) + \sum_{n=1}^{N} \ln p(t_n | \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

$$= \frac{M}{2} \ln \beta - \frac{1}{2} \ln |\mathbf{S}_0| - \frac{\beta}{2} (\mathbf{w} - \mathbf{m}_0)^{\mathrm{T}} \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0)$$

$$-b_0 \beta + (a_0 - 1) \ln \beta$$

$$+ \frac{N}{2} \ln \beta - \frac{\beta}{2} \sum_{n=1}^{N} {\{\mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n) - t_n\}^2 + \text{const.}}$$

Using the product rule, the posterior distribution can be written as  $p(\mathbf{w}, \beta | \mathbf{t}) = p(\mathbf{w} | \beta, \mathbf{t}) p(\beta | \mathbf{t})$ . Consider first the dependence on  $\mathbf{w}$ . We have

$$\ln p(\mathbf{w}|\beta, \mathbf{t}) = -\frac{\beta}{2} \mathbf{w}^{\mathrm{T}} \left[ \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} + \mathbf{S}_{0}^{-1} \right] \mathbf{w} + \mathbf{w}^{\mathrm{T}} \left[ \beta \mathbf{S}_{0}^{-1} \mathbf{m}_{0} + \beta \mathbf{\Phi}^{\mathrm{T}} \mathbf{t} \right] + \text{const.}$$

Thus we see that  $p(\mathbf{w}|\beta, \mathbf{t})$  is a Gaussian distribution with mean and covariance given by

$$\mathbf{m}_{N} = \mathbf{S}_{N} \left[ \mathbf{S}_{0}^{-1} \mathbf{m}_{0} + \mathbf{\Phi}^{\mathrm{T}} \mathbf{t} \right]$$
 (135)

$$\beta \mathbf{S}_{N}^{-1} = \beta \left( \mathbf{S}_{0}^{-1} + \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \right). \tag{136}$$

To find  $p(\beta|\mathbf{t})$  we first need to complete the square over  $\mathbf{w}$  to ensure that we pick up all terms involving  $\beta$  (any terms independent of  $\beta$  may be discarded since these will be absorbed into the normalization coefficient which itself will be found by inspection at the end). We also need to remember that a factor of  $(M/2) \ln \beta$  will be absorbed by the normalisation factor of  $p(\mathbf{w}|\beta,\mathbf{t})$ . Thus

$$\ln p(\beta|\mathbf{t}) = -\frac{\beta}{2}\mathbf{m}_0^{\mathrm{T}}\mathbf{S}_0^{-1}\mathbf{m}_0 + \frac{\beta}{2}\mathbf{m}_N^{\mathrm{T}}\mathbf{S}_N^{-1}\mathbf{m}_N$$
$$+ \frac{N}{2}\ln \beta - b_0\beta + (a_0 - 1)\ln \beta - \frac{\beta}{2}\sum_{n=1}^N t_n^2 + \text{const.}$$

We recognize this as the log of a Gamma distribution. Reading off the coefficients of  $\beta$  and  $\ln \beta$  we then have

$$a_N = a_0 + \frac{N}{2} (137)$$

$$b_N = b_0 + \frac{1}{2} \left( \mathbf{m}_0^{\mathrm{T}} \mathbf{S}_0^{-1} \mathbf{m}_0 - \mathbf{m}_N^{\mathrm{T}} \mathbf{S}_N^{-1} \mathbf{m}_N + \sum_{n=1}^N t_n^2 \right).$$
 (138)

**3.13** Following the line of presentation from Section 3.3.2, the predictive distribution is now given by

$$p(t|\mathbf{x}, \mathbf{t}) = \iint \mathcal{N}\left(t|\phi(\mathbf{x})^{\mathrm{T}}\mathbf{w}, \beta^{-1}\right) \mathcal{N}\left(\mathbf{w}|\mathbf{m}_{N}, \beta^{-1}\mathbf{S}_{N}\right) d\mathbf{w}$$

$$\operatorname{Gam}\left(\beta|a_{N}, b_{N}\right) d\beta \quad (139)$$

We begin by performing the integral over w. Identifying (2.113) with (3.49) and (2.114) with (3.8), using (3.3), such that

$$\mathbf{x} \Rightarrow \mathbf{w} \quad \boldsymbol{\mu} \Rightarrow \mathbf{m}_N \quad \boldsymbol{\Lambda}^{-1} \Rightarrow \mathbf{S}_N$$

$$\mathbf{y} \Rightarrow t \quad \mathbf{A} \Rightarrow \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} = \boldsymbol{\phi}^{\mathrm{T}} \quad \mathbf{b} \Rightarrow 0 \quad \mathbf{L}^{-1} \Rightarrow \boldsymbol{\beta}^{-1},$$

(2.115) and (136) give

$$p(t|\beta) = \mathcal{N}\left(t|\boldsymbol{\phi}^{\mathrm{T}}\mathbf{m}_{N}, \beta^{-1} + \boldsymbol{\phi}^{\mathrm{T}}\mathbf{S}_{N}\boldsymbol{\phi}\right)$$
$$= \mathcal{N}\left(t|\boldsymbol{\phi}^{\mathrm{T}}\mathbf{m}_{N}, \beta^{-1}\left(1 + \boldsymbol{\phi}^{\mathrm{T}}(\mathbf{S}_{0} + \boldsymbol{\phi}^{\mathrm{T}}\boldsymbol{\phi})^{-1}\boldsymbol{\phi}\right)\right).$$

Substituting this back into (139) we get

$$p(t|\mathbf{x}, \mathbf{X}, \mathbf{t}) = \int \mathcal{N}(t|\phi^{\mathrm{T}}\mathbf{m}_N, \beta^{-1}s) \operatorname{Gam}(\beta|a_N, b_N) d\beta,$$

where we have defined

$$s = 1 + \boldsymbol{\phi}^{\mathrm{T}} (\mathbf{S}_0 + \boldsymbol{\phi}^{\mathrm{T}} \boldsymbol{\phi})^{-1} \boldsymbol{\phi}.$$

We can now use (2.158)–(2.160) to obtain the final result:

$$p(t|\mathbf{x}, \mathbf{X}, \mathbf{t}) = \operatorname{St}(t|\mu, \lambda, \nu)$$

where

$$\mu = \boldsymbol{\phi}^{\mathrm{T}} \mathbf{m}_{N} \quad \lambda = \frac{a_{N}}{b_{N}} s^{-1} \quad \nu = 2a_{N}.$$

**3.14** For  $\alpha = 0$  the covariance matrix  $S_N$  becomes

$$\mathbf{S}_N = (\beta \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi})^{-1}. \tag{140}$$

#### **70** Solution 3.14

Let us define a new set of orthonormal basis functions given by linear combinations of the original basis functions so that

$$\psi(\mathbf{x}) = \mathbf{V}\phi(\mathbf{x}) \tag{141}$$

where V is an  $M \times M$  matrix. Since both the original and the new basis functions are linearly independent and span the same space, this matrix must be invertible and hence

$$\phi(\mathbf{x}) = \mathbf{V}^{-1}\psi(\mathbf{x}).$$

For the data set  $\{x_n\}$ , (141) and (3.16) give

$$\boldsymbol{\Psi} = \boldsymbol{\Phi} \mathbf{V}^T$$

and consequently

$$\Phi = \Psi V^{-T}$$

where  $V^{-T}$  denotes  $(V^{-1})^{T}$ . Orthonormality implies

$$\Psi^{\mathrm{T}}\Psi = \mathbf{I}.$$

Note that  $(\mathbf{V}^{-1})^{\mathrm{T}} = (\mathbf{V}^{\mathrm{T}})^{-1}$  as is easily verified. From (140), the covariance matrix then becomes

$$\mathbf{S}_N = \beta^{-1} (\mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi})^{-1} = \beta^{-1} (\mathbf{V}^{-\mathrm{T}} \mathbf{\Psi}^{\mathrm{T}} \mathbf{\Psi} \mathbf{V}^{-1})^{-1} = \beta^{-1} \mathbf{V}^{\mathrm{T}} \mathbf{V}.$$

Here we have used the orthonormality of the  $\psi_i(\mathbf{x})$ . Hence the equivalent kernel becomes

$$k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^{\mathrm{T}} \mathbf{S}_N \phi(\mathbf{x}') = \phi(\mathbf{x})^{\mathrm{T}} \mathbf{V}^{\mathrm{T}} \mathbf{V} \phi(\mathbf{x}') = \psi(\mathbf{x})^{\mathrm{T}} \psi(\mathbf{x}')$$

as required. From the orthonormality condition, and setting j = 1, it follows that

$$\sum_{n=1}^{N} \psi_i(\mathbf{x}_n) \psi_1(\mathbf{x}_n) = \sum_{n=1}^{N} \psi_i(\mathbf{x}_n) = \delta_{i1}$$

where we have used  $\psi_1(\mathbf{x}) = 1$ . Now consider the sum

$$\sum_{n=1}^{N} k(\mathbf{x}, \mathbf{x}_n) = \sum_{n=1}^{N} \boldsymbol{\psi}(\mathbf{x})^{\mathrm{T}} \boldsymbol{\psi}(\mathbf{x}_n) = \sum_{n=1}^{N} \sum_{i=1}^{M} \psi_i(\mathbf{x}) \psi_i(\mathbf{x}_n)$$
$$= \sum_{i=1}^{M} \psi_i(\mathbf{x}) \delta_{i1} = \psi_1(\mathbf{x}) = 1$$

which proves the summation constraint as required.

**3.15** This is easily shown by substituting the re-estimation formulae (3.92) and (3.95) into (3.82), giving

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^{\mathrm{T}} \mathbf{m}_N$$
$$= \frac{N - \gamma}{2} + \frac{\gamma}{2} = \frac{N}{2}.$$

**3.16** The likelihood function is a product of independent univariate Gaussians and so can be written as a joint Gaussian distribution over **t** with diagonal covariance matrix in the form

$$p(\mathbf{t}|\mathbf{w},\beta) = \mathcal{N}(\mathbf{t}|\mathbf{\Phi}\mathbf{w},\beta^{-1}\mathbf{I}_N). \tag{142}$$

Identifying (2.113) with the prior distribution  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$  and (2.114) with (142), such that

$$\mathbf{x} \Rightarrow \mathbf{w} \quad \boldsymbol{\mu} \Rightarrow \mathbf{0} \quad \boldsymbol{\Lambda}^{-1} \Rightarrow \alpha^{-1} \mathbf{I}_{M}$$
 $\mathbf{v} \Rightarrow \mathbf{t} \quad \mathbf{A} \Rightarrow \boldsymbol{\Phi} \quad \mathbf{b} \Rightarrow \mathbf{0} \quad \mathbf{L}^{-1} \Rightarrow \boldsymbol{\beta}^{-1} \mathbf{I}_{N}.$ 

(2.115) gives

$$p(\mathbf{t}|\alpha,\beta) = \mathcal{N}(\mathbf{t}|\mathbf{0},\beta^{-1}\mathbf{I}_N + \alpha^{-1}\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}}).$$

Taking the log we obtain

$$\ln p(\mathbf{t}|\alpha,\beta) = -\frac{N}{2}\ln(2\pi) - \frac{1}{2}\ln\left|\beta^{-1}\mathbf{I}_N + \alpha^{-1}\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}}\right| - \frac{1}{2}\mathbf{t}^{\mathrm{T}}\left(\beta^{-1}\mathbf{I}_N + \alpha^{-1}\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}}\right)\mathbf{t}. \quad (143)$$

Using the result (C.14) for the determinant we have

$$\begin{vmatrix} \beta^{-1} \mathbf{I}_N + \alpha^{-1} \mathbf{\Phi} \mathbf{\Phi}^{\mathrm{T}} | &= \beta^{-N} | \mathbf{I}_N + \beta \alpha^{-1} \mathbf{\Phi} \mathbf{\Phi}^{\mathrm{T}} | \\ &= \beta^{-N} | \mathbf{I}_M + \beta \alpha^{-1} \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} | \\ &= \beta^{-N} \alpha^{-M} | \alpha \mathbf{I}_M + \beta \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} | \\ &= \beta^{-N} \alpha^{-M} | \mathbf{A} | \end{aligned}$$

where we have used (3.81). Next consider the quadratic term in  $\mathbf{t}$  and make use of the identity (C.7) together with (3.81) and (3.84) to give

$$-\frac{1}{2}\mathbf{t} \left(\beta^{-1}\mathbf{I}_{N} + \alpha^{-1}\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}}\right)^{-1}\mathbf{t}$$

$$= -\frac{1}{2}\mathbf{t}^{\mathrm{T}} \left[\beta\mathbf{I}_{N} - \beta\mathbf{\Phi} \left(\alpha\mathbf{I}_{M} + \beta\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi}\right)^{-1}\mathbf{\Phi}^{\mathrm{T}}\beta\right]\mathbf{t}$$

$$= -\frac{\beta}{2}\mathbf{t}^{\mathrm{T}}\mathbf{t} + \frac{\beta^{2}}{2}\mathbf{t}^{\mathrm{T}}\mathbf{\Phi}\mathbf{A}^{-1}\mathbf{\Phi}^{\mathrm{T}}\mathbf{t}$$

欢迎关注公众号望机器学习与算法之道

where in the last step, we have exploited results from Solution 3.18. Substituting for the determinant and the quadratic term in (143) we obtain (3.86).

**3.17** Using (3.11), (3.12) and (3.52) together with the definition for the Gaussian, (2.43), we can rewrite (3.77) as follows:

$$p(\mathbf{t}|\alpha,\beta) = \int p(\mathbf{t}|\mathbf{w},\beta)p(\mathbf{w}|\alpha) d\mathbf{w}$$

$$= \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{M/2} \int \exp\left(-\beta E_D(\mathbf{w})\right) \exp\left(-\frac{\alpha}{2}\mathbf{w}^{\mathrm{T}}\mathbf{w}\right) d\mathbf{w}$$

$$= \left(\frac{\beta}{2\pi}\right)^{N/2} \left(\frac{\alpha}{2\pi}\right)^{M/2} \int \exp\left(-E(\mathbf{w})\right) d\mathbf{w},$$

where  $E(\mathbf{w})$  is defined by (3.79).

**3.18** We can rewrite (3.79)

$$\frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^{\mathrm{T}} \mathbf{w} 
= \frac{\beta}{2} (\mathbf{t}^{\mathrm{T}} \mathbf{t} - 2 \mathbf{t}^{\mathrm{T}} \mathbf{\Phi} \mathbf{w} + \mathbf{w}^{\mathrm{T}} \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \mathbf{w}) + \frac{\alpha}{2} \mathbf{w}^{\mathrm{T}} \mathbf{w} 
= \frac{1}{2} (\beta \mathbf{t}^{\mathrm{T}} \mathbf{t} - 2\beta \mathbf{t}^{\mathrm{T}} \mathbf{\Phi} \mathbf{w} + \mathbf{w}^{\mathrm{T}} \mathbf{A} \mathbf{w})$$

where, in the last line, we have used (3.81). We now use the tricks of adding  $\mathbf{0} = \mathbf{m}_{\mathbf{N}}^{\mathrm{T}} \mathbf{A} \mathbf{m}_{\mathbf{N}} - \mathbf{m}_{\mathbf{N}}^{\mathrm{T}} \mathbf{A} \mathbf{m}_{\mathbf{N}}$  and using  $\mathbf{I} = \mathbf{A}^{-1} \mathbf{A}$ , combined with (3.84), as follows:

$$\begin{split} &\frac{1}{2} \left( \beta \mathbf{t}^{\mathrm{T}} \mathbf{t} - 2\beta \mathbf{t}^{\mathrm{T}} \mathbf{\Phi} \mathbf{w} + \mathbf{w}^{\mathrm{T}} \mathbf{A} \mathbf{w} \right) \\ &= \frac{1}{2} \left( \beta \mathbf{t}^{\mathrm{T}} \mathbf{t} - 2\beta \mathbf{t}^{\mathrm{T}} \mathbf{\Phi} \mathbf{A}^{-1} \mathbf{A} \mathbf{w} + \mathbf{w}^{\mathrm{T}} \mathbf{A} \mathbf{w} \right) \\ &= \frac{1}{2} \left( \beta \mathbf{t}^{\mathrm{T}} \mathbf{t} - 2 \mathbf{m}_{N}^{\mathrm{T}} \mathbf{A} \mathbf{w} + \mathbf{w}^{\mathrm{T}} \mathbf{A} \mathbf{w} + \mathbf{m}_{N}^{\mathrm{T}} \mathbf{A} \mathbf{m}_{N} - \mathbf{m}_{N}^{\mathrm{T}} \mathbf{A} \mathbf{m}_{N} \right) \\ &= \frac{1}{2} \left( \beta \mathbf{t}^{\mathrm{T}} \mathbf{t} - \mathbf{m}_{N}^{\mathrm{T}} \mathbf{A} \mathbf{m}_{N} \right) + \frac{1}{2} (\mathbf{w} - \mathbf{m}_{N})^{\mathrm{T}} \mathbf{A} (\mathbf{w} - \mathbf{m}_{N}). \end{split}$$

Here the last term equals term the last term of (3.80) and so it remains to show that the first term equals the r.h.s. of (3.82). To do this, we use the same tricks again:

$$\begin{split} &\frac{1}{2}\left(\beta\mathbf{t}^{\mathrm{T}}\mathbf{t}-\mathbf{m}_{N}^{\mathrm{T}}\mathbf{A}\mathbf{m}_{N}\right)=\frac{1}{2}\left(\beta\mathbf{t}^{\mathrm{T}}\mathbf{t}-2\mathbf{m}_{N}^{\mathrm{T}}\mathbf{A}\mathbf{m}_{N}+\mathbf{m}_{N}^{\mathrm{T}}\mathbf{A}\mathbf{m}_{N}\right)\\ &=\frac{1}{2}\left(\beta\mathbf{t}^{\mathrm{T}}\mathbf{t}-2\mathbf{m}_{N}^{\mathrm{T}}\mathbf{A}\mathbf{A}^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{t}\beta+\mathbf{m}_{N}^{\mathrm{T}}\left(\alpha\mathbf{I}+\beta\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\right)\mathbf{m}_{N}\right)\\ &=\frac{1}{2}\left(\beta\mathbf{t}^{\mathrm{T}}\mathbf{t}-2\mathbf{m}_{N}^{\mathrm{T}}\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{t}\beta+\beta\mathbf{m}_{N}^{\mathrm{T}}\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\mathbf{m}_{N}+\alpha\mathbf{m}_{N}^{\mathrm{T}}\mathbf{m}_{N}\right)\\ &=\frac{1}{2}\left(\beta(\mathbf{t}-\boldsymbol{\Phi}\mathbf{m}_{N})^{\mathrm{T}}(\mathbf{t}-\boldsymbol{\Phi}\mathbf{m}_{N})+\alpha\mathbf{m}_{N}^{\mathrm{T}}\mathbf{m}_{N}\right)\\ &=\frac{\beta}{2}\left\|\mathbf{t}-\boldsymbol{\Phi}\mathbf{m}_{N}\right\|^{2}+\frac{\alpha}{2}\mathbf{m}_{N}^{\mathrm{T}}\mathbf{m}_{N}\end{split}$$

as required.

**3.19** From (3.80) we see that the integrand of (3.85) is an unnormalized Gaussian and hence integrates to the inverse of the corresponding normalizing constant, which can be read off from the r.h.s. of (2.43) as

$$(2\pi)^{M/2} |\mathbf{A}^{-1}|^{1/2}$$
.

Using (3.78), (3.85) and the properties of the logarithm, we get

$$\ln p(\mathbf{t}|\alpha,\beta) = \frac{M}{2} (\ln \alpha - \ln(2\pi)) + \frac{N}{2} (\ln \beta - \ln(2\pi)) + \ln \int \exp\{-E(\mathbf{w})\} d\mathbf{w}$$
$$= \frac{M}{2} (\ln \alpha - \ln(2\pi)) + \frac{N}{2} (\ln \beta - \ln(2\pi)) - E(\mathbf{m}_N) - \frac{1}{2} \ln |\mathbf{A}| + \frac{M}{2} \ln(2\pi)$$

which equals (3.86).

**3.20** We only need to consider the terms of (3.86) that depend on  $\alpha$ , which are the first, third and fourth terms.

Following the sequence of steps in Section 3.5.2, we start with the last of these terms,

$$-\frac{1}{2}\ln|\mathbf{A}|.$$

From (3.81), (3.87) and the fact that that eigenvectors  $\mathbf{u}_i$  are orthonormal (see also Appendix C), we find that the eigenvectors of  $\mathbf{A}$  to be  $\alpha + \lambda_i$ . We can then use (C.47) and the properties of the logarithm to take us from the left to the right side of (3.88).

The derivatives for the first and third term of (3.86) are more easily obtained using standard derivatives and (3.82), yielding

$$\frac{1}{2} \left( \frac{M}{\alpha} + \mathbf{m}_N^{\mathrm{T}} \mathbf{m}_N \right).$$

We combine these results into (3.89), from which we get (3.92) via (3.90). The expression for  $\gamma$  in (3.91) is obtained from (3.90) by substituting

$$\sum_{i}^{M} \frac{\lambda_i + \alpha}{\lambda_i + \alpha}$$

for M and re-arranging.

**3.21** The eigenvector equation for the  $M \times M$  real, symmetric matrix  $\mathbf A$  can be written as

$$\mathbf{A}\mathbf{u}_i = \eta_i \mathbf{u}_i$$

#### **74** Solution 3.21

where  $\{\mathbf{u}_i\}$  are a set of M orthonormal vectors, and the M eigenvalues  $\{\eta_i\}$  are all real. We first express the left hand side of (3.117) in terms of the eigenvalues of  $\mathbf{A}$ . The log of the determinant of  $\mathbf{A}$  can be written as

$$\ln |\mathbf{A}| = \ln \prod_{i=1}^{M} \eta_i = \sum_{i=1}^{M} \ln \eta_i.$$

Taking the derivative with respect to some scalar  $\alpha$  we obtain

$$\frac{\mathrm{d}}{\mathrm{d}\alpha}\ln|\mathbf{A}| = \sum_{i=1}^{M} \frac{1}{\eta_i} \frac{\mathrm{d}}{\mathrm{d}\alpha} \eta_i.$$

We now express the right hand side of (3.117) in terms of the eigenvector expansion and show that it takes the same form. First we note that  $\mathbf{A}$  can be expanded in terms of its own eigenvectors to give

$$\mathbf{A} = \sum_{i=1}^{M} \eta_i \mathbf{u}_i \mathbf{u}_i^{\mathrm{T}}$$

and similarly the inverse can be written as

$$\mathbf{A}^{-1} = \sum_{i=1}^{M} \frac{1}{\eta_i} \mathbf{u}_i \mathbf{u}_i^{\mathrm{T}}.$$

Thus we have

$$\operatorname{Tr}\left(\mathbf{A}^{-1} \frac{\mathrm{d}}{\mathrm{d}\alpha} \mathbf{A}\right) = \operatorname{Tr}\left(\sum_{i=1}^{M} \frac{1}{\eta_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \frac{\mathrm{d}}{\mathrm{d}\alpha} \sum_{j=1}^{M} \eta_{j} \mathbf{u}_{j} \mathbf{u}_{j}^{\mathrm{T}}\right)$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \frac{1}{\eta_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \left\{\sum_{j=1}^{M} \frac{\mathrm{d}\eta_{j}}{\mathrm{d}\alpha} \mathbf{u}_{j} \mathbf{u}_{j}^{\mathrm{T}} + \eta_{j} \left(\mathbf{b}_{j} \mathbf{u}_{j}^{\mathrm{T}} + \mathbf{u}_{j} \mathbf{b}_{j}^{\mathrm{T}}\right)\right\}\right)$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \frac{1}{\eta_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \sum_{j=1}^{M} \frac{\mathrm{d}\eta_{j}}{\mathrm{d}\alpha} \mathbf{u}_{j} \mathbf{u}_{j}^{\mathrm{T}}\right)$$

$$+ \operatorname{Tr}\left(\sum_{i=1}^{M} \frac{1}{\eta_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \sum_{j=1}^{M} \eta_{j} \left(\mathbf{b}_{j} \mathbf{u}_{j}^{\mathrm{T}} + \mathbf{u}_{j} \mathbf{b}_{j}^{\mathrm{T}}\right)\right)$$

$$(144)$$

where  $\mathbf{b}_{j} = d\mathbf{u}_{j}/d\alpha$ . Using the properties of the trace and the orthogonality of

eigenvectors, we can rewrite the second term as

$$\operatorname{Tr}\left(\sum_{i=1}^{M} \frac{1}{\eta_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \sum_{j=1}^{M} \eta_{j} \left(\mathbf{b}_{j} \mathbf{u}_{j}^{\mathrm{T}} + \mathbf{u}_{j} \mathbf{b}_{j}^{\mathrm{T}}\right)\right)$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \frac{1}{\eta_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \sum_{j=1}^{M} 2 \eta_{j} \mathbf{u}_{j} \mathbf{b}_{j}^{\mathrm{T}}\right)$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \sum_{j=1}^{M} \frac{2 \eta_{j}}{\eta_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \mathbf{u}_{j} \mathbf{b}_{j}^{\mathrm{T}}\right)$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \left(\mathbf{b}_{j} \mathbf{u}_{j}^{\mathrm{T}} + \mathbf{u}_{j} \mathbf{b}_{j}^{\mathrm{T}}\right)\right)$$

$$= \operatorname{Tr}\left(\frac{\mathrm{d}}{\mathrm{d}\alpha} \sum_{i}^{M} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}}\right).$$

However,

$$\sum_{i}^{M}\mathbf{u}_{i}\mathbf{u}_{i}^{\mathrm{T}}=\mathbf{I}$$

which is constant and thus its derivative w.r.t.  $\alpha$  will be zero and the second term in (144) vanishes.

For the first term in (144), we again use the properties of the trace and the orthognality of eigenvectors to obtain

$$\operatorname{Tr}\left(\mathbf{A}^{-1}\frac{\mathrm{d}}{\mathrm{d}\alpha}\mathbf{A}\right) = \sum_{i=1}^{M} \frac{1}{\eta_i} \frac{\mathrm{d}\eta_i}{\mathrm{d}\alpha}.$$

We have now shown that both the left and right hand sides of (3.117) take the same form when expressed in terms of the eigenvector expansion. Next, we use (3.117) to differentiate (3.86) w.r.t.  $\alpha$ , yielding

$$\frac{\mathrm{d}}{\mathrm{d}\alpha} \ln p(\mathbf{t}|\alpha\beta) = \frac{M}{2} \frac{1}{\alpha} - \frac{1}{2} \mathbf{m}_N^{\mathrm{T}} \mathbf{m}_N - \frac{1}{2} \mathrm{Tr} \left( \mathbf{A}^{-1} \frac{\mathrm{d}}{\mathrm{d}\alpha} \mathbf{A} \right)$$

$$= \frac{1}{2} \left( \frac{M}{\alpha} - \mathbf{m}_N^{\mathrm{T}} \mathbf{m}_N - \mathrm{Tr} \left( \mathbf{A}^{-1} \right) \right)$$

$$= \frac{1}{2} \left( \frac{M}{\alpha} - \mathbf{m}_N^{\mathrm{T}} \mathbf{m}_N - \sum_i \frac{1}{\lambda_i + \alpha} \right)$$

which we recognize as the r.h.s. of (3.89), from which (3.92) can be derived as detailed in Section 3.5.2, immediately following (3.89).

- **3.22** Using (3.82) and (3.93)—the derivation of latter is detailed in Section 3.5.2—we get the derivative of (3.86) w.r.t.  $\beta$  as the r.h.s. of (3.94). Rearranging this, collecting the  $\beta$ -dependent terms on one side of the equation and the remaining term on the other, we obtain (3.95).
- **3.23** From (3.10), (3.112) and the properties of the Gaussian and Gamma distributions (see Appendix B), we get

$$p(\mathbf{t}) = \iint p(\mathbf{t}|\mathbf{w}, \beta) p(\mathbf{w}|\beta) \, d\mathbf{w} p(\beta) \, d\beta$$

$$= \iint \left(\frac{\beta}{2\pi}\right)^{N/2} \exp\left\{-\frac{\beta}{2}(\mathbf{t} - \mathbf{\Phi}\mathbf{w})^{\mathrm{T}}(\mathbf{t} - \mathbf{\Phi}\mathbf{w})\right\}$$

$$\left(\frac{\beta}{2\pi}\right)^{M/2} |\mathbf{S}_{0}|^{-1/2} \exp\left\{-\frac{\beta}{2}(\mathbf{w} - \mathbf{m}_{0})^{\mathrm{T}}\mathbf{S}_{0}^{-1}(\mathbf{w} - \mathbf{m}_{0})\right\} \, d\mathbf{w}$$

$$\Gamma(a_{0})^{-1} b_{0}^{a_{0}} \beta^{a_{0}-1} \exp(-b_{0}\beta) \, d\beta$$

$$= \frac{b_{0}^{a_{0}}}{((2\pi)^{M+N}|\mathbf{S}_{0}|)^{1/2}} \iint \exp\left\{-\frac{\beta}{2}(\mathbf{t} - \mathbf{\Phi}\mathbf{w})^{\mathrm{T}}(\mathbf{t} - \mathbf{\Phi}\mathbf{w})\right\}$$

$$\exp\left\{-\frac{\beta}{2}(\mathbf{w} - \mathbf{m}_{0})^{\mathrm{T}}\mathbf{S}_{0}^{-1}(\mathbf{w} - \mathbf{m}_{0})\right\} \, d\mathbf{w}$$

$$\beta^{a_{0}-1} \beta^{N/2} \beta^{M/2} \exp(-b_{0}\beta) \, d\beta$$

$$= \frac{b_{0}^{a_{0}}}{((2\pi)^{M+N}|\mathbf{S}_{0}|)^{1/2}} \iint \exp\left\{-\frac{\beta}{2}(\mathbf{w} - \mathbf{m}_{N})^{\mathrm{T}}\mathbf{S}_{N}^{-1}(\mathbf{w} - \mathbf{m}_{N})\right\} \, d\mathbf{w}$$

$$\exp\left\{-\frac{\beta}{2}\left(\mathbf{t}^{\mathrm{T}}\mathbf{t} + \mathbf{m}_{0}^{\mathrm{T}}\mathbf{S}_{0}^{-1}\mathbf{m}_{0} - \mathbf{m}_{N}^{\mathrm{T}}\mathbf{S}_{N}^{-1}\mathbf{m}_{N}\right)\right\}$$

$$\beta^{a_{N}-1} \beta^{M/2} \exp(-b_{0}\beta) \, d\beta$$

where we have completed the square for the quadratic form in w, using

$$\mathbf{m}_{N} = \mathbf{S}_{N} \left[ \mathbf{S}_{0}^{-1} \mathbf{m}_{0} + \mathbf{\Phi}^{\mathrm{T}} \mathbf{t} \right]$$

$$\mathbf{S}_{N}^{-1} = \beta \left( \mathbf{S}_{0}^{-1} + \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \right)$$

$$a_{N} = a_{0} + \frac{N}{2}$$

$$b_{N} = b_{0} + \frac{1}{2} \left( \mathbf{m}_{0}^{\mathrm{T}} \mathbf{S}_{0}^{-1} \mathbf{m}_{0} - \mathbf{m}_{N}^{\mathrm{T}} \mathbf{S}_{N}^{-1} \mathbf{m}_{N} + \sum_{n=1}^{N} t_{n}^{2} \right).$$

Now we are ready to do the integration, first over w and then  $\beta$ , and re-arrange the

terms to obtain the desired result

$$p(\mathbf{t}) = \frac{b_0^{a_0}}{((2\pi)^{M+N}|\mathbf{S}_0|)^{1/2}} (2\pi)^{M/2} |\mathbf{S}_N|^{1/2} \int \beta^{a_N-1} \exp(-b_N \beta) \, \mathrm{d}\beta$$
$$= \frac{1}{(2\pi)^{N/2}} \frac{|\mathbf{S}_N|^{1/2}}{|\mathbf{S}_0|^{1/2}} \frac{b_0^{a_0}}{b_N^{a_N}} \frac{\Gamma(a_N)}{\Gamma(a_0)}.$$

**3.24** Substituting the r.h.s. of (3.10), (3.112) and (3.113) into (3.119), we get

$$p(\mathbf{t}) = \frac{\mathcal{N}(\mathbf{t}|\mathbf{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I}) \mathcal{N}(\mathbf{w}|\mathbf{m}_{0}, \beta^{-1}\mathbf{S}_{0}) \operatorname{Gam}(\beta|a_{0}, b_{0})}{\mathcal{N}(\mathbf{w}|\mathbf{m}_{N}, \beta^{-1}\mathbf{S}_{N}) \operatorname{Gam}(\beta|a_{N}, b_{N})}.$$
 (145)

Using the definitions of the Gaussian and Gamma distributions, we can write this as

$$\left(\frac{\beta}{2\pi}\right)^{N/2} \exp\left(-\frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi}\mathbf{w}\|^{2}\right) 
\left(\frac{\beta}{2\pi}\right)^{M/2} |\mathbf{S}_{0}|^{1/2} \exp\left(-\frac{\beta}{2} (\mathbf{w} - \mathbf{m}_{0})^{\mathrm{T}} \mathbf{S}_{0}^{-1} (\mathbf{w} - \mathbf{m}_{0})\right) 
\Gamma(a_{0})^{-1} b_{0}^{a_{0}} \beta^{a_{0}-1} \exp(-b_{0}\beta) 
\left\{\left(\frac{\beta}{2\pi}\right)^{M/2} |\mathbf{S}_{N}|^{1/2} \exp\left(-\frac{\beta}{2} (\mathbf{w} - \mathbf{m}_{N})^{\mathrm{T}} \mathbf{S}_{N}^{-1} (\mathbf{w} - \mathbf{m}_{N})\right) 
\Gamma(a_{N})^{-1} b_{N}^{a_{N}} \beta^{a_{N}-1} \exp(-b_{N}\beta)\right\}^{-1}. (146)$$

Concentrating on the factors corresponding to the denominator in (145), i.e. the fac-

tors inside  $\{...\}$ )<sup>-1</sup> in (146), we can use (135)–(138) to get

$$\mathcal{N}\left(\mathbf{w}|\mathbf{m}_{N}, \beta^{-1}\mathbf{S}_{N}\right) \operatorname{Gam}\left(\beta|a_{N}, b_{N}\right) \\
= \left(\frac{\beta}{2\pi}\right)^{M/2} |\mathbf{S}_{N}|^{1/2} \exp\left(-\frac{\beta}{2}\left(\mathbf{w}^{\mathrm{T}}\mathbf{S}_{N}^{-1}\mathbf{w} - \mathbf{w}^{\mathrm{T}}\mathbf{S}_{N}^{-1}\mathbf{m}_{N} - \mathbf{m}_{N}^{\mathrm{T}}\mathbf{S}_{N}^{-1}\mathbf{w}\right) \\
+\mathbf{m}_{N}^{\mathrm{T}}\mathbf{S}_{N}^{-1}\mathbf{m}_{N}\right) \Gamma(a_{N})^{-1}b_{N}^{a_{N}}\beta^{a_{N}-1} \exp(-b_{N}\beta) \\
= \left(\frac{\beta}{2\pi}\right)^{M/2} |\mathbf{S}_{N}|^{1/2} \exp\left(-\frac{\beta}{2}\left(\mathbf{w}^{\mathrm{T}}\mathbf{S}_{0}^{-1}\mathbf{w} + \mathbf{w}^{\mathrm{T}}\mathbf{\Phi}^{\mathrm{T}}\mathbf{\Phi}\mathbf{w} - \mathbf{w}^{\mathrm{T}}\mathbf{S}_{0}^{-1}\mathbf{m}_{0}\right) \\
-\mathbf{w}^{\mathrm{T}}\mathbf{\Phi}^{\mathrm{T}}\mathbf{t} - \mathbf{m}_{0}^{\mathrm{T}}\mathbf{S}_{N}^{-1}\mathbf{w} - \mathbf{t}^{\mathrm{T}}\mathbf{\Phi}\mathbf{w} + \mathbf{m}_{N}^{\mathrm{T}}\mathbf{S}_{N}^{-1}\mathbf{m}_{N}\right) \\
\Gamma(a_{N})^{-1}b_{N}^{a_{N}}\beta^{a_{0}+N/2-1} \\
\exp\left(-\left(b_{0} + \frac{1}{2}\left(\mathbf{m}_{0}^{\mathrm{T}}\mathbf{S}_{0}^{-1}\mathbf{m}_{0} - \mathbf{m}_{N}^{\mathrm{T}}\mathbf{S}_{N}^{-1}\mathbf{m}_{N} + \mathbf{t}^{\mathrm{T}}\mathbf{t}\right)\right)\beta\right) \\
= \left(\frac{\beta}{2\pi}\right)^{M/2} |\mathbf{S}_{N}|^{1/2} \exp\left(-\frac{\beta}{2}\left((\mathbf{w} - \mathbf{m}_{0})^{\mathrm{T}}\mathbf{S}_{0}(\mathbf{w} - \mathbf{m}_{0}) + \|\mathbf{t} - \mathbf{\Phi}\mathbf{w}\|^{2}\right)\right) \\
\Gamma(a_{N})^{-1}b_{N}^{a_{N}}\beta^{a_{N}+N/2-1} \exp(-b_{0}\beta).$$

Substituting this into (146), the exponential factors along with  $\beta^{a_0+N/2-1}(\beta/2\pi)^{M/2}$  cancel and we are left with (3.118).

# **Chapter 4** Linear Models for Classification

**4.1** Assume that the convex hulls of  $\{x_n\}$  and  $\{y_m\}$  intersect. Then there exist a point z such that

$$\mathbf{z} = \sum_{n} \alpha_n \mathbf{x}_n = \sum_{m} \beta_m \mathbf{y}_m$$

where  $\beta_m \geqslant 0$  for all m and  $\sum_m \beta_m = 1$ . If  $\{\mathbf{x}_n\}$  and  $\{\mathbf{y}_m\}$  also were to be linearly separable, we would have that

$$\widehat{\mathbf{w}}^{\mathrm{T}}\mathbf{z} + w_0 = \sum_n \alpha_n \widehat{\mathbf{w}}^{\mathrm{T}}\mathbf{x}_n + w_0 = \sum_n \alpha_n ($$

since  $\widehat{\mathbf{w}}^{\mathrm{T}}\mathbf{x}_n + w_0 > 0$  and the  $\{\alpha_n\}$  are all non-negative and sum to 1, but by the corresponding argument

$$\widehat{\mathbf{w}}^{\mathrm{T}}\mathbf{z} + w_0 = \sum_{m} \beta_m \widehat{\mathbf{w}}^{\mathrm{T}}\mathbf{y}_m + w_0 = \sum_{m} \beta_m (\widehat{\mathbf{w}}^{\mathrm{T}}\mathbf{y}_m + w_0) < 0,$$

which is a contradiction and hence  $\{x_n\}$  and  $\{y_m\}$  cannot be linearly separable if their convex hulls intersect.

If we instead assume that  $\{\mathbf{x}_n\}$  and  $\{\mathbf{y}_m\}$  are linearly separable and consider a point  $\mathbf{z}$  in the intersection of their convex hulls, the same contradiction arise. Thus no such point can exist and the intersection of the convex hulls of  $\{\mathbf{x}_n\}$  and  $\{\mathbf{y}_m\}$  must be empty.

**4.2** For the purpose of this exercise, we make the contribution of the bias weights explicit in (4.15), giving

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \operatorname{Tr} \left\{ (\mathbf{X} \mathbf{W} + \mathbf{1} \mathbf{w}_0^{\mathrm{T}} - \mathbf{T})^{\mathrm{T}} (\mathbf{X} \mathbf{W} + \mathbf{1} \mathbf{w}_0^{\mathrm{T}} - \mathbf{T}) \right\},$$
(147)

where  $\mathbf{w}_0$  is the column vector of bias weights (the top row of  $\widetilde{\mathbf{W}}$  transposed) and  $\mathbf{1}$  is a column vector of N ones.

We can take the derivative of (147) w.r.t.  $\mathbf{w}_0$ , giving

$$2N\mathbf{w}_0 + 2(\mathbf{XW} - \mathbf{T})^{\mathrm{T}}\mathbf{1}.$$

Setting this to zero, and solving for  $w_0$ , we obtain

$$\mathbf{w}_0 = \bar{\mathbf{t}} - \mathbf{W}^{\mathrm{T}} \bar{\mathbf{x}} \tag{148}$$

where

$$\bar{\mathbf{t}} = \frac{1}{N}\mathbf{T}^{\mathrm{T}}\mathbf{1} \quad \text{and} \quad \bar{\mathbf{x}} = \frac{1}{N}\mathbf{X}^{\mathrm{T}}\mathbf{1}.$$

If we substitute (148) into (147), we get

$$E_D(\mathbf{W}) = \frac{1}{2} \operatorname{Tr} \left\{ (\mathbf{X} \mathbf{W} + \overline{\mathbf{T}} - \overline{\mathbf{X}} \mathbf{W} - \mathbf{T})^{\mathrm{T}} (\mathbf{X} \mathbf{W} + \overline{\mathbf{T}} - \overline{\mathbf{X}} \mathbf{W} - \mathbf{T}) \right\},\,$$

where

$$\overline{T} = 1\overline{t}^{\mathrm{T}}$$
 and  $\overline{X} = 1\overline{x}^{\mathrm{T}}$ .

Setting the derivative of this w.r.t. W to zero we get

$$\mathbf{W} = (\widehat{\mathbf{X}}^{\mathrm{T}}\widehat{\mathbf{X}})^{-1}\widehat{\mathbf{X}}^{\mathrm{T}}\widehat{\mathbf{T}} = \widehat{\mathbf{X}}^{\dagger}\widehat{\mathbf{T}},$$

where we have defined  $\widehat{X} = X - \overline{X}$  and  $\widehat{T} = T - \overline{T}$ .

Now consider the prediction for a new input vector  $\mathbf{x}^*$ ,

$$\mathbf{y}(\mathbf{x}^{\star}) = \mathbf{W}^{\mathrm{T}}\mathbf{x}^{\star} + \mathbf{w}_{0}$$

$$= \mathbf{W}^{\mathrm{T}}\mathbf{x}^{\star} + \bar{\mathbf{t}} - \mathbf{W}^{\mathrm{T}}\bar{\mathbf{x}}$$

$$= \bar{\mathbf{t}} - \widehat{\mathbf{T}}^{\mathrm{T}} \left(\widehat{\mathbf{X}}^{\dagger}\right)^{\mathrm{T}} (\mathbf{x}^{\star} - \bar{\mathbf{x}}).$$
(149)

If we apply (4.157) to  $\bar{\mathbf{t}}$ , we get

$$\mathbf{a}^{\mathrm{T}}\bar{\mathbf{t}} = \frac{1}{N}\mathbf{a}^{\mathrm{T}}\mathbf{T}^{\mathrm{T}}\mathbf{1} = -b.$$

Therefore, applying (4.157) to (149), we obtain

$$\mathbf{a}^{\mathrm{T}}\mathbf{y}(\mathbf{x}^{\star}) = \mathbf{a}^{\mathrm{T}}\bar{\mathbf{t}} + \mathbf{a}^{\mathrm{T}}\widehat{\mathbf{T}}^{\mathrm{T}}\left(\widehat{\mathbf{X}}^{\dagger}\right)^{\mathrm{T}}\left(\mathbf{x}^{\star} - \bar{\mathbf{x}}\right)$$
$$= \mathbf{a}^{\mathrm{T}}\bar{\mathbf{t}} = -b,$$

since 
$$\mathbf{a}^{\mathrm{T}} \widehat{\mathbf{T}}^{\mathrm{T}} = \mathbf{a}^{\mathrm{T}} (\mathbf{T} - \overline{\mathbf{T}})^{\mathrm{T}} = b(\mathbf{1} - \mathbf{1})^{\mathrm{T}} = \mathbf{0}^{\mathrm{T}}$$
.

**4.3** When we consider several simultaneous constraints, (4.157) becomes

$$\mathbf{At}_n + \mathbf{b} = \mathbf{0},\tag{150}$$

where **A** is a matrix and **b** is a column vector such that each row of **A** and element of **b** correspond to one linear constraint.

If we apply (150) to (149), we obtain

$$egin{array}{lll} \mathbf{A}\mathbf{y}(\mathbf{x}^{\star}) &=& \mathbf{A}ar{\mathbf{t}} - \mathbf{A}\widehat{\mathbf{T}}^{\mathrm{T}} \left(\widehat{\mathbf{X}}^{\dagger}\right)^{\mathrm{T}} \left(\mathbf{x}^{\star} - ar{\mathbf{x}}\right) \ &=& \mathbf{A}ar{\mathbf{t}} = -\mathbf{b}. \end{array}$$

since 
$$\mathbf{A}\widehat{\mathbf{T}}^{\mathrm{T}} = \mathbf{A}(\mathbf{T} - \overline{\mathbf{T}})^{\mathrm{T}} = \mathbf{b}\mathbf{1}^{\mathrm{T}} - \mathbf{b}\mathbf{1}^{\mathrm{T}} = \mathbf{0}^{\mathrm{T}}$$
. Thus  $\mathbf{A}\mathbf{y}(\mathbf{x}^{\star}) + \mathbf{b} = \mathbf{0}$ .

**4.4 NOTE**: In the 1<sup>st</sup> printing of PRML, the text of the exercise refers equation (4.23) where it should refer to (4.22).

From (4.22) we can construct the Lagrangian function

$$L = \mathbf{w}^{\mathrm{T}}(\mathbf{m}_2 - \mathbf{m}_1) + \lambda \left( \mathbf{w}^{\mathrm{T}} \mathbf{w} - 1 \right).$$

Taking the gradient of L we obtain

$$\nabla L = \mathbf{m}_2 - \mathbf{m}_1 + 2\lambda \mathbf{w} \tag{151}$$

and setting this gradient to zero gives

$$\mathbf{w} = -\frac{1}{2\lambda}(\mathbf{m}_2 - \mathbf{m}_1)$$

form which it follows that  $\mathbf{w} \propto \mathbf{m}_2 - \mathbf{m}_1$ .

**4.5** Starting with the numerator on the r.h.s. of (4.25), we can use (4.23) and (4.27) to rewrite it as follows:

$$(m_2 - m_1)^2 = (\mathbf{w}^{\mathrm{T}}(\mathbf{m}_2 - \mathbf{m}_1))^2$$
$$= \mathbf{w}^{\mathrm{T}}(\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^{\mathrm{T}}\mathbf{w}$$
$$= \mathbf{w}^{\mathrm{T}}\mathbf{S}_{\mathrm{B}}\mathbf{w}. \tag{152}$$

Similarly, we can use (4.20), (4.23), (4.24), and (4.28) to rewrite the denominator of the r.h.s. of (4.25):

$$s_1^2 + s_2^2 = \sum_{n \in \mathcal{C}_1} (y_n - m_1)^2 + \sum_{k \in \mathcal{C}_2} (y_k - m_2)^2$$

$$= \sum_{n \in \mathcal{C}_1} (\mathbf{w}^{\mathrm{T}} (\mathbf{x}_n - \mathbf{m}_1))^2 + \sum_{k \in \mathcal{C}_2} (\mathbf{w}^{\mathrm{T}} (\mathbf{x}_k - \mathbf{m}_2))^2$$

$$= \sum_{n \in \mathcal{C}_1} \mathbf{w}^{\mathrm{T}} (\mathbf{x}_n - \mathbf{m}_1) (\mathbf{x}_n - \mathbf{m}_1)^{\mathrm{T}} \mathbf{w}$$

$$+ \sum_{k \in \mathcal{C}_2} \mathbf{w}^{\mathrm{T}} (\mathbf{x}_k - \mathbf{m}_2) (\mathbf{x}_k - \mathbf{m}_2)^{\mathrm{T}} \mathbf{w}$$

$$= \mathbf{w}^{\mathrm{T}} \mathbf{S}_{\mathrm{W}} \mathbf{w}. \tag{153}$$

Substituting (152) and (153) in (4.25) we obtain (4.26).

**4.6** Using (4.21) and (4.34) along with the chosen target coding scheme, we can re-write the l.h.s. of (4.33) as follows:

$$\sum_{n=1}^{N} (\mathbf{w}^{\mathrm{T}} \mathbf{x}_{n} - w_{0} - t_{n}) \mathbf{x}_{n} = \sum_{n=1}^{N} (\mathbf{w}^{\mathrm{T}} \mathbf{x}_{n} - \mathbf{w}^{\mathrm{T}} \mathbf{m} - t_{n}) \mathbf{x}_{n}$$

$$= \sum_{n=1}^{N} \{ (\mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - \mathbf{x}_{n} \mathbf{m}^{\mathrm{T}}) \mathbf{w} - \mathbf{x}_{n} t_{n} \}$$

$$= \sum_{n \in \mathcal{C}_{1}} \{ (\mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - \mathbf{x}_{n} \mathbf{m}^{\mathrm{T}}) \mathbf{w} - \mathbf{x}_{n} t_{n} \}$$

$$= \left( \sum_{n \in \mathcal{C}_{1}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - N_{1} \mathbf{m}_{1} \mathbf{m}^{\mathrm{T}} \right) \mathbf{w} - N_{1} \mathbf{m}_{1} \frac{N}{N_{1}}$$

$$= \left( \sum_{n \in \mathcal{C}_{2}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - N_{2} \mathbf{m}_{2} \mathbf{m}^{\mathrm{T}} \right) \mathbf{w} + N_{2} \mathbf{m}_{2} \frac{N}{N_{2}}$$

$$= \sum_{n \in \mathcal{C}_{1}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - N_{2} \mathbf{m}_{2} \mathbf{m}^{\mathrm{T}} \right) \mathbf{w} + N_{2} \mathbf{m}_{2} \frac{N}{N_{2}}$$

$$= \sum_{n \in \mathcal{C}_{2}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - N_{2} \mathbf{m}_{2} \mathbf{m}^{\mathrm{T}} \right) \mathbf{w} + N_{2} \mathbf{m}_{2} \frac{N}{N_{2}}$$

$$= \sum_{n \in \mathcal{C}_{2}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} + \sum_{n \in \mathcal{C}_{2}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - \left( N_{1} \mathbf{m}_{1} + N_{2} \mathbf{m}_{2} \right) \mathbf{m}^{\mathrm{T}} \right) \mathbf{w}$$

$$= \sum_{n \in \mathcal{C}_{1}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - N_{1} \mathbf{m}_{1} \mathbf{m}^{\mathrm{T}} \right) \mathbf{w} - \mathbf{x}_{n} t_{n} \mathbf{x}_{n}^{\mathrm{T}}$$

$$= \sum_{n \in \mathcal{C}_{2}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - N_{1} \mathbf{m}_{1} \mathbf{m}^{\mathrm{T}} \right) \mathbf{w} - \mathbf{x}_{n} t_{n}^{\mathrm{T}} \mathbf{w}$$

$$= \sum_{n \in \mathcal{C}_{2}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - N_{1} \mathbf{m}_{1} \mathbf{m}^{\mathrm{T}} \mathbf{w} - \mathbf{x}_{n} t_{n}^{\mathrm{T}} \mathbf{w}$$

$$= \sum_{n \in \mathcal{C}_{2}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - N_{1} \mathbf{m}_{1} \mathbf{m}^{\mathrm{T}} \mathbf{w} - \mathbf{x}_{n} t_{n}^{\mathrm{T}} \mathbf{w}$$

$$= \sum_{n \in \mathcal{C}_{2}} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} \mathbf$$

We then use the identity

$$\sum_{i \in \mathcal{C}_k} (\mathbf{x}_i - \mathbf{m}_k) (\mathbf{x}_i - \mathbf{m}_k)^{\mathrm{T}} = \sum_{i \in \mathcal{C}_k} (\mathbf{x}_i \mathbf{x}_i^{\mathrm{T}} - \mathbf{x}_i \mathbf{m}_k^{\mathrm{T}} - \mathbf{m}_k \mathbf{x}_i^{\mathrm{T}} + \mathbf{m}_k \mathbf{m}_k^{\mathrm{T}})$$
$$= \sum_{i \in \mathcal{C}_k} \mathbf{x}_i \mathbf{x}_i^{\mathrm{T}} - N_k \mathbf{m}_k \mathbf{m}_k^{\mathrm{T}}$$

together with (4.28) and (4.36) to rewrite (154) as

$$\begin{split} &\left(\mathbf{S}_{\mathrm{W}} + N_{1}\mathbf{m}_{1}\mathbf{m}_{1}^{\mathrm{T}} + N_{2}\mathbf{m}_{2}\mathbf{m}_{2}^{\mathrm{T}} \right. \\ &\left. - (N_{1}\mathbf{m}_{1} + N_{2}\mathbf{m}_{2})\frac{1}{N}(N_{1}\mathbf{m}_{1} + N_{2}\mathbf{m}_{2})\right)\mathbf{w} - N(\mathbf{m}_{1} - \mathbf{m}_{2}) \\ &= \left(\mathbf{S}_{\mathrm{W}} + \left(N_{1} - \frac{N_{1}^{2}}{N}\right)\mathbf{m}_{1}\mathbf{m}_{1}^{\mathrm{T}} - \frac{N_{1}N_{2}}{N}(\mathbf{m}_{1}\mathbf{m}_{2}^{\mathrm{T}} + \mathbf{m}_{2}\mathbf{m}_{1}) \right. \\ &\left. + \left(N_{2} - \frac{N_{2}^{2}}{N}\right)\mathbf{m}_{2}\mathbf{m}_{2}^{\mathrm{T}}\right)\mathbf{w} - N(\mathbf{m}_{1} - \mathbf{m}_{2}) \\ &= \left(\mathbf{S}_{\mathrm{W}} + \frac{(N_{1} + N_{2})N_{1} - N_{1}^{2}}{N}\mathbf{m}_{1}\mathbf{m}_{1}^{\mathrm{T}} - \frac{N_{1}N_{2}}{N}(\mathbf{m}_{1}\mathbf{m}_{2}^{\mathrm{T}} + \mathbf{m}_{2}\mathbf{m}_{1}) \right. \\ &\left. + \frac{(N_{1} + N_{2})N_{2} - N_{2}^{2}}{N}\mathbf{m}_{2}\mathbf{m}_{2}^{\mathrm{T}}\right)\mathbf{w} - N(\mathbf{m}_{1} - \mathbf{m}_{2}) \\ &= \left(\mathbf{S}_{\mathrm{W}} + \frac{N_{2}N_{1}}{N}\left(\mathbf{m}_{1}\mathbf{m}_{1}^{\mathrm{T}} - \mathbf{m}_{1}\mathbf{m}_{2}^{\mathrm{T}} - \mathbf{m}_{2}\mathbf{m}_{1} + \mathbf{m}_{2}\mathbf{m}_{2}^{\mathrm{T}}\right)\right)\mathbf{w} \\ &\left. - N(\mathbf{m}_{1} - \mathbf{m}_{2}) \\ &= \left(\mathbf{S}_{\mathrm{W}} + \frac{N_{2}N_{1}}{N}\mathbf{S}_{\mathrm{B}}\right)\mathbf{w} - N(\mathbf{m}_{1} - \mathbf{m}_{2}), \end{split}$$

where in the last line we also made use of (4.27). From (4.33), this must equal zero, and hence we obtain (4.37).

### **4.7** From (4.59) we have

$$1 - \sigma(a) = 1 - \frac{1}{1 + e^{-a}} = \frac{1 + e^{-a} - 1}{1 + e^{-a}}$$
$$= \frac{e^{-a}}{1 + e^{-a}} = \frac{1}{e^a + 1} = \sigma(-a).$$

The inverse of the logistic sigmoid is easily found as follows

$$y = \sigma(a) = \frac{1}{1 + e^{-a}}$$

$$\Rightarrow \frac{1}{y} - 1 = e^{-a}$$

$$\Rightarrow \ln\left\{\frac{1 - y}{y}\right\} = -a$$

$$\Rightarrow \ln\left\{\frac{y}{1 - y}\right\} = a = \sigma^{-1}(y).$$

**4.8** Substituting (4.64) into (4.58), we see that the normalizing constants cancel and we are left with

$$a = \ln \frac{\exp\left(-\frac{1}{2}\left(\mathbf{x} - \boldsymbol{\mu}_{1}\right)^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \boldsymbol{\mu}_{1}\right)\right) p(\mathcal{C}_{1})}{\exp\left(-\frac{1}{2}\left(\mathbf{x} - \boldsymbol{\mu}_{2}\right)^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \boldsymbol{\mu}_{2}\right)\right) p(\mathcal{C}_{2})}$$

$$= -\frac{1}{2} \left(\mathbf{x} \boldsymbol{\Sigma}^{\mathrm{T}} \mathbf{x} - \mathbf{x} \boldsymbol{\Sigma} \boldsymbol{\mu}_{1} - \boldsymbol{\mu}_{1}^{\mathrm{T}} \boldsymbol{\Sigma} \mathbf{x} + \boldsymbol{\mu}_{1}^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{\mu}_{1}$$

$$-\mathbf{x} \boldsymbol{\Sigma}^{\mathrm{T}} \mathbf{x} + \mathbf{x} \boldsymbol{\Sigma} \boldsymbol{\mu}_{2} + \boldsymbol{\mu}_{2}^{\mathrm{T}} \boldsymbol{\Sigma} \mathbf{x} - \boldsymbol{\mu}_{2}^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{\mu}_{2}\right) + \ln \frac{p(\mathcal{C}_{1})}{p(\mathcal{C}_{2})}$$

$$= (\boldsymbol{\mu}_{1} - \boldsymbol{\mu}_{2})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \left(\boldsymbol{\mu}_{1}^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{1} - \boldsymbol{\mu}_{2}^{\mathrm{T}} \boldsymbol{\Sigma} \boldsymbol{\mu}_{2}\right) + \ln \frac{p(\mathcal{C}_{1})}{p(\mathcal{C}_{2})}.$$

Substituting this into the rightmost form of (4.57) we obtain (4.65), with w and  $w_0$  given by (4.66) and (4.67), respectively.

**4.9** The likelihood function is given by

$$p(\{\phi_n, \mathbf{t}_n\} | \{\pi_k\}) = \prod_{n=1}^{N} \prod_{k=1}^{K} \{p(\phi_n | \mathcal{C}_k) \pi_k\}^{t_{nk}}$$

and taking the logarithm, we obtain

$$\ln p(\{\phi_n, \mathbf{t}_n\} | \{\pi_k\}) = \sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \{\ln p(\phi_n | \mathcal{C}_k) + \ln \pi_k\}.$$
 (155)

In order to maximize the log likelihood with respect to  $\pi_k$  we need to preserve the constraint  $\sum_k \pi_k = 1$ . This can be done by introducing a Lagrange multiplier  $\lambda$  and maximizing

$$\ln p\left(\{\phi_n, \mathbf{t}_n\} | \{\pi_k\}\right) + \lambda \left(\sum_{k=1}^K \pi_k - 1\right).$$

Setting the derivative with respect to  $\pi_k$  equal to zero, we obtain

$$\sum_{n=1}^{N} \frac{t_{nk}}{\pi_k} + \lambda = 0.$$

Re-arranging then gives

$$-\pi_k \lambda = \sum_{n=1}^N t_{nk} = N_k. \tag{156}$$

Summing both sides over k we find that  $\lambda = -N$ , and using this to eliminate  $\lambda$  we obtain (4.159).

**4.10** If we substitute (4.160) into (155) and then use the definition of the multivariate Gaussian, (2.43), we obtain

$$\ln p\left(\left\{\phi_{n}, \mathbf{t}_{n}\right\} \middle| \left\{\pi_{k}\right\}\right) = -\frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \left\{\ln |\mathbf{\Sigma}| + (\boldsymbol{\phi}_{n} - \boldsymbol{\mu}_{k})^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\boldsymbol{\phi} - \boldsymbol{\mu})\right\}, \quad (157)$$

where we have dropped terms independent of  $\{\mu_k\}$  and  $\Sigma$ .

Setting the derivative of the r.h.s. of (157) w.r.t.  $\mu_k$ , obtained by using (C.19), to zero, we get

$$\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \boldsymbol{\Sigma}^{-1} (\boldsymbol{\phi}_n - \boldsymbol{\mu}_k) = 0.$$

Making use of (156), we can re-arrange this to obtain (4.161).

Rewriting the r.h.s. of (157) as

$$-\frac{1}{2}b\sum_{n=1}^{N}\sum_{k=1}^{K}t_{nk}\left\{\ln|\mathbf{\Sigma}|+\operatorname{Tr}\left[\mathbf{\Sigma}^{-1}(\boldsymbol{\phi}_{n}-\boldsymbol{\mu}_{k})(\boldsymbol{\phi}-\boldsymbol{\mu}_{k})^{\mathrm{T}}\right]\right\},$$

we can use (C.24) and (C.28) to calculate the derivative w.r.t.  $\Sigma^{-1}$ . Setting this to zero we obtain

$$\frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{T} t_{nk} \left\{ \mathbf{\Sigma} - (\boldsymbol{\phi}_{n} - \boldsymbol{\mu}_{n})(\boldsymbol{\phi}_{n} - \boldsymbol{\mu}_{k})^{\mathrm{T}} \right\} = 0.$$

Again making use of (156), we can re-arrange this to obtain (4.162), with  $S_k$  given by (4.163).

Note that, as in Exercise 2.34, we do not enforce that  $\Sigma$  should be symmetric, but simply note that the solution is automatically symmetric.

85

**4.11** The generative model for  $\phi$  corresponding to the chosen coding scheme is given by

$$p\left(\boldsymbol{\phi}\mid\mathcal{C}_{k}\right)=\prod_{m=1}^{M}p\left(\boldsymbol{\phi}_{m}\mid\mathcal{C}_{k}\right)$$

where

$$p\left(\boldsymbol{\phi}_{m} \mid \mathcal{C}_{k}\right) = \prod_{l=1}^{L} \mu_{kml}^{\phi_{ml}},$$

where in turn  $\{\mu_{kml}\}$  are the parameters of the multinomial models for  $\phi$ . Substituting this into (4.63) we see that

$$a_{k} = \ln p \left(\phi \mid \mathcal{C}_{k}\right) p \left(\mathcal{C}_{k}\right)$$

$$= \ln p \left(\mathcal{C}_{k}\right) + \sum_{m=1}^{M} \ln p \left(\phi_{m} \mid \mathcal{C}_{k}\right)$$

$$= \ln p \left(\mathcal{C}_{k}\right) + \sum_{m=1}^{M} \sum_{l=1}^{L} \phi_{ml} \ln \mu_{kml},$$

which is linear in  $\phi_{ml}$ .

**4.12** Differentiating (4.59) we obtain

$$\frac{d\sigma}{da} = \frac{e^{-a}}{(1+e^{-a})^2} 
= \sigma(a) \left\{ \frac{e^{-a}}{1+e^{-a}} \right\} 
= \sigma(a) \left\{ \frac{1+e^{-a}}{1+e^{-a}} - \frac{1}{1+e^{-a}} \right\} 
= \sigma(a) (1-\sigma(a)).$$

**4.13** We start by computing the derivative of (4.90) w.r.t.  $y_n$ 

$$\frac{\partial E}{\partial y_n} = \frac{1 - t_n}{1 - y_n} - \frac{t_n}{y_n} \tag{158}$$

$$= \frac{y_n (1 - t_n) - t_n (1 - y_n)}{y_n (1 - y_n)}$$

$$= \frac{y_n - y_n t_n - t_n + y_n t_n}{y_n (1 - y_n)}$$

$$= \frac{y_n - t_n}{y_n (1 - y_n)}.$$
(159)

From (4.88), we see that

$$\frac{\partial y_n}{\partial a_n} = \frac{\partial \sigma(a_n)}{\partial a_n} = \sigma(a_n) \left( 1 - \sigma(a_n) \right) = y_n (1 - y_n). \tag{161}$$

Finally, we have

$$\nabla a_n = \phi_n \tag{162}$$

where  $\nabla$  denotes the gradient with respect to w. Combining (160), (161) and (162) using the chain rule, we obtain

$$\nabla E = \sum_{n=1}^{N} \frac{\partial E}{\partial y_n} \frac{\partial y_n}{\partial a_n} \nabla a_n$$
$$= \sum_{n=1}^{N} (y_n - t_n) \phi_n$$

as required.

**4.14** If the data set is linearly separable, any decision boundary separating the two classes will have the property

$$\mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}_n \left\{ \begin{array}{l} \geqslant 0 & \text{if } t_n = 1, \\ < 0 & \text{otherwise.} \end{array} \right.$$

Moreover, from (4.90) we see that the negative log-likelihood will be minimized (i.e., the likelihood maximized) when  $y_n = \sigma(\mathbf{w}_T \phi_n) = t_n$  for all n. This will be the case when the sigmoid function is saturated, which occurs when its argument,  $\mathbf{w}^T \phi$ , goes to  $\pm \infty$ , i.e., when the magnitude of  $\mathbf{w}$  goes to infinity.

**4.15 NOTE**: In PRML, "concave" should be "convex" on the last line of the exercise.

Assuming that the argument to the sigmoid function (4.87) is finite, the diagonal elements of  ${\bf R}$  will be strictly positive. Then

$$\mathbf{v}^{\mathrm{T}}\mathbf{\Phi}^{\mathrm{T}}\mathbf{R}\mathbf{\Phi}\mathbf{v} = \left(\mathbf{v}^{\mathrm{T}}\mathbf{\Phi}^{\mathrm{T}}\mathbf{R}^{1/2}\right)\left(\mathbf{R}^{1/2}\mathbf{\Phi}\mathbf{v}\right) = \left\|\mathbf{R}^{1/2}\mathbf{\Phi}\mathbf{v}\right\|^{2} > 0$$

where  $\mathbf{R}^{1/2}$  is a diagonal matrix with elements  $(y_n(1-y_n))^{1/2}$ , and thus  $\mathbf{\Phi}^T \mathbf{R} \mathbf{\Phi}$  is positive definite.

Now consider a Taylor expansion of  $E(\mathbf{w})$  around a minima,  $\mathbf{w}^*$ ,

$$E(\mathbf{w}) = E(\mathbf{w}^*) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^{\mathrm{T}} \mathbf{H} (\mathbf{w} - \mathbf{w}^*)$$

where the linear term has vanished since  $\mathbf{w}^{\star}$  is a minimum. Now let

$$\mathbf{w} = \mathbf{w}^* + \lambda \mathbf{v}$$

where  $\mathbf{v}$  is an arbitrary, non-zero vector in the weight space and consider

$$\frac{\partial^2 E}{\partial \lambda^2} = \mathbf{v}^{\mathrm{T}} \mathbf{H} \mathbf{v} > 0.$$

This shows that  $E(\mathbf{w})$  is convex. Moreover, at the minimum of  $E(\mathbf{w})$ ,

$$\mathbf{H}\left(\mathbf{w} - \mathbf{w}^{\star}\right) = 0$$

and since  ${\bf H}$  is positive definite,  ${\bf H}^{-1}$  exists and  ${\bf w}={\bf w}^{\star}$  must be the unique minimum.

**4.16** If the values of the  $\{t_n\}$  were known then each data point for which  $t_n=1$  would contribute  $p(t_n=1|\phi(\mathbf{x}_n))$  to the log likelihood, and each point for which  $t_n=0$  would contribute  $1-p(t_n=1|\phi(\mathbf{x}_n))$  to the log likelihood. A data point whose probability of having  $t_n=1$  is given by  $\pi_n$  will therefore contribute

$$\pi_n p(t_n = 1 | \phi(\mathbf{x}_n)) + (1 - \pi_n)(1 - p(t_n = 1 | \phi(\mathbf{x}_n)))$$

and so the overall log likelihood for the data set is given by

$$\sum_{n=1}^{N} \pi_n \ln p \left( t_n = 1 \mid \phi(\mathbf{x}_n) \right) + (1 - \pi_n) \ln \left( 1 - p \left( t_n = 1 \mid \phi(\mathbf{x}_n) \right) \right). \tag{163}$$

This can also be viewed from a sampling perspective by imagining sampling the value of each  $t_n$  some number M times, with probability of  $t_n = 1$  given by  $\pi_n$ , and then constructing the likelihood function for this expanded data set, and dividing by M. In the limit  $M \to \infty$  we recover (163).

**4.17** From (4.104) we have

$$\frac{\partial y_k}{\partial a_k} = \frac{e^{a_k}}{\sum_i e^{a_i}} - \left(\frac{e^{a_k}}{\sum_i e^{a_i}}\right)^2 = y_k (1 - y_k),$$

$$\frac{\partial y_k}{\partial a_j} = -\frac{e^{a_k} e^{a_j}}{\left(\sum_i e^{a_i}\right)^2} = -y_k y_j, \qquad j \neq k.$$

Combining these results we obtain (4.106).

**4.18 NOTE**: In the 1<sup>st</sup> printing of PRML, the text of the exercise refers equation (4.91) where it should refer to (4.106).

From (4.108) we have

$$\frac{\partial E}{\partial y_{nk}} = -\frac{t_{nk}}{y_{nk}}.$$

If we combine this with (4.106) using the chain rule, we get

$$\frac{\partial E}{\partial a_{nj}} = \sum_{k=1}^{K} \frac{\partial E}{\partial y_{nk}} \frac{\partial y_{nk}}{\partial a_{nj}}$$

$$= -\sum_{k=1}^{K} \frac{t_{nk}}{y_{nk}} y_{nk} (I_{kj} - y_{nj})$$

$$= y_{nj} - t_{nj},$$

where we have used that  $\forall n : \sum_{k} t_{nk} = 1$ .

If we combine this with (162), again using the chain rule, we obtain (4.109).

**4.19** Using the cross-entropy error function (4.90), and following Exercise 4.13, we have

$$\frac{\partial E}{\partial y_n} = \frac{y_n - t_n}{y_n(1 - y_n)}. (164)$$

Also

$$\nabla a_n = \phi_n. \tag{165}$$

From (4.115) and (4.116) we have

$$\frac{\partial y_n}{\partial a_n} = \frac{\partial \Phi(a_n)}{\partial a_n} = \frac{1}{\sqrt{2\pi}} e^{-a_n^2}.$$
 (166)

Combining (164), (165) and (166), we get

$$\nabla E = \sum_{n=1}^{N} \frac{\partial E}{\partial y_n} \frac{\partial y_n}{\partial a_n} \nabla a_n = \sum_{n=1}^{N} \frac{y_n - t_n}{y_n (1 - y_n)} \frac{1}{\sqrt{2\pi}} e^{-a_n^2} \phi_n.$$
 (167)

In order to find the expression for the Hessian, it is is convenient to first determine

$$\frac{\partial}{\partial y_n} \frac{y_n - t_n}{y_n (1 - y_n)} = \frac{y_n (1 - y_n)}{y_n^2 (1 - y_n)^2} - \frac{(y_n - t_n)(1 - 2y_n)}{y_n^2 (1 - y_n)^2} 
= \frac{y_n^2 + t_n - 2y_n t_n}{y_n^2 (1 - y_n)^2}.$$
(168)

Then using (165)–(168) we have

$$\nabla \nabla E = \sum_{n=1}^{N} \left\{ \frac{\partial}{\partial y_n} \left[ \frac{y_n - t_n}{y_n (1 - y_n)} \right] \frac{1}{\sqrt{2\pi}} e^{-a_n^2} \phi_n \nabla y_n \right.$$

$$\left. + \frac{y_n - t_n}{y_n (1 - y_n)} \frac{1}{\sqrt{2\pi}} e^{-a_n^2} (-2a_n) \phi_n \nabla a_n \right\}$$

$$= \sum_{n=1}^{N} \left( \frac{y_n^2 + t_n - 2y_n t_n}{y_n (1 - y_n)} \frac{1}{\sqrt{2\pi}} e^{-a_n^2} - 2a_n (y_n - t_n) \right) \frac{e^{-2a_n^2} \phi_n \phi_n^{\mathrm{T}}}{\sqrt{2\pi} y_n (1 - y_n)}.$$

**4.20 NOTE**: In the 1<sup>st</sup> printing of PRML, equation (4.110) contains an incorrect leading minus sign ('–') on the right hand side.

We first write out the components of the  $MK \times MK$  Hessian matrix in the form

$$\frac{\partial^2 E}{\partial w_{ki} \partial w_{jl}} = \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi_{ni} \phi_{nl}.$$

To keep the notation uncluttered, consider just one term in the summation over n and show that this is positive semi-definite. The sum over n will then also be positive semi-definite. Consider an arbitrary vector of dimension MK with elements  $u_{ki}$ . Then

$$\mathbf{u}^{\mathrm{T}}\mathbf{H}\mathbf{u} = \sum_{i,j,k,l} u_{ki} y_k (I_{kj} - y_j) \phi_i \phi_l u_{jl}$$

$$= \sum_{j,k} b_j y_k (I_{kj} - y_j) b_k$$

$$= \sum_k y_k b_k^2 - \left(\sum_k b_k y_k\right)^2$$

where

$$b_k = \sum_i u_{ki} \phi_{ni}.$$

We now note that the quantities  $y_k$  satisfy  $0 \leqslant y_k \leqslant 1$  and  $\sum_k y_k = 1$ . Furthermore, the function  $f(b) = b^2$  is a concave function. We can therefore apply Jensen's inequality to give

$$\sum_{k} y_k b_k^2 = \sum_{k} y_k f(b_k) \geqslant f\left(\sum_{k} y_k b_k\right) = \left(\sum_{k} y_k b_k\right)^2$$

and hence

$$\mathbf{u}^{\mathrm{T}}\mathbf{H}\mathbf{u} \geqslant 0.$$

Note that the equality will never arise for finite values of  $a_k$  where  $a_k$  is the set of arguments to the softmax function. However, the Hessian can be positive *semi*-definite since the basis vectors  $\phi_{ni}$  could be such as to have zero dot product for a linear subspace of vectors  $u_{ki}$ . In this case the minimum of the error function would comprise a continuum of solutions all having the same value of the error function.

**4.21 NOTE**: In PRML, (4.116) should read

$$\Phi(a) = \frac{1}{2} \left\{ 1 + \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \right\}.$$

Note that  $\Phi$  should be  $\Phi$  (i.e. not bold) on the l.h.s.

We consider the two cases where  $a \ge 0$  and a < 0 separately. In the first case, we can use (2.42) to rewrite (4.114) as

$$\Phi(a) = \int_{-\infty}^{0} \mathcal{N}(\theta|0,1) d\theta + \int_{0}^{a} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\theta^{2}}{2}\right) d\theta$$
$$= \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \int_{0}^{a/\sqrt{2}} \exp\left(-u^{2}\right) \sqrt{2} du$$
$$= \frac{1}{2} \left\{ 1 + \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \right\},$$

where, in the last line, we have used (4.115).

When a < 0, the symmetry of the Gaussian distribution gives

$$\Phi(a) = 1 - \Phi(-a).$$

Combining this with the above result, we get

$$\Phi(a) = 1 - \frac{1}{2} \left\{ 1 + \operatorname{erf}\left(-\frac{a}{\sqrt{2}}\right) \right\}$$
$$= \frac{1}{2} \left\{ 1 + \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \right\},$$

where we have used the fact that the erf function is is anti-symmetric, i.e., erf(-a) = -erf(a).

**4.22** Starting from (4.136), using (4.135), we have

$$p(\mathcal{D}) = \int p(\mathcal{D} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

$$\simeq p(\mathcal{D} \mid \boldsymbol{\theta}_{MAP}) p(\boldsymbol{\theta}_{MAP})$$

$$\int \exp\left(-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP}) \mathbf{A}^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})\right) d\boldsymbol{\theta}$$

$$= p(\mathcal{D} \mid \boldsymbol{\theta}_{MAP}) p(\boldsymbol{\theta}_{MAP}) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}},$$

where  $\bf A$  is given by (4.138). Taking the logarithm of this yields (4.137).

**4.23 NOTE**: In the  $1^{st}$  printing of PRML, the text of the exercise contains a typographical error. Following the equation, it should say that **H** is the matrix of second derivatives of the *negative* log likelihood.

The BIC approximation can be viewed as a large N approximation to the log model evidence. From (4.138), we have

$$\mathbf{A} = -\nabla \nabla \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) p(\boldsymbol{\theta}_{\text{MAP}})$$
$$= \mathbf{H} - \nabla \nabla \ln p(\boldsymbol{\theta}_{\text{MAP}})$$

91

and if  $p(\theta) = \mathcal{N}(\theta|\mathbf{m}, \mathbf{V}_0)$ , this becomes

$$\mathbf{A} = \mathbf{H} + \mathbf{V}_0^{-1}.$$

If we assume that the prior is broad, or equivalently that the number of data points is large, we can neglect the term  $\mathbf{V}_0^{-1}$  compared to  $\mathbf{H}$ . Using this result, (4.137) can be rewritten in the form

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) - \frac{1}{2}(\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m})\mathbf{V}_{0}^{-1}(\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m}) - \frac{1}{2}\ln|\mathbf{H}| + \text{const}$$
(169)

as required. Note that the phrasing of the question is misleading, since the assumption of a broad prior, or of large N, is required in order to derive this form, as well as in the subsequent simplification.

We now again invoke the broad prior assumption, allowing us to neglect the second term on the right hand side of (169) relative to the first term.

Since we assume i.i.d. data,  $\mathbf{H} = -\nabla\nabla \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}})$  consists of a sum of terms, one term for each datum, and we can consider the following approximation:

$$\mathbf{H} = \sum_{n=1}^{N} \mathbf{H}_n = N\widehat{\mathbf{H}}$$

where  $\mathbf{H}_n$  is the contribution from the  $n^{\mathrm{th}}$  data point and

$$\widehat{\mathbf{H}} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{H}_{n}.$$

Combining this with the properties of the determinant, we have

$$\ln |\mathbf{H}| = \ln |N\widehat{\mathbf{H}}| = \ln \left( N^M |\widehat{\mathbf{H}}| \right) = M \ln N + \ln |\widehat{\mathbf{H}}|$$

where M is the dimensionality of  $\boldsymbol{\theta}$ . Note that we are assuming that  $\widehat{\mathbf{H}}$  has full rank M. Finally, using this result together (169), we obtain (4.139) by dropping the  $\ln |\widehat{\mathbf{H}}|$  since this O(1) compared to  $\ln N$ .

**4.24** Consider a rotation of the coordinate axes of the M-dimensional vector  $\mathbf{w}$  such that  $\mathbf{w} = (w_{\parallel}, \mathbf{w}_{\perp})$  where  $\mathbf{w}^{\mathrm{T}} \boldsymbol{\phi} = w_{\parallel} \|\boldsymbol{\phi}\|$ , and  $\mathbf{w}_{\perp}$  is a vector of length M-1. We then have

$$\int \sigma(\mathbf{w}^{\mathrm{T}}\boldsymbol{\phi})q(\mathbf{w}) \, d\mathbf{w} = \iint \sigma\left(w_{\parallel}\|\boldsymbol{\phi}\|\right) q(\mathbf{w}_{\perp}|w_{\parallel})q(w_{\parallel}) \, dw_{\parallel} \, d\mathbf{w}_{\perp}$$
$$= \int \sigma(w_{\parallel}\|\boldsymbol{\phi}\|)q(w_{\parallel}) \, dw_{\parallel}.$$

 $=\int \sigma(w_{\parallel}\|\phi\|)q(w_{\parallel})\,\mathrm{d}w_{\parallel}.$  Note that the joint distribution  $q(w_{\parallel})$  is Gaussian Theoretic the marginal distribution  $q(w_{\parallel})$  is also Gaussian and can be found using the standard results presented in

Section 2.3.2. Denoting the unit vector

$$\mathbf{e} = \frac{1}{\|\boldsymbol{\phi}\|} \boldsymbol{\phi}$$

we have

$$q(w_{\parallel}) = \mathcal{N}(w_{\parallel}|\mathbf{e}^{\mathrm{T}}\mathbf{m}_{N}, \mathbf{e}^{\mathrm{T}}\mathbf{S}_{N}\mathbf{e}).$$

Defining  $a=w_{\parallel}\|\phi\|$  we see that the distribution of a is given by a simple re-scaling of the Gaussian, so that

$$q(a) = \mathcal{N}(a|\boldsymbol{\phi}^{\mathrm{T}}\mathbf{m}_{N}, \boldsymbol{\phi}^{\mathrm{T}}\mathbf{S}_{N}\boldsymbol{\phi})$$

where we have used  $\|\phi\|\mathbf{e} = \phi$ . Thus we obtain (4.151) with  $\mu_a$  given by (4.149) and  $\sigma_a^2$  given by (4.150).

**4.25** From (4.88) we have that

$$\frac{\mathrm{d}\sigma}{\mathrm{d}a}\Big|_{a=0} = \sigma(0)(1-\sigma(0))$$

$$= \frac{1}{2}\left(1-\frac{1}{2}\right) = \frac{1}{4}.$$
(170)

Since the derivative of a cumulative distribution function is simply the corresponding density function, (4.114) gives

$$\frac{\mathrm{d}\Phi(\lambda a)}{\mathrm{d}a}\bigg|_{a=0} = \lambda \mathcal{N}(0|0,1)$$
$$= \lambda \frac{1}{\sqrt{2\pi}}.$$

Setting this equal to (170), we see that

$$\lambda = \frac{\sqrt{2\pi}}{4}$$
 or equivalently  $\lambda^2 = \frac{\pi}{8}$ .

This is illustrated in Figure 4.9.

**4.26** First of all consider the derivative of the right hand side with respect to  $\mu$ , making use of the definition of the probit function, giving

$$\left(\frac{1}{2\pi}\right)^{1/2} \exp\left\{-\frac{\mu^2}{2(\lambda^{-2}+\sigma^2)}\right\} \frac{1}{(\lambda^{-2}+\sigma^2)^{1/2}}.$$

Now make the change of variable  $a = \mu + \sigma z$ , so that the left hand side of (4.152) becomes

$$\int_{-\infty}^{\infty} \Phi(\lambda \mu + \lambda \sigma z) \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2}z^2\right\} \sigma dz$$

where we have substituted for the Gaussian distribution. Now differentiate with respect to  $\mu$ , making use of the definition of the probit function, giving

$$\frac{1}{2\pi} \int_{-\infty}^{\infty} \exp\left\{-\frac{1}{2}z^2 - \frac{\lambda^2}{2}(\mu + \sigma z)^2\right\} \sigma \,dz.$$

The integral over z takes the standard Gaussian form and can be evaluated analytically by making use of the standard result for the normalization coefficient of a Gaussian distribution. To do this we first complete the square in the exponent

$$\begin{split} &-\frac{1}{2}z^2 - \frac{\lambda^2}{2}(\mu + \sigma z)^2 \\ &= &-\frac{1}{2}z^2(1 + \lambda^2\sigma^2) - z\lambda^2\mu\sigma - \frac{1}{2}\lambda^2\mu^2 \\ &= &-\frac{1}{2}\left[z + \lambda^2\mu\sigma(1 + \lambda^2\sigma^2)^{-1}\right]^2(1 + \lambda^2\sigma^2) + \frac{1}{2}\frac{\lambda^4\mu^2\sigma^2}{(1 + \lambda^2\sigma^2)} - \frac{1}{2}\lambda^2\mu^2. \end{split}$$

Integrating over z then gives the following result for the derivative of the left hand side

$$\frac{1}{(2\pi)^{1/2}} \frac{1}{(1+\lambda^2\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2}\lambda^2\mu^2 + \frac{1}{2}\frac{\lambda^4\mu^2\sigma^2}{(1+\lambda^2\sigma^2)}\right\} 
= \frac{1}{(2\pi)^{1/2}} \frac{1}{(1+\lambda^2\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2}\frac{\lambda^2\mu^2}{(1+\lambda^2\sigma^2)}\right\}.$$

Thus the derivatives of the left and right hand sides of (4.152) with respect to  $\mu$  are equal. It follows that the left and right hand sides are equal up to a function of  $\sigma^2$  and  $\lambda$ . Taking the limit  $\mu \to -\infty$  the left and right hand sides both go to zero, showing that the constant of integration must also be zero.

## **Chapter 5 Neural Networks**

**5.1** NOTE: In the 1<sup>st</sup> printing of PRML, the text of this exercise contains a typographical error. On line 2,  $g(\cdot)$  should be replaced by  $h(\cdot)$ .

See Solution 3.1.

**5.2** The likelihood function for an i.i.d. data set,  $\{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\}$ , under the conditional distribution (5.16) is given by

$$\prod_{n=1}^{N} \mathcal{N}\left(\mathbf{t}_{n}|\mathbf{y}(\mathbf{x}_{n},\mathbf{w}),\beta^{-1}\mathbf{I}\right).$$

If we take the logarithm of this, using (2.43), we get

$$\sum_{n=1}^{N} \ln \mathcal{N} \left( \mathbf{t}_{n} | \mathbf{y}(\mathbf{x}_{n}, \mathbf{w}), \beta^{-1} \mathbf{I} \right)$$

$$= -\frac{1}{2} \sum_{n=1}^{N} \left( \mathbf{t}_{n} - \mathbf{y}(\mathbf{x}_{n}, \mathbf{w}) \right)^{\mathrm{T}} (\beta \mathbf{I}) \left( \mathbf{t}_{n} - \mathbf{y}(\mathbf{x}_{n}, \mathbf{w}) \right) + \text{const}$$

$$= -\frac{\beta}{2} \sum_{n=1}^{N} \| \mathbf{t}_{n} - \mathbf{y}(\mathbf{x}_{n}, \mathbf{w}) \|^{2} + \text{const},$$

where 'const' comprises terms which are independent of w. The first term on the right hand side is proportional to the negative of (5.11) and hence maximizing the log-likelihood is equivalent to minimizing the sum-of-squares error.

#### **5.3** In this case, the likelihood function becomes

$$p(\mathbf{T}|\mathbf{X}, \mathbf{w}, \mathbf{\Sigma}) = \prod_{n=1}^{N} \mathcal{N}(\mathbf{t}_{n}|\mathbf{y}(\mathbf{x}_{n}, \mathbf{w}), \mathbf{\Sigma}),$$

with the corresponding log-likelihood function

$$\ln p(\mathbf{T}|\mathbf{X}, \mathbf{w}, \mathbf{\Sigma})$$

$$= -\frac{N}{2} \left( \ln |\mathbf{\Sigma}| + K \ln(2\pi) \right) - \frac{1}{2} \sum_{n=1}^{N} (\mathbf{t}_n - \mathbf{y}_n)^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{y}_n), \quad (171)$$

where  $y_n = y(x_n, w)$  and K is the dimensionality of y and t.

If we first treat  $\Sigma$  as fixed and known, we can drop terms that are independent of w from (171), and by changing the sign we get the error function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} (\mathbf{t}_n - \mathbf{y}_n)^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{y}_n).$$

If we consider maximizing (171) w.r.t.  $\Sigma$ , the terms that need to be kept are

$$-\frac{N}{2}\ln|\mathbf{\Sigma}| - \frac{1}{2}\sum_{n=1}^{N}(\mathbf{t}_{n} - \mathbf{y}_{n})^{\mathrm{T}}\mathbf{\Sigma}^{-1}(\mathbf{t}_{n} - \mathbf{y}_{n}).$$

By rewriting the second term we get

$$-\frac{N}{2}\ln|\mathbf{\Sigma}| - \frac{1}{2}\mathrm{Tr}\left[\mathbf{\Sigma}^{-1}\sum_{n=1}^{N}(\mathbf{t}_{n} - \mathbf{y}_{n})(\mathbf{t}_{n} - \mathbf{y}_{n})^{\mathrm{T}}\right].$$

Using results from Appendix C, we can maximize this by setting the derivative w.r.t.  $\Sigma^{-1}$  to zero, yielding

$$\Sigma = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{t}_n - \mathbf{y}_n) (\mathbf{t}_n - \mathbf{y}_n)^{\mathrm{T}}.$$

Thus the optimal value for  $\Sigma$  depends on w through  $y_n$ .

A possible way to address this mutual dependency between w and  $\Sigma$  when it comes to optimization, is to adopt an iterative scheme, alternating between updates of w and  $\Sigma$  until some convergence criterion is reached.

**5.4** Let  $t \in \{0,1\}$  denote the data set label and let  $k \in \{0,1\}$  denote the true class label. We want the network output to have the interpretation  $y(\mathbf{x}, \mathbf{w}) = p(k = 1 | \mathbf{x})$ . From the rules of probability we have

$$p(t=1|\mathbf{x}) = \sum_{k=0}^{1} p(t=1|k)p(k|\mathbf{x}) = (1-\epsilon)y(\mathbf{x}, \mathbf{w}) + \epsilon(1-y(\mathbf{x}, \mathbf{w})).$$

The conditional probability of the data label is then

$$p(t|\mathbf{x}) = p(t = 1|\mathbf{x})^{t} (1 - p(t = 1|\mathbf{x})^{1-t}.$$

Forming the likelihood and taking the negative logarithm we then obtain the error function in the form

$$E(\mathbf{w}) = -\sum_{n=1}^{N} \left\{ t_n \ln \left[ (1 - \epsilon) y(\mathbf{x}_n, \mathbf{w}) + \epsilon (1 - y(\mathbf{x}_n, \mathbf{w})) \right] + (1 - t_n) \ln \left[ 1 - (1 - \epsilon) y(\mathbf{x}_n, \mathbf{w}) - \epsilon (1 - y(\mathbf{x}_n, \mathbf{w})) \right] \right\}.$$

See also Solution 4.16.

**5.5** For the given interpretation of  $y_k(\mathbf{x}, \mathbf{w})$ , the conditional distribution of the target vector for a multiclass neural network is

$$p(\mathbf{t}|\mathbf{w}_1,\ldots,\mathbf{w}_K) = \prod_{k=1}^K y_k^{t_k}.$$

Thus, for a data set of N points, the likelihood function will be

$$p(\mathbf{T}|\mathbf{w}_1,\ldots,\mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}.$$

Taking the negative logarithm in order to derive an error function we obtain (5.24) as required. Note that this is the same result as for the multiclass logistic regression model, given by (4.108).

**5.6** Differentiating (5.21) with respect to the activation  $a_n$  corresponding to a particular data point n, we obtain

$$\frac{\partial E}{\partial a_n} = -t_n \frac{1}{y_n} \frac{\partial y_n}{\partial a_n} + (1 - t_n) \frac{1}{1 - y_n} \frac{\partial y_n}{\partial a_n}.$$
 (172)

From (4.88), we have

$$\frac{\partial y_n}{\partial a_n} = y_n (1 - y_n). \tag{173}$$

Substituting (173) into (172), we get

$$\frac{\partial E}{\partial a_n} = -t_n \frac{y_n(1-y_n)}{y_n} + (1-t_n) \frac{y_n(1-y_n)}{(1-y_n)}$$
$$= y_n - t_n$$

as required.

- **5.7** See Solution 4.17.
- **5.8** From (5.59), using standard derivatives, we get

$$\frac{\mathrm{d}\tanh}{\mathrm{d}a} = \frac{e^{a}}{e^{a} + e^{-a}} - \frac{e^{a}(e^{a} - e^{-a})}{(e^{a} + e^{-a})^{2}} + \frac{e^{-a}}{e^{a} + e^{-a}} + \frac{e^{-a}(e^{a} - e^{-a})}{(e^{a} + e^{-a})^{2}}$$

$$= \frac{e^{a} + e^{-a}}{e^{a} + e^{-a}} + \frac{1 - e^{2a} - e^{-2a} + 1}{(e^{a} + e^{-a})^{2}}$$

$$= 1 - \frac{e^{2a} - 2 + e^{-2a}}{(e^{a} + e^{-a})^{2}}$$

$$= 1 - \frac{(e^{a} - e^{-a})(e^{a} - e^{-a})}{(e^{a} + e^{-a})(e^{a} + e^{-a})}$$

$$= 1 - \tanh^{2}(a)$$

**5.9** This simply corresponds to a scaling and shifting of the binary outputs, which directly gives the activation function, using the notation from (5.19), in the form

$$y = 2\sigma(a) - 1.$$

The corresponding error function can be constructed from (5.21) by applying the inverse transform to  $y_n$  and  $t_n$ , yielding

$$E(\mathbf{w}) = -\sum_{n=1}^{N} \frac{1+t_n}{2} \ln \frac{1+y_n}{2} + \left(1 - \frac{1+t_n}{2}\right) \ln \left(1 - \frac{1+y_n}{2}\right)$$
$$= -\frac{1}{2} \sum_{n=1}^{N} \left\{ (1+t_n) \ln(1+y_n) + (1-t_n) \ln(1-y_n) \right\} + N \ln 2$$

97

where the last term can be dropped, since it is independent of w.

To find the corresponding activation function we simply apply the linear transformation to the logistic sigmoid given by (5.19), which gives

$$y(a) = 2\sigma(a) - 1 = \frac{2}{1 + e^{-a}} - 1$$
$$= \frac{1 - e^{-a}}{1 + e^{-a}} = \frac{e^{a/2} - e^{-a/2}}{e^{a/2} + e^{-a/2}}$$
$$= \tanh(a/2).$$

**5.10** From (5.33) and (5.35) we have

$$\mathbf{u}_i^{\mathrm{T}} \mathbf{H} \mathbf{u}_i = \mathbf{u}_i^{\mathrm{T}} \lambda_i \mathbf{u}_i = \lambda_i.$$

Assume that **H** is positive definite, so that (5.37) holds. Then by setting  $\mathbf{v} = \mathbf{u}_i$  it follows that

$$\lambda_i = \mathbf{u}_i^{\mathrm{T}} \mathbf{H} \mathbf{u}_i > 0 \tag{174}$$

for all values of i. Thus, if **H** is positive definite, all of its eigenvalues will be positive.

Conversely, assume that (174) holds. Then, for any vector,  $\mathbf{v}$ , we can make use of (5.38) to give

$$\mathbf{v}^{\mathrm{T}}\mathbf{H}\mathbf{v} = \left(\sum_{i} c_{i} \mathbf{u}_{i}\right)^{\mathrm{T}} \mathbf{H} \left(\sum_{j} c_{j} \mathbf{u}_{j}\right)$$

$$= \left(\sum_{i} c_{i} \mathbf{u}_{i}\right)^{\mathrm{T}} \left(\sum_{j} \lambda_{j} c_{j} \mathbf{u}_{j}\right)$$

$$= \sum_{i} \lambda_{i} c_{i}^{2} > 0$$

where we have used (5.33) and (5.34) along with (174). Thus, if all of the eigenvalues are positive, the Hessian matrix will be positive definite.

**5.11 NOTE**: In PRML, Equation (5.32) contains a typographical error: = should be  $\simeq$ .

We start by making the change of variable given by (5.35) which allows the error function to be written in the form (5.36). Setting the value of the error function  $E(\mathbf{w})$  to a constant value C we obtain

$$E(\mathbf{w}^*) + \frac{1}{2} \sum_{i} \lambda_i \alpha_i^2 = C.$$

Re-arranging gives

$$\sum_{i} \lambda_{i} \alpha_{i}^{2} = 2C - 2E(\mathbf{w}^{*}) = \widetilde{C}$$

where  $\widetilde{C}$  is also a constant. This is the equation for an ellipse whose axes are aligned with the coordinates described by the variables  $\{\alpha_i\}$ . The length of axis j is found by setting  $\alpha_i = 0$  for all  $i \neq j$ , and solving for  $\alpha_j$  giving

$$\alpha_j = \left(\frac{\widetilde{C}}{\lambda_j}\right)^{1/2}$$

which is inversely proportional to the square root of the corresponding eigenvalue.

**5.12 NOTE**: See note in Solution 5.11.

From (5.37) we see that, if **H** is positive definite, then the second term in (5.32) will be positive whenever  $(\mathbf{w} - \mathbf{w}^*)$  is non-zero. Thus the smallest value which  $E(\mathbf{w})$  can take is  $E(\mathbf{w}^*)$ , and so  $\mathbf{w}^*$  is the minimum of  $E(\mathbf{w})$ .

Conversely, if  $\mathbf{w}^*$  is the minimum of  $E(\mathbf{w})$ , then, for any vector  $\mathbf{w} \neq \mathbf{w}^*$ ,  $E(\mathbf{w}) > E(\mathbf{w}^*)$ . This will only be the case if the second term of (5.32) is positive for all values of  $\mathbf{w} \neq \mathbf{w}^*$  (since the first term is independent of  $\mathbf{w}$ ). Since  $\mathbf{w} - \mathbf{w}^*$  can be set to any vector of real numbers, it follows from the definition (5.37) that  $\mathbf{H}$  must be positive definite.

**5.13** From exercise 2.21 we know that a  $W \times W$  matrix has W(W+1)/2 independent elements. Add to that the W elements of the gradient vector  $\mathbf{b}$  and we get

$$\frac{W(W+1)}{2} + W = \frac{W(W+1) + 2W}{2} = \frac{W^2 + 3W}{2} = \frac{W(W+3)}{2}.$$

**5.14** We are interested in determining how the correction term

$$\delta = E'(w_{ij}) - \frac{E(w_{ij} + \epsilon) - E(w_{ij} - \epsilon)}{2\epsilon}$$
(175)

depend on  $\epsilon$ .

Using Taylor expansions, we can rewrite the numerator of the first term of (175) as

$$E(w_{ij}) + \epsilon E'(w_{ij}) + \frac{\epsilon^2}{2} E''(w_{ij}) + O(\epsilon^3)$$
$$-E(w_{ij}) + \epsilon E'(w_{ij}) - \frac{\epsilon^2}{2} E''(w_{ij}) + O(\epsilon^3) = 2\epsilon E'(w_{ij}) + O(\epsilon^3).$$

Note that the  $\epsilon^2$ -terms cancel. Substituting this into (175) we get,

$$\delta = \frac{2\epsilon E'(w_{ij}) + O(\epsilon^3)}{2\epsilon} - E'(w_{ij}) = O(\epsilon^2).$$

**5.15** The alternative forward propagation scheme takes the first line of (5.73) as its starting point. However, rather than proceeding with a 'recursive' definition of  $\partial y_k/\partial a_j$ , we instead make use of a corresponding definition for  $\partial a_j/\partial x_i$ . More formally

$$J_{ki} = \frac{\partial y_k}{\partial x_i} = \sum_j \frac{\partial y_k}{\partial a_j} \frac{\partial a_j}{\partial x_i}$$

where  $\partial y_k/\partial a_j$  is defined by (5.75), (5.76) or simply as  $\delta_{kj}$ , for the case of linear output units. We define  $\partial a_j/\partial x_i = w_{ji}$  if  $a_j$  is in the first hidden layer and otherwise

$$\frac{\partial a_j}{\partial x_i} = \sum_{l} \frac{\partial a_j}{\partial a_l} \frac{\partial a_l}{\partial x_i} \tag{176}$$

where

$$\frac{\partial a_j}{\partial a_l} = w_{jl} h'(a_l). \tag{177}$$

Thus we can evaluate  $J_{ki}$  by forward propagating  $\partial a_j/\partial x_i$ , with initial value  $w_{ij}$ , alongside  $a_j$ , using (176) and (177).

**5.16** The multivariate form of (5.82) is

$$E = \frac{1}{2} \sum_{n=1}^{N} (\mathbf{y}_n - \mathbf{t}_n)^{\mathrm{T}} (\mathbf{y}_n - \mathbf{t}_n).$$

The elements of the first and second derivatives then become

$$\frac{\partial E}{\partial w_i} = \sum_{n=1}^{N} (\mathbf{y}_n - \mathbf{t}_n)^{\mathrm{T}} \frac{\partial \mathbf{y}_n}{\partial w_i}$$

and

$$\frac{\partial^2 E}{\partial w_i \partial w_j} = \sum_{n=1}^N \left\{ \frac{\partial \mathbf{y}_n}{\partial w_j}^{\mathrm{T}} \frac{\partial \mathbf{y}_n}{\partial w_i} + (\mathbf{y}_n - \mathbf{t}_n)^{\mathrm{T}} \frac{\partial^2 \mathbf{y}_n}{\partial w_j \partial w_i} \right\}.$$

As for the univariate case, we again assume that the second term of the second derivative vanishes and we are left with

$$\mathbf{H} = \sum_{n=1}^{N} \mathbf{B}_n \mathbf{B}_n^{\mathrm{T}},$$

where  $\mathbf{B}_n$  is a  $W \times K$  matrix, K being the dimensionality of  $\mathbf{y}_n$ , with elements

$$(\mathbf{B}_n)_{lk} = \frac{\partial y_{nk}}{\partial w_l}.$$

**5.17** Taking the second derivatives of (5.193) with respect to two weights  $w_r$  and  $w_s$  we obtain

$$\frac{\partial^{2} E}{\partial w_{r} \partial w_{s}} = \sum_{k} \int \left\{ \frac{\partial y_{k}}{\partial w_{r}} \frac{\partial y_{k}}{\partial w_{s}} \right\} p(\mathbf{x}) d\mathbf{x} 
+ \sum_{k} \int \left\{ \frac{\partial^{2} y_{k}}{\partial w_{r} \partial w_{s}} (y_{k}(\mathbf{x}) - \mathbb{E}_{t_{k}}[t_{k}|\mathbf{x}]) \right\} p(\mathbf{x}) d\mathbf{x}. (178)$$

Using the result (1.89) that the outputs  $y_k(\mathbf{x})$  of the trained network represent the conditional averages of the target data, we see that the second term in (178) vanishes. The Hessian is therefore given by an integral of terms involving only the products of first derivatives. For a finite data set, we can write this result in the form

$$\frac{\partial^2 E}{\partial w_r \partial w_s} = \frac{1}{N} \sum_{n=1}^{N} \sum_{k} \frac{\partial y_k^n}{\partial w_r} \frac{\partial y_k^n}{\partial w_s}$$

which is identical with (5.84) up to a scaling factor.

**5.18** If we introduce skip layer weights, **U**, into the model described in Section 5.3.2, this will only affect the last of the forward propagation equations, (5.64), which becomes

$$y_k = \sum_{j=0}^{M} w_{kj}^{(2)} z_j + \sum_{i=1}^{D} u_{ki} x_i.$$

Note that there is no need to include the input bias. The derivative w.r.t.  $u_{ki}$  can be expressed using the output  $\{\delta_k\}$  of (5.65),

$$\frac{\partial E}{\partial u_{ki}} = \delta_k x_i.$$

**5.19** If we take the gradient of (5.21) with respect to w, we obtain

$$\nabla E(\mathbf{w}) = \sum_{n=1}^{N} \frac{\partial E}{\partial a_n} \nabla a_n = \sum_{n=1}^{N} (y_n - t_n) \nabla a_n,$$

where we have used the result proved earlier in the solution to Exercise 5.6. Taking the second derivatives we have

$$\nabla \nabla E(\mathbf{w}) = \sum_{n=1}^{N} \left\{ \frac{\partial y_n}{\partial a_n} \nabla a_n \nabla a_n + (y_n - t_n) \nabla \nabla a_n \right\}.$$

Dropping the last term and using the result (4.88) for the derivative of the logistic sigmoid function, proved in the solution to Exercise 4.12, we finally get

$$\nabla \nabla E(\mathbf{w}) \simeq \sum_{n=1}^{N} y_n (1 - y_n) \nabla a_n \nabla a_n = \sum_{n=1}^{N} y_n (1 - y_n) \mathbf{b}_n \mathbf{b}_n^{\mathrm{T}}$$

where  $\mathbf{b}_n \equiv \nabla a_n$ .

**5.20** Using the chain rule, we can write the first derivative of (5.24) as

$$\frac{\partial E}{\partial w_i} = \sum_{n=1}^{N} \sum_{k=1}^{K} \frac{\partial E}{\partial a_{nk}} \frac{\partial a_{nk}}{\partial w_i}.$$
 (179)

From Exercise 5.7, we know that

$$\frac{\partial E}{\partial a_{nk}} = y_{nk} - t_{nk}.$$

Using this and (4.106), we can get the derivative of (179) w.r.t.  $w_i$  as

$$\frac{\partial^2 E}{\partial w_i \, \partial w_j} = \sum_{n=1}^N \sum_{k=1}^K \left( \sum_{l=1}^K y_{nk} (I_{kl} - y_{nl}) \, \frac{\partial a_{nk}}{\partial w_i} \, \frac{\partial a_{nl}}{\partial w_j} + (y_{nk} - t_{nk}) \frac{\partial^2 a_{nk}}{\partial w_i \, \partial w_j} \right).$$

For a trained model, the network outputs will approximate the conditional class probabilities and so the last term inside the parenthesis will vanish in the limit of a large data set, leaving us with

$$(\mathbf{H})_{ij} \simeq \sum_{n=1}^{N} \sum_{k=1}^{K} \sum_{l=1}^{K} y_{nk} (I_{kl} - y_{nl}) \frac{\partial a_{nk}}{\partial w_i} \frac{\partial a_{nl}}{\partial w_j}.$$

**5.21 NOTE**: In PRML, the text in the exercise could be misunderstood; a clearer formulation is: "Extend the expression (5.86) for the outer product approximation of the Hessian matrix to the case of K > 1 output units. Hence, derive a form that allows (5.87) to be used to incorporate sequentially contributions from individual outputs as well as individual patterns. This, together with the identity (5.88), will allow the use of (5.89) for finding the inverse of the Hessian by sequentially incorporating contributions from individual outputs and patterns."

From (5.44) and (5.46), we see that the multivariate form of (5.82) is

$$E = \frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{K} (y_{nk} - t_{nk})^{2}.$$

Consequently, the multivariate form of (5.86) is given by

$$\mathbf{H}_{NK} = \sum_{n=1}^{N} \sum_{k=1}^{K} \mathbf{b}_{nk} \mathbf{b}_{nk}^{\mathrm{T}}$$

$$\tag{180}$$

where  $\mathbf{b}_{nk} \equiv \nabla a_{nk} = \nabla y_{nk}$ . The double index indicate that we will now iterate over outputs as well as patterns in the sequential build-up of the Hessian. However, in terms of the end result, there is no real need to attribute terms in this sum to specific outputs or specific patterns. Thus, by changing the indexation in (180), we can write it

$$\mathbf{H}_{J} = \sum_{j=1}^{J} \mathbf{c}_{j} \mathbf{c}_{j}^{\mathrm{T}}$$

$$\tag{181}$$

where J = NK and

欢迎关注公众号a机器常为与算法之道 $k(j) = (j-1) \circ K + 1$ 

with  $\oslash$  and  $\odot$  denoting integer division and remainder, respectively. The advantage of the indextion in (181) is that now we have a single indexed sum and so we can use (5.87)–(5.89) as they stand, just replacing  $\mathbf{b}_L$  with  $\mathbf{c}_L$ , letting L run from 0 to J-1.

**5.22 NOTE**: The first printing of PRML contained typographical errors in equation (5.95). On the r.h.s.,  $H_{kk'}$  should be  $M_{kk'}$ . Moreover, the indices j and j' should be swapped on the r.h.s.

Using the chain rule together with (5.48) and (5.92), we have

$$\frac{\partial E_n}{\partial w_{kj}^{(2)}} = \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial w_{kj}^{(2)}}$$

$$= \delta_k z_j \tag{182}$$

Thus,

$$\frac{\partial^2 E_n}{\partial w_{kj}^{(2)} \partial w_{k'j'}^{(2)}} = \frac{\partial \delta_k z_j}{\partial w_{k'j'}^{(2)}}$$

and since  $z_i$  is independent of the second layer weights,

$$\frac{\partial^{2} E_{n}}{\partial w_{kj}^{(2)} \partial w_{k'j'}^{(2)}} = z_{j} \frac{\partial \delta_{k}}{\partial w_{k'j'}^{(2)}} 
= z_{j} \frac{\partial^{2} E_{n}}{\partial a_{k} \partial a_{k'}} \frac{\partial a_{k}}{\partial w_{k'j'}^{(2)}} 
= z_{j} z_{j'} M_{kk'},$$

where we again have used the chain rule together with (5.48) and (5.92).

If both weights are in the first layer, we again used the chain rule, this time together with (5.48), (5.55) and (5.56), to get

$$\frac{\partial E_n}{\partial w_{ji}^{(1)}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}^{(1)}} 
= x_i \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} 
= x_i h'(a_j) \sum_k w_{kj}^{(2)} \delta_k.$$

Thus we have

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{j'i'}^{(1)}} = \frac{\partial}{\partial w_{j'i'}} \left( x_i h'(a_j) \sum_k w_{kj}^{(2)} \delta_k \right).$$

Now we note that  $x_i$  and  $w_{kj}^{(2)}$  do not depend on  $w_{j'i'}^{(1)}$ , while  $h'(a_j)$  is only affected in the case where j=j'. Using these observations together with (5.48), we get

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{j'i'}^{(1)}} = x_i x_{i'} h''(a_j) I_{jj'} \sum_k w_{kj}^{(2)} \delta_k + x_i h'(a_j) \sum_k w_{kj}^{(2)} \frac{\partial \delta_k}{\partial w_{j'i'}^{(1)}}.$$
 (183)

From (5.48), (5.55), (5.56), (5.92) and the chain rule, we have

$$\frac{\partial \delta_{k}}{\partial w_{j'i'}^{(1)}} = \sum_{k'} \frac{\partial^{2} E_{n}}{\partial a_{k} \partial a_{k'}} \frac{\partial a_{k'}}{\partial a_{j'}} \frac{\partial a_{j'}}{\partial w_{j'i'}^{(1)}}$$

$$= x_{i'}h'(a_{j}) \sum_{k'} w_{k'j'}^{(2)} M_{kk'}.$$
(184)

Substituting this back into (183), we obtain (5.94).

Finally, from (182) we have

$$\frac{\partial^2 E_n}{\partial w_{ji}^{(1)} \partial w_{kj'}^{(2)}} = \frac{\partial \delta_k z_{j'}}{\partial w_{ji}^{(1)}}.$$

Using (184), we get

$$\frac{\partial^{2} E_{n}}{\partial w_{ij}^{(1)} \partial w_{kj'}^{(2)}} = z_{j'} x_{i} h'(a_{j}) \sum_{k'} w_{k'j}^{(2)} M_{kk'} + \delta_{k} I_{jj'} h'(a_{j}) x_{i}$$

$$= x_{i} h'(a_{j}) \left( \delta_{k} I_{jj'} + \sum_{k'} w_{k'j}^{(2)} M_{kk'} \right).$$

**5.23** If we introduce skip layer weights into the model discussed in Section 5.4.5, three new cases are added to three already covered in Exercise 5.22.

The first derivative w.r.t. skip layer weight  $u_{ki}$  can be written

$$\frac{\partial E_n}{\partial u_{ki}} = \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial u_{ki}} = \frac{\partial E_n}{\partial a_k} x_i. \tag{185}$$

Using this, we can consider the first new case, where both weights are in the skip layer,

$$\frac{\partial^2 E_n}{\partial u_{ki} \partial u_{k'i'}} = \frac{\partial^2 E_n}{\partial a_k \partial a_{k'}} \frac{\partial a_{k'}}{\partial u_{k'i'}} x_i$$
$$= M_{kk'} x_i x_{i'},$$

where we have also used (5.92).

When one weight is in the skip layer and the other weight is in the hidden-to-output layer, we can use (185), (5.48) and (5.92) to get

$$\frac{\partial^2 E_n}{\partial u_{ki} \partial w_{k'j}^{(2)}} = \frac{\partial^2 E_n}{\partial a_k \partial a_{k'}} \frac{\partial a_{k'}}{\partial w_{k'j}^{(2)}} x_i$$

$$= M_{kk'} z_j x_i.$$

Finally, if one weight is a skip layer weight and the other is in the input-to-hidden layer, (185), (5.48), (5.55), (5.56) and (5.92) together give

$$\frac{\partial^{2} E_{n}}{\partial u_{ki} \partial w_{ji'}^{(1)}} = \frac{\partial}{\partial w_{ji'}^{(1)}} \left( \frac{\partial E_{n}}{\partial a_{k}} x_{i} \right) 
= \sum_{k'} \frac{\partial^{2} E_{n}}{\partial a_{k} \partial a_{k'}} \frac{\partial a_{k'}}{\partial w_{ji'}^{(1)}} x_{i} 
= x_{i} x_{i'} h'(a_{j}) \sum_{k'} M_{kk'} w_{k'j}^{(2)}.$$

**5.24** With the transformed inputs, weights and biases, (5.113) becomes

$$z_j = h\left(\sum_i \widetilde{w}_{ji}\widetilde{x}_i + \widetilde{w}_{j0}\right).$$

Using (5.115)–(5.117), we can rewrite the argument of  $h(\cdot)$  on the r.h.s. as

$$\sum_{i} \frac{1}{a} w_{ji} (ax_{i} + b) + w_{j0} - \frac{b}{a} \sum_{i} w_{ji}$$

$$= \sum_{i} w_{ji} x_{i} + \frac{b}{a} \sum_{i} w_{ji} + w_{j0} - \frac{b}{a} \sum_{i} w_{ji}$$

$$= \sum_{i} w_{ji} x_{i} + w_{j0}.$$

Similarly, with the transformed outputs, weights and biases, (5.114) becomes

$$\widetilde{y}_k = \sum_{j} \widetilde{w}_{kj} z_j + \widetilde{w}_{k0}.$$

Using (5.118)–(5.120), we can rewrite this as

$$cy_k + d = \sum_k cw_{kj}z_j + cw_{k0} + d$$
$$= c\left(\sum_i w_{kj}z_j + w_{k0}\right) + d.$$

By subtracting d and subsequently dividing by c on both sides, we recover (5.114) in its original form.

**5.25** The gradient of (5.195) is given

$$\nabla E = \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

and hence update formula (5.196) becomes

$$\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \rho \mathbf{H} (\mathbf{w}^{(\tau-1)} - \mathbf{w}^{\star}).$$

Pre-multiplying both sides with  $\mathbf{u}_i^{\mathrm{T}}$  we get

$$w_{j}^{(\tau)} = \mathbf{u}_{j}^{\mathrm{T}} \mathbf{w}^{(\tau)}$$

$$= \mathbf{u}_{j}^{\mathrm{T}} \mathbf{w}^{(\tau-1)} - \rho \mathbf{u}_{j}^{\mathrm{T}} \mathbf{H} (\mathbf{w}^{(\tau-1)} - \mathbf{w}^{\star})$$

$$= w_{j}^{(\tau-1)} - \rho \eta_{j} \mathbf{u}_{j}^{\mathrm{T}} (\mathbf{w} - \mathbf{w}^{\star})$$

$$= w_{j}^{(\tau-1)} - \rho \eta_{j} (w_{j}^{(\tau-1)} - w_{j}^{\star}),$$
(187)

where we have used (5.198). To show that

$$w_i^{(\tau)} = \{1 - (1 - \rho \eta_j)^{\tau}\} w_i^{\star}$$

for  $\tau = 1, 2, ...$ , we can use proof by induction. For  $\tau = 1$ , we recall that  $\mathbf{w}^{(0)} = \mathbf{0}$  and insert this into (187), giving

$$w_{j}^{(1)} = w_{j}^{(0)} - \rho \eta_{j} (w_{j}^{(0)} - w_{j}^{\star})$$

$$= \rho \eta_{j} w_{j}^{\star}$$

$$= \{1 - (1 - \rho \eta_{j})\} w_{j}^{\star}.$$

Now we assume that the result holds for  $\tau = N - 1$  and then make use of (187)

$$\begin{split} w_j^{(N)} &= w_j^{(N-1)} - \rho \eta_j (w_j^{(N-1)} - w_j^{\star}) \\ &= w_j^{(N-1)} (1 - \rho \eta_j) + \rho \eta_j w_j^{\star} \\ &= \left\{ 1 - (1 - \rho \eta_j)^{N-1} \right\} w_j^{\star} (1 - \rho \eta_j) + \rho \eta_j w_j^{\star} \\ &= \left\{ (1 - \rho \eta_j) - (1 - \rho \eta_j)^N \right\} w_j^{\star} + \rho \eta_j w_j^{\star} \\ &= \left\{ 1 - (1 - \rho \eta_j)^N \right\} w_j^{\star} \end{split}$$

as required

Provided that  $|1-\rho\eta_j|<1$  then we have  $(1-\rho\eta_j)^{\tau}\to 0$  as  $\tau\to\infty$ , and hence  $\left\{1-(1-\rho\eta_j)^N\right\}\to 1$  and  $\mathbf{w}^{(\tau)}\to\mathbf{w}^{\star}$ .

If  $\tau$  is finite but  $\eta_j \gg (\rho \tau)^{-1}$ ,  $\tau$  must still be large, since  $\eta_j \rho \tau \gg 1$ , even though  $|1 - \rho \eta_j| < 1$ . If  $\tau$  is large, it follows from the argument above that  $w_j^{(\tau)} \simeq w_j^{\star}$ .

If, on the other hand,  $\eta_j \ll (\rho \tau)^{-1}$ , this means that  $\rho \eta_j$  must be small, since  $\rho \eta_j \tau \ll 1$  and  $\tau$  is an integer greater than or equal to one. If we expand,

$$(1 - \rho \eta_j)^{\tau} = 1 - \tau \rho \eta_j + O(\rho \eta_j^2)$$

and insert this into (5.197), we get

$$|w_{j}^{(\tau)}| = |\{1 - (1 - \rho \eta_{j})^{\tau}\} w_{j}^{\star}|$$

$$= |\{1 - (1 - \tau \rho \eta_{j} + O(\rho \eta_{j}^{2}))\} w_{j}^{\star}|$$

$$\simeq \tau \rho \eta_{j} |w_{j}^{\star}| \ll |w_{j}^{\star}|$$

Recall that in Section 3.5.3 we showed that when the regularization parameter (called  $\alpha$  in that section) is much larger than one of the eigenvalues (called  $\lambda_j$  in that section) then the corresponding parameter value  $w_i$  will be close to zero. Conversely, when  $\alpha$  is much smaller than  $\lambda_i$  then  $w_i$  will be close to its maximum likelihood value. Thus  $\alpha$  is playing an analogous role to  $\rho\tau$ .

**5.26 NOTE**: In PRML, equation (5.201) should read

$$\Omega_n = \frac{1}{2} \sum_k \left( \mathcal{G} y_k \right)^2 \bigg|_{\mathbf{x}_n}.$$

In this solution, we will indicate dependency on  $\mathbf{x}_n$  with a subscript n on relevant symbols.

Substituting the r.h.s. of (5.202) into (5.201) and then using (5.70), we get

$$\Omega_n = \frac{1}{2} \sum_{k} \left( \sum_{i} \tau_{ni} \frac{\partial y_{nk}}{\partial x_{ni}} \right)^2$$
 (188)

$$= \frac{1}{2} \sum_{k} \left( \sum_{i} \tau_{ni} J_{nki} \right)^{2} \tag{189}$$

where  $J_{nki}$  denoted  $J_{ki}$  evaluated at  $\mathbf{x}_n$ . Summing (189) over n, we get (5.128).

By applying  $\mathcal{G}$  from (5.202) to the equations in (5.203) and making use of (5.205) we obtain (5.204). From this, we see that  $\beta_{nl}$  can be written in terms of  $\alpha_{ni}$ , which in turn can be written as functions of  $\beta_{ni}$  from the previous layer. For the input layer, using (5.204) and (5.205), we get

$$\beta_{nj} = \sum_{i} w_{ji} \alpha_{ni}$$

$$= \sum_{i} w_{ji} \mathcal{G} x_{ni}$$

$$= \sum_{i} w_{ji} \sum_{i'} \tau_{ni'} \frac{\partial x_{ni}}{\partial x_{ni'}}$$

$$= \sum_{i} w_{ji} \tau_{ni}.$$
(190)

Thus we see that, starting from (190),  $\tau_n$  is propagated forward by subsequent application of the equations in (5.204), yielding the  $\beta_{nl}$  for the output layer, from which  $\Omega_n$  can be computed using (5.201),

$$\Omega_n = \frac{1}{2} \sum_k (\mathcal{G}y_{nk})^2 = \frac{1}{2} \sum_k \alpha_{nk}^2.$$

Considering  $\partial \Omega_n/\partial w_{rs}$ , we start from (5.201) and make use of the chain rule, together with (5.52), (5.205) and (5.207), to obtain

$$\frac{\partial \Omega_n}{\partial w_{rs}} = \sum_k (\mathcal{G}y_{nk}) \mathcal{G} (\delta_{nkr} z_{ns})$$

$$= \sum_k \alpha_{nk} (\phi_{nkr} z_{ns} + \delta_{nkr} \alpha_{ns}).$$

The backpropagation formula for computing  $\delta_{nkr}$  follows from (5.74), which is used in computing the Jacobian matrix, and is given by

$$\delta_{nkr} = h'(a_{nr}) \sum_{l} w_{lr} \delta_{nkl}.$$

Using this together with (5.205) and (5.207), we can obtain backpropagation equations for  $\phi_{nkr}$ ,

$$\phi_{nkr} = \mathcal{G}\delta_{nkr}$$

$$= \mathcal{G}\left(h'(a_{nr})\sum_{l}w_{lr}\delta_{nkl}\right)$$

$$= h''(a_{nr})\beta_{nr}\sum_{l}w_{lr}\delta_{nkl} + h'(a_{nr})\sum_{l}w_{lr}\phi_{nkl}.$$

**5.27** If  $s(x, \xi) = x + \xi$ , then

$$\frac{\partial s_k}{\partial \xi_i} = I_{ki}$$
, i.e.,  $\frac{\partial \mathbf{s}}{\partial \boldsymbol{\xi}} = \mathbf{I}$ ,

and since the first order derivative is constant, there are no higher order derivatives. We now make use of this result to obtain the derivatives of y w.r.t.  $\xi_i$ :

$$\frac{\partial y}{\partial \xi_i} = \sum_k \frac{\partial y}{\partial s_k} \frac{\partial s_k}{\partial \xi_i} = \frac{\partial y}{\partial s_i} = b_i$$

$$\frac{\partial y}{\partial \xi_i \partial \xi_j} = \frac{\partial b_i}{\partial \xi_j} = \sum_k \frac{\partial b_i}{\partial s_k} \frac{\partial s_k}{\partial \xi_j} = \frac{\partial b_i}{\partial s_j} = B_{ij}$$

Using these results, we can write the expansion of  $\widetilde{E}$  as follows:

$$\begin{split} \widetilde{E} &= \frac{1}{2} \iiint \{y(\mathbf{x}) - t\}^2 p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) \, \mathrm{d}\boldsymbol{\xi} \, \mathrm{d}\mathbf{x} \, \mathrm{d}t \\ &+ \iiint \{y(\mathbf{x}) - t\} \mathbf{b}^\mathrm{T} \boldsymbol{\xi} p(\boldsymbol{\xi}) p(t|\mathbf{x}) p(\mathbf{x}) \, \mathrm{d}\boldsymbol{\xi} \, \mathrm{d}\mathbf{x} \, \mathrm{d}t \\ &+ \frac{1}{2} \iiint \boldsymbol{\xi}^\mathrm{T} \left( \{y(\mathbf{x}) - t\} \mathbf{B} + \mathbf{b} \mathbf{b}^\mathrm{T} \right) \boldsymbol{\xi} p(\boldsymbol{\xi}) p(t|\mathbf{x}) p(\mathbf{x}) \, \mathrm{d}\boldsymbol{\xi} \, \mathrm{d}\mathbf{x} \, \mathrm{d}t. \end{split}$$

The middle term will again disappear, since  $\mathbb{E}[\xi] = \mathbf{0}$  and thus we can write  $\widetilde{E}$  on the form of (5.131) with

$$\Omega = \frac{1}{2} \iiint \boldsymbol{\xi}^{\mathrm{T}} \left( \{ y(\mathbf{x}) - t \} \mathbf{B} + \mathbf{b} \mathbf{b}^{\mathrm{T}} \right) \boldsymbol{\xi} p(\boldsymbol{\xi}) p(t|\mathbf{x}) p(\mathbf{x}) \, \mathrm{d} \boldsymbol{\xi} \, \mathrm{d} \mathbf{x} \, \mathrm{d} t.$$

Again the first term within the parenthesis vanishes to leading order in  $\xi$  and we are left with

$$\Omega \simeq \frac{1}{2} \iint \boldsymbol{\xi}^{T} \left( \mathbf{b} \mathbf{b}^{T} \right) \boldsymbol{\xi} p(\boldsymbol{\xi}) p(\mathbf{x}) \, \mathrm{d} \boldsymbol{\xi} \, \mathrm{d} \mathbf{x}$$

$$= \frac{1}{2} \iint \operatorname{Trace} \left[ \left( \boldsymbol{\xi} \boldsymbol{\xi}^{T} \right) \left( \mathbf{b} \mathbf{b}^{T} \right) \right] p(\boldsymbol{\xi}) p(\mathbf{x}) \, \mathrm{d} \boldsymbol{\xi} \, \mathrm{d} \mathbf{x}$$

$$= \frac{1}{2} \int \operatorname{Trace} \left[ \mathbf{I} \left( \mathbf{b} \mathbf{b}^{T} \right) \right] p(\mathbf{x}) \, \mathrm{d} \mathbf{x}$$

$$= \frac{1}{2} \int \mathbf{b}^{T} \mathbf{b} p(\mathbf{x}) \, \mathrm{d} \mathbf{x} = \frac{1}{2} \int \|\nabla y(\mathbf{x})\|^{2} p(\mathbf{x}) \, \mathrm{d} \mathbf{x},$$

where we used the fact that  $\mathbb{E}[\boldsymbol{\xi}\boldsymbol{\xi}^{\mathrm{T}}]=\mathbf{I}.$ 

**5.28** The modifications only affect derivatives with respect to weights in the convolutional layer. The units within a feature map (indexed m) have different inputs, but all share a common weight vector,  $\mathbf{w}^{(m)}$ . Thus, errors  $\delta^{(m)}$  from all units within a feature map will contribute to the derivatives of the corresponding weight vector. In this situation, (5.50) becomes

$$\frac{\partial E_n}{\partial w_i^{(m)}} = \sum_j \frac{\partial E_n}{\partial a_j^{(m)}} \frac{\partial a_j^{(m)}}{\partial w_i^{(m)}} = \sum_j \delta_j^{(m)} z_{ji}^{(m)}.$$

Here  $a_j^{(m)}$  denotes the activation of the  $j^{\rm th}$  unit in the  $m^{\rm th}$  feature map, whereas  $w_i^{(m)}$  denotes the  $i^{\rm th}$  element of the corresponding feature vector and, finally,  $z_{ji}^{(m)}$  denotes the  $i^{\rm th}$  input for the  $j^{\rm th}$  unit in the  $m^{\rm th}$  feature map; the latter may be an actual input or the output of a preceding layer.

Note that  $\delta_j^{(m)} = \partial E_n/\partial a_j^{(m)}$  will typically be computed recursively from the  $\delta s$  of the units in the following layer, using (5.55). If there are layer(s) preceding the

convolutional layer, the standard backward propagation equations will apply; the weights in the convolutional layer can be treated as if they were independent parameters, for the purpose of computing the  $\delta s$  for the preceding layer's units.

**5.29** This is easily verified by taking the derivative of (5.138), using (1.46) and standard derivatives, yielding

$$\frac{\partial \Omega}{\partial w_i} = \frac{1}{\sum_k \pi_k \mathcal{N}(w_i | \mu_k, \sigma_k^2)} \sum_j \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \frac{(w_i - \mu_j)}{\sigma^2}.$$

Combining this with (5.139) and (5.140), we immediately obtain the second term of (5.141).

**5.30** Since the  $\mu_i$ s only appear in the regularization term,  $\Omega(\mathbf{w})$ , from (5.139) we have

$$\frac{\partial \widetilde{E}}{\partial \mu_j} = \lambda \frac{\partial \Omega}{\partial \mu_j}.$$
 (191)

Using (2.42), (5.138) and (5.140) and standard rules for differentiation, we can calculate the derivative of  $\Omega(\mathbf{w})$  as follows:

$$\frac{\partial \Omega}{\partial \mu_j} = -\sum_i \frac{1}{\sum_{j'} \pi_{j'} \mathcal{N}\left(w_i | \mu_{j'}, \sigma_{j'}^2\right)} \pi_j \mathcal{N}\left(w_i | \mu_j, \sigma_j^2\right) \frac{w_i - \mu_j}{\sigma_j^2} 
= -\sum_i \gamma_j(\mathbf{w}_i) \frac{w_i - \mu_j}{\sigma_j^2}.$$

Combining this with (191), we get (5.142).

**5.31** Following the same line of argument as in Solution 5.30, we need the derivative of  $\Omega(\mathbf{w})$  w.r.t.  $\sigma_j$ . Again using (2.42), (5.138) and (5.140) and standard rules for differentiation, we find this to be

$$\frac{\partial\Omega}{\partial\sigma_{j}} = -\sum_{i} \frac{1}{\sum_{j'} \pi_{j'} \mathcal{N}\left(w_{i} | \mu_{j'}, \sigma_{j'}^{2}\right)} \pi_{j} \frac{1}{(2\pi)^{1/2}} \left\{ -\frac{1}{\sigma_{j}^{2}} \exp\left(-\frac{(w_{i} - \mu_{j})^{2}}{2\sigma_{j}^{2}}\right) + \frac{1}{\sigma_{j}} \exp\left(-\frac{(w_{i} - \mu_{j})^{2}}{2\sigma_{j}^{2}}\right) \frac{(w_{i} - \mu_{j})^{2}}{\sigma_{j}^{3}} \right\} 
= \sum_{i} \gamma_{j}(w_{i}) \left\{ \frac{1}{\sigma_{j}} - \frac{(w_{i} - \mu_{j})^{2}}{\sigma_{j}^{3}} \right\}.$$

Combining this with (191), we get (5.143).

**5.32 NOTE**: In the first printing of PRML, there is a leading  $\lambda$  missing on the r.h.s. of equation (5.147). Moreover, in the text of the exercise (last line), the equation of the constraint to be used should read " $\sum_{k} \gamma_k(w_i) = 1$  for all i".

Equation (5.208) follows from (5.146) in exactly the same way that (4.106) follows from (4.104) in Solution 4.17.

Just as in Solutions 5.30 and 5.31,  $\eta_j$  only affect  $\widetilde{E}$  through  $\Omega(\mathbf{w})$ . However,  $\eta_j$  will affect  $\pi_k$  for all values of k (not just j=k). Thus we have

$$\frac{\partial \Omega}{\partial \eta_j} = \sum_k \frac{\partial \Omega}{\partial \pi_k} \frac{\partial \pi_k}{\partial \eta_j}.$$
 (192)

From (5.138) and (5.140), we get

$$\frac{\partial \Omega}{\partial \pi_k} = -\sum_i \frac{\gamma_k(w_i)}{\pi_k}.$$

Substituting this and (5.208) into (192) yields

$$\frac{\partial \Omega}{\partial \eta_j} = \frac{\partial \widetilde{E}}{\partial \eta_j} = -\sum_k \sum_i \frac{\gamma_k(w_i)}{\pi_k} \left\{ \delta_{jk} \pi_j - \pi_j \pi_k \right\}$$
$$= \sum_i \left\{ \pi_j - \gamma_j(w_i) \right\},$$

where we have used the fact that  $\sum_k \gamma_k(w_i) = 1$  for all i.

**5.33** From standard trigometric rules we get the position of the end of the first arm,

$$(x_1^{(1)}, x_2^{(1)}) = (L_1 \cos(\theta_1), L_1 \sin(\theta_1)).$$

Similarly, the position of the end of the second arm relative to the end of the first arm is given by the corresponding equation, with an angle offset of  $\pi$  (see Figure 5.18), which equals a change of sign

$$\begin{pmatrix} x_1^{(2)}, x_2^{(2)} \end{pmatrix} = (L_2 \cos(\theta_1 + \theta_2 - \pi), L_1 \sin(\theta_1 + \theta_2 - \pi))$$

$$= -(L_2 \cos(\theta_1 + \theta_2), L_2 \sin(\theta_1 + \theta_2)).$$

Putting this together, we must also taken into account that  $\theta_2$  is measured relative to the first arm and so we get the position of the end of the second arm relative to the attachment point of the first arm as

$$(x_1, x_2) = (L_1 \cos(\theta_1) - L_2 \cos(\theta_1 + \theta_2), L_1 \sin(\theta_1) - L_2 \sin(\theta_1 + \theta_2)).$$

**5.34 NOTE**: In the 1<sup>st</sup> printing of PRML, the l.h.s. of (5.154) should be replaced with  $\gamma_{nk} = \gamma_k(\mathbf{t}_n|\mathbf{x}_n)$ . Accordingly, in (5.155) and (5.156),  $\gamma_k$  should be replaced by  $\gamma_{nk}$  and in (5.156),  $t_l$  should be  $t_{nl}$ .

We start by using the chain rule to write

$$\frac{\partial E_n}{\partial a_k^{\pi}} = \sum_{j=1}^K \frac{\partial E_n}{\partial \pi_j} \frac{\partial \pi_j}{\partial a_k^{\pi}}.$$
 (193)

111

Note that because of the coupling between outputs caused by the softmax activation function, the dependence on the activation of a single output unit involves all the output units.

For the first factor inside the sum on the r.h.s. of (193), standard derivatives applied to the  $n^{\rm th}$  term of (5.153) gives

$$\frac{\partial E_n}{\partial \pi_j} = -\frac{\mathcal{N}_{nj}}{\sum_{l=1}^K \pi_l \mathcal{N}_{nl}} = -\frac{\gamma_{nj}}{\pi_j}.$$
 (194)

For the for the second factor, we have from (4.106) that

$$\frac{\partial \pi_j}{\partial a_k^{\pi}} = \pi_j (I_{jk} - \pi_k). \tag{195}$$

Combining (193), (194) and (195), we get

$$\frac{\partial E_n}{\partial a_k^{\pi}} = -\sum_{j=1}^K \frac{\gamma_{nj}}{\pi_j} \pi_j (I_{jk} - \pi_k)$$

$$= -\sum_{j=1}^K \gamma_{nj} (I_{jk} - \pi_k) = -\gamma_{nk} + \sum_{j=1}^K \gamma_{nj} \pi_k = \pi_k - \gamma_{nk},$$

where we have used the fact that, by (5.154),  $\sum_{j=1}^{K} \gamma_{nj} = 1$  for all n.

**5.35 NOTE**: See Solution 5.34.

From (5.152) we have

$$a_{kl}^{\mu} = \mu_{kl}$$

and thus

$$\frac{\partial E_n}{\partial a_{kl}^{\mu}} = \frac{\partial E_n}{\partial \mu_{kl}}.$$

From (2.43), (5.153) and (5.154), we get

$$\frac{\partial E_n}{\partial \mu_{kl}} = -\frac{\pi_k \mathcal{N}_{nk}}{\sum_{k'} \pi_{k'} \mathcal{N}_{nk'}} \frac{t_{nl} - \mu_{kl}}{\sigma_k^2(\mathbf{x}_n)}$$
$$= \gamma_{nk} \left( \mathbf{t}_n | \mathbf{x}_n \right) \frac{\mu_{kl} - t_{nl}}{\sigma_k^2(\mathbf{x}_n)}.$$

**5.36 NOTE**: In the 1<sup>st</sup> printing of PRML, equation (5.157) is incorrect and the correct equation appears at the end of this solution; see also Solution 5.34.

From (5.151) and (5.153), we see that   
欢迎关注公众号@机器。学验和数据,算法之道   

$$\frac{\partial a_k}{\partial a_k}$$
 (196)

where, from (5.151),

$$\frac{\partial \sigma_k}{\partial a_k^{\sigma}} = \sigma_k. \tag{197}$$

From (2.43), (5.153) and (5.154), we get

$$\frac{\partial E_n}{\partial \sigma_k} = -\frac{1}{\sum_{k'} \mathcal{N}_{nk'}} \left(\frac{L}{2\pi}\right)^{L/2} \left\{ -\frac{L}{\sigma^{L+1}} \exp\left(-\frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{2\sigma_k^2}\right) + \frac{1}{\sigma^L} \exp\left(-\frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{2\sigma_k^2}\right) \frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{\sigma_k^3} \right\} \\
= \gamma_{nk} \left(\frac{L}{\sigma_k} - \frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{\sigma_k^3}\right).$$

Combining this with (196) and (197), we get

$$\frac{\partial E_n}{\partial a_k^{\sigma}} = \gamma_{nk} \left( L - \frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{\sigma_k^2} \right).$$

### **5.37** From (2.59) and (5.148) we have

$$\mathbb{E}\left[\mathbf{t}|\mathbf{x}\right] = \int \mathbf{t}p\left(\mathbf{t}|\mathbf{x}\right) d\mathbf{t}$$

$$= \int \mathbf{t} \sum_{k=1}^{K} \pi_k(\mathbf{x}) \mathcal{N}\left(\mathbf{t}|\boldsymbol{\mu}_k(\mathbf{x}), \sigma_k^2(\mathbf{x})\right) d\mathbf{t}$$

$$= \sum_{k=1}^{K} \pi_k(\mathbf{x}) \int \mathbf{t} \mathcal{N}\left(\mathbf{t}|\boldsymbol{\mu}_k(\mathbf{x}), \sigma_k^2(\mathbf{x})\right) d\mathbf{t}$$

$$= \sum_{k=1}^{K} \pi_k(\mathbf{x}) \boldsymbol{\mu}_k(\mathbf{x}).$$

We now introduce the shorthand notation

$$\overline{\mathbf{t}}_k = \boldsymbol{\mu}_k(\mathbf{x}) \quad ext{and} \quad \overline{\mathbf{t}} = \sum_{k=1}^K \pi_k(\mathbf{x}) \overline{\mathbf{t}}_k.$$

Using this together with (2.59), (2.62), (5.148) and (5.158), we get

$$s^{2}(\mathbf{x}) = \mathbb{E}\left[\|\mathbf{t} - \mathbb{E}\left[\mathbf{t}|\mathbf{x}\right]\|^{2}|\mathbf{x}\right] = \int \|\mathbf{t} - \overline{\mathbf{t}}\|^{2} p\left(\mathbf{t}|\mathbf{x}\right) d\mathbf{t}$$

$$= \int \left(\mathbf{t}^{T}\mathbf{t} - \mathbf{t}^{T}\overline{\mathbf{t}} - \overline{\mathbf{t}}^{T}\mathbf{t} + \overline{\mathbf{t}}^{T}\overline{\mathbf{t}}\right) \sum_{k=1}^{K} \pi_{k} \mathcal{N}\left(\mathbf{t}|\boldsymbol{\mu}_{k}(\mathbf{x}), \sigma_{k}^{2}(\mathbf{x})\right) d\mathbf{t}$$

$$= \sum_{k=1}^{K} \pi_{k}(\mathbf{x}) \left\{\sigma_{k}^{2} + \overline{\mathbf{t}}_{k}^{T}\overline{\mathbf{t}}_{k} - \overline{\mathbf{t}}_{k}^{T}\overline{\mathbf{t}} - \overline{\mathbf{t}}^{T}\overline{\mathbf{t}}_{k} + \overline{\mathbf{t}}^{T}\overline{\mathbf{t}}\right\}$$

$$= \sum_{k=1}^{K} \pi_{k}(\mathbf{x}) \left\{\sigma_{k}^{2} + \|\overline{\mathbf{t}}_{k} - \overline{\mathbf{t}}\|^{2}\right\}$$

$$= \sum_{k=1}^{K} \pi_{k}(\mathbf{x}) \left\{\sigma_{k}^{2} + \|\overline{\mathbf{t}}_{k} - \overline{\mathbf{t}}\|^{2}\right\}$$

$$= \sum_{k=1}^{K} \pi_{k}(\mathbf{x}) \left\{\sigma_{k}^{2} + \|\boldsymbol{\mu}_{k}(\mathbf{x}) - \sum_{l}^{K} \pi_{l} \boldsymbol{\mu}_{l}(\mathbf{x})\|^{2}\right\}.$$

**5.38** Making the following substitions from the r.h.s. of (5.167) and (5.171),

$$\mathbf{x} \Rightarrow \mathbf{w} \quad \boldsymbol{\mu} \Rightarrow \mathbf{w}_{\mathrm{MAP}} \quad \boldsymbol{\Lambda}^{-1} \Rightarrow \boldsymbol{A}^{-1}$$

$$\mathbf{y} \Rightarrow t \quad \mathbf{A} \Rightarrow \mathbf{g}^{\mathrm{T}} \quad \mathbf{b} \Rightarrow y(\mathbf{x}, \mathbf{w}_{\mathrm{MAP}}) - \mathbf{g}^{\mathrm{T}} \mathbf{w}_{\mathrm{MAP}} \quad \mathbf{L}^{-1} \Rightarrow \beta^{-1},$$

in (2.113) and (2.114), (2.115) becomes

$$p(t) = \mathcal{N} \left( t | \mathbf{g}^{\mathrm{T}} \mathbf{w}_{\mathrm{MAP}} + y(\mathbf{x}, \mathbf{w}_{\mathrm{MAP}}) - \mathbf{g}^{\mathrm{T}} \mathbf{w}_{\mathrm{MAP}}, \beta^{-1} + \mathbf{g}^{\mathrm{T}} \mathbf{A}^{-1} \mathbf{g} \right)$$
$$= \mathcal{N} \left( t | y(\mathbf{x}, \mathbf{w}_{\mathrm{MAP}}), \sigma^{2} \right),$$

where  $\sigma^2$  is defined by (5.173).

**5.39** Using (4.135), we can approximate (5.174) as

$$p(\mathcal{D}|\alpha,\beta) \simeq p(\mathcal{D}|\mathbf{w}_{\text{MAP}},\beta)p(\mathbf{w}_{\text{MAP}}|\alpha)$$
$$\int \exp\left\{-\frac{1}{2}\left(\mathbf{w} - \mathbf{w}_{\text{MAP}}\right)^{\text{T}}\mathbf{A}\left(\mathbf{w} - \mathbf{w}_{\text{MAP}}\right)\right\} d\mathbf{w},$$

where **A** is given by (5.166), as  $p(\mathcal{D}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)$  is proportional to  $p(\mathbf{w}|\mathcal{D}, \alpha, \beta)$ . Using (4.135), (5.162) and (5.163), we can rewrite this as

$$p(\mathcal{D}|\alpha,\beta) \simeq \prod_{n=1}^{N} \mathcal{N}(t_n|y(\mathbf{x}_n,\mathbf{w}_{\text{MAP}}),\beta^{-1}) \mathcal{N}(\mathbf{w}_{\text{MAP}}|\mathbf{0},\alpha^{-1}\mathbf{I}) \frac{(2\pi)^{W/2}}{|\mathbf{A}|^{1/2}}.$$

Taking the logarithm of both sides and then using (2.42) and (2.43), we obtain the desired result.

#### 114 Solutions 5.40–6.1

**5.40** For a K-class neural network, the likelihood function is given by

$$\prod_{n}^{N} \prod_{k}^{K} y_{k}(\mathbf{x}_{n}, \mathbf{w})^{t_{nk}}$$

and the corresponding error function is given by (5.24).

Again we would use a Laplace approximation for the posterior distribution over the weights, but the corresponding Hessian matrix, **H**, in (5.166), would now be derived from (5.24). Similarly, (5.24), would replace the binary cross entropy error term in the regularized error function (5.184).

The predictive distribution for a new pattern would again have to be approximated, since the resulting marginalization cannot be done analytically. However, in contrast to the two-class problem, there is no obvious candidate for this approximation, although Gibbs (1997) discusses various alternatives.

**5.41 NOTE**: In PRML, the final "const" term in Equation (5.183) should be ommitted.

This solutions is similar to Solution 5.39, with the difference that the log-likelihood term is now given by (5.181). Again using (4.135), the corresponding approximation of the marginal likelihood becomes

$$p(\mathcal{D}|\alpha) \simeq p(\mathcal{D}|\mathbf{w}_{\text{MAP}})p(\mathbf{w}_{\text{MAP}}|\alpha)$$
$$\int \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\text{T}}\mathbf{A}(\mathbf{w} - \mathbf{w}_{\text{MAP}})\right) d\mathbf{w}, \quad (198)$$

where now

$$\mathbf{A} = -\nabla\nabla \ln p(\mathcal{D}|\mathbf{w}) = \mathbf{H} + \alpha \mathbf{I}.$$

Performing the integral in (198) using (4.135) and then taking the logarithm on, we get (5.183).

# **Chapter 6** Kernel Methods

**6.1** We first of all note that  $J(\mathbf{a})$  depends on a only through the form  $\mathbf{K}\mathbf{a}$ . Since typically the number N of data points is greater than the number M of basis functions, the matrix  $\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}}$  will be rank deficient. There will then be M eigenvectors of  $\mathbf{K}$  having non-zero eigenvalues, and N-M eigenvectors with eigenvalue zero. We can then decompose  $\mathbf{a} = \mathbf{a}_{\parallel} + \mathbf{a}_{\perp}$  where  $\mathbf{a}_{\parallel}^{\mathrm{T}}\mathbf{a}_{\perp} = 0$  and  $\mathbf{K}\mathbf{a}_{\perp} = \mathbf{0}$ . Thus the value of  $\mathbf{a}_{\perp}$  is not determined by  $J(\mathbf{a})$ . We can remove the ambiguity by setting  $\mathbf{a}_{\perp} = \mathbf{0}$ , or equivalently by adding a regularizer term

$$rac{\epsilon}{2}\mathbf{a}_{\perp}^{\mathrm{T}}\mathbf{a}_{\perp}$$

to  $J(\mathbf{a})$  where  $\epsilon$  is a small positive constant. Then  $\mathbf{a} = \mathbf{a}_{\parallel}$  where  $\mathbf{a}_{\parallel}$  lies in the span of  $\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}}$  and hence can be written as a linear combination of the columns of  $\mathbf{\Phi}$ , so that in component notation

$$a_n = \sum_{i=1}^M u_i \phi_i(\mathbf{x}_n)$$

or equivalently in vector notation

$$\mathbf{a} = \mathbf{\Phi}\mathbf{u}.\tag{199}$$

Substituting (199) into (6.7) we obtain

$$J(\mathbf{u}) = \frac{1}{2} (\mathbf{K} \mathbf{\Phi} \mathbf{u} - \mathbf{t})^{\mathrm{T}} (\mathbf{K} \mathbf{\Phi} \mathbf{u} - \mathbf{t}) + \frac{\lambda}{2} \mathbf{u}^{\mathrm{T}} \mathbf{\Phi}^{\mathrm{T}} \mathbf{K} \mathbf{\Phi} \mathbf{u}$$
$$= \frac{1}{2} (\mathbf{\Phi} \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \mathbf{u} - \mathbf{t})^{\mathrm{T}} (\mathbf{\Phi} \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \mathbf{u} - \mathbf{t}) + \frac{\lambda}{2} \mathbf{u}^{\mathrm{T}} \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \mathbf{u} \quad (200)$$

Since the matrix  $\Phi^T\Phi$  has full rank we can define an equivalent parametrization given by

$$\mathbf{w} = \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \mathbf{u}$$

and substituting this into (200) we recover the original regularized error function (6.2).

**6.2** Starting with an initial weight vector  $\mathbf{w} = \mathbf{0}$  the Perceptron learning algorithm increments  $\mathbf{w}$  with vectors  $t_n \phi(\mathbf{x}_n)$  where n indexes a pattern which is misclassified by the current model. The resulting weight vector therefore comprises a linear combination of vectors of the form  $t_n \phi(\mathbf{x}_n)$  which we can represent in the form

$$\mathbf{w} = \sum_{n=1}^{N} \alpha_n t_n \phi(\mathbf{x}_n)$$
 (201)

where  $\alpha_n$  is an integer specifying the number of times that pattern n was used to update w during training. The corresponding predictions made by the trained Perceptron are therefore given by

$$y(\mathbf{x}) = \operatorname{sign} (\mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}))$$

$$= \operatorname{sign} \left( \sum_{n=1}^{N} \alpha_n t_n \boldsymbol{\phi}(\mathbf{x}_n)^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}) \right)$$

$$= \operatorname{sign} \left( \sum_{n=1}^{N} \alpha_n t_n k(\mathbf{x}_n, \mathbf{x}) \right).$$

Thus the predictive function of the Perceptron has been expressed purely in terms of the kernel function. The learning algorithm of the Perceptron can similarly be written as

$$\alpha_n \to \alpha_n + 1$$

### 116 Solutions 6.3–6.5

for patterns which are misclassified, in other words patterns which satisfy

$$t_n\left(\mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}(\mathbf{x}_n)\right) \geqslant 0.$$

Using (201) together with  $\alpha_n \geqslant 0$ , this can be written in terms of the kernel function in the form

$$t_n\left(\sum_{m=1}^N k(\mathbf{x}_m, \mathbf{x}_n)\right) \geqslant 0$$

and so the learning algorithm depends only on the elements of the Gram matrix.

**6.3** The distance criterion for the nearest neighbour classifier can be expressed in terms of the kernel as follows

$$D(\mathbf{x}, \mathbf{x}_n) = \|\mathbf{x} - \mathbf{x}_n\|^2$$
$$= \mathbf{x}^{\mathrm{T}} \mathbf{x} + \mathbf{x}_n^{\mathrm{T}} \mathbf{x}_n - 2\mathbf{x}^{\mathrm{T}} \mathbf{x}_n$$
$$= k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}, \mathbf{x}_n)$$

where  $k(\mathbf{x}, \mathbf{x}_n) = \mathbf{x}^T \mathbf{x}_n$ . We then obtain a non-linear kernel classifier by replacing the linear kernel with some other choice of kernel function.

**6.4** An example of such a matrix is

$$\begin{pmatrix} 2 & -2 \\ -3 & 4 \end{pmatrix}$$
.

We can verify this by calculating the determinant of

$$\left(\begin{array}{cc} 2-\lambda & -2 \\ -3 & 4-\lambda \end{array}\right),\,$$

setting the resulting expression equal to zero and solve for the eigenvalues  $\lambda$ , yielding

$$\lambda_1 \simeq 5.65$$
 and  $\lambda_2 \simeq 0.35$ ,

which are both positive.

**6.5** The results (6.13) and (6.14) are easily proved by using (6.1) which defines the kernel in terms of the scalar product between the feature vectors for two input vectors. If  $k_1(\mathbf{x}, \mathbf{x}')$  is a valid kernel then there must exist a feature vector  $\phi(\mathbf{x})$  such that

$$k_1(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}').$$

It follows that

$$ck_1(\mathbf{x}, \mathbf{x}') = \mathbf{u}(\mathbf{x})^T \mathbf{u}(\mathbf{x}')$$

where

$$\mathbf{u}(\mathbf{x}) = c^{1/2} \boldsymbol{\phi}(\mathbf{x})$$

and so  $ck_1(\mathbf{x}, \mathbf{x}')$  can be expressed as the scalar product of feature vectors, and hence is a valid kernel.

Similarly, for (6.14) we can write

$$f(\mathbf{x})k_1(\mathbf{x},\mathbf{x}')f(\mathbf{x}') = \mathbf{v}(\mathbf{x})^{\mathrm{T}}\mathbf{v}(\mathbf{x}')$$

where we have defined

$$\mathbf{v}(\mathbf{x}) = f(\mathbf{x})\phi(\mathbf{x}).$$

Again, we see that  $f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$  can be expressed as the scalar product of feature vectors, and hence is a valid kernel.

Alternatively, these results can be proved be appealing to the general result that the Gram matrix,  $\mathbf{K}$ , whose elements are given by  $k(\mathbf{x}_n, \mathbf{x}_m)$ , should be positive semidefinite for all possible choices of the set  $\{\mathbf{x}_n\}$ , by following a similar argument to Solution 6.7 below.

**6.6** Equation (6.15) follows from (6.13), (6.17) and (6.18).

For (6.16), we express the exponential as a power series, yielding

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$$
  
=  $\sum_{m=0}^{\infty} \frac{(k_1(\mathbf{x}, \mathbf{x}'))^m}{m!}$ .

Since this is a polynomial in  $k_1(\mathbf{x}, \mathbf{x}')$  with positive coefficients, (6.16) follows from (6.15).

**6.7** (6.17) is most easily proved by making use of the result, discussed on page 295, that a necessary and sufficient condition for a function  $k(\mathbf{x}, \mathbf{x}')$  to be a valid kernel is that the Gram matrix  $\mathbf{K}$ , whose elements are given by  $k(\mathbf{x}_n, \mathbf{x}_m)$ , should be positive semidefinite for all possible choices of the set  $\{\mathbf{x}_n\}$ . A matrix  $\mathbf{K}$  is positive semidefinite if, and only if,

$$\mathbf{a}^{\mathrm{T}}\mathbf{K}\mathbf{a}\geqslant 0$$

for any choice of the vector  $\mathbf{a}$ . Let  $\mathbf{K}_1$  be the Gram matrix for  $k_1(\mathbf{x}, \mathbf{x}')$  and let  $\mathbf{K}_2$  be the Gram matrix for  $k_2(\mathbf{x}, \mathbf{x}')$ . Then

$$\mathbf{a}^T(\mathbf{K}_1 + \mathbf{K}_2)\mathbf{a} = \mathbf{a}^T\mathbf{K}_1\mathbf{a} + \mathbf{a}^T\mathbf{K}_2\mathbf{a} \geqslant 0$$

where we have used the fact that  $\mathbf{K}_1$  and  $\mathbf{K}_2$  are positive semi-definite matrices, together with the fact that the sum of two non-negative numbers will itself be non-negative. Thus, (6.17) defines a valid kernel.

To prove (6.18), we take the approach adopted in Solution 6.5. Since we know that  $k_1(\mathbf{x}, \mathbf{x}')$  and  $k_2(\mathbf{x}, \mathbf{x}')$  are valid kernels, we know that there exist mappings  $\phi(\mathbf{x})$  and  $\psi(\mathbf{x})$  such that

$$k_1(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^{\mathrm{T}} \phi(\mathbf{x}')$$
 and  $k_2(\mathbf{x}, \mathbf{x}') = \psi(\mathbf{x})^{\mathrm{T}} \psi(\mathbf{x}')$ .

Hence

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$$

$$= \phi(\mathbf{x})^{\mathrm{T}} \phi(\mathbf{x}') \psi(\mathbf{x})^{\mathrm{T}} \psi(\mathbf{x}')$$

$$= \sum_{m=1}^{M} \phi_m(\mathbf{x}) \phi_m(\mathbf{x}') \sum_{n=1}^{N} \psi_n(\mathbf{x}) \psi_n(\mathbf{x}')$$

$$= \sum_{m=1}^{M} \sum_{n=1}^{N} \phi_m(\mathbf{x}) \phi_m(\mathbf{x}') \psi_n(\mathbf{x}) \psi_n(\mathbf{x}')$$

$$= \sum_{k=1}^{K} \varphi_k(\mathbf{x}) \varphi_k(\mathbf{x}')$$

$$= \varphi(\mathbf{x})^{\mathrm{T}} \varphi(\mathbf{x}'),$$

where K = MN and

$$\varphi_k(\mathbf{x}) = \phi_{((k-1) \otimes N)+1}(\mathbf{x}) \psi_{((k-1) \otimes N)+1}(\mathbf{x}),$$

where in turn  $\oslash$  and  $\odot$  denote integer division and remainder, respectively.

**6.8** If we consider the Gram matrix, **K**, corresponding to the l.h.s. of (6.19), we have

$$(\mathbf{K})_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = k_3 \left( \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \right) = (\mathbf{K}_3)_{ij}$$

where  $\mathbf{K}_3$  is the Gram matrix corresponding to  $k_3(\cdot,\cdot)$ . Since  $k_3(\cdot,\cdot)$  is a valid kernel

$$\mathbf{u}^{\mathrm{T}}\mathbf{K}\mathbf{u} = \mathbf{u}^{\mathrm{T}}\mathbf{K}_{2}\mathbf{u} \geq 0.$$

For (6.20), let  $\mathbf{K} = \mathbf{X}^{\mathrm{T}} \mathbf{A} \mathbf{X}$ , so that  $(\mathbf{K})_{ij} = \mathbf{x}_i^{\mathrm{T}} \mathbf{A} \mathbf{x}_j$ , and consider

$$\mathbf{u}^{\mathrm{T}}\mathbf{K}\mathbf{u} = \mathbf{u}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}}\mathbf{A}\mathbf{X}\mathbf{u}$$
$$= \mathbf{v}^{\mathrm{T}}\mathbf{A}\mathbf{v} \geqslant 0$$

where ,  $\mathbf{v} = \mathbf{X}\mathbf{u}$  and we have used that  $\mathbf{A}$  is positive semidefinite.

- **6.9** Equations (6.21) and (6.22) are special cases of (6.17) and (6.18), respectively, where  $k_a(\cdot,\cdot)$  and  $k_b(\cdot,\cdot)$  only depend on particular elements in their argument vectors. Thus (6.21) and (6.22) follow from the more general results.
- **6.10** Any solution of a linear learning machine based on this kernel must take the form

$$y(\mathbf{x}) = \sum_{n=1}^{N} \alpha_n k(\mathbf{x}_n, \mathbf{x}) = \left(\sum_{n=1}^{N} \alpha_n f(\mathbf{x}_n)\right) f(\mathbf{x}) = Cf(\mathbf{x}).$$

**6.11** As discussed in Solution 6.6, the exponential kernel (6.16) can be written as an infinite sum of terms, each of which can itself be written as an inner product of feature vectors, according to (6.15). Thus, by concatenating the feature vectors of the indvidual terms in that sum, we can write this as an inner product of infinite dimension feature vectors. More formally,

$$\exp \left(\mathbf{x}^{\mathrm{T}}\mathbf{x}'/\sigma^{2}\right) = \sum_{m=0}^{\infty} \phi_{m}(\mathbf{x})^{\mathrm{T}} \phi_{0}(\mathbf{x}')$$
$$= \psi(\mathbf{x})^{\mathrm{T}} \psi(\mathbf{x}')$$

where  $\psi(\mathbf{x})^{\mathrm{T}} = [\phi_0(\mathbf{x})^{\mathrm{T}}, \phi_1(\mathbf{x})^{\mathrm{T}}, \ldots]$ . Hence, we can write (6.23) as

$$k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\varphi}(\mathbf{x})^{\mathrm{T}} \boldsymbol{\varphi}(\mathbf{x}')$$

where

$$oldsymbol{arphi}(\mathbf{x}) = \exp\left(rac{\mathbf{x}^{\mathrm{T}}\mathbf{x}}{\sigma^{2}}
ight) oldsymbol{\psi}(\mathbf{x}).$$

**6.12 NOTE**: In the 1<sup>st</sup> printing of PRML, there is an error in the text relating to this exercise. Immediately following (6.27), it says: |A| denotes the number of *subsets* in A; it should have said: |A| denotes the number of *elements* in A.

Since A may be equal to D (the subset relation was not defined to be strict),  $\phi(D)$  must be defined. This will map to a vector of  $2^{|D|}$  1s, one for each possible subset of D, including D itself as well as the empty set. For  $A \subset D$ ,  $\phi(A)$  will have 1s in all positions that correspond to subsets of A and 0s in all other positions. Therefore,  $\phi(A_1)^{\mathrm{T}}\phi(A_2)$  will count the number of subsets shared by  $A_1$  and  $A_2$ . However, this can just as well be obtained by counting the number of elements in the intersection of  $A_1$  and  $A_2$ , and then raising 2 to this number, which is exactly what (6.27) does.

**6.13** In the case of the transformed parameter  $\psi(\theta)$ , we have

$$\mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) = \mathbf{M}\mathbf{g}_{\boldsymbol{\eta}} \tag{202}$$

where M is a matrix with elements

$$M_{ij} = \frac{\partial \psi_i}{\partial \theta_j}$$

(recall that  $\psi(\theta)$  is assumed to be differentiable) and

$$\mathbf{g}_{\psi} = \nabla_{\psi} \ln p \left( \mathbf{x} | \boldsymbol{\psi}(\boldsymbol{\theta}) \right).$$

The Fisher information matrix then becomes

$$\mathbf{F} = \mathbb{E}_{\mathbf{x}} \left[ \mathbf{M} \mathbf{g}_{\psi} \mathbf{g}_{\psi}^{\mathrm{T}} \mathbf{M}^{\mathrm{T}} \right]$$
$$= \mathbf{M} \mathbb{E}_{\mathbf{x}} \left[ \mathbf{g}_{\psi} \mathbf{g}_{\psi}^{\mathrm{T}} \right] \mathbf{M}^{\mathrm{T}}. \tag{203}$$

Substituting (202) and (203) into (6.33), we get

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{g}_{\psi}^{\mathrm{T}} \mathbf{M}^{\mathrm{T}} \left( \mathbf{M} \mathbb{E}_{\mathbf{x}} \left[ \mathbf{g}_{\psi} \mathbf{g}_{\psi}^{\mathrm{T}} \right] \mathbf{M}^{\mathrm{T}} \right)^{-1} \mathbf{M} \mathbf{g}_{\psi}$$

$$= \mathbf{g}_{\psi}^{\mathrm{T}} \mathbf{M}^{\mathrm{T}} \left( \mathbf{M}^{\mathrm{T}} \right)^{-1} \mathbb{E}_{\mathbf{x}} \left[ \mathbf{g}_{\psi} \mathbf{g}_{\psi}^{\mathrm{T}} \right]^{-1} \mathbf{M}^{-1} \mathbf{M} \mathbf{g}_{\psi}$$

$$= \mathbf{g}_{\psi}^{\mathrm{T}} \mathbb{E}_{\mathbf{x}} \left[ \mathbf{g}_{\psi} \mathbf{g}_{\psi}^{\mathrm{T}} \right]^{-1} \mathbf{g}_{\psi}, \qquad (204)$$

where we have used (C.3) and the fact that  $\psi(\theta)$  is assumed to be invertible. Since  $\theta$  was simply replaced by  $\psi(\theta)$ , (204) corresponds to the original form of (6.33).

**6.14** In order to evaluate the Fisher kernel for the Gaussian we first note that the covariance is assumed to be fixed, and hence the parameters comprise only the elements of the mean  $\mu$ . The first step is to evaluate the Fisher score defined by (6.32). From the definition (2.43) of the Gaussian we have

$$g(\mu, \mathbf{x}) = \nabla_{\mu} \ln \mathcal{N}(\mathbf{x} | \mu, \mathbf{S}) = \mathbf{S}^{-1}(\mathbf{x} - \mu).$$

Next we evaluate the Fisher information matrix using the definition (6.34), giving

$$\mathbf{F} = \mathbb{E}_{\mathbf{x}} \left[ \mathbf{g}(\boldsymbol{\mu}, \mathbf{x}) \mathbf{g}(\boldsymbol{\mu}, \mathbf{x})^T \right] = \mathbf{S}^{-1} \mathbb{E}_{\mathbf{x}} \left[ (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^T \right] \mathbf{S}^{-1}.$$

Here the expectation is with respect to the original Gaussian distribution, and so we can use the standard result

$$\mathbb{E}_{\mathbf{x}}\left[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}}\right] = \mathbf{S}$$

from which we obtain

$$\mathbf{F} = \mathbf{S}^{-1}.$$

Thus the Fisher kernel is given by

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \mathbf{S}^{-1} (\mathbf{x}' - \boldsymbol{\mu}),$$

which we note is just the squared Mahalanobis distance.

**6.15** The determinant for the  $2 \times 2$  Gram matrix

$$\left(\begin{array}{cc} k(x_1, x_1) & k(x_1, x_2) \\ k(x_2, x_1) & k(x_2, x_2) \end{array}\right)$$

equals

$$k(x_1, x_1)k(x_2, x_2) - k(x_1, x_2)^2$$

where we have used the fact that  $k(x_1, x_2) = k(x_2, x_1)$ . Then (6.96) follows directly from the fact that this determinant must be non-negative for a positive semidefinite matrix.

**6.16 NOTE**: In the 1<sup>st</sup> printing of PRML, a detail is missing in this exercise; the text "where  $\mathbf{w}_{\perp}^{\mathrm{T}} \phi(\mathbf{x}_n) = 0$  for all n," should be inserted at the beginning of the line immediately following equation (6.98).

We start by rewriting (6.98) as

$$\mathbf{w} = \mathbf{w}_{\parallel} + \mathbf{w}_{\perp} \tag{205}$$

where

$$\mathbf{w}_{\parallel} = \sum_{n=1}^{N} \alpha_n \boldsymbol{\phi}(\mathbf{x}_n).$$

Note that since  $\mathbf{w}_{\perp}^{\mathrm{T}} \phi(\mathbf{x}_n) = 0$  for all n,

$$\mathbf{w}_{\perp}^{\mathrm{T}}\mathbf{w}_{\parallel} = 0. \tag{206}$$

Using (205) and (206) together with the fact that  $\mathbf{w}_{\perp}^{\mathrm{T}} \phi(\mathbf{x}_n) = 0$  for all n, we can rewrite (6.97) as

$$J(\mathbf{w}) = f\left((\mathbf{w}_{\parallel} + \mathbf{w}_{\perp})^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_{1}), \dots, (\mathbf{w}_{\parallel} + \mathbf{w}_{\perp})^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_{N})\right) + g\left((\mathbf{w}_{\parallel} + \mathbf{w}_{\perp})^{\mathrm{T}} (\mathbf{w}_{\parallel} + \mathbf{w}_{\perp})\right) = f\left(\mathbf{w}_{\parallel}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_{1}), \dots, \mathbf{w}_{\parallel}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_{N})\right) + g\left(\mathbf{w}_{\parallel}^{\mathrm{T}} \mathbf{w}_{\parallel} + \mathbf{w}_{\perp}^{\mathrm{T}} \mathbf{w}_{\perp}\right).$$

Since  $g(\cdot)$  is monotonically increasing, it will have its minimum w.r.t.  $\mathbf{w}_{\perp}$  at  $\mathbf{w}_{\perp} = \mathbf{0}$ , in which case

$$\mathbf{w} = \mathbf{w}_{\parallel} = \sum_{n=1}^{N} \alpha_n \phi(\mathbf{x}_n)$$

as desired.

**6.17 NOTE**: In the 1<sup>st</sup> printing of PRML, there are typographical errors in the text relating to this exercise. In the sentence following immediately after (6.39),  $f(\mathbf{x})$  should be replaced by  $y(\mathbf{x})$ . Also, on the l.h.s. of (6.40),  $y(\mathbf{x}_n)$  should be replaced by  $y(\mathbf{x})$ . There were also errors in Appendix D, which might cause confusion; please consult the errata on the PRML website.

Following the discussion in Appendix D we give a first-principles derivation of the solution. First consider a variation in the function  $y(\mathbf{x})$  of the form

$$y(\mathbf{x}) \to y(\mathbf{x}) + \epsilon \eta(\mathbf{x}).$$

Substituting into (6.39) we obtain

$$E[y + \epsilon \eta] = \frac{1}{2} \sum_{n=1}^{N} \int \left\{ y(\mathbf{x}_n + \boldsymbol{\xi}) + \epsilon \eta(\mathbf{x}_n + \boldsymbol{\xi}) - t_n \right\}^2 \nu(\boldsymbol{\xi}) \, d\boldsymbol{\xi}.$$

Now we expand in powers of  $\epsilon$  and set the coefficient of  $\epsilon$ , which corresponds to the functional first derivative, equal to zero, giving

This must hold for every choice of the variation function  $\eta(\mathbf{x})$ . Thus we can choose

$$\eta(\mathbf{x}) = \delta(\mathbf{x} - \mathbf{z})$$

where  $\delta(\cdot)$  is the Dirac delta function. This allows us to evaluate the integral over  $\xi$  giving

$$\sum_{n=1}^{N} \int \left\{ y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n \right\} \delta(\mathbf{x}_n + \boldsymbol{\xi} - \mathbf{z}) \nu(\boldsymbol{\xi}) \, d\boldsymbol{\xi} = \sum_{n=1}^{N} \left\{ y(\mathbf{z}) - t_n \right\} \nu(\mathbf{z} - \mathbf{x}_n).$$

Substituting this back into (207) and rearranging we then obtain the required result (6.40).

#### **6.18** From the product rule we have

$$p(t|x) = \frac{p(t,x)}{p(x)}.$$

With p(t, x) given by (6.42) and

$$f(x - x_n, t - t_n) = \mathcal{N}\left([x - x_n, t - t_n]^{\mathrm{T}}|\mathbf{0}, \sigma^2 \mathbf{I}\right)$$

this becomes

$$p(t|x) = \frac{\sum_{n=1}^{N} \mathcal{N}\left([x - x_n, t - t_n]^{\mathrm{T}} | \mathbf{0}, \sigma^2 \mathbf{I}\right)}{\int \sum_{m=1}^{N} \mathcal{N}\left([x - x_m, t - t_m]^{\mathrm{T}} | \mathbf{0}, \sigma^2 \mathbf{I}\right) dt}$$
$$= \frac{\sum_{n=1}^{N} \mathcal{N}\left(x - x_n | 0, \sigma^2\right) \mathcal{N}\left(t - t_n | 0, \sigma^2\right)}{\sum_{m=1}^{N} \mathcal{N}\left(x - x_m | 0, \sigma^2\right)}.$$

From (6.46), (6.47), the definition of f(x,t) and the properties of the Gaussian distribution, we can rewrite this as

$$p(t|x) = \sum_{n=1}^{N} k(x, x_n) \mathcal{N} \left( t - t_n | 0, \sigma^2 \right)$$
$$= \sum_{n=1}^{N} k(x, x_n) \mathcal{N} \left( t | t_n, \sigma^2 \right)$$
(208)

where

$$k(x, x_n) = \frac{\mathcal{N}(x - x_n | 0, \sigma^2)}{\sum_{m=1}^{N} \mathcal{N}(x - x_m | 0, \sigma^2)}.$$

We see that this a Gaussian mixture model where  $k(x, x_n)$  play the role of input dependent mixing coefficients.

Using (208) it is straightforward to calculate various expectations:

$$\mathbb{E}[t|x] = \int t \, p(t|x) \, \mathrm{d}t$$

$$= \int t \sum_{n=1}^{N} k(x, x_n) \mathcal{N}\left(t|t_n, \sigma^2\right) \, \mathrm{d}t$$

$$= \sum_{n=1}^{N} k(x, x_n) \int t \, \mathcal{N}\left(t|t_n, \sigma^2\right) \, \mathrm{d}t$$

$$= \sum_{n=1}^{N} k(x, x_n) \, t_n$$

and

$$\operatorname{var}[t|x] = \mathbb{E}\left[\left(t - \mathbb{E}[t|x]\right)^{2}\right]$$

$$= \int \left(t - \mathbb{E}[t|x]\right)^{2} p(t|x) dt$$

$$= \sum_{n=1}^{N} k(x, x_{n}) \int \left(t - \mathbb{E}[t|x]\right)^{2} \mathcal{N}\left(t|t_{n}, \sigma^{2}\right) dt$$

$$= \sum_{n=1}^{N} k(x, x_{n}) \left(\sigma^{2} + t_{n}^{2} - 2t_{n} \mathbb{E}[t|x] + \mathbb{E}[t|x]^{2}\right)$$

$$= \sigma^{2} - \mathbb{E}[t|x]^{2} + \sum_{n=1}^{N} k(x, x_{n}) t_{n}^{2}.$$

**6.19** Changing variables to  $\mathbf{z}_n = \mathbf{x}_n - \boldsymbol{\xi}_n$  we obtain

$$E = \frac{1}{2} \sum_{n=1}^{N} \int [y(\mathbf{z}_n) - t_n]^2 g(\mathbf{x}_n - \mathbf{z}_n) d\mathbf{z}_n.$$

If we set the functional derivative of E with respect to the function  $y(\mathbf{x})$ , for some general value of  $\mathbf{x}$ , to zero using the calculus of variations (see Appendix D) we have

$$\frac{\delta E}{\delta \mathbf{y}(\mathbf{x})} = \sum_{n=1}^{N} \int [y(\mathbf{z}_n) - t_n] g(\mathbf{x}_n - \mathbf{z}_n) \delta(\mathbf{x} - \mathbf{z}_n) d\mathbf{z}_n$$
$$= \sum_{n=1}^{N} [y(\mathbf{x}) - t_n] g(\mathbf{x}_n - \mathbf{x}) = 0.$$

Solving for  $y(\mathbf{x})$  we obtain

$$y(\mathbf{x}) = \sum_{n=1}^{N} k(\mathbf{x}, \mathbf{x}_n) t_n$$
 (209)

where we have defined

$$k(\mathbf{x}, \mathbf{x}_n) = \frac{g(\mathbf{x}_n - \mathbf{x})}{\sum_n g(\mathbf{x}_n - \mathbf{x})}.$$

This an expansion in kernel functions, where the kernels satisfy the summation constraint  $\sum_n k(\mathbf{x}, \mathbf{x}_n) = 1$ .

**6.20** Given the joint distribution (6.64), we can identify  $t_{N+1}$  with  $\mathbf{x}_a$  and  $\mathbf{t}$  with  $\mathbf{x}_b$  in (2.65). Note that this means that we are prepending rather than appending  $t_{N+1}$  to  $\mathbf{t}$  and  $\mathbf{C}_{N+1}$  therefore gets redefined as

$$\mathbf{C}_{N+1} = \left( \begin{array}{cc} c & \mathbf{k}^{\mathrm{T}} \\ \mathbf{k} & \mathbf{C}_{N} \end{array} \right).$$

It then follows that

$$\mu_a = 0$$
  $\mu_b = \mathbf{0}$   $\mathbf{x}_b = \mathbf{t}$ 
 $\Sigma_{aa} = c$   $\Sigma_{bb} = \mathbf{C}_N$   $\Sigma_{ab} = \mathbf{\Sigma}_{ba}^{\mathrm{T}} = \mathbf{k}^{\mathrm{T}}$ 

in (2.81) and (2.82), from which (6.66) and (6.67) follows directly.

**6.21** Both the Gaussian process and the linear regression model give rise to Gaussian predictive distributions  $p(t_{N+1}|\mathbf{x}_{N+1})$  so we simply need to show that these have the same mean and variance. To do this we make use of the expression (6.54) for the kernel function defined in terms of the basis functions. Using (6.62) the covariance matrix  $\mathbf{C}_N$  then takes the form

$$\mathbf{C}_N = \frac{1}{\alpha} \mathbf{\Phi} \mathbf{\Phi}^{\mathrm{T}} + \beta^{-1} \mathbf{I}_N \tag{210}$$

where  $\Phi$  is the design matrix with elements  $\Phi_{nk} = \phi_k(\mathbf{x}_n)$ , and  $\mathbf{I}_N$  denotes the  $N \times N$  unit matrix. Consider first the mean of the Gaussian process predictive distribution, which from (210), (6.54), (6.66) and the definitions in the text preceding (6.66) is given by

$$m_{N+1} = \alpha^{-1} \boldsymbol{\phi}(\mathbf{x}_{N+1})^{\mathrm{T}} \boldsymbol{\Phi}^{\mathrm{T}} \left( \alpha^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^{\mathrm{T}} + \beta^{-1} \mathbf{I}_{N} \right)^{-1} \mathbf{t}.$$

We now make use of the matrix identity (C.6) to give

$$\boldsymbol{\Phi}^{\mathrm{T}} \left( \boldsymbol{\alpha}^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^{\mathrm{T}} + \boldsymbol{\beta}^{-1} \mathbf{I}_{N} \right)^{-1} = \alpha \boldsymbol{\beta} \left( \boldsymbol{\beta} \boldsymbol{\Phi}^{\mathrm{T}} \boldsymbol{\Phi} + \alpha \mathbf{I}_{M} \right)^{-1} \boldsymbol{\Phi}^{\mathrm{T}} = \alpha \boldsymbol{\beta} \mathbf{S}_{N} \boldsymbol{\Phi}^{\mathrm{T}}.$$

Thus the mean becomes

$$m_{N+1} = \beta \boldsymbol{\phi}(\mathbf{x}_{N+1})^{\mathrm{T}} \mathbf{S}_N \boldsymbol{\Phi}^{\mathrm{T}} \mathbf{t}$$

which we recognize as the mean of the predictive distribution for the linear regression model given by (3.58) with  $\mathbf{m}_N$  defined by (3.53) and  $\mathbf{S}_N$  defined by (3.54).

For the variance we similarly substitute the expression (210) for the kernel function into the Gaussian process variance given by (6.67) and then use (6.54) and the definitions in the text preceding (6.66) to obtain

$$\sigma_{N+1}^{2}(\mathbf{x}_{N+1}) = \alpha^{-1}\phi(\mathbf{x}_{N+1})^{\mathrm{T}}\phi(\mathbf{x}_{N+1}) + \beta^{-1}$$

$$-\alpha^{-2}\phi(\mathbf{x}_{N+1})^{\mathrm{T}}\mathbf{\Phi}^{\mathrm{T}}\left(\alpha^{-1}\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}} + \beta^{-1}\mathbf{I}_{N}\right)^{-1}\mathbf{\Phi}\phi(\mathbf{x}_{N+1})$$

$$= \beta^{-1} + \phi(\mathbf{x}_{N+1})^{\mathrm{T}}\left(\alpha^{-1}\mathbf{I}_{M}\right)$$

$$-\alpha^{-2}\mathbf{\Phi}^{\mathrm{T}}\left(\alpha^{-1}\mathbf{\Phi}\mathbf{\Phi}^{\mathrm{T}} + \beta^{-1}\mathbf{I}_{N}\right)^{-1}\mathbf{\Phi}\phi(\mathbf{x}_{N+1}). \tag{211}$$

We now make use of the matrix identity (C.7) to give

$$\alpha^{-1}\mathbf{I}_{M} - \alpha^{-1}\mathbf{I}_{M}\boldsymbol{\Phi}^{\mathrm{T}} \left(\boldsymbol{\Phi}(\alpha^{-1}\mathbf{I}_{M})\boldsymbol{\Phi}^{\mathrm{T}} + \beta^{-1}\mathbf{I}_{N}\right)^{-1}\boldsymbol{\Phi}\alpha^{-1}\mathbf{I}_{M}$$
$$= \left(\alpha\mathbf{I} + \beta\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\right)^{-1} = \mathbf{S}_{N},$$

where we have also used (3.54). Substituting this in (211), we obtain

$$\sigma_N^2(\mathbf{x}_{N+1}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x}_{N+1})^{\mathrm{T}} \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}_{N+1})$$

as derived for the linear regression model in Section 3.3.2.

### **6.22** From (6.61) we have

$$p\left(\left[\begin{array}{c} \mathbf{t}_{1...N} \\ \mathbf{t}_{N+1...N+L} \end{array}\right]\right) = \mathcal{N}\left(\left[\begin{array}{c} \mathbf{t}_{1...N} \\ \mathbf{t}_{N+1...N+L} \end{array}\right] \left|\begin{array}{c} \mathbf{0}, \mathbf{C} \end{array}\right)$$

with C specified by (6.62).

For our purposes, it is useful to consider the following partition<sup>2</sup> of C:

$$\mathbf{C} = \left( \begin{array}{cc} \mathbf{C}_{bb} & \mathbf{C}_{ba} \\ \mathbf{C}_{ab} & \mathbf{C}_{aa} \end{array} \right),$$

where  $C_{aa}$  corresponds to  $\mathbf{t}_{N+1...N+L}$  and  $C_{bb}$  corresponds to  $\mathbf{t}_{1...N}$ . We can use this together with (2.94)–(2.97) and (6.61) to obtain the conditional distribution

$$p(\mathbf{t}_{N+1...N+L}|\mathbf{t}_{1...N}) = \mathcal{N}\left(\mathbf{t}_{N+1...N+L}|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}^{-1}\right)$$
(212)

where, from (2.78)–(2.80),

$$\Lambda_{aa}^{-1} = \mathbf{C}_{aa} - \mathbf{C}_{ab} \mathbf{C}_{bb}^{-1} \mathbf{C}_{ba}$$

$$\Lambda_{ab} = -\Lambda_{aa} \mathbf{C}_{ab} \mathbf{C}_{bb}^{-1}$$
(213)

 $<sup>^{2}</sup>$ The indexing and ordering of this partition have been chosen to match the indexing used in (2.94)–(2.97) as well as the ordering of elements used in the single variate case, as seen in (6.64)–(6.65).

and

$$\boldsymbol{\mu}_{a|b} = -\boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} \mathbf{t}_{1...N} = \mathbf{C}_{ab} \mathbf{C}_{bb}^{-1} \mathbf{t}_{1...N}. \tag{214}$$

Restricting (212) to a single test target, we obtain the corresponding marginal distribution, where  $C_{aa}$ ,  $C_{ba}$  and  $C_{bb}$  correspond to c, k and  $C_N$  in (6.65), respectively. Making the matching substitutions in (213) and (214), we see that they equal (6.67) and (6.66), respectively.

**6.23 NOTE**: In the 1<sup>st</sup> printing of PRML, a typographical mistake appears in the text of the exercise at line three, where it should say "... a training set of input vectors  $\mathbf{x}_1, \ldots, \mathbf{x}_N$ ".

If we assume that the target variables,  $t_1, \ldots, t_D$ , are independent given the input vector,  $\mathbf{x}$ , this extension is straightforward.

Using analogous notation to the univariate case,

$$p(\mathbf{t}_{N+1}|\mathbf{T}) = \mathcal{N}(\mathbf{t}_{N+1}|\mathbf{m}(\mathbf{x}_{N+1}), \sigma(\mathbf{x}_{N+1})\mathbf{I}),$$

where  ${f T}$  is a N imes D matrix with the vectors  ${f t}_1^{
m T}, \ldots, {f t}_N^{
m T}$  as its rows,

$$\mathbf{m}(\mathbf{x}_{N+1})^{\mathrm{T}} = \mathbf{k}^{\mathrm{T}} \mathbf{C}_N \mathbf{T}$$

and  $\sigma(\mathbf{x}_{N+1})$  is given by (6.67). Note that  $\mathbf{C}_N$ , which only depend on the input vectors, is the same in the uni- and multivariate models.

**6.24** Since the diagonal elements of a diagonal matrix are also the eigenvalues of the matrix, **W** is positive definite (see Appendix C). Alternatively, for an arbitrary, non-zero vector **x**,

$$\mathbf{x}^{\mathrm{T}}\mathbf{W}\mathbf{x} = \sum_{i} x_{i}^{2} W_{ii} > 0.$$

If  $\mathbf{x}^T \mathbf{W} \mathbf{x} > 0$  and  $\mathbf{x}^T \mathbf{V} \mathbf{x} > 0$  for an arbitrary, non-zero vector  $\mathbf{x}$ , then

$$\mathbf{x}^{\mathrm{T}}(\mathbf{W} + \mathbf{V})\mathbf{x} = \mathbf{x}^{\mathrm{T}}\mathbf{W}\mathbf{x} + \mathbf{x}^{\mathrm{T}}\mathbf{V}\mathbf{x} > 0.$$

**6.25** Substituting the gradient and the Hessian into the Newton-Raphson formula we obtain

$$\mathbf{a}_{N}^{\text{new}} = \mathbf{a}_{N} + (\mathbf{C}_{N}^{-1} + \mathbf{W}_{N})^{-1} \left[ \mathbf{t}_{N} - \boldsymbol{\sigma}_{N} - \mathbf{C}_{N}^{-1} \mathbf{a}_{N} \right]$$
$$= (\mathbf{C}_{N}^{-1} + \mathbf{W}_{N})^{-1} \left[ \mathbf{t}_{N} - \boldsymbol{\sigma}_{N} + \mathbf{W}_{N} \mathbf{a}_{N} \right]$$
$$= \mathbf{C}_{N} (\mathbf{I} + \mathbf{W}_{N} \mathbf{C}_{N})^{-1} \left[ \mathbf{t}_{N} - \boldsymbol{\sigma}_{N} + \mathbf{W}_{N} \mathbf{a}_{N} \right]$$

**6.26** Using (2.115) the mean of the posterior distribution  $p(a_{N+1}|\mathbf{t}_N)$  is given by

$$\mathbf{k}^{\mathrm{T}} \mathbf{C}_{N}^{-1} \mathbf{a}_{N}^{*}$$
.

Combining this with the condition

$$\mathbf{C}_N^{-1}\mathbf{a}_N^* = \mathbf{t}_N - \boldsymbol{\sigma}_N$$

satisfied by  $\mathbf{a}_{N}^{*}$  we obtain (6.87).

Similarly, from (2.115) the variance of the posterior distribution  $p(a_{N+1}|\mathbf{t}_N)$  is given by

$$\operatorname{var}[a_{N+1}|\mathbf{t}_{N}] = c - \mathbf{k}^{\mathrm{T}} \mathbf{C}_{N}^{-1} \mathbf{k} + \mathbf{k}^{\mathrm{T}} \mathbf{C}_{N}^{-1} \mathbf{C}_{N} (\mathbf{I} + \mathbf{W}_{N} \mathbf{C}_{N})^{-1} \mathbf{C}_{N}^{-1} \mathbf{k}$$

$$= c - \mathbf{k}^{\mathrm{T}} \mathbf{C}_{N}^{-1} \left[ \mathbf{I} - (\mathbf{C}_{N}^{-1} + \mathbf{W}_{N})^{-1} \mathbf{C}_{N}^{-1} \right] \mathbf{k}$$

$$= c - \mathbf{k}^{\mathrm{T}} \mathbf{C}_{N}^{-1} (\mathbf{C}_{N}^{-1} + \mathbf{W}_{N})^{-1} \mathbf{W}_{N} \mathbf{k}$$

$$= c - \mathbf{k}^{\mathrm{T}} (\mathbf{W}_{N}^{-1} + \mathbf{C}_{N})^{-1} \mathbf{k}$$

as required.

**6.27** Using (4.135), (6.80) and (6.85), we can approximate (6.89) as follows:

$$p(\mathbf{t}_{N}|\boldsymbol{\theta}) = \int p(\mathbf{t}_{N}|\mathbf{a}_{N})p(\mathbf{a}_{N}|\boldsymbol{\theta}) d\mathbf{a}_{N}$$

$$\simeq p(\mathbf{t}_{N}|\mathbf{a}_{N}^{\star})p(\mathbf{a}_{N}^{\star}|\boldsymbol{\theta})$$

$$\int \exp\left\{-\frac{1}{2}(\mathbf{a}_{N} - \mathbf{a}_{N}^{\star})^{\mathrm{T}}\mathbf{H}(\mathbf{a}_{N} - \mathbf{a}_{N}^{\star})\right\} d\mathbf{a}_{N}$$

$$= \exp\left(\Psi(\mathbf{a}_{N}^{\star})\right)\frac{(2\pi)^{N/2}}{|\mathbf{H}|^{1/2}}.$$

Taking the logarithm, we obtain (6.90).

To derive (6.91), we gather the terms from (6.90) that involve  $C_N$ , yielding

$$-\frac{1}{2} \left( \mathbf{a}_{N}^{\star \mathrm{T}} \mathbf{C}_{N}^{-1} \mathbf{a}_{N}^{\star} + \ln |\mathbf{C}_{N}| + \ln |\mathbf{W}_{N} + \mathbf{C}_{N}^{-1}| \right)$$

$$= -\frac{1}{2} \mathbf{a}_{N}^{\star \mathrm{T}} \mathbf{C}_{N}^{-1} \mathbf{a}_{N}^{\star} - \frac{1}{2} \ln |\mathbf{C}_{N} \mathbf{W}_{N} + \mathbf{I}|.$$

Applying (C.21) and (C.22) to the first and second terms, respectively, we get (6.91). Applying (C.22) to the l.h.s. of (6.92), we get

$$-\frac{1}{2}\sum_{n=1}^{N} \frac{\partial \ln |\mathbf{W}_{N} + \mathbf{C}_{N}^{-1}|}{\partial a_{n}^{\star}} \frac{\partial a_{n}^{\star}}{\partial \theta_{j}} = -\frac{1}{2}\sum_{n=1}^{N} \operatorname{Tr}\left(\left(\mathbf{W}_{N} + \mathbf{C}_{N}^{-1}\right)^{-1} \frac{\partial \mathbf{W}}{\partial a_{n}^{\star}}\right) \frac{\partial a_{n}^{\star}}{\partial \theta_{j}}$$
$$= -\frac{1}{2}\sum_{n=1}^{N} \operatorname{Tr}\left(\left(\mathbf{C}_{N}\mathbf{W}_{N} + \mathbf{I}\right)^{-1} \mathbf{C}_{N} \frac{\partial \mathbf{W}}{\partial a_{n}^{\star}}\right) \frac{\partial a_{n}^{\star}}{\partial \theta_{j}}. \tag{215}$$

Using the definition of W together with (4.88), we have

$$\frac{\mathrm{d}W_{nn}}{\mathrm{d}a_n^{\star}} = \frac{\mathrm{d}\sigma_n^{\star}(1 - \sigma_n^{\star})}{\mathrm{d}a_n^{\star}}$$
$$= \sigma_n^{\star}(1 - \sigma_n^{\star})^2 - \sigma_n^{\star 2}(1 - \sigma_n^{\star})$$
$$= \sigma_n^{\star}(1 - \sigma_n^{\star})(1 - 2\sigma_n^{\star})$$

and substituting this into (215) we the r.h.s. of (6.92).

Gathering all the terms in (6.93) involving  $\partial a_n^{\star}/\partial \theta_i$  on one side, we get

$$(\mathbf{I} + \mathbf{C}_N \mathbf{W}_N) \frac{\partial a_n^*}{\partial \theta_j} = \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N).$$

Left-multiplying both sides with  $(\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1}$ , we obtain (6.94).

# **Chapter 7** Sparse Kernel Machines

### **7.1** From Bayes' theorem we have

$$p(t|\mathbf{x}) \propto p(\mathbf{x}|t)p(t)$$

where, from (2.249),

$$p(\mathbf{x}|t) = \frac{1}{N_t} \sum_{n=1}^{N} \frac{1}{Z_k} k(\mathbf{x}, \mathbf{x}_n) \delta(t, t_n).$$

Here  $N_t$  is the number of input vectors with label t (+1 or -1) and  $N=N_{+1}+N_{-1}$ .  $\delta(t,t_n)$  equals 1 if  $t=t_n$  and 0 otherwise.  $Z_k$  is the normalisation constant for the kernel. The minimum misclassification-rate is achieved if, for each new input vector,  $\tilde{\mathbf{x}}$ , we chose  $\tilde{t}$  to maximise  $p(\tilde{t}|\tilde{\mathbf{x}})$ . With equal class priors, this is equivalent to maximizing  $p(\tilde{\mathbf{x}}|\tilde{t})$  and thus

$$\tilde{t} = \left\{ \begin{array}{ll} +1 & \text{iff } \frac{1}{N_{+1}} \sum_{i:t_i = +1} k(\tilde{\mathbf{x}}, \mathbf{x}_i) \geqslant \frac{1}{N_{-1}} \sum_{j:t_j = -1} k(\tilde{\mathbf{x}}, \mathbf{x}_j) \\ -1 & \text{otherwise.} \end{array} \right.$$

Here we have dropped the factor  $1/Z_k$  since it only acts as a common scaling factor. Using the encoding scheme for the label, this classification rule can be written in the more compact form

$$\tilde{t} = \operatorname{sign}\left(\sum_{n=1}^{N} \frac{t_n}{N_{t_n}} k(\tilde{\mathbf{x}}, \mathbf{x}_n)\right).$$

Now we take  $k(\mathbf{x}, \mathbf{x}_n) = \mathbf{x}^T \mathbf{x}_n$ , which results in the kernel density

$$p(\mathbf{x}|t=+1) = \frac{1}{N_{+1}} \sum_{n:t=-1} \mathbf{x}^{\mathrm{T}} \mathbf{x}_{n} = \mathbf{x}^{\mathrm{T}} \bar{\mathbf{x}}^{+}.$$

Here, the sum in the middle experssion runs over all vectors  $\mathbf{x}_n$  for which  $t_n = +1$  and  $\bar{\mathbf{x}}^+$  denotes the mean of these vectors, with the corresponding definition for the negative class. Note that this density is improper, since it cannot be normalized.

However, we can still compare likelihoods under this density, resulting in the classification rule

$$\tilde{t} = \begin{cases} +1 & \text{if } \tilde{\mathbf{x}}^{\mathrm{T}} \bar{\mathbf{x}}^{+} \geqslant \tilde{\mathbf{x}}^{\mathrm{T}} \bar{\mathbf{x}}^{-}, \\ -1 & \text{otherwise.} \end{cases}$$

The same argument would of course also apply in the feature space  $\phi(\mathbf{x})$ .

- **7.2** Consider multiplying both sides of (7.5) by  $\gamma > 0$ . Accordingly, we would then replace all occurences of w and b in (7.3) with  $\gamma$ w and  $\gamma$ b, respectively. However, as discussed in the text following (7.3), its solution w.r.t. w and b is invariant to a common scaling factor and hence would remain unchanged.
- **7.3** Given a data set of two data points,  $\mathbf{x}_1 \in \mathcal{C}_+$   $(t_1 = +1)$  and  $\mathbf{x}_2 \in \mathcal{C}_ (t_2 = -1)$ , the maximum margin hyperplane is determined by solving (7.6) subject to the constraints

$$\mathbf{w}^{\mathrm{T}}\mathbf{x}_{1} + b = +1 \tag{216}$$

$$\mathbf{w}^{\mathrm{T}}\mathbf{x}_{2} + b = -1. \tag{217}$$

We do this by introducing Lagrange multipliers  $\lambda$  and  $\eta$ , and solving

$$\underset{\mathbf{w},b}{\operatorname{arg\,min}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \lambda \left( \mathbf{w}^{\mathrm{T}} \mathbf{x}_1 + b - 1 \right) + \eta \left( \mathbf{w}^{\mathrm{T}} \mathbf{x}_2 + b + 1 \right) \right\}.$$

Taking the derivative of this w.r.t. w and b and setting the results to zero, we obtain

$$0 = \mathbf{w} + \lambda \mathbf{x}_1 + \eta \mathbf{x}_2 \tag{218}$$

$$0 = \lambda + \eta. \tag{219}$$

Equation (219) immediately gives  $\lambda = -\eta$ , which together with (218) give

$$\mathbf{w} = \lambda \left( \mathbf{x}_1 - \mathbf{x}_2 \right). \tag{220}$$

For b, we first rearrange and sum (216) and (217) to obtain

$$2b = -\mathbf{w}^{\mathrm{T}} \left( \mathbf{x}_1 + \mathbf{x}_2 \right).$$

Using (220), we can rewrite this as

$$b = -\frac{\lambda}{2} (\mathbf{x}_1 - \mathbf{x}_2)^{\mathrm{T}} (\mathbf{x}_1 + \mathbf{x}_2)$$
$$= -\frac{\lambda}{2} (\mathbf{x}_1^{\mathrm{T}} \mathbf{x}_1 - \mathbf{x}_2^{\mathrm{T}} \mathbf{x}_2).$$

Note that the Lagrange multiplier  $\lambda$  remains undetermined, which reflects the inherent indeterminacy in the magnitude of  $\mathbf{w}$  and b.

**7.4** From Figure 4.1 and (7.4), we see that the value of the margin

$$\rho = \frac{1}{\|\mathbf{w}\|} \quad \text{and so} \quad \frac{1}{\rho^2} = \|\mathbf{w}\|^2.$$

From (7.16) we see that, for the maximum margin solution, the second term of (7.7) vanishes and so we have

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2.$$

Using this together with (7.8), the dual (7.10) can be written as

$$\frac{1}{2} \|\mathbf{w}\|^2 = \sum_{n=1}^{N} a_n - \frac{1}{2} \|\mathbf{w}\|^2,$$

from which the desired result follows.

- **7.5** These properties follow directly from the results obtained in the solution to the previous exercise, 7.4.
- **7.6** If  $p(t = 1|y) = \sigma(y)$ , then

$$p(t = -1|y) = 1 - p(t = 1|y) = 1 - \sigma(y) = \sigma(-y),$$

where we have used (4.60). Thus, given i.i.d. data  $\mathcal{D} = \{(t_1, \mathbf{x}_n), \dots, (t_N, \mathbf{x}_N)\}$ , we can write the corresponding likelihood as

$$p(\mathcal{D}) = \prod_{t_n=1} \sigma(y_n) \prod_{t_{n'}=-1} \sigma(-y_{n'}) = \prod_{n=1}^{N} \sigma(t_n y_n),$$

where  $y_n = y(\mathbf{x}_n)$ , as given by (7.1). Taking the negative logarithm of this, we get

$$-\ln p(\mathcal{D}) = -\ln \prod_{n=1}^{N} \sigma(t_n y_n)$$

$$= \sum_{n=1}^{N} \ln \sigma(t_n y_n)$$

$$= \sum_{n=1}^{N} \ln(1 + \exp(-t_n y_n)),$$

where we have used (4.59). Combining this with the regularization term  $\lambda \|\mathbf{w}\|^2$ , we obtain (7.47).

**7.7** We start by rewriting (7.56) as

$$L = \sum_{n=1}^{N} C\xi_n + \sum_{n=1}^{N} C\widehat{\xi}_n + \frac{1}{2} \mathbf{w}^{\mathrm{T}} \mathbf{w} - \sum_{n=1}^{N} (\mu_n \xi_n + \widehat{\mu}_n \widehat{\xi}_n)$$
$$- \sum_{n=1}^{N} a_n (\epsilon + \xi_n + \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n) + b - t_n)$$
$$- \sum_{n=1}^{N} \widehat{a}_n (\epsilon + \widehat{\xi}_n - \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n) - b + t_n),$$

where we have used (7.1). We now use (7.1), (7.57), (7.59) and (7.60) to rewrite this as

$$L = \sum_{n=1}^{N} (a_n + \mu_n) \xi_n + \sum_{n=1}^{N} (\widehat{a}_n + \widehat{\mu}_n) \widehat{\xi}_n$$

$$+ \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} (a_n - \widehat{a}_n) (a_m - \widehat{a}_m) \phi(\mathbf{x}_n)^{\mathrm{T}} \phi(\mathbf{x}_m) - \sum_{n=1}^{N} (\mu_n \xi_n + \widehat{\mu}_n \widehat{\xi}_n)$$

$$- \sum_{n=1}^{N} (a_n \xi_n + \widehat{a}_n \widehat{\xi}_n) - \epsilon \sum_{n=1}^{N} (a_n + \widehat{a}_n) + \sum_{n=1}^{N} (a_n - \widehat{a}_n) t_n$$

$$- \sum_{n=1}^{N} \sum_{m=1}^{N} (a_n - \widehat{a}_n) (a_m - \widehat{a}_m) \phi(\mathbf{x}_n)^{\mathrm{T}} \phi(\mathbf{x}_m) - b \sum_{n=1}^{N} (a_n - \widehat{a}_n).$$

If we now eliminate terms that cancel out and use (7.58) to eliminate the last term, what we are left with equals the r.h.s. of (7.61).

**7.8** This follows from (7.67) and (7.68), which in turn follow from the KKT conditions, (E.9)–(E.11), for  $\mu_n$ ,  $\xi_n$ ,  $\widehat{\mu}_n$  and  $\widehat{\xi}_n$ , and the results obtained in (7.59) and (7.60). For example, for  $\mu_n$  and  $\xi_n$ , the KKT conditions are

$$\begin{aligned}
\xi_n &\geqslant 0 \\
\mu_n &\geqslant 0 \\
\mu_n \xi_n &= 0
\end{aligned} \tag{221}$$

and from (7.59) we have that

$$\mu_n = C - a_n. \tag{222}$$

Combining (221) and (222), we get (7.67); similar reasoning for  $\widehat{\mu}_n$  and  $\widehat{\xi}_n$  lead to (7.68).



$$\mathbf{v} \Rightarrow \mathbf{t} \quad \mathbf{A} \Rightarrow \mathbf{\Phi} \quad \mathbf{b} \Rightarrow \mathbf{0} \quad \mathbf{L} \Rightarrow \beta \mathbf{I}.$$

in (2.113) and (2.114), upon which the desired result follows from (2.116) and (2.117).

**7.10** We first note that this result is given immediately from (2.113)–(2.115), but the task set in the exercise was to practice the technique of completing the square. In this solution and that of Exercise 7.12, we broadly follow the presentation in Section 3.5.1. Using (7.79) and (7.80), we can write (7.84) in a form similar to (3.78)

$$p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \left(\frac{\beta}{2\pi}\right)^{N/2} \frac{1}{(2\pi)^{N/2}} \prod_{i=1}^{M} \alpha_i \int \exp\left\{-E(\mathbf{w})\right\} d\mathbf{w}$$
 (223)

where

$$E(\mathbf{w}) = \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{w}\|^2 + \frac{1}{2} \mathbf{w}^{\mathrm{T}} \mathbf{A} \mathbf{w}$$

and  $\mathbf{A} = \operatorname{diag}(\boldsymbol{\alpha})$ .

Completing the square over w, we get

$$E(\mathbf{w}) = \frac{1}{2} (\mathbf{w} - \mathbf{m})^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{w} - \mathbf{m}) + E(\mathbf{t})$$
 (224)

where m and  $\Sigma$  are given by (7.82) and (7.83), respectively, and

$$E(\mathbf{t}) = \frac{1}{2} \left( \beta \mathbf{t}^{\mathrm{T}} \mathbf{t} - \mathbf{m}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \mathbf{m} \right). \tag{225}$$

Using (224), we can evaluate the integral in (223) to obtain

$$\int \exp\left\{-E(\mathbf{w})\right\} d\mathbf{w} = \exp\left\{-E(\mathbf{t})\right\} (2\pi)^{M/2} |\mathbf{\Sigma}|^{1/2}.$$
 (226)

Considering this as a function of  $\mathbf{t}$  we see from (7.83), that we only need to deal with the factor  $\exp\{-E(\mathbf{t})\}$ . Using (7.82), (7.83), (C.7) and (7.86), we can re-write (225) as follows

$$E(\mathbf{t}) = \frac{1}{2} (\beta \mathbf{t}^{\mathrm{T}} \mathbf{t} - \mathbf{m}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \mathbf{m})$$

$$= \frac{1}{2} (\beta \mathbf{t}^{\mathrm{T}} \mathbf{t} - \beta \mathbf{t}^{\mathrm{T}} \mathbf{\Phi} \mathbf{\Sigma} \mathbf{\Sigma}^{-1} \mathbf{\Sigma} \mathbf{\Phi}^{\mathrm{T}} \mathbf{t} \beta)$$

$$= \frac{1}{2} \mathbf{t}^{\mathrm{T}} (\beta \mathbf{I} - \beta \mathbf{\Phi} \mathbf{\Sigma} \mathbf{\Phi}^{\mathrm{T}} \beta) \mathbf{t}$$

$$= \frac{1}{2} \mathbf{t}^{\mathrm{T}} (\beta \mathbf{I} - \beta \mathbf{\Phi} (\mathbf{A} + \beta \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi})^{-1} \mathbf{\Phi}^{\mathrm{T}} \beta) \mathbf{t}$$

$$= \frac{1}{2} \mathbf{t}^{\mathrm{T}} (\beta^{-1} \mathbf{I} + \mathbf{\Phi} \mathbf{A}^{-1} \mathbf{\Phi}^{\mathrm{T}})^{-1} \mathbf{t}$$

$$= \frac{1}{2} \mathbf{t}^{\mathrm{T}} \mathbf{C}^{-1} \mathbf{t}.$$

This gives us the last term on the r.h.s. of (7.85); the two preceding terms are given implicitly, as they form the normalization constant for the posterior Gaussian distribution  $p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ .

- **7.11** If we make the same substitutions as in Exercise 7.9, the desired result follows from (2.115).
- **7.12** Using the results (223)–(226) from Solution 7.10, we can write (7.85) in the form of (3.86):

$$\ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{N}{2} \ln \boldsymbol{\beta} + \frac{1}{2} \sum_{i}^{N} \ln \alpha_{i} - E(\mathbf{t}) - \frac{1}{2} \ln |\mathbf{\Sigma}| - \frac{N}{2} \ln(2\pi). \quad (227)$$

By making use of (225) and (7.83) together with (C.22), we can take the derivatives of this w.r.t  $\alpha_i$ , yielding

$$\frac{\partial}{\partial \alpha_i} \ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2\alpha_i} - \frac{1}{2} \Sigma_{ii} - \frac{1}{2} m_i^2.$$
 (228)

Setting this to zero and re-arranging, we obtain

$$\alpha_i = \frac{1 - \alpha_i \Sigma_{ii}}{m_i^2} = \frac{\gamma_i}{m_i^2},$$

where we have used (7.89). Similarly, for  $\beta$  we see that

$$\frac{\partial}{\partial \beta} \ln p(\mathbf{t}|\mathbf{X}, \boldsymbol{\alpha}, \beta) = \frac{1}{2} \left( \frac{N}{\beta} - \|\mathbf{t} - \boldsymbol{\Phi}\mathbf{m}\|^2 - \text{Tr} \left[ \boldsymbol{\Sigma} \boldsymbol{\Phi}^{\mathrm{T}} \boldsymbol{\Phi} \right] \right). \tag{229}$$

Using (7.83), we can rewrite the argument of the trace operator as

$$\Sigma \Phi^{T} \Phi = \Sigma \Phi^{T} \Phi + \beta^{-1} \Sigma \mathbf{A} - \beta^{-1} \Sigma \mathbf{A}$$

$$= \Sigma (\Phi^{T} \Phi \beta + \mathbf{A}) \beta^{-1} - \beta^{-1} \Sigma \mathbf{A}$$

$$= (\mathbf{A} + \beta \Phi^{T} \Phi)^{-1} (\Phi^{T} \Phi \beta + \mathbf{A}) \beta^{-1} - \beta^{-1} \Sigma \mathbf{A}$$

$$= (\mathbf{I} - \mathbf{A} \Sigma) \beta^{-1}.$$
(230)

Here the first factor on the r.h.s. of the last line equals (7.89) written in matrix form. We can use this to set (229) equal to zero and then re-arrange to obtain (7.88).

**7.13** We start by introducing prior distributions over  $\alpha$  and  $\beta$ ,

$$p(\alpha_i) = \operatorname{Gam}(\alpha_i | a_{\alpha 0}, b_{\beta 0}), i = 1, \dots, N,$$
  
 $p(\beta) = \operatorname{Gam}(\beta | a_{\beta 0}, b_{\beta 0}).$ 

Note that we use an independent, common prior for all  $\alpha_i$ . We can then combine this with (7.84) to obtain

$$p(\boldsymbol{\alpha}, \beta, \mathbf{t} | \mathbf{X}) = p(\mathbf{t} | \mathbf{X}, \boldsymbol{\alpha}, \beta) p(\boldsymbol{\alpha}) p(\beta).$$

Rather than maximizing the r.h.s. directly, we first take the logarithm, which enables us to use results from Solution 7.12. Using (227) and (B.26), we get

$$\ln p(\boldsymbol{\alpha}, \beta, \mathbf{t} | \mathbf{X}) = \frac{N}{2} \ln \beta + \frac{1}{2} \sum_{i=1}^{N} \ln \alpha_{i} - E(\mathbf{t}) - \frac{1}{2} \ln |\boldsymbol{\Sigma}| - \frac{N}{2} \ln(2\pi)$$
$$-N \ln \Gamma(a_{\alpha 0})^{-1} + N a_{\alpha 0} \ln b_{\alpha 0} + \sum_{i=1}^{N} ((a_{\alpha 0} - 1) \ln \alpha_{i} - b_{\alpha 0} \alpha_{i})$$
$$- \ln \Gamma(a_{\beta 0})^{-1} + a_{\beta 0} \ln b_{\beta 0} + (a_{\beta 0} - 1) \ln \beta - b_{\beta 0} \beta.$$

Using (228), we obtain the derivative of this w.r.t.  $\alpha_i$  as

$$\frac{\partial}{\partial \alpha_i} \ln p(\boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{t} | \mathbf{X}) = \frac{1}{2\alpha_i} - \frac{1}{2} \Sigma_{ii} - \frac{1}{2} m_i^2 + \frac{a_{\alpha 0} - 1}{\alpha_i} - b_{\alpha 0}.$$

Setting this to zero and rearranging (cf. Solution 7.12) we obtain

$$\alpha_i^{\text{new}} = \frac{\gamma_i + 2a_{\alpha 0} - 2}{m_i^2 - 2b_{\alpha 0}},$$

where we have used (7.89).

For  $\beta$ , we can use (229) together with (B.26) to get

$$\frac{\partial}{\partial \beta} \ln p(\boldsymbol{\alpha}, \beta, \mathbf{t} | \mathbf{X}) = \frac{1}{2} \left( \frac{N}{\beta} - \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}\|^2 - \text{Tr} \left[ \mathbf{\Sigma} \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \right] \right) + \frac{a_{\beta 0} - 1}{\beta} - b_{\beta 0}.$$

Setting this equal to zero and using (7.89) and (230), we get

$$\frac{1}{\beta^{\text{new}}} = \frac{\|\mathbf{t} - \mathbf{\Phi}\mathbf{m}\|^2 + 2b_{\beta 0}}{a_{\beta 0} + 2 + N - \sum_i \gamma_i}.$$

**7.14** If we make the following substitions from (7.81) into (2.113),

$$\mathbf{x} \Rightarrow \mathbf{w} \quad \boldsymbol{\mu} \Rightarrow \mathbf{m} \quad \boldsymbol{\Lambda}^{-1} \Rightarrow \boldsymbol{\Sigma},$$

and from (7.76) and (7.77) into (2.114)

$$\mathbf{y} \Rightarrow t \quad \mathbf{A} \Rightarrow \phi(\mathbf{x})^{\mathrm{T}} \quad \mathbf{b} \Rightarrow \mathbf{0} \quad \mathbf{L} \Rightarrow \beta^{\star} \mathbf{I},$$

(7.90) and (7.91) can be read off directly from (2.115).

**7.15** Using (7.94), (7.95) and (7.97)–(7.99), we can rewrite (7.85) as follows

$$\begin{split} \ln p(\mathbf{t}|\mathbf{X},\boldsymbol{\alpha},\boldsymbol{\beta}) &= -\frac{1}{2} \bigg\{ N \ln(2\pi) + \ln |\mathbf{C}_{-i}| |1 + \alpha_i^{-1} \boldsymbol{\varphi}_i^{\mathrm{T}} \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i| \\ &+ \mathbf{t}^{\mathrm{T}} \left( \mathbf{C}_{-i}^{-1} - \frac{\mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^{\mathrm{T}} \mathbf{C}_{-i}^{-1}}{\alpha_i + \boldsymbol{\varphi}_i^{\mathrm{T}} \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i} \right) \mathbf{t} \bigg\} \\ &= -\frac{1}{2} \left\{ N \ln(2\pi) + \ln |\mathbf{C}_{-i}| + \mathbf{t}^{\mathrm{T}} \mathbf{C}_{-i}^{-1} \mathbf{t} \right\} \\ &+ \frac{1}{2} \left[ -\ln |1 + \alpha_i^{-1} \boldsymbol{\varphi}_i^{\mathrm{T}} \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i| + \mathbf{t}^{\mathrm{T}} \frac{\mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^{\mathrm{T}} \mathbf{C}_{-i}^{-1}}{\alpha_i + \boldsymbol{\varphi}_i^{\mathrm{T}} \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i} \mathbf{t} \right] \\ &= L(\alpha_{-i}) + \frac{1}{2} \left[ \ln \alpha_i - \ln(\alpha_i + s_i) + \frac{q_i^2}{\alpha_i + s_i} \right] \\ &= L(\alpha_{-i}) + \lambda(\alpha_i) \end{split}$$

**7.16** If we differentiate (7.97) twice w.r.t.  $\alpha_i$ , we get

$$\frac{\mathrm{d}^2 \lambda}{\mathrm{d}\alpha_i^2} = -\frac{1}{2} \left( \frac{1}{\alpha_i^2} + \frac{1}{(\alpha_i + s_i)^2} \right).$$

This second derivative must be negative and thus the solution given by (7.101) corresponds to a maximum.

**7.17** Using (7.83), (7.86) and (C.7), we have

$$\mathbf{C}^{-1} = \beta \mathbf{I} - \beta^2 \mathbf{\Phi} \left( \mathbf{A} + \beta \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \right)^{-1} \mathbf{\Phi}^{\mathrm{T}} = \beta \mathbf{I} - \beta^2 \mathbf{\Phi} \mathbf{\Sigma} \mathbf{\Phi}^{\mathrm{T}}.$$

Substituting this into (7.102) and (7.103), we immediately obtain (7.106) and (7.107), respectively.

**7.18** As the RVM can be regarded as a regularized logistic regression model, we can follow the sequence of steps used to derive (4.91) in Exercise 4.13 to derive the first term of the r.h.s. of (7.110), whereas the second term follows from standard matrix derivatives (see Appendix C). Note however, that in Exercise 4.13 we are dealing with the *negative* log-likelhood.

To derive (7.111), we make use of (161) and (162) from Exercise 4.13. If we write the first term of the r.h.s. of (7.110) in component form we get

$$\frac{\partial}{\partial w_j} \sum_{n=1}^{N} (t_n - y_n) \phi_{ni} = -\sum_{n=1}^{N} \frac{\partial y_n}{\partial a_n} \frac{\partial a_n}{\partial w_j} \phi_{ni}$$
$$= -\sum_{n=1}^{N} y_n (1 - y_n) \phi_{nj} \phi_{ni},$$

which, written in matrix form, equals the first term inside the parenthesis on the r.h.s. of (7.111). The second term again follows from standard matrix derivatives.

**7.19 NOTE**: In the 1<sup>st</sup> printing of PRML, on line 1 of the text of this exercise, "approximate log marginal" should be "approximate marginal".

We start by taking the logarithm of (7.114), which, omitting terms that do not depend on  $\alpha$ , leaves us with

$$\ln p(\mathbf{w}^{\star}|\boldsymbol{\alpha}) + \frac{1}{2}\ln |\boldsymbol{\Sigma}| = -\frac{1}{2} \left( \ln |\boldsymbol{\Sigma}^{-1}| + \sum_{i} (w_{i}^{\star})^{2} \alpha_{i} - \ln \alpha_{i} \right),$$

where we have used (7.80). Making use of (7.113) and (C.22), we can differentiate this to obtain (7.115), from which we get (7.116) by using  $\gamma_i = 1 - \alpha_i \Sigma_{ii}$ .

# **Chapter 8** Graphical Models

**8.1** We want to show that, for (8.5),

$$\sum_{x_1} \dots \sum_{x_K} p(\mathbf{x}) = \sum_{x_1} \dots \sum_{x_K} \prod_{k=1}^K p(x_k | pa_k) = 1.$$

We assume that the nodes in the graph has been numbered such that  $x_1$  is the root node and no arrows lead from a higher numbered node to a lower numbered node. We can then marginalize over the nodes in reverse order, starting with  $x_K$ 

$$\sum_{x_1} \dots \sum_{x_K} p(\mathbf{x}) = \sum_{x_1} \dots \sum_{x_K} p(x_K | \mathbf{pa}_K) \prod_{k=1}^{K-1} p(x_k | \mathbf{pa}_k)$$
$$= \sum_{x_1} \dots \sum_{x_{K-1}} \prod_{k=1}^{K-1} p(x_k | \mathbf{pa}_k),$$

since each of the conditional distributions is assumed to be correctly normalized and none of the other variables depend on  $x_K$ . Repeating this process K-2 times we are left with

$$\sum_{x_1} p(x_1|\emptyset) = 1.$$

**8.2** Consider a directed graph in which the nodes of the graph are numbered such that are no edges going from a node to a lower numbered node. If there exists a directed cycle in the graph then the subset of nodes belonging to this directed cycle must also satisfy the same numbering property. If we traverse the cycle in the direction of the edges the node numbers cannot be monotonically increasing since we must end up back at the starting node. It follows that the cycle cannot be a directed cycle.

**Table 1** Comparison of the distribution p(a,b) with the product of marginals p(a)p(b) showing that these are not equal for the given joint distribution p(a,b,c).

	a		ŀ	)	p(a,b)
	0		(	)	336.000
			1		264.000
	1		(	)	256.000
	1	_	1	-	144.000
(	a	i	b		p(a)p(b)
(	)	(	О	(	355200.000
(	)		1		244800.000
	1 0		236800.000		
	1		1		163200.000

**8.3** The distribution p(a, b) is found by summing the complete joint distribution p(a, b, c) over the states of c so that

$$p(a,b) = \sum_{c \in \{0,1\}} p(a,b,c)$$

and similarly the marginal distributions p(a) and p(b) are given by

$$p(a) = \sum_{b \in \{0,1\}} \sum_{c \in \{0,1\}} p(a,b,c) \text{ and } p(b) = \sum_{a \in \{0,1\}} \sum_{c \in \{0,1\}} p(a,b,c).$$
 (231)

Table 1 shows the joint distribution p(a,b) as well as the product of marginals p(a)p(b), demonstrating that these are not equal for the specified distribution.

The conditional distribution p(a, b|c) is obtained by conditioning on the value of c and normalizing

$$p(a,b|c) = \frac{p(a,b,c)}{\sum_{a \in \{0,1\}} \sum_{b \in \{0,1\}} p(a,b,c)}.$$

Similarly for the conditionals p(a|c) and p(b|c) we have

$$p(a|c) = \frac{\sum_{b \in \{0,1\}} p(a,b,c)}{\sum_{a \in \{0,1\}} \sum_{b \in \{0,1\}} p(a,b,c)}$$

and

$$p(b|c) = \frac{\sum_{a \in \{0,1\}} p(a,b,c)}{\sum_{a \in \{0,1\}} \sum_{b \in \{0,1\}} p(a,b,c)}.$$
 (232)

Table 2 compares the conditional distribution p(a,b|c) with the product of marginals p(a|c)p(b|c), showing that these are equal for the given joint distribution p(a,b,c) for both c=0 and c=1.

**8.4** In the previous exercise we have already computed p(a) in (231) and p(b|c) in (232). There remains to compute p(c|a) which is done using

$$p(c|a) = \frac{\sum_{b \in \{0,1\}} p(a,b,c)}{\sum_{b \in \{0,1\}} \sum_{c \in \{0,1\}} p(a,b,c)}.$$

**Table 2** Comparison of the conditional distribution p(a,b|c) with the product of marginals p(a|c)p(b|c) showing that these are equal for the given distribution.

a	b	c	p(a,b c)
0	0	0	0.400
0	1	0	0.100
1	0	0	0.400
1	1	0	0.100
0	0	1	0.277
0	1	1	0.415
1	0	1	0.123
1	1	1	0.185

a	b	c	p(a c)p(b c)
0	0	0	0.400
0	1	0	0.100
1	0	0	0.400
1	1	0	0.100
0	0	1	0.277
0	1	1	0.415
1	0	1	0.123
1	1	1	0.185

The required distributions are given in Table 3.

**Table 3** Tables of p(a), p(c|a) and p(b|c) evaluated by marginalizing and conditioning the joint distribution of Table 8.2.

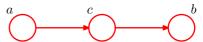
a	p(a)
0	600.000
1	400.000

c	a	p(c a)
0	0	0.400
1	0	0.600
0	1	0.600
1	1	0.400

	b	c	p(b c)
	0	0	0.800
	1	0	0.200
ĺ	0	1	0.400
ĺ	1	1	0.600

Multiplying the three distributions together we recover the joint distribution p(a,b,c) given in Table 8.2, thereby allowing us to verify the validity of the decomposition p(a,b,c)=p(a)p(c|a)p(b|c) for this particular joint distribution. We can express this decomposition using the graph shown in Figure 4.

**Figure 4** Directed graph representing the joint distribution *a* given in Table 8.2.



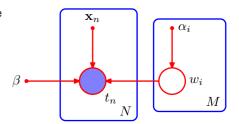
**8.5** NOTE: In PRML, Equation (7.79) contains a typographical error:  $p(t_n|\mathbf{x}_n, \mathbf{w}, \beta^{-1})$  should be  $p(t_n|\mathbf{x}_n, \mathbf{w}, \beta)$ . This correction is provided for completeness only; it does not affect this solution.

The solution is given in Figure 5.

**8.6 NOTE**: In PRML, the text of the exercise should be slightly altered; please consult the PRML errata.

In order to interpret (8.104) suppose initially that  $\mu_0=0$  and that  $\mu_i=1-\epsilon$  where  $\epsilon\ll 1$  for  $i=1,\ldots,K$ . We see that, if all of the  $x_i=0$  then  $p(y=1|x_1,\ldots,x_K)=0$  while if L of the  $x_i=1$  then  $p(y=1|x_1,\ldots,x_K)=1-\epsilon^L$  which is close to 1. For  $\epsilon\to 0$  this represents the logical OR function in which y=1 if one or more of the  $x_i=1$ , and y=0 otherwise. More generally, if just one of the  $x_i=1$  with all remaining  $x_{j\neq i}=0$  then  $p(y=1|x_1,\ldots,x_K)=\mu_i$  and so we can interpret  $\mu_i$  as the probability of y=1 given that only this one  $x_i=1$ . We can similarly interpret  $\mu_0$  as the probability of y=1 when all of the  $x_i=0$ . An example of the application of this model would be in medical diagnosis in which y=1

Figure 5 The graphical representation of the relevance vector machine (RVM); Solution 8.5.



represents the presence or absence of a symptom, and each of the  $x_i$  represents the presence or absence of some disease. For the  $i^{\rm th}$  disease there is a probability  $\mu_i$  that it will give rise to the symptom. There is also a background probability  $\mu_0$  that the symptom will be observed even in the absence of disease. In practice we might observe that the symptom is indeed present (so that y=1) and we wish to infer the posterior probability for each disease. We can do this using Bayes' theorem once we have defined prior probabilities  $p(x_i)$  for the diseases.

#### **8.7** Starting with $\mu$ , (8.11) and (8.15) directly gives

$$\mu_1 = \sum_{j \in \emptyset} w_{1j} \mathbb{E}[x_j] + b_1 = b_1,$$

$$\mu_2 = \sum_{j \in \{x_1\}} w_{2j} \mathbb{E}[x_j] + b_2 = w_{21}b_1 + b_2$$

and

$$\mu_3 = \sum_{j \in \{x_2\}} w_{3j} \mathbb{E}[x_j] + b_3 = w_{32}(w_{21}b_1 + b_2) + b_3.$$

Similarly for  $\Sigma$ , using (8.11) and (8.16), we get

$$\begin{aligned} & \operatorname{cov}[x_1, x_1] = \sum_{k \in \emptyset} w_{1j} \operatorname{cov}[x_1, x_k] + I_{11} v_1 = v_1, \\ & \operatorname{cov}[x_1, x_2] = \sum_{k \in \{x_1\}} w_{2j} \operatorname{cov}[x_1, x_k] + I_{12} v_2 = w_{21} v_1, \\ & \operatorname{cov}[x_1, x_3] = \sum_{k \in \{x_2\}} w_{3j} \operatorname{cov}[x_1, x_k] + I_{13} v_3 = w_{32} w_{21} v_1, \\ & \operatorname{cov}[x_2, x_2] = \sum_{k \in \{x_1\}} w_{2j} \operatorname{cov}[x_2, x_k] + I_{22} v_2 = w_{21}^2 v_1 + v_2, \\ & \operatorname{cov}[x_2, x_3] = \sum_{k \in \{x_2\}} w_{3j} \operatorname{cov}[x_2, x_k] + I_{23} v_3 = w_{32} (w_{21}^2 v_1 + v_2) \end{aligned}$$

and

$$\operatorname{cov}[x_3, x_3] = \sum_{k \in \{x_2\}} w_{3j} \operatorname{cov}[x_3, x_k] + I_{33} v_3 = w_{32}^2 (w_{21}^2 v_1 + v_2) + v_3,$$

where the symmetry of  $\Sigma$  gives the below diagonal elements.

**8.8**  $a \perp \!\!\!\perp b, c \mid d$  can be written as

$$p(a, b, c|d) = p(a|d)p(b, c|d).$$

Summing (or integrating) both sides with respect to c, we obtain

$$p(a, b|d) = p(a|d)p(b|d)$$
 or  $a \perp \!\!\!\perp b \mid d$ ,

as desired.

- **8.9** Consider Figure 8.26. In order to apply the d-separation criterion we need to consider all possible paths from the central node  $x_i$  to all possible nodes external to the Markov blanket. There are three possible categories of such paths. First, consider paths via the parent nodes. Since the link from the parent node to the node  $x_i$  has its tail connected to the parent node, it follows that for any such path the parent node must be either tail-to-tail or head-to-tail with respect to the path. Thus the observation of the parent node will block any such path. Second consider paths via one of the child nodes of node  $x_i$  which do not pass directly through any of the co-parents. By definition such paths must pass to a child of the child node and hence will be head-to-tail with respect to the child node and so will be blocked. The third and final category of path passes via a child node of  $x_i$  and then a co-parent node. This path will be head-to-head with respect to the observed child node and hence will not be blocked by the observed child node. However, this path will either tail-totail or head-to-tail with respect to the co-parent node and hence observation of the co-parent will block this path. We therefore see that all possible paths leaving node  $x_i$  will be blocked and so the distribution of  $x_i$ , conditioned on the variables in the Markov blanket, will be independent of all of the remaining variables in the graph.
- **8.10** From Figure 8.54, we see that

$$p(a, b, c, d) = p(a)p(b)p(c|a, b)p(d|c).$$

Following the examples in Section 8.2.1, we see that

$$p(a,b) = \sum_{c} \sum_{d} p(a,b,c,d)$$
$$= p(a)p(b) \sum_{c} p(c|a,b) \sum_{d} p(d|c)$$
$$= p(a)p(b).$$

Similarly,

$$p(a,b|d) = \frac{\sum_{c} p(a,b,c,d)}{\sum_{a} \sum_{b} \sum_{c} p(a,b,c,d)}$$
$$= \frac{p(d|a,b)p(a)p(b)}{p(d)}$$
$$\neq p(a|d)p(b|d)$$

in general. Note that this result could also be obtained directly from the graph in Figure 8.54 by using d-separation, discussed in Section 8.2.2.

**8.11** The described situation correspond to the graph shown in Figure 8.54 with a=B, b=F, c=G and d=D (cf. Figure 8.21). To evaluate the probability that the tank is empty given the driver's report that the gauge reads zero, we use Bayes' theorem

$$p(F = 0|D = 0) = \frac{p(D = 0|F = 0)p(F = 0)}{p(D = 0)}.$$

To evaluate p(D=0|F=0), we marginalize over B and G,

$$p(D=0|F=0) = \sum_{B,G} p(D=0|G)p(G|B,F=0)p(B) = 0.748$$
 (233)

and to evaluate p(D=0), we marginalize also over F,

$$p(D=0) = \sum_{B,G,F} p(D=0|G)p(G|B,F)p(B)p(F) = 0.352.$$
 (234)

Combining these results with p(F = 0), we get

$$p(F = 0|D = 0) = 0.213.$$

Note that this is slightly lower than the probability obtained in (8.32), reflecting the fact that the driver is not completely reliable.

If we now also observe B=0, we longer marginalize over B in (233) and (234), but instead keep it fixed at its observed value, yielding

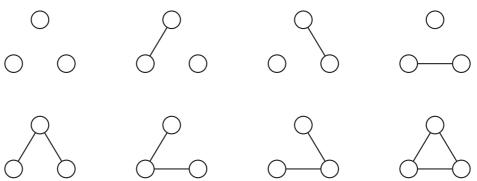
$$p(F = 0|D = 0, B = 0) = 0.110$$

which is again lower than what we obtained with a direct observation of the fuel gauge in (8.33). More importantly, in both cases the value is lower than before we observed B=0, since this observation provides an alternative explanation why the gauge should read zero; see also discussion following (8.33).

In an undirected graph of M nodes there could peterially be a link between each pair of nodes. The number of distinct graphs is then a raised to the power of the number of potential links. To evaluate the number of distinct links, note that there

### **Solutions 8.13–8.15**

are M nodes each of which could have a link to any of the other M-1 nodes, making a total of M(M-1) links. However, each link is counted twice since, in an undirected graph, a link from node a to node b is equivalent to a link from node b to node a. The number of distinct potential links is therefore M(M-1)/2 and so the number of distinct graphs is  $2^{M(M-1)/2}$ . The set of 8 possible graphs over three nodes is shown in Figure 6.



**Figure 6** The set of 8 distinct undirected graphs which can be constructed over M=3 nodes.

**8.13** The change in energy is

$$E(x_j = +1) - E(x_j = -1) = 2h - 2\beta \sum_{i \in ne(j)} x_i - 2\eta y_j$$

where ne(j) denotes the nodes which are neighbours of  $x_j$ .

- **8.14** The most probable configuration corresponds to the configuration with the lowest energy. Since  $\eta$  is a positive constant (and  $h=\beta=0$ ) and  $x_i,y_i\in\{-1,+1\}$ , this will be obtained when  $x_i=y_i$  for all  $i=1,\ldots,D$ .
- **8.15** The marginal distribution  $p(x_{n-1}, x_n)$  is obtained by marginalizing the joint distribution  $p(\mathbf{x})$  over all variables except  $x_{n-1}$  and  $x_n$ ,

$$p(x_{n-1}, x_n) = \sum_{x_1} \dots \sum_{x_{n-2}} \sum_{x_{n+1}} \dots \sum_{x_N} p(\mathbf{x}).$$

This is analogous to the marginal distribution for a single variable, given by (8.50).

Following the same steps as in the single variable case described in Section 8.4.1,

we arrive at a modified form of (8.52),

$$p(x_n) = \frac{1}{Z}$$

$$\underbrace{\left[\sum_{x_{n-2}} \psi_{n-2,n-1}(x_{n-2}, x_{n-1}) \cdots \left[\sum_{x_1} \psi_{1,2}(x_1, x_2)\right] \cdots \right]}_{\mu_{\alpha}(x_{n-1})} \psi_{n-1,n}(x_{n-1}, x_n)$$

$$\underbrace{\left[\sum_{x_{n+1}} \psi_{n,n+1}(x_n, x_{n+1}) \cdots \left[\sum_{x_N} \psi_{N-1,N}(x_{N-1}, x_N)\right] \cdots \right]}_{\mu_{\beta}(x_n)},$$

from which (8.58) immediately follows.

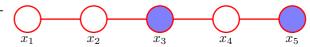
**8.16** Observing  $\mathbf{x}_N = \widehat{\mathbf{x}}_N$  will only change the initial expression (message) for the  $\beta$ -recursion, which now becomes

$$\mu_{\beta}(\mathbf{x}_{N-1}) = \psi_{N-1,N}(\mathbf{x}_{N-1}, \widehat{\mathbf{x}}_N).$$

Note that there is no summation over  $\mathbf{x}_N$ .  $p(\mathbf{x}_n)$  can then be evaluated using (8.54)–(8.57) for all n = 1, ..., N - 1.

**8.17** With N=5 and  $x_3$  and  $x_5$  observed, the graph from Figure 8.38 will look like in Figure 7. This graph is undirected, but from Figure 8.32 we see that the equivalent

Figure 7 The graph discussed in Solution 8.17.



directed graph can be obtained by simply directing all the edges from left to right. (**NOTE**: In PRML, the labels of the two rightmost nodes in Figure 8.32b should be interchanged to be the same as in Figure 8.32a.) In this directed graph, the edges on the path from  $x_2$  to  $x_5$  meet head-to-tail at  $x_3$  and since  $x_3$  is observed, by d-separation  $x_2 \perp \!\!\! \perp x_5 | x_3$ ; note that we would have obtained the same result if we had chosen to direct the arrows from right to left. Alternatively, we could have obtained this result using graph separation in undirected graphs, illustrated in Figure 8.27.

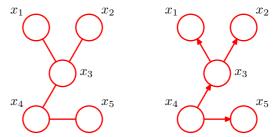
From (8.54), we have

$$p(x_2) = \frac{1}{7}\mu_{\alpha}(x_2)\mu_{\beta}(x_2). \tag{235}$$

 $\mu_{\alpha}(x_2)$  is given by (8.56), while for  $\mu_{\beta}(x_2)$ , (8.57) gives

$$\mu_{\beta}(x_2) = \sum_{x_3} \psi_{2,3}(x_2, x_3) \mu_{\beta}(x_3)$$
$$= \psi_{2,3}(x_2, \widehat{x}_3) \mu_{\beta}(\widehat{x}_3)$$

Figure 8 The graph on the left is an undirected tree. If we pick  $x_4$  to be the root node and direct all the edges in the graph to point from the root to the leaf nodes  $(x_1, x_2)$  and  $(x_1, x_2)$  and  $(x_2, x_3)$ , we obtain the directed tree shown on the right.



since  $x_3$  is observed and we denote the observed value  $\hat{x}_3$ . Thus, any influence that  $x_5$  might have on  $\mu_{\beta}(\hat{x}_3)$  will be in terms of a scaling factor that is indepedent of  $x_2$  and which will be absorbed into the normalization constant Z in (235) and so

$$p(x_2|x_3, x_5) = p(x_2|x_3).$$

**8.18** The joint probability distribution over the variables in a general directed graphical model is given by (8.5). In the particular case of a tree, each node has a single parent, so  $pa_k$  will be a singleton for each node, k, except for the root node for which it will empty. Thus, the joint probability distribution for a tree will be similar to the joint probability distribution over a chain, (8.44), with the difference that the same variable may occur to the right of the conditioning bar in several conditional probability distributions, rather than just one (in other words, although each node can only have one parent, it can have several children). Hence, the argument in Section 8.3.4, by which (8.44) is re-written as (8.45), can also be applied to probability distributions over trees. The result is a Markov random field model where each potential function corresponds to one conditional probability distribution in the directed tree. The prior for the root node, e.g.  $p(x_1)$  in (8.44), can again be incorporated in one of the potential functions associated with the root node or, alternatively, can be incorporated as a single node potential.

This transformation can also be applied in the other direction. Given an undirected tree, we pick a node arbitrarily as the root. Since the graph is a tree, there is a unique path between every pair of nodes, so, starting at root and working outwards, we can direct all the edges in the graph to point from the root to the leaf nodes. An example is given in Figure 8. Since every edge in the tree correspond to a two-node potential function, by normalizing this appropriately, we obtain a conditional probability distribution for the child given the parent.

Since there is a unique path between every pair of nodes in an undirected tree, once we have chosen the root node, the remainder of the resulting directed tree is given. Hence, from an undirected tree with N nodes, we can construct N different directed trees, one for each choice of root node.

**8.19** If we convert the chain model discussed in Section 8.4.1 into a factor graph, each potential function in (8.49) will become a factor. Under this factor graph model,  $p(x_n)$  is given by (8.63) as

$$p(x_n) = \mu_{f_{n-1}, n \to x_n}(x_n) \mu_{f_{n-n+1} \to x_n}(x_n)$$
(236)

where we have adopted the indexing of potential functions from (8.49) to index the factors. From (8.64)–(8.66), we see that

$$\mu_{f_{n-1,n}\to x_n}(x_n) = \sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, x_n) \mu_{x_{n-1}\to f_{n-1,n}}(x_{n-1})$$
 (237)

and

$$\mu_{f_{n,n+1}\to x_n}(x_n) = \sum_{x_{n+1}} \psi_{n,n+1}(x_n, x_{n+1}) \mu_{x_{n+1}\to f_{n,n+1}}(x_{n+1}).$$
 (238)

From (8.69), we further see that

$$\mu_{x_{n-1} \to f_{n-1,n}}(x_{n-1}) = \mu_{f_{n-2,n-1} \to x_{n-1}}(x_{n-1})$$

and

$$\mu_{x_{n+1}\to f_{n,n+1}}(x_{n+1}) = \mu_{f_{n+1,n+2}\to x_{n+1}}(x_{n+1}).$$

Substituting these into (237) and (238), respectively, we get

$$\mu_{f_{n-1,n}\to x_n}(x_n) = \sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, x_n) \mu_{f_{n-2,n-1}\to x_{n-1}}(x_{n-1})$$
 (239)

and

$$\mu_{f_{n,n+1}\to x_n}(x_n) = \sum_{x_{n+1}} \psi_{n,n+1}(x_n, x_{n+1}) \mu_{f_{n+1,n+2}\to x_{n+1}}(x_{n+1}).$$
 (240)

Since the messages are uniquely identified by the index of their arguments and whether the corresponding factor comes before or after the argument node in the chain, we can rename the messages as

$$\mu_{f_{n-2,n-1} \to x_{n-1}}(x_{n-1}) = \mu_{\alpha}(x_{n-1})$$

and

$$\mu_{f_{n+1,n+2} \to x_{n+1}}(x_{n+1}) = \mu_{\beta}(x_{n+1}).$$

Applying these name changes to both sides of (239) and (240), respectively, we recover (8.55) and (8.57), and from these and (236) we obtain (8.54); the normalization constant 1/Z can be easily computed by summing the (unnormalized) r.h.s. of (8.54). Note that the end nodes of the chain are variable nodes which send unit messages to their respective neighbouring factors (cf. (8.56)).

**8.20** We do the induction over the size of the tree and we grow the tree one node at a time while, at the same time, we update the message passing schedule. Note that we can build up any tree this way.

For a single root node, the required condition holds trivially true, since there are no messages to be passed. We then assume that it holds for a tree with N nodes. In the induction step we add a new leaf node to such a tree. This new leaf node need not

to wait for any messages from other nodes in order to send its outgoing message and so it can be scheduled to send it first, before any other messages are sent. Its parent node will receive this message, whereafter the message propagation will follow the schedule for the original tree with N nodes, for which the condition is assumed to hold.

For the propagation of the outward messages from the root back to the leaves, we first follow the propagation schedule for the original tree with N nodes, for which the condition is assumed to hold. When this has completed, the parent of the new leaf node will be ready to send its outgoing message to the new leaf node, thereby completing the propagation for the tree with N+1 nodes.

**8.21 NOTE**: In the 1<sup>st</sup> printing of PRML, this exercise contains a typographical error. On line 2,  $f_x(\mathbf{x}_s)$  should be  $f_s(\mathbf{x}_s)$ .

To compute  $p(\mathbf{x}_s)$ , we marginalize  $p(\mathbf{x})$  over all other variables, analogously to (8.61),

$$p(\mathbf{x}_s) = \sum_{\mathbf{x} \setminus \mathbf{x}_s} p(\mathbf{x}).$$

Using (8.59) and the defintion of  $F_s(x, X_s)$  that followed (8.62), we can write this as

$$p(\mathbf{x}_s) = \sum_{\mathbf{x} \setminus \mathbf{x}_s} f_s(\mathbf{x}_s) \prod_{i \in \text{ne}(f_s)} \prod_{j \in \text{ne}(x_i) \setminus f_s} F_j(x_i, X_{ij})$$

$$= f_s(\mathbf{x}_s) \prod_{i \in \text{ne}(f_s)} \sum_{\mathbf{x} \setminus \mathbf{x}_s} \prod_{j \in \text{ne}(x_i) \setminus f_s} F_j(x_i, X_{ij})$$

$$= f_s(\mathbf{x}_s) \prod_{i \in \text{ne}(f_s)} \mu_{x_i \to f_s}(x_i),$$

where in the last step, we used (8.67) and (8.68). Note that the marginalization over the different sub-trees rooted in the neighbours of  $f_s$  would only run over variables in the respective sub-trees.

**8.22** Let  $X_a$  denote the set of variable nodes in the connected subgraph of interest and  $X_b$  the remaining variable nodes in the full graph. To compute the joint distribution over the variables in  $X_a$ , we need to marginalize  $p(\mathbf{x})$  over  $X_b$ ,

$$p(X_a) = \sum_{X_b} p(\mathbf{x}).$$

We can use the sum-product algorithm to perform this marginalization efficiently, in the same way that we used it to marginalize over all variables but  $x_n$  when computing  $p(x_n)$ . Following the same steps as in the single variable case (see Section 8.4.4),

we can write can write  $p(X_a)$  in a form corresponding to (8.63),

$$p(X_a) = \prod_{s_a} f_{s_a}(X_{s_a}) \prod_{s \in ne} \sum_{X_s} F_s(x_s, X_s)$$

$$= \prod_{s_a} f_{s_a}(X_{s_a}) \prod_{s \in ne} \chi_a \mu_{f_s \to x_s}(x_s).$$
(241)

Here,  $s_a$  indexes factors that only depend on variables in  $X_a$  and so  $X_{s_a} \subseteq X_a$  for all values of  $s_a$ ; s indexes factors that connect  $X_a$  and  $X_b$  and hence also the corresponding nodes,  $x_s \in X_a$ .  $X_s \subseteq X_b$  denotes the variable nodes connected to  $x_s$  via factor  $f_s$ . The messages  $\mu_{f_s \to x_s}(x_s)$  can be computed using the sum-product algorithm, starting from the leaf nodes in, or connected to nodes in,  $X_b$ . Note that the density in (241) may require normalization, which will involve summing the r.h.s. of (241) over all possible combination of values for  $X_a$ .

**8.23** This follows from the fact that the message that a node,  $x_i$ , will send to a factor  $f_s$ , consists of the product of all other messages received by  $x_i$ . From (8.63) and (8.69), we have

$$p(x_i) = \prod_{s \in ne(x_i)} \mu_{f_s \to x_i}(x_i)$$

$$= \mu_{f_s \to x_i}(x_i) \prod_{t \in ne(x_i) \setminus f_s} \mu_{f_t \to x_i}(x_i)$$

$$= \mu_{f_s \to x_i}(x_i) \mu_{x_i \to f_s}(x_i).$$

**8.24 NOTE**: In PRML, this exercise contains a typographical error. On the last line,  $f(\mathbf{x}_s)$  should be  $f_s(\mathbf{x}_s)$ .

See Solution 8.21.

**8.25 NOTE**: In the 1<sup>st</sup> printing of PRML, equation (8.86) contains a typographical error. On the third line, the second summation should sum over  $x_3$ , not  $x_2$ . Furthermore, in equation (8.79), " $\mu_{x_2 \to f_b}$ " (no argument) should be " $\mu_{x_2 \to f_b}(x_2)$ ".

Starting from (8.63), using (8.73), (8.77) and (8.81)–(8.83), we get

$$\widetilde{p}(x_1) = \mu_{f_a \to x_1}(x_1)$$

$$= \sum_{x_2} f_a(x_1, x_2) \mu_{x_2 \to f_a}(x_2)$$

$$= \sum_{x_2} f_a(x_1, x_2) \mu_{f_b \to x_2}(x_2) \mu_{f_c \to x_2}(x_2)$$

$$= \sum_{x_2} f_a(x_1, x_2) \sum_{x_3} f_b(x_2, x_3) \sum_{x_4} f_c(x_2, x_4)$$

$$= \sum_{x_2} \sum_{x_3} \sum_{x_4} f_a(x_1, x_2) f_b(x_2, x_3) f_c(x_2, x_4)$$

$$= \sum_{x_2} \sum_{x_3} \sum_{x_4} \widetilde{p}(\mathbf{x}).$$

Similarly, starting from (8.63), using (8.73), (8.75) and (8.77)–(8.79), we get

$$\widetilde{p}(x_3) = \mu_{f_b \to x_3}(x_3)$$

$$= \sum_{x_2} f_b(x_2, x_3) \mu_{x_2 \to f_b}(x_2)$$

$$= \sum_{x_2} f_b(x_2, x_3) \mu_{f_a \to x_2}(x_2) \mu_{f_c \to x_2}(x_2)$$

$$= \sum_{x_2} f_b(x_2, x_3) \sum_{x_1} f_a(x_1, x_2) \sum_{x_4} f_c(x_2, x_4)$$

$$= \sum_{x_1} \sum_{x_2} \sum_{x_4} f_a(x_1, x_2) f_b(x_2, x_3) f_c(x_2, x_4)$$

$$= \sum_{x_1} \sum_{x_2} \sum_{x_4} \widetilde{p}(\mathbf{x}).$$

Finally, starting from (8.72), using (8.73), (8.74), (8.77), (8.81) and (8.82), we get

$$\begin{split} \widetilde{p}(x_1, x_2) &= f_a(x_1, x_2) \mu_{x_1 \to f_a}(x_1) \mu_{x_2 \to f_a}(x_2) \\ &= f_a(x_1, x_2) \mu_{f_b \to x_2}(x_2) \mu_{f_c \to x_2}(x_2) \\ &= f_a(x_1, x_2) \sum_{x_3} f_b(x_2, x_3) \sum_{x_4} f_b(x_2, x_4) \\ &= \sum_{x_3} \sum_{x_4} f_a(x_1, x_2) f_b(x_2, x_3) f_b(x_2, x_4) \\ &= \sum_{x_3} \sum_{x_4} \widetilde{p}(\mathbf{x}). \end{split}$$

**8.26** We start by using the product and sum rules to write

$$p(x_a, x_b) = p(x_b|x_a)p(x_a) = \sum_{\mathbf{x}_{\backslash ab}} p(\mathbf{x})$$
 (242)

where  $\mathbf{x}_{\backslash ab}$  denote the set of all all variables in the graph except  $x_a$  and  $x_b$ .

We can use the sum-product algorithm from Section 8.4.4 to first evaluate  $p(x_a)$ , by marginalizing over all other variables (including  $x_b$ ). Next we successively fix  $x_a$  at all its allowed values and for each value, we use the sum-product algorithm to evaluate  $p(x_b|x_a)$ , by marginalizing over all variables except  $x_b$  and  $x_a$ , the latter of which will only appear in the formulae at its current, fixed value. Finally, we use (242) to evaluate the joint distribution  $p(x_a, x_b)$ .

**8.27** An example is given by

for which  $\widehat{x} = 2$  and  $\widehat{y} = 2$ .

**8.28** If a graph has one or more cycles, there exists at least one set of nodes and edges such that, starting from an arbitrary node in the set, we can visit all the nodes in the set and return to the starting node, without traversing any edge more than once.

Consider one particular such cycle. When one of the nodes  $n_1$  in the cycle sends a message to one of its neighbours  $n_2$  in the cycle, this causes a pending messages on the edge to the next node  $n_3$  in that cycle. Thus sending a pending message along an edge in the cycle always generates a pending message on the next edge in that cycle. Since this is true for every node in the cycle it follows that there will always exist at least one pending message in the graph.

**8.29** We show this by induction over the number of nodes in the tree-structured factor graph.

First consider a graph with two nodes, in which case only two messages will be sent across the single edge, one in each direction. None of these messages will induce any pending messages and so the algorithm terminates.

We then assume that for a factor graph with N nodes, there will be no pending messages after a finite number of messages have been sent. Given such a graph, we can construct a new graph with N+1 nodes by adding a new node. This new node will have a single edge to the original graph (since the graph must remain a tree) and so if this new node receives a message on this edge, it will induce no pending messages. A message sent from the new node will trigger propagation of messages in the original graph with N nodes, but by assumption, after a finite number of messages have been sent, there will be no pending messages and the algorithm will terminate.

## **Chapter 9 Mixture Models and EM**

**9.1** Since both the E- and the M-step minimise the distortion measure (9.1), the algorithm will never change from a particular assignment of data points to prototypes, unless the new assignment has a lower value for (9.1).

Since there is a finite number of possible assignments, each with a corresponding unique minimum of (9.1) w.r.t. the prototypes,  $\{\mu_k\}$ , the K-means algorithm will converge after a finite number of steps, when no re-assignment of data points to prototypes will result in a decrease of (9.1). When no-reassignment takes place, there also will not be any change in  $\{\mu_k\}$ .

**9.2** Taking the derivative of (9.1), which in this case only involves  $\mathbf{x}_n$ , w.r.t.  $\boldsymbol{\mu}_k$ , we get

$$\frac{\partial J}{\partial \mu_k} = -2r_{nk}(\mathbf{x}_n - \boldsymbol{\mu}_k) = z(\boldsymbol{\mu}_k).$$

Substituting this into (2.129), with  $\mu_k$  replacing  $\theta$ , we get

$$\boldsymbol{\mu}_k^{ ext{new}} = \boldsymbol{\mu}_k^{ ext{old}} + \eta_n (\mathbf{x}_n - \boldsymbol{\mu}_k^{ ext{old}})$$

where by (9.2),  $\mu_k^{\text{old}}$  will be the prototype nearest to  $\mathbf{x}_n$  and the factor of 2 has been absorbed into  $\eta_n$ .

**9.3** From (9.10) and (9.11), we have

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) = \sum_{\mathbf{z}} \prod_{k=1}^{K} (\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_k}.$$

Exploiting the 1-of-K representation for z, we can re-write the r.h.s. as

$$\sum_{j=1}^K \prod_{k=1}^K \left( \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)^{I_{kj}} = \sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

where  $I_{kj} = 1$  if k = j and 0 otherwise.

**9.4** From Bayes' theorem we have

$$p(\boldsymbol{\theta}|\mathbf{X}) = \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{X})}.$$

To maximize this w.r.t.  $\theta$ , we only need to consider the numerator on the r.h.s. and we shall find it more convenient to operate with the logarithm of this expression,

$$ln p(\mathbf{X}|\boldsymbol{\theta}) + ln p(\boldsymbol{\theta}) \tag{243}$$

where we recognize the first term as the l.h.s. of (9.29). Thus we follow the steps in Section 9.3 in dealing with the latent variables, **Z**. Note that the second term in

(243) does not involve **Z** and will not affect the corresponding E-step, which hence gives (9.30). In the M-step, however, we are maximizing over  $\theta$  and so we need to include the second term of (243), yielding

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{\mathrm{old}}) + \ln p(\boldsymbol{\theta}).$$

- **9.5** Consider any two of the latent variable nodes, which we denote  $\mathbf{z}_l$  and  $\mathbf{z}_m$ . We wish to determine whether these variables are independent, conditioned on the observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$  and on the parameters  $\boldsymbol{\mu}, \boldsymbol{\Sigma}$  and  $\boldsymbol{\pi}$ . To do this we consider every possible path from  $\mathbf{z}_l$  to  $\mathbf{z}_m$ . The plate denotes that there are N separate copies of the notes  $\mathbf{z}_n$  and  $\mathbf{x}_n$ . Thus the only paths which connect  $\mathbf{z}_l$  and  $\mathbf{z}_m$  are those which go via one of the parameter nodes  $\boldsymbol{\mu}, \boldsymbol{\Sigma}$  or  $\boldsymbol{\pi}$ . Since we are conditioning on these parameters they represent observed nodes. Furthermore, any path through one of these parameter nodes must be tail-to-tail at the parameter node, and hence all such paths are blocked. Thus  $\mathbf{z}_l$  and  $\mathbf{z}_m$  are independent, and since this is true for any pair of such nodes it follows that the posterior distribution factorizes over the data set.
- **9.6** In this case, the expected complete-data log likelihood function becomes

$$\mathbb{E}_{\mathbf{Z}}\left[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})\right] = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \left\{\ln \pi_k + \ln \mathcal{N}\left(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}\right)\right\}$$

where  $\gamma(z_{nk})$  is defined in (9.16). Differentiating this w.r.t.  $\Sigma^{-1}$ , using (C.24) and (C.28), we get

$$\frac{N}{2} \mathbf{\Sigma} - \frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \left( \mathbf{x}_{n} - \boldsymbol{\mu}_{k} \right) \left( \mathbf{x}_{n} - \boldsymbol{\mu}_{k} \right)^{\mathrm{T}}$$

where we have also used that  $\sum_{k=1}^{K} \gamma(z_{nk}) = 1$  for all n. Setting this equal to zero and rearranging, we obtain

$$\Sigma = \frac{1}{N} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^{\mathrm{T}}.$$

**9.7** Consider first the optimization with respect to the parameters  $\{\mu_k, \Sigma_k\}$ . For this we can ignore the terms in (9.36) which depend on  $\ln \pi_k$ . We note that, for each data point n, the quantities  $z_{nk}$  are all zero except for a particular element which equals one. We can therefore partition the data set into K groups, denoted  $\mathbf{X}_k$ , such that all the data points  $\mathbf{x}_n$  assigned to component k are in group  $\mathbf{X}_k$ . The complete-data log likelihood function can then be written



This represents the sum of K independent terms, one for each component in the mixture. When we maximize this term with respect to  $\mu_k$  and  $\Sigma_k$  we will simply be fitting the  $k^{\rm th}$  component to the data set  $\mathbf{X}_k$ , for which we will obtain the usual maximum likelihood results for a single Gaussian, as discussed in Chapter 2.

For the mixing coefficients we need only consider the terms in  $\ln \pi_k$  in (9.36), but we must introduce a Lagrange multiplier to handle the constraint  $\sum_k \pi_k = 1$ . Thus we maximize

$$\sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} \ln \pi_k + \lambda \left( \sum_{k=1}^{K} \pi_k - 1 \right)$$

which gives

$$0 = \sum_{n=1}^{N} \frac{z_{nk}}{\pi_k} + \lambda.$$

Multiplying through by  $\pi_k$  and summing over k we obtain  $\lambda = -N$ , from which we have

$$\pi_k = \frac{1}{N} \sum_{n=1}^{N} z_{nk} = \frac{N_k}{N}$$

where  $N_k$  is the number of data points in group  $\mathbf{X}_k$ .

**9.8** Using (2.43), we can write the r.h.s. of (9.40) as

$$-\frac{1}{2}\sum_{n=1}^{N}\sum_{j=1}^{K}\gamma(z_{nj})(\mathbf{x}_{n}-\boldsymbol{\mu}_{j})^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}_{n}-\boldsymbol{\mu}_{j})+\mathrm{const.},$$

where 'const.' summarizes terms independent of  $\mu_j$  (for all j). Taking the derivative of this w.r.t.  $\mu_k$ , we get

$$-\sum_{n=1}^{N} \gamma(z_{nk}) \left( \mathbf{\Sigma}^{-1} \boldsymbol{\mu}_{k} - \mathbf{\Sigma}^{-1} \mathbf{x}_{n} \right),\,$$

and setting this to zero and rearranging, we obtain (9.17).

**9.9** If we differentiate (9.40) w.r.t.  $\Sigma_k^{-1}$ , while keeping the  $\gamma(z_{nk})$  fixed, we get

$$\frac{\partial}{\partial \mathbf{\Sigma}^{-1}} \mathbb{E}_{\mathbf{Z}} \left[ \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) \right] = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \frac{1}{2} \left( \mathbf{\Sigma}_{k} - (x_{n} - \boldsymbol{\mu}_{k})(x_{n} - \boldsymbol{\mu}_{k})^{\mathrm{T}} \right)$$

where we have used (C.28). Setting this equal to zero and rearranging, we obtain (9.19).

Appendix E For  $\pi_k$ , we add a Lagrange multiplier term to (9.40) to enforce the constraint

$$\sum_{k=1}^{K} \pi_k = 1$$

yielding

$$\mathbb{E}_{\mathbf{Z}}\left[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})\right] + \lambda \left(\sum_{k=1}^{K} \pi_k - 1\right).$$

Differentiating this w.r.t.  $\pi_k$ , we get

$$\sum_{n=1}^{N} \gamma(z_{nk}) \frac{1}{\pi_k} + \lambda = \frac{N_k}{\pi_k} + \lambda$$

where we have used (9.18). Setting this equal to zero and rearranging, we get

$$N_k = -\pi_k \lambda$$
.

Summing both sides over k, making use of (9.9), we see that  $-\lambda = N$  and thus

$$\pi_k = \frac{N_k}{N}.$$

**9.10** For the mixture model the joint distribution can be written

$$p(\mathbf{x}_a, \mathbf{x}_b) = \sum_{k=1}^{K} \pi_k p(\mathbf{x}_a, \mathbf{x}_b | k).$$

We can find the conditional density  $p(\mathbf{x}_b|\mathbf{x}_a)$  by making use of the relation

$$p(\mathbf{x}_b|\mathbf{x}_a) = \frac{p(\mathbf{x}_a, \mathbf{x}_b)}{p(\mathbf{x}_a)}.$$

For mixture model the marginal density of  $x_a$  is given by

$$p(\mathbf{x}_a) = \sum_{k=1}^K \pi_k p(\mathbf{x}_a|k)$$

where

$$p(\mathbf{x}_a|k) = \int p(\mathbf{x}_a, \mathbf{x}_b|k) \, d\mathbf{x}_b.$$

Thus we can write the conditional density in the form

$$p(\mathbf{x}_b|\mathbf{x}_a) = \frac{\sum_{k=1}^K \pi_k p(\mathbf{x}_a, \mathbf{x}_b|k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_a|j)}.$$

Now we decompose the numerator using

$$p(\mathbf{x}_a, \mathbf{x}_b|k) = p(\mathbf{x}_b|\mathbf{x}_a, k)p(\mathbf{x}_a|k)$$

which allows us finally to write the conditional density as a mixture model of the form

$$p(\mathbf{x}_b|\mathbf{x}_a) = \sum_{k=1}^{K} \lambda_k p(\mathbf{x}_b|\mathbf{x}_a, k)$$
 (244)

where the mixture coefficients are given by

$$\lambda_k \equiv p(k|\mathbf{x}_a) = \frac{\pi_k p(\mathbf{x}_a|k)}{\sum_j \pi_j p(\mathbf{x}_a|j)}$$
(245)

and  $p(\mathbf{x}_b|\mathbf{x}_a, k)$  is the conditional for component k.

**9.11** As discussed in Section 9.3.2,  $\gamma(z_{nk}) \to r_{nk}$  as  $\epsilon \to 0$ .  $\Sigma_k = \epsilon \mathbf{I}$  for all k and are no longer free parameters.  $\pi_k$  will equal the proportion of data points assigned to cluster k and assuming reasonable initialization of  $\pi$  and  $\{\mu_k\}$ ,  $\pi_k$  will remain strictly positive. In this situation, we can maximize (9.40) w.r.t.  $\{\mu_k\}$  independently of  $\pi$ , leaving us with

$$\sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \ln \mathcal{N}\left(\mathbf{x}_{n} | \boldsymbol{\mu}_{k}, \epsilon \mathbf{I}\right) = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \left(-\frac{1}{2\epsilon} \|\mathbf{x}_{n} - \boldsymbol{\mu}_{k}\|^{2}\right) + \text{const.}$$

which equal the negative of (9.1) upto a scaling factor (which is independent of  $\{\mu_k\}$ ).

**9.12** Since the expectation of a sum is the sum of the expectations we have

$$\mathbb{E}[\mathbf{x}] = \sum_{k=1}^K \pi_k \mathbb{E}_k[\mathbf{x}] = \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k$$

where  $\mathbb{E}_k[\mathbf{x}]$  denotes the expectation of  $\mathbf{x}$  under the distribution  $p(\mathbf{x}|k)$ . To find the covariance we use the general relation

$$\mathrm{cov}[\mathbf{x}] = \mathbb{E}[\mathbf{x}\mathbf{x}^\mathrm{T}] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^\mathrm{T}$$

to give

$$cov[\mathbf{x}] = \mathbb{E}[\mathbf{x}\mathbf{x}^{\mathrm{T}}] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^{\mathrm{T}}$$

$$= \sum_{k=1}^{K} \pi_{k} \mathbb{E}_{k}[\mathbf{x}\mathbf{x}^{\mathrm{T}}] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^{\mathrm{T}}$$

$$= \sum_{k=1}^{K} \pi_{k} \left\{ \mathbf{\Sigma}_{k} + \boldsymbol{\mu}_{k} \boldsymbol{\mu}_{k}^{\mathrm{T}} \right\} - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^{\mathrm{T}}.$$

9.13 The expectation of x under the mixture distribution is given by

$$\mathbb{E}[\mathbf{x}] = \sum_{k=1}^{K} \pi_k \mathbb{E}_k[\mathbf{x}] = \sum_{k=1}^{K} \pi_k \boldsymbol{\mu}_k.$$

Now we make use of (9.58) and (9.59) to give

$$\mathbb{E}[\mathbf{x}] = \sum_{k=1}^{K} \pi_k \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) \mathbf{x}_n$$

$$= \sum_{n=1}^{N} \mathbf{x}_n \frac{1}{N} \sum_{k=1}^{K} \gamma(z_{nk})$$

$$= \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n$$

$$= \overline{\mathbf{x}}$$

where we have used  $\pi_k = N_k/N$ , and the fact that  $\gamma(z_{nk})$  are posterior probabilities and hence  $\sum_k \gamma(z_{nk}) = 1$ .

Now suppose we initialize a mixture of Bernoulli distributions by setting the means to a common value  $\mu_k = \widehat{\mu}$  for  $k = 1, \dots, K$  and then run the EM algorithm. In the E-step we first compute the responsibilities which will be given by

$$\gamma(z_{nk}) = \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n | \boldsymbol{\mu}_j)} = \frac{\pi_k}{\sum_{j=1}^K \pi_j} = \pi_k$$

and are therefore independent of n. In the subsequent M-step the revised means are given by

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

$$= \frac{1}{N_k} \pi_k \sum_{n=1}^N \mathbf{x}_n$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

$$= \overline{\mathbf{x}}$$

where again we have made use of  $\pi_k = N_k/N$ . Note that since these are again the same for all k it follows from the previous discussion that the responsibilities on the

next E-step will again be given by  $\gamma(z_{nk}) = \pi_k$  and hence will be unchanged. The revised mixing coefficients are given by

$$\frac{1}{N} \sum_{n=1}^{N} \gamma(z_{nk}) = \pi_k$$

and so are also unchanged. Thus the EM algorithm has converged and no further changes will take place with subsequent E and M steps. Note that this is a degenerate solution in which all of the components of the mixture are identical, and so this distribution is equivalent to a single multivariate Bernoulli distribution.

**9.14** Forming the product of (9.52) and (9.53), we get

$$\prod_{k=1}^{K} p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \prod_{j=1}^{K} \pi_j^{z_j} = \prod_{k=1}^{K} \left( p(\mathbf{x}|\boldsymbol{\mu}_k) \pi_k \right)^{z_k}.$$

If we marginalize this over z, we get

$$\sum_{\mathbf{z}} \prod_{k=1}^{K} (p(\mathbf{x}|\boldsymbol{\mu}_k) \pi_k)^{z_k} = \sum_{j=1}^{K} \prod_{k=1}^{K} (p(\mathbf{x}|\boldsymbol{\mu}_k) \pi_k)^{I_{jk}}$$
$$= \sum_{j=1}^{K} \pi_j p(\mathbf{x}|\boldsymbol{\mu}_j)$$

where we have exploited the 1-of-K coding scheme used for z.

**9.15** This is easily shown by calculating the derivatives of (9.55), setting them to zero and solve for  $\mu_{ki}$ . Using standard derivatives, we get

$$\frac{\partial}{\partial \mu_{ki}} \mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\pi})] = \sum_{n=1}^{N} \gamma(z_{nk}) \left( \frac{x_{ni}}{\mu_{ki}} - \frac{1 - x_{ni}}{1 - \mu_{ki}} \right) \\
= \frac{\sum_{n} \gamma(z_{nk}) x_{ni} - \sum_{n} \gamma(z_{nk}) \mu_{ki}}{\mu_{ki} (1 - \mu_{ki})}$$

Setting this to zero and solving for  $\mu_{ki}$ , we get

$$\mu_{ki} = \frac{\sum_{n} \gamma(z_{nk}) x_{ni}}{\sum_{n} \gamma(z_{nk})},$$

which equals (9.59) when written in vector form.

**9.16** This is identical with the maximization w.r.t.  $\pi_k$  in the Gaussian mixture model, detailed in the second half of Solution 9.9.

- **9.17** This follows directly from the equation for the incomplete log-likelihood, (9.51). The largest value that the argument to the logarithm on the r.h.s. of (9.51) can have is 1, since  $\forall n, k : 0 \leq p(\mathbf{x}_n | \boldsymbol{\mu}_k) \leq 1, \ 0 \leq \pi_k \leq 1 \ \text{and} \ \sum_k^K \pi_k = 1$ . Therefore, the maximum value for  $\ln p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\pi})$  equals 0.
- **9.18** From Solution 9.4, which dealt with MAP estimation for a general mixture model, we know that the E-step will remain unchanged. In the M-step we maximize

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) + \ln p(\boldsymbol{\theta})$$

which in the case of the given model becomes,

$$\sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^{D} \left[ x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki}) \right] \right\}$$

$$+ \sum_{j=1}^{K} \sum_{i'=1}^{D} \left\{ (a_j - 1) \ln \mu_{ji'} + (b_j - 1) \ln(1 - \mu_{ji'}) \right\} + \sum_{l=1}^{K} (\alpha_l - 1) \ln \pi_l$$
 (246)

where we have used (9.55), (2.13) and (2.38), and we have dropped terms independent of  $\{\mu_k\}$  and  $\pi$ . Note that we have assumed that each parameter  $\mu_{ki}$  has the same prior for each i, but this can differ for different components k.

Differentiating (246) w.r.t.  $\mu_{ki}$  yields

$$\sum_{n=1}^{N} \gamma(z_{nk}) \left\{ \frac{x_{ni}}{\mu_{ki}} - \frac{1 - x_{ni}}{1 - \mu_{ki}} \right\} + \frac{a_k}{\mu_{ki}} - \frac{1 - b_k}{1 - \mu_{ki}}$$

$$= \frac{N_k \overline{x}_{ki} + a - 1}{\mu_{ki}} - \frac{N_k - N_k \overline{x}_{ki} + b - 1}{1 - \mu_{ki}}$$

where  $N_k$  is given by (9.57) and  $\overline{x}_{ki}$  is the  $i^{\text{th}}$  element of  $\overline{\mathbf{x}}$  defined in (9.58). Setting this equal to zero and rearranging, we get

$$\mu_{ki} = \frac{N_k \overline{x}_{ki} + a - 1}{N_k + a - 1 + b - 1}.$$
(247)

Note that if  $a_k = b_k = 1$  for all k, this reduces to the standard maximum likelihood result. Also, as N becomes large, (247) will approach the maximum likelihood result.

When maximizing w.r.t.  $\pi_k$ , we need to enforce the constraint  $\sum_k \pi_k = 1$ , which we do by adding a Lagrange multiplier term to (246). Dropping terms independent of  $\pi$  we are left with

$$\sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \ln \pi_k + \sum_{l=1}^{K} (\alpha_l - 1) \ln \pi_l + \lambda \left( \sum_{j=1}^{K} \pi_j - 1 \right).$$

### Appendix E

Differentiating this w.r.t.  $\pi_k$ , we get

$$\frac{N_k + \alpha_k - 1}{\pi_k} + \lambda$$

and setting this equal to zero and rearranging, we have

$$N_k + \alpha_k - 1 = -\lambda \pi_k.$$

Summing both sides over k, using  $\sum_k \pi_k = 1$ , we see that  $-\lambda = N + \alpha_0 - K$ , where  $\alpha_0$  is given by (2.39), and thus

$$\pi_k = \frac{N_k + \alpha_k - 1}{N + \alpha_0 - K}. (248)$$

Also in this case, if  $\alpha_k=1$  for all k, we recover the maximum likelihood result exactly. Similarly, as N gets large, (248) will approach the maximum likelihood result.

**9.19** As usual we introduce a latent variable  $\mathbf{z}_n$  corresponding to each observation. The conditional distribution of the observed data set, given the latent variables, is then

$$p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}) = \prod_{n=1}^{N} p(\mathbf{x}_n | \boldsymbol{\mu}_k)^{z_{nk}}.$$

Similarly, the distribution of the latent variables is given by

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

The expected value of the complete-data log likelihood function is given by

$$\sum_{n=1}^{N} \sum_{k=1}^{k} \gamma(z_{nk}) \left\{ \ln \pi_k + \sum_{i=1}^{D} \sum_{j=1}^{M} x_{nij} \ln \mu_{kij} \right\}$$

where as usual we have defined responsibilities given by

$$\gamma(z_{nk}) = \mathbb{E}[z_{nk}] = \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k)}{\sum_{j=1}^{M} \pi_j p(\mathbf{x}_n | \boldsymbol{\mu}_j)}.$$

These represent the E-step equations.

To derive the M-step equations we add to the expected complete-data log likelihood function a set of Lagrange multiplier terms given by

$$\lambda \left( \sum_{k=1}^{K} \pi_k - 1 \right) + \sum_{k=1}^{K} \sum_{i=1}^{D} \eta_{ki} \left( \sum_{j=1}^{M} \mu_{kij} - 1 \right)$$

to enforce the constraint  $\sum_k \pi_k = 1$  as well as the set of constraints

$$\sum_{j=1}^{M} \mu_{kij} = 1$$

for all values of i and k. Maximizing with respect to the mixing coefficients  $\pi_k$ , and eliminating the Lagrange multiplier  $\lambda$  in the usual way, we obtain

$$\pi_k = \frac{N_k}{N}$$

where we have defined

$$N_k = \sum_{n=1}^{N} \gamma(z_{nk}).$$

Similarly maximizing with respect to the parameters  $\mu_{kij}$ , and again eliminating the Lagrange multipliers, we obtain

$$\mu_{kij} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma(z_{nk}) x_{nij}.$$

This is an intuitively reasonable result which says that the value of  $\mu_{kij}$  for component k is given by the fraction of those counts assigned to component k which have non-zero values of the corresponding elements i and j.

**9.20** If we take the derivatives of (9.62) w.r.t.  $\alpha$ , we get

$$\frac{\partial}{\partial \alpha} \mathbb{E}\left[\ln p(\mathbf{t}, \mathbf{w} | \alpha, \beta)\right] = \frac{M}{2} \frac{1}{\alpha} - \frac{1}{2} \mathbb{E}\left[\mathbf{w}^{\mathrm{T}} \mathbf{w}\right].$$

Setting this equal to zero and re-arranging, we obtain (9.63).

**9.21** Taking the derivative of (9.62) w.r.t.  $\beta$ , we obtain

$$\frac{\partial}{\partial \beta} \mathbb{E}\left[\ln p(\mathbf{t}, \mathbf{w} | \alpha, \beta)\right] = \frac{N}{2} \frac{1}{\beta} - \frac{1}{2} \sum_{n=1}^{N} \mathbb{E}\left[\left(t_n - \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}_n\right)^2\right]. \tag{249}$$

From (3.49)–(3.51), we see that

$$\mathbb{E}\left[(t_n - \mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}_n)^2\right] = \mathbb{E}\left[t_n^2 - 2t_n\mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}_n + \mathrm{Tr}[\boldsymbol{\phi}_n\boldsymbol{\phi}_n^{\mathrm{T}}\mathbf{w}\mathbf{w}^{\mathrm{T}}]\right]$$

$$= t_n^2 - 2t_n\mathbf{m}_N^{\mathrm{T}}\boldsymbol{\phi}_n + \mathrm{Tr}\left[\boldsymbol{\phi}_n\boldsymbol{\phi}_n^{\mathrm{T}}(\mathbf{m}_N\mathbf{m}_N^{\mathrm{T}} + \mathbf{S}_N)\right]$$

$$= (t_n - \mathbf{m}_N^{\mathrm{T}}\boldsymbol{\phi}_n)^2 + \mathrm{Tr}\left[\boldsymbol{\phi}_n\boldsymbol{\phi}_n^{\mathrm{T}}\mathbf{S}_N\right].$$

Substituting this into (249) and rearranging, we obtain

$$\frac{1}{\beta} = \frac{1}{N} \left( \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N\|^2 + \text{Tr} \left[ \mathbf{\Phi}^{\text{T}} \mathbf{\Phi} \mathbf{S}_N \right] \right).$$

**9.22 NOTE**: In PRML, a pair of braces is missing from (9.66), which should read

$$\mathbb{E}_{\mathbf{w}} \left[ \ln \left\{ p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\alpha) \right\} \right].$$

Moreover  $m_N$  should be m in the numerator on the r.h.s. of (9.68).

Using (7.76)–(7.83) and associated definitions, we can rewrite (9.66) as

$$\mathbb{E}_{\mathbf{w}} \left[ \ln \mathcal{N} \left( \mathbf{t} | \mathbf{\Phi} \mathbf{w}, \beta^{-1} \mathbf{I} \right) + \ln \mathcal{N} \left( \mathbf{w} | \mathbf{0}, \mathbf{A}^{-1} \right) \right]$$

$$= \frac{1}{2} \mathbb{E}_{\mathbf{w}} \left[ N \ln \beta - \beta || \mathbf{t} - \mathbf{\Phi} \mathbf{w} ||^{2} + \sum_{i=1}^{M} \ln \alpha_{i} - \text{Tr} \left[ \mathbf{A} \mathbf{w} \mathbf{w}^{T} \right] \right] + \text{const}$$

$$= \frac{1}{2} \left( N \ln \beta - \beta \left( || \mathbf{t} - \mathbf{\Phi} \mathbf{m} ||^{2} + \text{Tr} \left[ \mathbf{\Phi}^{T} \mathbf{\Phi} \mathbf{\Sigma} \right] \right) + \sum_{i=1}^{M} \ln \alpha_{i} - \text{Tr} \left[ \mathbf{A} \left( \mathbf{m} \mathbf{m}^{T} + \mathbf{\Sigma} \right) \right] \right) + \text{const}.$$
(250)

Differentiating this w.r.t.  $\alpha_i$ , using (C.23), and setting the result equal to zero, we get

$$\frac{1}{2}\frac{1}{\alpha_i} - \frac{1}{2}\left(m_i^2 + \Sigma_{ii}\right) = 0$$

which we can rearrange to obtain (9.67).

Differentiating (250) w.r.t.  $\beta$  and setting the result equal to zero we get

$$\frac{N}{2} \frac{1}{\beta} - \frac{1}{2} \left( \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}\|^2 + \text{Tr} \left[ \mathbf{\Phi}^{T} \mathbf{\Phi} \mathbf{\Sigma} \right] \right) = 0.$$
 (251)

Using (7.83), (C.6) and (C.7) together with the fact that **A** is diagonal, we can rewrite  $\Phi^{T}\Phi\Sigma$  as follows:

$$\begin{split} \boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\boldsymbol{\Sigma} &= \boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\mathbf{A}^{-1}\left(\mathbf{I} + \beta\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\mathbf{A}^{-1}\right)^{-1} \\ &= \boldsymbol{\Phi}^{\mathrm{T}}\left(\mathbf{I} + \beta\boldsymbol{\Phi}\mathbf{A}^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\right)^{-1}\boldsymbol{\Phi}\mathbf{A}^{-1} \\ &= \beta\left(\mathbf{I} - \mathbf{I} + \boldsymbol{\Phi}^{\mathrm{T}}\left(\beta^{-1}\mathbf{I} + \boldsymbol{\Phi}\mathbf{A}^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\right)^{-1}\boldsymbol{\Phi}\mathbf{A}^{-1}\right) \\ &= \beta\left(\mathbf{I} - \mathbf{A}\left(\mathbf{A}^{-1} + \mathbf{A}^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\left(\beta^{-1}\mathbf{I} + \boldsymbol{\Phi}\mathbf{A}^{-1}\boldsymbol{\Phi}^{\mathrm{T}}\right)^{-1}\boldsymbol{\Phi}\mathbf{A}^{-1}\right)\right) \\ &= \beta\left(\mathbf{I} - \mathbf{A}\left(\mathbf{A} + \beta\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\right)^{-1}\right) = \beta\left(\mathbf{I} - \mathbf{A}\boldsymbol{\Sigma}\right). \end{split}$$

Using this together with (7.89), we obtain (9.68) from (251).

**9.23 NOTE**: In the 1<sup>st</sup> printing of PRML, the task set in this exercise is to show that the two sets of re-estimation equations are formally equivalent, without any restriction. However, it really should be restricted to stationary points of the objective function.

Considering the case when the optimization has converged, we can start with  $\alpha_i$ , as defined by (7.87), and use (7.89) to re-write this as

$$\alpha_i^{\star} = \frac{1 - \alpha_i^{\star} \Sigma_{ii}}{m_N^2},$$

where  $\alpha_i^\star=\alpha_i^{\rm new}=\alpha_i$  is the value reached at convergence. We can re-write this as

$$\alpha_i^{\star}(m_i^2 + \Sigma_{ii}) = 1$$

which is easily re-written as (9.67).

For  $\beta$ , we start from (9.68), which we re-write as

$$\frac{1}{\beta^{\star}} = \frac{\|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N\|^2}{N} + \frac{\sum_i \gamma_i}{\beta^{\star} N}.$$

As in the  $\alpha$ -case,  $\beta^\star=\beta^{\rm new}=\beta$  is the value reached at convergence. We can re-write this as

$$\frac{1}{\beta^{\star}} \left( N - \sum_{i} \gamma_{i} \right) = \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_{N}\|^{2},$$

which can easily be re-written as (7.88).

- **9.24** This is analogous to Solution 10.1, with the integrals replaced by sums.
- **9.25** This follows from the fact that the Kullback-Leibler divergence, KL(q||p), is at its minimum, 0, when q and p are identical. This means that

$$\frac{\partial}{\partial \boldsymbol{\theta}} \mathrm{KL}(q \| p) = \mathbf{0},$$

since  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$  depends on  $\boldsymbol{\theta}$ . Therefore, if we compute the gradient of both sides of (9.70) w.r.t.  $\boldsymbol{\theta}$ , the contribution from the second term on the r.h.s. will be  $\mathbf{0}$ , and so the gradient of the first term must equal that of the l.h.s.

**9.26** From (9.18) we get

$$N_k^{\text{old}} = \sum_n \gamma^{\text{old}}(z_{nk}). \tag{252}$$

We get  $N_k^{\text{new}}$  by recomputing the responsibilities,  $\gamma(z_{mk})$ , for a specific data point,  $\mathbf{x}_m$ , yielding

$$N_k^{\text{new}} = \sum_{n \neq m} \gamma^{\text{old}}(z_{nk}) + \gamma^{\text{new}}(z_{mk}).$$
 (253)

Combining this with (252), we get (9.79).

Similarly, from (9.17) we have

# 欢迎关注公众号@机器学习点集法之道

and recomputing the responsibilities,  $\gamma(z_{mk})$ , we get

$$\mu_{k}^{\text{new}} = \frac{1}{N_{k}^{\text{new}}} \left( \sum_{n \neq m} \gamma^{\text{old}}(z_{nk}) \mathbf{x}_{n} + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_{m} \right)$$

$$= \frac{1}{N_{k}^{\text{new}}} \left( N_{k}^{\text{old}} \boldsymbol{\mu}_{k}^{\text{old}} - \gamma^{\text{old}}(z_{mk}) \mathbf{x}_{m} + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_{m} \right)$$

$$= \frac{1}{N_{k}^{\text{new}}} \left( \left( N_{k}^{\text{new}} - \gamma^{\text{new}}(z_{mk}) + \gamma^{\text{old}}(z_{mk}) \right) \boldsymbol{\mu}_{k}^{\text{old}} \right)$$

$$- \gamma^{\text{old}}(z_{mk}) \mathbf{x}_{m} + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_{m}$$

$$= \boldsymbol{\mu}_{k}^{\text{old}} + \left( \frac{\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N_{k}^{\text{new}}} \right) (\mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{old}}),$$

where we have used (9.79).

**9.27** Following the treatment of  $\mu_k$  in Solution 9.26, (9.19) gives

$$\boldsymbol{\Sigma}_k^{\text{old}} = \frac{1}{N_k^{\text{old}}} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{old}})^{\text{T}}$$

where  $N_k^{\rm old}$  is given by (252). Recomputing the responsibilities  $\gamma(z_{mk})$ , and using

$$\begin{split} \boldsymbol{\Sigma}_{k}^{\text{new}} &= \frac{1}{N_{k}^{\text{new}}} \left( \sum_{n \neq m} \gamma^{\text{old}}(z_{nk}) \left( \mathbf{x}_{n} - \boldsymbol{\mu}_{k}^{\text{old}} \right) \left( \mathbf{x}_{n} - \boldsymbol{\mu}_{k}^{\text{old}} \right)^{\text{T}} \right. \\ &+ \gamma^{\text{new}}(z_{mk}) \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{new}} \right) \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{new}} \right)^{\text{T}} \right) \\ &= \frac{1}{N_{k}^{\text{new}}} \left( N_{k}^{\text{old}} \boldsymbol{\Sigma}_{k}^{\text{old}} - \gamma^{\text{old}}(z_{mk}) \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{old}} \right) \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{old}} \right)^{\text{T}} \right. \\ &+ \gamma^{\text{new}}(z_{mk}) \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{new}} \right) \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{new}} \right)^{\text{T}} \right) \\ &= \boldsymbol{\Sigma}_{k}^{\text{old}} - \frac{\gamma^{\text{old}}(z_{mk})}{N_{k}^{\text{new}}} \left( \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{old}} \right) \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{old}} \right)^{\text{T}} - \boldsymbol{\Sigma}_{k}^{\text{old}} \right) \\ &+ \frac{\gamma^{\text{new}}(z_{mk})}{N_{k}^{\text{new}}} \left( \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{new}} \right) \left( \mathbf{x}_{m} - \boldsymbol{\mu}_{k}^{\text{new}} \right)^{\text{T}} - \boldsymbol{\Sigma}_{k}^{\text{old}} \right) \end{split}$$

where we have also used (9.79).

For  $\pi_k$ , (9.22) gives

$$\pi_k^{\text{old}} = \frac{N_k^{\text{old}}}{N} = \frac{1}{N} \sum_{n=1}^N \gamma^{\text{old}}(z_{nk})$$

and thus recomputing  $\gamma(z_{nk})$  we get

$$\pi_k^{\text{new}} = \frac{1}{N} \left( \sum_{n \neq m}^N \gamma^{\text{old}}(z_{nk}) + \gamma^{\text{new}}(z_{mk}) \right)$$

$$= \frac{1}{N} \left( N \pi_k^{\text{old}} - \gamma^{\text{old}}(z_{mk}) + \gamma^{\text{new}}(z_{mk}) \right)$$

$$= \pi_k^{\text{old}} - \frac{\gamma^{\text{old}}(z_{mk})}{N} + \frac{\gamma^{\text{new}}(z_{mk})}{N}.$$

## **Chapter 10** Approximate Inference

**10.1** Starting from (10.3), we use the product rule together with (10.4) to get

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

$$= \int q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X} \mid \mathbf{Z}) p(\mathbf{X})}{q(\mathbf{Z})} \right\} d\mathbf{Z}$$

$$= \int q(\mathbf{Z}) \left( \ln \left\{ \frac{p(\mathbf{X} \mid \mathbf{Z})}{q(\mathbf{Z})} \right\} + \ln p(\mathbf{X}) \right) d\mathbf{Z}$$

$$= -KL(q \mid\mid p) + \ln p(\mathbf{X}).$$

Rearranging this, we immediately get (10.2).

**10.2** By substituting  $\mathbb{E}[z_1] = m_1 = \mu_1$  and  $\mathbb{E}[z_2] = m_2 = \mu_2$  in (10.13) and (10.15), respectively, we see that both equations are satisfied and so this is a solution.

To show that it is indeed the only solution when  $p(\mathbf{z})$  is non-singular, we first substitute  $\mathbb{E}[z_1] = m_1$  and  $\mathbb{E}[z_2] = m_2$  in (10.13) and (10.15), respectively. Next, we substitute the r.h.s. of (10.13) for  $m_1$  in (10.15), yielding

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21} \left( \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} \left( m_2 - \mu_2 \right) - \mu_1 \right)$$
  
=  $\mu_2 - \Lambda_{22}^{-1} \Lambda_{21} \Lambda_{11}^{-1} \Lambda_{12} \left( m_2 - \mu_2 \right)$ 

which we can rewrite as

$$m_2 \left( 1 - \Lambda_{22}^{-1} \Lambda_{21} \Lambda_{11}^{-1} \Lambda_{12} \right) = \mu_2 \left( 1 - \Lambda_{22}^{-1} \Lambda_{21} \Lambda_{11}^{-1} \Lambda_{12} \right).$$

Thus, unless  $\Lambda_{22}^{-1}\Lambda_{21}\Lambda_{11}^{-1}\Lambda_{12}=1$ , the solution  $\mu_2=m_2$  is unique. If  $p(\mathbf{z})$  is non-singular,

$$|\mathbf{\Lambda}| = \Lambda_{11}\Lambda_{22} - \Lambda_{21}\Lambda_{12} \neq 0$$

which we can rewrite as

$$\Lambda_{11}^{-1}\Lambda_{22}^{-1}\Lambda_{21}\Lambda_{12} \neq 1$$

as desired. Since  $\mu_2=m_2$  is the unique solution to (10.15),  $\mu_1=m_1$  is the unique solution to (10.13).

**10.3** Starting from (10.16) and optimizing w.r.t.  $q_i(\mathbf{Z}_i)$ , we get

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

$$= -\int \left( p(\mathbf{Z}) \ln q_{j}(\mathbf{Z}_{j}) + p(\mathbf{Z}) \sum_{i \neq j} \ln q_{i}(\mathbf{Z}_{i}) \right) d\mathbf{Z} + \text{const.}$$

$$= -\int p(\mathbf{Z}) \ln q_{j}(\mathbf{Z}_{j}) d\mathbf{Z} + \text{const.}$$

$$= -\int \ln q_{j}(\mathbf{Z}_{j}) \left[ \int p(\mathbf{Z}) \prod_{i \neq j} d\mathbf{Z}_{i} \right] d\mathbf{Z}_{j} + \text{const.}$$

$$= -\int F_{j}(\mathbf{Z}_{j}) \ln q_{j}(\mathbf{Z}_{j}) d\mathbf{Z}_{j} + \text{const.},$$

where terms independent of  $q_{j}\left(\mathbf{Z}_{j}\right)$  have been absorbed into the constant term and we have defined

$$F_{j}(\mathbf{Z}_{j}) = \int p(\mathbf{Z}) \prod_{i \neq j} d\mathbf{Z}_{i}.$$

We use a Lagrange multiplier to ensure that  $q_j(\mathbf{Z}_j)$  integrates to one, yielding

$$-\int F_{j}(\mathbf{Z}_{j}) \ln q_{j}(\mathbf{Z}_{j}) d\mathbf{Z}_{j} + \lambda \left( \int q_{j}(\mathbf{Z}_{j}) d\mathbf{Z}_{j} - 1 \right).$$

Using the results from Appendix D, we then take the functional derivative of this w.r.t.  $q_i$  and set this to zero, to obtain

$$-\frac{F_j(\mathbf{Z}_j)}{q_j(\mathbf{Z}_j)} + \lambda = 0.$$

From this, we see that

$$\lambda q_i(\mathbf{Z}_i) = F_i(\mathbf{Z}_i).$$

Integrating both sides over  $\mathbf{Z}_{j}$ , we see that, since  $q_{j}(\mathbf{Z}_{j})$  must integrate to one,

$$\lambda = \int F_{j}(\mathbf{Z}_{j}) \, d\mathbf{Z}_{j} = \int \left[ \int p\left(\mathbf{Z}\right) \prod_{i \neq j} \, d\mathbf{Z}_{i} \right] \, d\mathbf{Z}_{j} = 1,$$

and thus

$$q_{j}(\mathbf{Z}_{j}) = F_{j}(\mathbf{Z}_{j}) = \int p(\mathbf{Z}) \prod_{i \neq j} d\mathbf{Z}_{i}.$$

**10.4** The Kullback-Leibler divergence takes the form

$$KL(p||q) = -\int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} + \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x}.$$

Substituting the Gaussian for  $q(\mathbf{x})$  we obtain

$$KL(p||q) = -\int p(\mathbf{x}) \left\{ -\frac{1}{2} \ln |\mathbf{\Sigma}| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} d\mathbf{x} + \text{const.}$$

$$= \frac{1}{2} \left\{ \ln |\mathbf{\Sigma}| + \text{Tr} \left( \mathbf{\Sigma}^{-1} \mathbb{E} \left[ (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \right] \right) \right\} + \text{const.}$$

$$= \frac{1}{2} \left\{ \ln |\mathbf{\Sigma}| + \boldsymbol{\mu}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \boldsymbol{\mu} - 2\boldsymbol{\mu}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \mathbb{E}[\mathbf{x}] + \text{Tr} \left( \mathbf{\Sigma}^{-1} \mathbb{E} \left[ \mathbf{x} \mathbf{x}^{\mathrm{T}} \right] \right) \right\} + \text{const.}$$

$$(254)$$

Differentiating this w.r.t.  $\mu$ , using results from Appendix C, and setting the result to zero, we see that

$$\mu = \mathbb{E}[\mathbf{x}]. \tag{255}$$

Similarly, differentiating (254) w.r.t.  $\Sigma^{-1}$ , again using results from Appendix C and also making use of (255) and (1.42), we see that

$$\Sigma = \mathbb{E}\left[\mathbf{x}\mathbf{x}^{\mathrm{T}}\right] - \mu\mu^{\mathrm{T}} = \mathrm{cov}[\mathbf{x}].$$

**10.5** We assume that  $q(\mathbf{Z}) = q(\mathbf{z})q(\boldsymbol{\theta})$  and so we can optimize w.r.t.  $q(\mathbf{z})$  and  $q(\boldsymbol{\theta})$  independently.

For  $q(\mathbf{z})$ , this is equivalent to minimizing the Kullback-Leibler divergence, (10.4), which here becomes

$$\mathrm{KL}(q \parallel p) = -\iint q\left(\boldsymbol{\theta}\right) q\left(\mathbf{z}\right) \ln \frac{p\left(\mathbf{z}, \boldsymbol{\theta} \mid \mathbf{X}\right)}{q\left(\mathbf{z}\right) q\left(\boldsymbol{\theta}\right)} \, \mathrm{d}\mathbf{z} \, \mathrm{d}\boldsymbol{\theta}.$$

For the particular chosen form of  $q(\theta)$ , this is equivalent to

$$KL(q \parallel p) = -\int q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \boldsymbol{\theta}_0 \mid \mathbf{X})}{q(\mathbf{z})} d\mathbf{z} + \text{const.}$$

$$= -\int q(\mathbf{z}) \ln \frac{p(\mathbf{z} \mid \boldsymbol{\theta}_0, \mathbf{X}) p(\boldsymbol{\theta}_0 \mid \mathbf{X})}{q(\mathbf{z})} d\mathbf{z} + \text{const.}$$

$$= -\int q(\mathbf{z}) \ln \frac{p(\mathbf{z} \mid \boldsymbol{\theta}_0, \mathbf{X})}{q(\mathbf{z})} d\mathbf{z} + \text{const.},$$

where const accumulates all terms independent of  $q(\mathbf{z})$ . This KL divergence is minimized when  $q(\mathbf{z}) = p(\mathbf{z}|\boldsymbol{\theta}_0, \mathbf{X})$ , which corresponds exactly to the E-step of the EM algorithm.

To determine  $q(\theta)$ , we consider

$$\int q(\boldsymbol{\theta}) \int q(\mathbf{z}) \ln \frac{p(\mathbf{X}, \boldsymbol{\theta}, \mathbf{z})}{q(\boldsymbol{\theta}) q(\mathbf{z})} d\mathbf{z} d\boldsymbol{\theta} 
= \int q(\boldsymbol{\theta}) \mathbb{E}_{q(\mathbf{z})} [\ln p(\mathbf{X}, \boldsymbol{\theta}, \mathbf{z})] d\boldsymbol{\theta} - \int q(\boldsymbol{\theta}) \ln q(\boldsymbol{\theta}) d\boldsymbol{\theta} + \text{const.}$$

where the last term summarizes terms independent of  $q(\theta)$ . Since  $q(\theta)$  is constrained to be a point density, the contribution from the entropy term (which formally diverges) will be constant and independent of  $\theta_0$ . Thus, the optimization problem is reduced to maximizing expected complete log posterior distribution

$$\mathbb{E}_{q(\mathbf{z})} \left[ \ln p \left( \mathbf{X}, \boldsymbol{\theta}_0, \mathbf{z} \right) \right],$$

w.r.t.  $\theta_0$ , which is equivalent to the M-step of the EM algorithm.

**10.6** We start by rewriting (10.19) as

$$D(p||q) = \frac{4}{1 - \alpha^2} \left( 1 - \int p(x)^{\gamma_p} q(x)^{\gamma_q} dx \right)$$
 (256)

so that

$$\gamma_p = \frac{1+\alpha}{2} \quad \text{and} \quad \gamma_q = \frac{1-\alpha}{2}.$$
(257)

We note that

$$\lim_{\alpha \to 1} \gamma_q = 0 \tag{258}$$

$$\lim_{\alpha \to 1} \gamma_q = 0$$

$$\lim_{\alpha \to 1} \gamma_p = 1$$

$$1 - \gamma_p = \gamma_q.$$
(258)
(259)

$$1 - \gamma_p = \gamma_q. \tag{260}$$

Based on observation (258), we make a Maclaurin expansion of  $q(x)^{\gamma_q}$  in  $\gamma_q$  as follows

$$q^{\gamma_q} = \exp\left(\gamma_q \ln q\right) = 1 + \gamma_q \ln q + O\left(\gamma_q^2\right) \tag{261}$$

where q is a shorthand notation for q(x). Similarly, based on (259), we make a Taylor expansion of  $p(x)^{\gamma_p}$  in  $\gamma_p$  around 1,

$$p^{\gamma_p} = \exp(\gamma_p \ln p)$$

$$= p - (1 - \gamma_p)p \ln p + O((\gamma_p - 1)^2)$$

$$= p - \gamma_q p \ln p + O(\gamma_q^2)$$
(262)

where we have used (260) and we have adopted corresponding shorthand notation for p(x).

Using (261) and (262), we can rewrite (256) as

$$D(p||q) = \frac{4}{1-\alpha^2} \left( 1 - \int \left[ p - \gamma_q p \ln p + O\left(\gamma_q^2\right) \right] \left[ 1 + \gamma_q \ln q + O\left(\gamma_q^2\right) \right] dx \right)$$
$$= \frac{4}{1-\alpha^2} \left( 1 - \int p + \gamma_q \left( p \ln q - p \ln p \right) dx + O\left(\gamma_q^2\right) \right)$$
(263)

where  $O\left(\gamma_q^2\right)$  account for all higher order terms. From (257) we have

$$\frac{4}{1-\alpha^2}\gamma_q = \frac{2(1-\alpha)}{1-\alpha^2} = \frac{2}{(1+\alpha)}$$
$$\frac{4}{1-\alpha^2}\gamma_q^2 = \frac{(1-\alpha)^2}{1-\alpha^2} = \frac{1-\alpha}{(1+\alpha)}$$

and thus

$$\lim_{\alpha \to 1} \frac{4}{1 - \alpha^2} \gamma_q = 1$$

$$\lim_{\alpha \to 1} \frac{4}{1 - \alpha^2} \gamma_q^2 = 0.$$

Using these results together with (259), and (263), we see that

$$\lim_{\alpha \to 1} D(p||q) = -\int p(\ln q - \ln p) \, \mathrm{d}x = \mathrm{KL}(p||q).$$

The proof that  $\alpha \to -1$  yields  $\mathrm{KL}(q\|p)$  is analogous.

#### **10.7 NOTE**: See note in Solution 10.9.

We take the  $\mu$ -dependent term from the last line of (10.25) as our starting point. We can rewrite this as follows

$$-\frac{\mathbb{E}[\tau]}{2} \left\{ \lambda_0 (\mu - \mu_0)^2 + \sum_{n=1}^N (x_n - \mu)^2 \right\}$$

$$= -\frac{\mathbb{E}[\tau]}{2} \left\{ (\lambda_0 + N) \mu^2 + \sum_{n=1}^N x_n^2 - 2\mu (\lambda_0 \mu_0 + N\overline{x}) \right\}$$

$$= -\frac{\mathbb{E}[\tau]}{2} \left\{ (\lambda_0 + N) \left( \mu - \frac{\lambda_0 \mu_0 + N\overline{x}}{\lambda_0 + N} \right)^2 + \sum_{n=1}^N x_n^2 - \frac{(\lambda_0 m u_0 + N\overline{x})^2}{\lambda_0 + N} \right\}$$

where in the last step we have completed the square over  $\mu$ . The last two terms are independent of  $\mu$  and hence can be dropped. Taking the exponential of the remainder, we obtain an unnormalized Gaussian with mean and precision given by (10.26) and (10.27), respectively.

For the posterior over  $\tau$ , we take the last two lines of (10.28) as our starting point. Gathering terms that multiply  $\tau$  and  $\ln \tau$  into two groups, we can rewrite this as

$$\left(a_0 + \frac{N+1}{2} - 1\right) \ln \tau - \left(b_0 + \frac{1}{2} \mathbb{E} \left[ \sum_{n=1}^{N} (x - \mu)^2 + \lambda_0 (\mu - \mu_0) \right] \right) \tau + \text{const.}$$

Taking the exponential of this we get an unnormalized Gamma distribution with shape and inverse scale parameters given by (10.29) and (10.30), respectively.

**10.8 NOTE**: See note in Solution 10.9.

If we substitute the r.h.s. of (10.29) and (10.30) for a and b, respectively, in (B.27) and (B.28), we get

$$\mathbb{E}[\tau] = \frac{2a_0 + N + 1}{2b_0 + \mathbb{E}\left[\lambda_0(\mu - \mu_0) + \sum_{n=1}^{N} (x_n - \mu)^2\right]}$$

$$\text{var}[\tau] = \frac{2a_0 + N + 1}{2\left(b_0 + \frac{1}{2}\mathbb{E}\left[\lambda_0(\mu - \mu_0) + \sum_{n=1}^{N} (x_n - \mu)^2\right]\right)^2}$$

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

From this we see directly that

$$\lim_{N \to \infty} \mathbb{E}[\tau] = \frac{N}{\mathbb{E}\left[\sum_{n=1}^{N} (x_n - \mu)^2\right]}$$
$$\lim_{N \to \infty} \text{var}[\tau] = 0$$

as long as the data set is not singular.

**10.9 NOTE**: In the  $1^{st}$  printing of PRML, an extra term of 1/2 should be added to the r.h.s. of (10.29), with consequential changes to (10.31) and (10.33), which should read

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

and

$$\frac{1}{\mathbb{E}[\tau]} = (\overline{x^2} - \overline{x}^2) = \frac{1}{N} \sum_{n=1}^{N} (x_n - \overline{x})^2$$

respectively.

Assuming  $a_0 = b_0 = \lambda_0 = 0$ , (10.29), (10.30) and (B.27) give

$$\frac{1}{\mathbb{E}[\tau]} = \frac{1}{N+1} \sum_{n=1}^{N} (x_n - \mu)^2$$
$$= \frac{1}{N+1} \sum_{n=1}^{N} (x_n^2 - 2x_n\mu + \mu^2)$$

Taking the expectation of this under  $q(\mu)$ , making use of (10.32), we get

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

which we can rearrange to obtain (10.33).

**10.10 NOTE**: In the 1<sup>st</sup> printing of PRML, there are errors that affect this exercise.  $\mathcal{L}_m$  used in (10.34) and (10.35) should really be  $\mathcal{L}$ , whereas  $\mathcal{L}_m$  used in (10.36) is given in Solution 10.11 below.

This completely analogous to Solution 10.1. Starting from (10.35), we can use the product rule to get,

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

$$= \sum_{m} \sum_{\mathbf{Z}} q(\mathbf{Z}|m) q(m) \ln \left\{ \frac{p(\mathbf{Z}, m|\mathbf{X}) p(\mathbf{X})}{q(\mathbf{Z}|m) q(m)} \right\}$$

$$= \sum_{m} \sum_{\mathbf{Z}} q(\mathbf{Z}|m) q(m) \ln \left\{ \frac{p(\mathbf{Z}, m|\mathbf{X})}{q(\mathbf{Z}|m) q(m)} \right\} + \ln p(\mathbf{X}).$$

Rearranging this, we obtain (10.34).

#### 170 Solutions 10.11–10.13

**10.11 NOTE**: Consult note preceding Solution 10.10 for some relevant corrections. We start by rewriting the lower bound as follows

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

$$= \sum_{m} \sum_{\mathbf{Z}} q(\mathbf{Z}|m)q(m) \left\{ \ln p(\mathbf{Z}, \mathbf{X}|m) + \ln p(m) - \ln q(\mathbf{Z}|m) - \ln q(m) \right\}$$

$$= \sum_{m} q(m) \left( \ln p(m) - \ln q(m) + \sum_{\mathbf{Z}} q(\mathbf{Z}|m) \left\{ \ln p(\mathbf{Z}, \mathbf{X}|m) - \ln q(\mathbf{Z}|m) \right\} \right)$$

$$= \sum_{m} q(m) \left\{ \ln (p(m) \exp\{\mathcal{L}_{m}\}) - \ln q(m) \right\}, \tag{264}$$

where

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

We recognize (264) as the negative KL divergence between q(m) and the (not necessarily normalized) distribution  $p(m)\exp\{\mathcal{L}_m\}$ . This will be maximized when the KL divergence is minimized, which will be the case when

$$q(m) \propto p(m) \exp{\{\mathcal{L}_m\}}.$$

- **10.12** This derivation is given in detail in Section 10.2.1, starting with the paragraph containing (10.43) (page 476) and ending with (10.49).
- 10.13 In order to derive the optimal solution for  $q(\mu_k, \Lambda_k)$  we start with the result (10.54) and keep only those term which depend on  $\mu_k$  or  $\Lambda_k$  to give

$$\ln q^{\star}(\boldsymbol{\mu}_{k}, \boldsymbol{\Lambda}_{k}) = \ln \mathcal{N}\left(\boldsymbol{\mu}_{k} | \mathbf{m}_{0}, (\beta_{0} \boldsymbol{\Lambda}_{k})^{-1}\right) + \ln \mathcal{W}(\boldsymbol{\Lambda}_{k} | \mathbf{W}_{0}, \nu_{0})$$

$$+ \sum_{n=1}^{N} \mathbb{E}[z_{nk}] \ln \mathcal{N}\left(\mathbf{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Lambda}_{k}^{-1}\right) + \text{const.}$$

$$= -\frac{\beta_{0}}{2} (\boldsymbol{\mu}_{k} - \mathbf{m}_{0})^{T} \boldsymbol{\Lambda}_{k} (\boldsymbol{\mu}_{k} - \mathbf{m}_{0}) + \frac{1}{2} \ln |\boldsymbol{\Lambda}_{k}| - \frac{1}{2} \text{Tr}\left(\boldsymbol{\Lambda}_{k} \mathbf{W}_{0}^{-1}\right)$$

$$+ \frac{(\nu_{0} - D - 1)}{2} \ln |\boldsymbol{\Lambda}_{k}| - \frac{1}{2} \sum_{n=1}^{N} \mathbb{E}[z_{nk}] (\mathbf{x}_{n} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Lambda}_{k} (\mathbf{x}_{n} - \boldsymbol{\mu}_{k})$$

$$+ \frac{1}{2} \left(\sum_{n=1}^{N} \mathbb{E}[z_{nk}]\right) \ln |\boldsymbol{\Lambda}_{k}| + \text{const.}$$
(265)

Using the product rule of probability, we can express  $\ln q^*(\mu_k, \Lambda_k)$  as  $\ln q^*(\mu_k | \Lambda_k) + \ln q^*(\Lambda_k)$ . Let us first of all identify the distribution for  $\mu_k$ . To do this we need only consider terms on the right hand side of (265) which depend on  $\mu_k$ , giving

$$\begin{aligned} & \ln q^{\star}(\boldsymbol{\mu}_{k}|\boldsymbol{\Lambda}_{k}) \\ & = & -\frac{1}{2}\boldsymbol{\mu}_{k}^{\mathrm{T}}\left[\beta_{0} + \sum_{n=1}^{N}\mathbb{E}[z_{nk}]\right]\boldsymbol{\Lambda}_{k}\boldsymbol{\mu}_{k} + \boldsymbol{\mu}_{k}^{\mathrm{T}}\boldsymbol{\Lambda}_{k}\left[\beta_{0}\mathbf{m}_{0} + \sum_{n=1}^{N}\mathbb{E}[z_{nk}]\mathbf{x}_{n}\right] \\ & + \mathrm{const.} \\ & = & -\frac{1}{2}\boldsymbol{\mu}_{k}^{\mathrm{T}}\left[\beta_{0} + N_{k}\right]\boldsymbol{\Lambda}_{k}\boldsymbol{\mu}_{k} + \boldsymbol{\mu}_{k}^{\mathrm{T}}\boldsymbol{\Lambda}_{k}\left[\beta_{0}\mathbf{m}_{0} + N_{k}\overline{\mathbf{x}}_{k}\right] + \mathrm{const.} \end{aligned}$$

where we have made use of (10.51) and (10.52). Thus we see that  $\ln q^*(\boldsymbol{\mu}_k|\boldsymbol{\Lambda}_k)$  depends quadratically on  $\boldsymbol{\mu}_k$  and hence  $q^*(\boldsymbol{\mu}_k|\boldsymbol{\Lambda}_k)$  is a Gaussian distribution. Completing the square in the usual way allows us to determine the mean and precision of this Gaussian, giving

$$q^{\star}(\boldsymbol{\mu}_k|\boldsymbol{\Lambda}_k) = \mathcal{N}(\boldsymbol{\mu}_k|\mathbf{m}_k, \beta_k \boldsymbol{\Lambda}_k)$$
 (266)

where

$$\beta_k = \beta_0 + N_k$$

$$\mathbf{m}_k = \frac{1}{\beta_k} (\beta_0 \mathbf{m}_0 + N_k \overline{\mathbf{x}}_k).$$

Next we determine the form of  $q^*(\Lambda_k)$  by making use of the relation

$$\ln q^{\star}(\mathbf{\Lambda}_k) = \ln q^{\star}(\boldsymbol{\mu}_k, \mathbf{\Lambda}_k) - \ln q^{\star}(\boldsymbol{\mu}_k|\mathbf{\Lambda}_k).$$

On the right hand side of this relation we substitute for  $\ln q^{\star}(\mu_k, \Lambda_k)$  using (265), and we substitute for  $\ln q^{\star}(\mu_k | \Lambda_k)$  using the result (266). Keeping only those terms which depend on  $\Lambda_k$  we obtain

$$\ln q^{\star}(\mathbf{\Lambda}_{k}) = -\frac{\beta_{0}}{2}(\boldsymbol{\mu}_{k} - \mathbf{m}_{0})^{\mathrm{T}} \mathbf{\Lambda}_{k}(\boldsymbol{\mu}_{k} - \mathbf{m}_{0}) + \frac{1}{2} \ln |\mathbf{\Lambda}_{k}| - \frac{1}{2} \mathrm{Tr} \left(\mathbf{\Lambda}_{k} \mathbf{W}_{0}^{-1}\right)$$

$$+ \frac{(\nu_{0} - D - 1)}{2} \ln |\mathbf{\Lambda}_{k}| - \frac{1}{2} \sum_{n=1}^{N} \mathbb{E}[z_{nk}] (\mathbf{x}_{n} - \boldsymbol{\mu}_{k})^{\mathrm{T}} \mathbf{\Lambda}_{k} (\mathbf{x}_{n} - \boldsymbol{\mu}_{k})$$

$$+ \frac{1}{2} \left(\sum_{n=1}^{N} \mathbb{E}[z_{nk}]\right) \ln |\mathbf{\Lambda}_{k}| + \frac{\beta_{k}}{2} (\boldsymbol{\mu}_{k} - \mathbf{m}_{k})^{\mathrm{T}} \mathbf{\Lambda}_{k} (\boldsymbol{\mu}_{k} - \mathbf{m}_{k})$$

欢迎关注公众号a机器第二点算法之。i+ const.

Here we have defined

$$\mathbf{W}_{k}^{-1} = \mathbf{W}_{0}^{-1} + \beta_{0}(\boldsymbol{\mu}_{k} - \mathbf{m}_{0})(\boldsymbol{\mu}_{k} - \mathbf{m}_{0})^{\mathrm{T}} + \sum_{n=1}^{N} \mathbb{E}[z_{nk}](\mathbf{x}_{n} - \boldsymbol{\mu}_{k})(\mathbf{x}_{n} - \boldsymbol{\mu}_{k})^{\mathrm{T}}$$

$$-\beta_{k}(\boldsymbol{\mu}_{k} - \mathbf{m}_{k})(\boldsymbol{\mu}_{k} - \mathbf{m}_{k})^{\mathrm{T}}$$

$$= \mathbf{W}_{0}^{-1} + N_{k}\mathbf{S}_{k} + \frac{\beta_{0}N_{k}}{\beta_{0} + N_{k}}(\overline{\mathbf{x}}_{k} - \mathbf{m}_{0})(\overline{\mathbf{x}}_{k} - \mathbf{m}_{0})^{\mathrm{T}}$$

$$\nu_{k} = \nu_{0} + \sum_{n=1}^{N} \mathbb{E}[z_{nk}]$$

$$= \nu_{0} + N_{k},$$

$$(267)$$

where we have made use of the result

$$\sum_{n=1}^{N} \mathbb{E}[z_{nk}] \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} = \sum_{n=1}^{N} \mathbb{E}[z_{nk}] (\mathbf{x}_{n} - \overline{\mathbf{x}}_{k}) (\mathbf{x}_{n} - \overline{\mathbf{x}}_{k})^{\mathrm{T}} + N_{k} \overline{\mathbf{x}}_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}}$$

$$= N_{k} \mathbf{S}_{k} + N_{k} \overline{\mathbf{x}}_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}}$$
(268)

and we have made use of (10.53). Note that the terms involving  $\mu_k$  have cancelled out in (267) as we expect since  $q^*(\Lambda_k)$  is independent of  $\mu_k$ .

Thus we see that  $q^*(\Lambda_k)$  is a Wishart distribution of the form

$$q^{\star}(\mathbf{\Lambda}_k) = \mathcal{W}(\mathbf{\Lambda}_k | \mathbf{W}_k, \nu_k).$$

**10.14** We can express the required expectation as an integration with respect to the variational posterior distribution  $q^*(\mu_k, \Lambda_k) = q^*(\mu_k | \Lambda_k) q^*(\Lambda_k)$ . Thus we have

$$\mathbb{E}_{\boldsymbol{\mu}_{k},\boldsymbol{\Lambda}_{k}}\left[(\mathbf{x}_{n}-\boldsymbol{\mu}_{k})^{\mathrm{T}}\boldsymbol{\Lambda}_{k}(\mathbf{x}_{n}-\boldsymbol{\mu}_{k})\right]$$

$$=\iint \operatorname{Tr}\left\{\boldsymbol{\Lambda}_{k}(\mathbf{x}_{n}-\boldsymbol{\mu}_{k})(\mathbf{x}_{n}-\boldsymbol{\mu}_{k})^{\mathrm{T}}\right\}q^{\star}(\boldsymbol{\mu}_{k}|\boldsymbol{\Lambda}_{k})q^{\star}(\boldsymbol{\Lambda}_{k})\,\mathrm{d}\boldsymbol{\mu}_{k}\,\mathrm{d}\boldsymbol{\Lambda}_{k}.$$

Next we use the result  $q^{\star}(\boldsymbol{\mu}_k|\boldsymbol{\Lambda}_k) = \mathcal{N}(\boldsymbol{\mu}_k|\mathbf{m}_k,\beta_k\boldsymbol{\Lambda}_k)$  to perform the integration over  $\boldsymbol{\mu}_k$  using the standard expressions for expectations under a Gaussian distribution, giving

$$\mathbb{E}[\boldsymbol{\mu}_k] = \mathbf{m}_k$$

$$\mathbb{E}[\boldsymbol{\mu}_k \boldsymbol{\mu}_k^{\mathrm{T}}] = \mathbf{m}_k \mathbf{m}_k^{\mathrm{T}} + \beta_k^{-1} \boldsymbol{\Lambda}_k^{-1}$$

from which we obtain the expectation with respect to  $\mu_k$  in the form

$$\begin{split} & \mathbb{E}_{\boldsymbol{\mu}_k} \left[ (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^{\mathrm{T}} \right] \\ & = & \mathbf{x}_n \mathbf{x}_n^{\mathrm{T}} - \mathbf{x}_n \mathbf{m}_k^{\mathrm{T}} - \mathbf{m}_k \mathbf{x}_n^{\mathrm{T}} + \mathbf{m}_k \mathbf{m}_k^{\mathrm{T}} + \beta_k^{-1} \boldsymbol{\Lambda}_k^{-1} \\ & = & (\mathbf{x}_n - \mathbf{m}_k) (\mathbf{x}_n - \mathbf{m}_k)^{\mathrm{T}} + \beta_k^{-1} \boldsymbol{\Lambda}_k^{-1}. \end{split}$$

Finally, taking the expectation with respect to  $\Lambda_k$  we have

$$\begin{split} & \mathbb{E}_{\boldsymbol{\mu}_{k},\boldsymbol{\Lambda}_{k}}\left[(\mathbf{x}_{n}-\boldsymbol{\mu}_{k})^{\mathrm{T}}\boldsymbol{\Lambda}_{k}(\mathbf{x}_{n}-\boldsymbol{\mu}_{k})\right] \\ & = \int \mathrm{Tr}\left\{\boldsymbol{\Lambda}_{k}\left[(\mathbf{x}_{n}-\mathbf{m}_{k})(\mathbf{x}_{n}-\mathbf{m}_{k})^{\mathrm{T}}+\boldsymbol{\beta}_{k}^{-1}\boldsymbol{\Lambda}_{k}^{-1}\right]\right\}\boldsymbol{q}^{\star}(\boldsymbol{\Lambda}_{k})\,\mathrm{d}\boldsymbol{\Lambda}_{k} \\ & = \int \left\{(\mathbf{x}_{n}-\mathbf{m}_{k})^{\mathrm{T}}\boldsymbol{\Lambda}_{k}(\mathbf{x}_{n}-\mathbf{m}_{k})+D\boldsymbol{\beta}_{k}^{-1}\right\}\boldsymbol{q}^{\star}(\boldsymbol{\Lambda}_{k})\,\mathrm{d}\boldsymbol{\Lambda}_{k} \\ & = D\boldsymbol{\beta}_{k}^{-1}+\nu_{k}(\mathbf{x}_{n}-\mathbf{m}_{k})^{\mathrm{T}}\mathbf{W}_{k}(\mathbf{x}_{n}-\mathbf{m}_{k}) \end{split}$$

as required. Here we have used  $q^{\star}(\mathbf{\Lambda}_k) = \mathcal{W}(\mathbf{\Lambda}_k | \mathbf{W}_k, \nu_k)$ , together with the standard result for the expectation under a Wishart distribution to give  $\mathbb{E}[\mathbf{\Lambda}_k] = \nu_k \mathbf{W}_k$ .

- **10.15** By substituting (10.58) into (B.17) and then using (B.24) together with the fact that  $\sum_k N_k = N$ , we obtain (10.69).
- **10.16** To derive (10.71) we make use of (10.38) to give

$$\mathbb{E}[\ln p(D|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] = \frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{K} \mathbb{E}[z_{nk}] \left\{ \mathbb{E}[\ln |\boldsymbol{\Lambda}_k|] - \mathbb{E}[(\mathbf{x}_n - \boldsymbol{\mu}_k)\boldsymbol{\Lambda}_k(\mathbf{x}_n - \boldsymbol{\mu}_k)] - D\ln(2\pi) \right\}.$$

We now use  $\mathbb{E}[z_{nk}] = r_{nk}$  together with (10.64) and the definition of  $\widetilde{\Lambda}_k$  given by (10.65) to give

$$\mathbb{E}[\ln p(D|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Lambda})] = \frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \{\ln \widetilde{\Lambda}_k - D\beta_k^{-1} - \nu_k (\mathbf{x}_n - \mathbf{m}_k)^{\mathrm{T}} \mathbf{W}_k (\mathbf{x}_n - \mathbf{m}_k) - D \ln(2\pi) \}.$$

Now we use the definitions (10.51) to (10.53) together with the result (268) to give (10.71).

We can derive (10.72) simply by taking the logarithm of  $p(\mathbf{z}|\boldsymbol{\pi})$  given by (10.37)

$$\mathbb{E}[\ln p(\mathbf{z}|\boldsymbol{\pi})] = \sum_{n=1}^{N} \sum_{k=1}^{K} \mathbb{E}[z_{nk}] \mathbb{E}[\ln \pi_k]$$

and then making use of  $\mathbb{E}[z_{nk}] = r_{nk}$  together with the definition of  $\widetilde{\pi}_k$  given by (10.65).

**10.17** The result (10.73) is obtained by using the definition of  $p(\pi)$  given by (10.39) together with the definition of  $\widetilde{\pi}_k$  given by (10.66).

For the result (10.74) we start with the definition of the prior  $p(\mu, \Lambda)$  given by (10.40) to give

$$\mathbb{E}[\ln p(\boldsymbol{\mu}, \boldsymbol{\Lambda})] = \frac{1}{2} \sum_{k=1}^{K} \left\{ D \ln \beta_0 - D \ln(2\pi) + \mathbb{E}[\ln |\boldsymbol{\Lambda}_k|] - \beta_0 \mathbb{E}[(\boldsymbol{\mu}_k - \mathbf{m}_0)^{\mathrm{T}} \boldsymbol{\Lambda}_k (\boldsymbol{\mu}_k - \mathbf{m}_0)] \right\} + K \ln B(\mathbf{W}_0, \nu_0) + \sum_{k=1}^{K} \left\{ \frac{(\nu_0 - D - 1)}{2} \mathbb{E}[\ln |\boldsymbol{\Lambda}_k|] - \frac{1}{2} \mathrm{Tr}(\mathbf{W}_0^{-1} \mathbb{E}[\boldsymbol{\Lambda}_k]) \right\}.$$

Now consider the term  $\mathbb{E}[(\boldsymbol{\mu}_k - \mathbf{m}_0)^T \boldsymbol{\Lambda}_k (\boldsymbol{\mu}_k - \mathbf{m}_0)]$ . To evaluate this expression we first perform the expectation with respect to  $q^*(\boldsymbol{\mu}_k | \boldsymbol{\Lambda}_k)$  then the subsequently perform the expectation with respect to  $q^*(\boldsymbol{\Lambda}_k)$ . Using the standard results for moments under a Gaussian we have

$$\begin{array}{rcl} \mathbb{E}[\boldsymbol{\mu}_k] & = & \mathbf{m}_k \\ \mathbb{E}[\boldsymbol{\mu}_k \boldsymbol{\mu}_k^{\mathrm{T}}] & = & \mathbf{m}_k \mathbf{m}_k^{\mathrm{T}} + \beta_k^{-1} \boldsymbol{\Lambda}_k^{-1} \end{array}$$

and hence

$$\begin{split} &\mathbb{E}_{\boldsymbol{\mu}_{k}\boldsymbol{\Lambda}_{k}}[(\boldsymbol{\mu}_{k}-\mathbf{m}_{0})^{\mathrm{T}}\boldsymbol{\Lambda}_{k}(\boldsymbol{\mu}_{k}-\mathbf{m}_{0})] = \mathrm{Tr}\left(\mathbb{E}_{\boldsymbol{\mu}_{k}\boldsymbol{\Lambda}_{k}}\left[\boldsymbol{\Lambda}_{k}(\boldsymbol{\mu}_{k}-\mathbf{m}_{0})(\boldsymbol{\mu}_{k}-\mathbf{m}_{0})^{\mathrm{T}}\right]\right) \\ &= \mathrm{Tr}\left(\mathbb{E}_{\boldsymbol{\Lambda}_{k}}\left[\boldsymbol{\Lambda}_{k}(\boldsymbol{\beta}_{k}^{-1}\boldsymbol{\Lambda}_{k}^{-1}+\mathbf{m}_{k}\mathbf{m}_{k}^{\mathrm{T}}-\mathbf{m}_{0}\mathbf{m}_{k}^{\mathrm{T}}-\mathbf{m}_{k}\mathbf{m}_{0}^{\mathrm{T}}+\mathbf{m}_{0}\mathbf{m}_{0}^{\mathrm{T}})\right]\right) \\ &= K\boldsymbol{\beta}_{k}^{-1}+(\mathbf{m}_{k}-\mathbf{m}_{0})^{\mathrm{T}}\mathbb{E}[\boldsymbol{\Lambda}_{k}](\mathbf{m}_{k}-\mathbf{m}_{0}). \end{split}$$

Now we use (B.80) to give  $\mathbb{E}[\mathbf{\Lambda}_k] = \nu_k \mathbf{W}_k$  and  $\mathbb{E}[\ln \mathbf{\Lambda}_k] = \ln \widetilde{\Lambda}_k$  from (10.65) to give (10.74).

For (10.75) we take use the result (10.48) for  $q^*(\mathbf{z})$  to give

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

and using  $\mathbb{E}[z_{nk}] = r_{nk}$  we obtain (10.75).

The solution (10.76) for  $\mathbb{E}[\ln q(\pi)]$  is simply the negative entropy of the corresponding Dirichlet distribution (10.57) and is obtained from (B.22).

Finally, we need the entropy of the Gaussian-Wishart distribution  $q(\boldsymbol{\mu}, \boldsymbol{\Lambda})$ . First of all we note that this distribution factorizes into a product of factors  $q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$  and the entropy of the product is the sum of the entropies of the individual terms, as is easily verified. Next we write

$$\ln q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) = \ln q(\boldsymbol{\mu}_k | \boldsymbol{\Lambda}_k) + \ln q(\boldsymbol{\Lambda}_k).$$

Consider first the quantity  $\mathbb{E}[\ln q(\mu_k|\Lambda_k)]$ . Taking the expectation first with respect to  $\mu_k$  we can make use of the standard result (B.41) for the entropy of a Gaussian to

give

$$\mathbb{E}_{\boldsymbol{\mu}_{k}\boldsymbol{\Lambda}_{k}}[\ln q(\boldsymbol{\mu}_{k}|\boldsymbol{\Lambda}_{k})] = \mathbb{E}_{\boldsymbol{\Lambda}_{k}}\left[\frac{1}{2}\ln|\boldsymbol{\Lambda}_{k}| + \frac{D}{2}\left(\ln\beta_{k} - 1 - \ln(2\pi)\right)\right]$$
$$= \frac{1}{2}\ln\widetilde{\boldsymbol{\Lambda}}_{k} + \frac{D}{2}\left(\ln\beta_{k} - 1 - \ln(2\pi)\right).$$

The term  $\mathbb{E}[\ln q(\mathbf{\Lambda}_k)]$  is simply the negative entropy of a Wishart distribution, which we write as  $-H[q(\mathbf{\Lambda}_k)]$ .

**10.18** We start with  $\beta_k$ , which appears in (10.71), (10.74) and (10.77). Using these, we can differentiate (10.70) w.r.t.  $\beta_k^{-1}$ , to get

$$\frac{\partial \mathcal{L}}{\partial \beta_k^{-1}} = \frac{D}{2} \left( -N_k - \beta_0 + \beta_k \right).$$

Setting this equal to zero and rearranging the terms, we obtain (10.60). We then consider  $\mathbf{m}_k$ , which appears in the quadratic terms of (10.71) and (10.74). Thus differentiation of (10.70) w.r.t.  $\mathbf{m}_k$  gives

$$\frac{\partial \mathcal{L}}{\partial \mathbf{m}_{k}} = -N_{k} \nu_{k} \left( \mathbf{W}_{k} \mathbf{m}_{k} - \mathbf{W}_{k} \overline{\mathbf{x}}_{k} \right) - \beta_{0} \nu_{k} \left( \mathbf{W}_{k} \mathbf{m}_{k} - \mathbf{W}_{k} \mathbf{m}_{0} \right).$$

Setting this equal to zero, using (10.60) and rearranging the terms, we obtain (10.61). Next we tackle  $\{\mathbf{W}_k, \nu_k\}$ . Here we need to perform a joint optimization w.r.t.  $\mathbf{W}_k$  and  $\nu_k$  for each  $k=1,\ldots,K$ . Like  $\beta_k$ ,  $\mathbf{W}_k$  and  $\nu_k$  appear in (10.71), (10.74) and (10.77). Using these, we can rewrite the r.h.s. of (10.70) as

$$\frac{1}{2} \sum_{k}^{K} \left( N_{k} \ln \widetilde{\Lambda}_{k} - N_{k} \nu_{k} \left\{ \operatorname{Tr} \left( \mathbf{S}_{k} \mathbf{W}_{k} \right) + \operatorname{Tr} \left( \mathbf{W}_{k} \left( \overline{\mathbf{x}}_{k} - \mathbf{m}_{k} \right) \left( \overline{\mathbf{x}}_{k} - \mathbf{m}_{k} \right)^{\mathrm{T}} \right) \right\} \\
+ \ln \widetilde{\Lambda}_{k} - \beta_{0} \nu_{k} \left( \mathbf{m}_{k} - \mathbf{m}_{0} \right)^{\mathrm{T}} \mathbf{W}_{k} \left( \mathbf{m}_{k} - \mathbf{m}_{0} \right) + \left( \nu_{0} - D - 1 \right) \ln \widetilde{\Lambda}_{k} \\
- \nu_{k} \operatorname{Tr} \left( \mathbf{W}_{0}^{-1} \mathbf{W}_{k} \right) - \ln \widetilde{\Lambda}_{k} + 2 \operatorname{H} \left[ q(\mathbf{\Lambda}_{k}) \right] \tag{269}$$

where we have dropped terms independent of  $\{\mathbf{W}_k, \nu_k\}$ ,  $\ln \widetilde{\Lambda}_k$  is given by (10.65),

$$H[q(\mathbf{\Lambda}_k)] = -\ln B(\mathbf{W}_k, \nu_k) - \frac{\nu_k - D - 1}{2} \ln \widetilde{\Lambda}_k \frac{\nu_k D}{2}$$
 (270)

where we have used (10.65) and (B.81), and, from (B.79),

$$\ln B(\mathbf{W}_k, \nu_k) = \frac{\nu_k}{2} \ln |\mathbf{W}_k| - \frac{\nu_k D}{2} - \sum_{i=1}^D \ln \Gamma\left(\frac{\nu_k + 1 - i}{2}\right). \tag{271}$$

Restricting attention to a single component, k, and making use of (270), (269) gives

$$\frac{1}{2}\left(N_k + \nu_0 - \nu_k\right) \ln \widetilde{\Lambda}_k - \frac{\nu_k}{2} \operatorname{Tr}\left(\mathbf{W}_k \mathbf{F}_k\right) - \ln B(\mathbf{W}_k, \nu_k) + \frac{\nu_k D}{2}$$
(272)

where

$$\mathbf{F}_{k} = \mathbf{W}_{0}^{-1} + N_{k} \mathbf{S}_{k} + N_{k} (\overline{\mathbf{x}}_{k} - \mathbf{m}_{k}) (\overline{\mathbf{x}}_{k} - \mathbf{m}_{k})^{\mathrm{T}} + \beta_{0} (\mathbf{m}_{k} - \mathbf{m}_{0}) (\mathbf{m}_{k} - \mathbf{m}_{0})^{\mathrm{T}} = \mathbf{W}_{0}^{-1} + N_{k} \mathbf{S}_{k} + \frac{N_{k} \beta_{0}}{N_{k} + \beta_{0}} (\overline{\mathbf{x}}_{k} - \mathbf{m}_{0}) (\overline{\mathbf{x}}_{k} - \mathbf{m}_{0})^{\mathrm{T}}$$
(273)

as we shall show below. Differentiating (272) w.r.t.  $\nu_k$ , making use of (271) and (10.65), and setting the result to zero, we get

$$0 = \frac{1}{2} \left( (N_k + \nu_0 - \nu_k) \frac{\mathrm{d} \ln \widetilde{\Lambda}_k}{\mathrm{d}\nu_k} - \ln \widetilde{\Lambda}_k - \mathrm{Tr} \left( \mathbf{W}_k \mathbf{F}_k \right) \right.$$
$$\left. + \ln |\mathbf{W}_k| + D \ln 2 + \sum_{i=1}^{D} \ln \Gamma \left( \frac{\nu_k + 1 - i}{2} \right) + D \right)$$
$$= \frac{1}{2} \left( (N_k + \nu_0 - \nu_k) \frac{\mathrm{d} \ln \widetilde{\Lambda}_k}{\mathrm{d}\nu_k} - \mathrm{Tr} \left( \mathbf{W}_k \mathbf{F}_k \right) + D \right). \tag{274}$$

Similarly, differentiating (272) w.r.t.  $W_k$ , making use of (271), (273) and (10.65), and setting the result to zero, we get

$$0 = \frac{1}{2} \left( (N_k + \nu_0 - \nu_k) \mathbf{W}_k^{-1} - \mathbf{F}_k + \mathbf{W}_k^{-1} \right)$$

$$= \frac{1}{2} \left( (N_k + \nu_0 - \nu_k) \mathbf{W}_k^{-1} - \mathbf{W}_0^{-1} - N_k \mathbf{S}_k - \frac{N_k \beta_0}{N_k + \beta_0} (\overline{\mathbf{x}}_k - \mathbf{m}_0) (\overline{\mathbf{x}}_k - \mathbf{m}_0)^{\mathrm{T}} + \mathbf{W}_k^{-1} \right)$$
(275)

We see that the simultaneous equations (274) and (275) are satisfied if and only if

$$0 = N_k + \nu_0 - \nu_k$$
  

$$0 = \mathbf{W}_0^{-1} + N_k \mathbf{S}_k + \frac{N_k \beta_0}{N_k + \beta_0} (\overline{\mathbf{x}}_k - \mathbf{m}_0) (\overline{\mathbf{x}}_k - \mathbf{m}_0)^{\mathrm{T}} - \mathbf{W}_k^{-1}$$

from which (10.63) and (10.62) follow, respectively. However, we still have to derive (273). From (10.60) and (10.61), derived above, we have

$$\mathbf{m}_k = \frac{\beta_0 \mathbf{m}_0 + N_k \overline{\mathbf{x}}_k}{\beta_0 + N_k}$$

and using this, we get

$$N_{k} (\overline{\mathbf{x}}_{k} - \mathbf{m}_{k}) (\overline{\mathbf{x}}_{k} - \mathbf{m}_{k})^{\mathrm{T}} + \beta_{0} (\mathbf{m}_{k} - \mathbf{m}_{0}) (\mathbf{m}_{k} - \mathbf{m}_{0})^{\mathrm{T}} =$$

$$N_{k} \overline{\mathbf{x}}_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}} - N_{k} \overline{\mathbf{x}}_{k} \frac{(\beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k})^{\mathrm{T}}}{\beta_{0} + N_{k}} - \frac{\beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k}}{\beta_{0} + N_{k}} N_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}}$$

$$+ \frac{N_{k} (\beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k}) (\beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k})^{\mathrm{T}}}{(\beta_{0} + N_{k})^{2}} + \frac{\beta_{0} (\beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k}) (\beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k})^{\mathrm{T}}}{(\beta_{0} + N_{k})^{2}}$$

$$- \beta_{0} \mathbf{m}_{0} \frac{(\beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k})^{\mathrm{T}}}{\beta_{0} + N_{k}} - \frac{\beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k}}{\beta_{0} \mathbf{m}_{0}^{\mathrm{T}} + \beta_{0} \mathbf{m}_{0} \mathbf{m}_{0}^{\mathrm{T}}}.$$

We now gather the coefficients of the terms on the r.h.s. as follows.

Coefficients of  $\overline{\mathbf{x}}_k \overline{\mathbf{x}}_k^{\mathrm{T}}$ :

$$\begin{split} N_k &- \frac{N_k^2}{\beta_0 + N_k} - \frac{N_k^2}{\beta_0 + N_k} + \frac{N_k^3}{(\beta_0 + N_k)^2} + \frac{\beta_0 N_k^2}{(\beta_0 + N_k)^2} \\ &= N_k - \frac{N_k^2}{\beta_0 + N_k} - \frac{N_k^2}{\beta_0 + N_k} + \frac{N_k^3}{(\beta_0 + N_k)^2} + \frac{\beta_0 N_k^2}{(\beta_0 + N_k)^2} \\ &= \frac{N_k \beta_0^2 + N_k^2 \beta_0 + N_k^2 \beta_0 + N_k^3 - 2(N_k^2 \beta_0 + N_k^3) + N_k^3 + \beta_0 N_k^2}{(\beta_0 + N_k)^2} \\ &= \frac{N_k \beta_0^2 + \beta_0 N_k^2}{(\beta_0 + N_k)^2} = \frac{N_k \beta_0 (\beta_0 + N_k)}{(\beta_0 + N_k)^2} = \frac{N_k \beta_0}{\beta_0 + N_k} \end{split}$$

Coefficients of  $\overline{\mathbf{x}}_k \mathbf{m}_0^{\mathrm{T}}$  and  $\mathbf{m}_0 \overline{\mathbf{x}}_k^{\mathrm{T}}$  (these are identical):

$$-\frac{N_k \beta_0}{\beta_0 + N_k} + \frac{\beta_0 N_k^2}{(\beta_0 + N_k)^2} + \frac{\beta_0^2 N_k}{(\beta_0 + N_k)^2} - \frac{N_k \beta_0}{\beta_0 + N_k}$$
$$= \frac{N_k \beta_0}{\beta_0 + N_k} - \frac{2N_k \beta_0}{\beta_0 + N_k} = -\frac{N_k \beta_0}{\beta_0 + N_k}$$

Coefficients of  $\mathbf{m}_0 \mathbf{m}_0^{\mathrm{T}}$ :

$$\begin{split} &\frac{N_k \beta_0^2}{(\beta_0 + N_k)^2} + \frac{\beta_0^3}{(\beta_0 + N_k)^2} - \frac{2\beta_0^2}{\beta_0 + N_k} + \beta_0 \\ &= \frac{\beta_0^2 (N_k + \beta_0)}{(\beta_0 + N_k)^2} - \frac{2\beta_0^2}{\beta_0 + N_k} + \beta_0 \\ &= \frac{\beta_0^2 - 2\beta_0^2 + \beta_0^2 + N_k \beta_0}{\beta_0 + N_k} = \frac{N_k \beta_0}{\beta_0 + N_k} \end{split}$$

Thus

$$N_{k} (\overline{\mathbf{x}}_{k} - \mathbf{m}_{k}) (\overline{\mathbf{x}}_{k} - \mathbf{m}_{k})^{\mathrm{T}} + \beta_{0} (\mathbf{m}_{k} - \mathbf{m}_{0}) (\mathbf{m}_{k} - \mathbf{m}_{0})^{\mathrm{T}}$$

$$= \frac{N_{k} \beta_{0}}{N_{k} + \beta_{0}} (\overline{\mathbf{x}}_{k} - \mathbf{m}_{0}) (\overline{\mathbf{x}}_{k} - \mathbf{m}_{0})^{\mathrm{T}}$$
(276)

as desired.

Now we turn our attention to  $\alpha$ , which appear in (10.72) and (10.73), through (10.66), and (10.76). Using these together with (B.23) and (B.24), we can differentiate (10.70) w.r.t.  $\alpha_k$  and set it equal to zero, yielding

$$\frac{\partial \mathcal{L}}{\partial \alpha_{k}} = [N_{k} + (\alpha_{0} - 1) - (\alpha_{k} - 1)] \frac{\partial \ln \widehat{\pi}_{k}}{\partial \alpha_{k}} - \ln \widehat{\pi}_{k} - \frac{\partial \ln C(\alpha)}{\partial \alpha_{k}}$$

$$= [N_{k} + (\alpha_{0} - 1) - (\alpha_{k} - 1)] \left\{ \psi_{1}(\alpha_{k}) - \psi_{1}(\widehat{\alpha}) \frac{\partial \widehat{\alpha}}{\partial \alpha_{k}} \right\}$$

$$+ \psi(\widehat{\alpha}) - \psi(\alpha_{k}) - \psi(\widehat{\alpha}) \frac{\partial \widehat{\alpha}}{\partial \alpha_{k}} + \psi(\alpha_{k})$$

$$= [N_{k} + (\alpha_{0} - 1) - (\alpha_{k} - 1)] \left\{ \psi_{1}(\alpha_{k}) - \psi_{1}(\widehat{\alpha}) \right\} = 0 \tag{277}$$

where  $\psi(\cdot)$  and  $\psi_1(\cdot)$  are di- and trigamma functions, respectively. If we assume that  $\alpha_0 > 0$ , (10.58) must hold for (277) to hold, since the trigamma function is strictly positive and monotoncally decreasing for arguments greater than zero.

Finally, we maximize (10.70) w.r.t.  $r_{nk}$ , subject to the constraints  $\sum_k r_{nk} = 1$  for all  $n = 1, \ldots, N$ . Note that  $r_{nk}$  not only appears in (10.72) and (10.75), but also in (10.71) through  $N_k$ ,  $\overline{\mathbf{x}}_k$  and  $\mathbf{S}_k$ , and so we must substitute using (10.51), (10.52) and (10.53), respectively. To simplify subsequent calculations, we start by considering the last two terms inside the braces of (10.71), which we write together as

$$\frac{1}{2} \sum_{k=1}^{K} \nu_k \operatorname{Tr}\left(\mathbf{W}_k \mathbf{Q}_k\right) \tag{278}$$

where, using (10.51), (10.52) and (10.53),

$$\mathbf{Q}_{k} = \sum_{n=1}^{N} r_{nk} \left( \mathbf{x}_{n} - \overline{\mathbf{x}}_{k} \right) \left( \mathbf{x}_{n} - \overline{\mathbf{x}}_{k} \right)^{\mathrm{T}} + N_{k} \left( \overline{\mathbf{x}}_{k} - \mathbf{m}_{k} \right) \left( \overline{\mathbf{x}}_{k} - \mathbf{m}_{k} \right)^{\mathrm{T}}$$

$$= \sum_{n=1}^{N} r_{nk} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathrm{T}} - 2N_{k} \overline{\mathbf{x}}_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}} + N_{k} \overline{\mathbf{x}}_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}}$$

$$+ N_{k} \overline{\mathbf{x}}_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}} - N_{k} \mathbf{m}_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}} - N_{k} \overline{\mathbf{x}}_{k} \mathbf{m}_{k}^{\mathrm{T}} + N_{k} \mathbf{m}_{k} \mathbf{m}_{k}^{\mathrm{T}}$$

$$= \sum_{n=1}^{N} r_{nk} \left( \mathbf{x}_{n} - \mathbf{m}_{k} \right) \left( \mathbf{x}_{n} - \mathbf{m}_{k} \right)^{\mathrm{T}}. \tag{279}$$

Using (10.51), (278) and (279), we can now consider all terms in (10.71) which

depend on  $r_{kn}$  and add the appropriate Lagrange multiplier terms, yielding

$$\frac{1}{2} \sum_{k=1}^{K} \sum_{n=1}^{N} r_{nk} \left( \ln \widetilde{\Lambda}_k - D\beta_k^{-1} \right) \\
- \frac{1}{2} \sum_{k=1}^{K} \sum_{n=1}^{N} r_{nk} \nu_k \left( \mathbf{x}_n - \mathbf{m}_k \right)^{\mathrm{T}} \mathbf{W}_k \left( \mathbf{x}_n - \mathbf{m}_k \right) \\
+ \sum_{k=1}^{K} \sum_{n=1}^{N} r_{nk} \ln \widetilde{\pi}_k - \sum_{k=1}^{K} \sum_{n=1}^{N} r_{nk} \ln r_{nk} + \sum_{n=1}^{N} \lambda_n \left( 1 - \sum_{k=1}^{K} r_{nk} \right).$$

Taking the derivative of this w.r.t.  $r_{kn}$  and setting it equal to zero we obtain

$$0 = \frac{1}{2} \ln \widetilde{\Lambda}_k - \frac{D}{2\beta_k} - \frac{1}{2} \nu_k (\mathbf{x}_n - \mathbf{m}_k)^{\mathrm{T}} \mathbf{W}_k (\mathbf{x}_n - \mathbf{m}_k) + \ln \widetilde{\pi}_k - \ln r_{nk} - 1 - \lambda_n.$$

Moving  $\ln r_{nk}$  to the l.h.s. and exponentiating both sides, we see that for each n

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

which is in agreement with (10.67); the normalized form is then given by (10.49).

**10.19** We start by performing the integration over  $\pi$  in (10.80), making use of the result

$$\mathbb{E}[\pi_k] = \frac{\alpha_k}{\widehat{\alpha}}$$

to give

$$p(\widehat{\mathbf{x}}|D) = \sum_{k=1}^K \frac{\alpha_k}{\widehat{\alpha}} \iint \mathcal{N}(\widehat{\mathbf{x}}|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) \, \mathrm{d}\boldsymbol{\mu}_k \, \mathrm{d}\boldsymbol{\Lambda}_k.$$

The variational posterior distribution over  $\mu$  and  $\Lambda$  is given from (10.59) by

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

Using this result we next perform the integration over  $\mu_k$ . This can be done explicitly by completing the square in the exponential in the usual way, or we can simply appeal to the general result (2.109) and (2.110) for the linear-Gaussian model from Chapter 2 to give

$$\int \mathcal{N}(\widehat{\mathbf{x}}|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1}) \mathcal{N}\left(\boldsymbol{\mu}_k|\mathbf{m}_k, (\beta_k \boldsymbol{\Lambda}_k)^{-1}\right) d\boldsymbol{\mu}_k = \mathcal{N}\left(\widehat{\mathbf{x}}|\mathbf{m}_k, (1 + \beta_k^{-1}) \boldsymbol{\Lambda}_k^{-1}\right).$$

Thus we have

$$p(\widehat{\mathbf{x}}|D) = \sum_{k=1}^{K} \frac{\alpha_k}{\widehat{\alpha}} \int \mathcal{N}\left(\widehat{\mathbf{x}}|\mathbf{m}_k, \left(1 + \beta_k^{-1}\right) \mathbf{\Lambda}_k^{-1}\right) \, \mathcal{W}(\mathbf{\Lambda}_k|\mathbf{W}_k, \nu_k) \, \mathrm{d}\mathbf{\Lambda}_k.$$

The final integration over  $\Lambda_k$  is the convolution of a Wishart with a Gaussian. Omitting multiplicative constants which are independent of  $\hat{\mathbf{x}}$  we have

$$\begin{split} \int \mathcal{N}\left(\widehat{\mathbf{x}}|\mathbf{m}_{k}, \left(1+\beta_{k}^{-1}\right) \mathbf{\Lambda}_{k}^{-1}\right) \, \mathcal{W}(\mathbf{\Lambda}_{k}|\mathbf{W}_{k}, \nu_{k}) \, \mathrm{d}\mathbf{\Lambda}_{k} \\ &\propto \int |\mathbf{\Lambda}_{k}|^{1/2 + (\nu_{k} - D - 1)/2} \exp\left\{-\frac{1}{2\left(1+\beta_{k}^{-1}\right)} \mathrm{Tr}\left[\mathbf{\Lambda}_{k}(\widehat{\mathbf{x}} - \mathbf{m}_{k})(\widehat{\mathbf{x}} - \mathbf{m}_{k})^{\mathrm{T}}\right] \right. \\ &\left. -\frac{1}{2} \mathrm{Tr}\left[\mathbf{\Lambda}_{k} \mathbf{W}_{k}^{-1}\right]\right\} \mathrm{d}\mathbf{\Lambda}_{k}. \end{split}$$

We can now perform this integration by observing that the argument of the integral is an un-normalized Wishart distribution (unsurprisingly since the Wishart is the conjugate prior for the precision of a Gaussian) and so we can write down the result of this integration, up to an overall constant, by using the known normalization coefficient for the Wishart, given by (B.79). Thus we have

$$\int \mathcal{N}\left(\widehat{\mathbf{x}}|\mathbf{m}_{k}, \left(1 + \beta_{k}^{-1}\right) \mathbf{\Lambda}_{k}^{-1}\right) \mathcal{W}(\mathbf{\Lambda}_{k}|\mathbf{W}_{k}, \nu_{k}) d\mathbf{\Lambda}_{k}$$

$$\propto \left|\mathbf{W}_{k}^{-1} + \frac{1}{\left(1 + \beta_{k}^{-1}\right)} (\widehat{\mathbf{x}} - \mathbf{m}_{k}) (\widehat{\mathbf{x}} - \mathbf{m}_{k})^{\mathrm{T}} \right|^{-(\nu_{k} + 1)/2}$$

$$\propto \left|\mathbf{I} + \frac{1}{\left(1 + \beta_{k}^{-1}\right)} \mathbf{W}_{k} (\widehat{\mathbf{x}} - \mathbf{m}_{k}) (\widehat{\mathbf{x}} - \mathbf{m}_{k})^{\mathrm{T}} \right|^{-(\nu_{k} + 1)/2}$$

where we have omitted factors independent of  $\hat{\mathbf{x}}$  since we are only interested in the functional dependence on  $\hat{\mathbf{x}}$ , and we have made use of the result  $|\mathbf{A}\mathbf{B}| = |\mathbf{A}||\mathbf{B}|$  and omitted an overall factor of  $|\mathbf{W}_k^{-1}|$ . Next we use the identity

$$\left|\mathbf{I} + \mathbf{a}\mathbf{b}^{\mathrm{T}}\right| = (1 + \mathbf{a}^{\mathrm{T}}\mathbf{b})$$

where a and b are D-dimensional vectors and I is the  $D \times D$  unit matrix, to give

$$\int \mathcal{N}\left(\widehat{\mathbf{x}}|\mathbf{m}_{k}, \left(1 + \beta_{k}^{-1}\right) \mathbf{\Lambda}_{k}^{-1}\right) \mathcal{W}(\mathbf{\Lambda}_{k}|\mathbf{W}_{k}, \nu_{k}) d\mathbf{\Lambda}_{k}$$

$$\propto \left\{1 + \frac{1}{\left(1 + \beta_{k}^{-1}\right)} (\widehat{\mathbf{x}} - \mathbf{m}_{k})^{\mathrm{T}} \mathbf{W}_{k} (\widehat{\mathbf{x}} - \mathbf{m}_{k})\right\}^{-(\nu_{k} + 1)/2}.$$

We recognize this result as being a Student distribution, and by comparison with the standard form (B.68) for the Student we see that it has mean  $\mathbf{m}_k$ , precision given by (10.82) and degrees of freedom parameter  $\nu_k + 1 - D$ . We can re-instate the normalization coefficient using the standard form for the Student distribution given in Appendix B.

**10.20** Consider first the posterior distribution over the precision of component k given by

$$q^{\star}(\mathbf{\Lambda}_k) = \mathcal{W}(\mathbf{\Lambda}_k | \mathbf{W}_k, \nu_k).$$

From (10.63) we see that for large N we have  $\nu_k \to N_k$ , and similarly from (10.62) we see that  $\mathbf{W}_k \to N_k^{-1} \mathbf{S}_k^{-1}$ . Thus the mean of the distribution over  $\mathbf{\Lambda}_k$ , given by  $\mathrm{E}[\mathbf{\Lambda}_k] = \nu_k \mathbf{W}_k \to \mathbf{S}_k^{-1}$  which is the maximum likelihood value (this assumes that the quantities  $r_{nk}$  reduce to the corresponding EM values, which is indeed the case as we shall show shortly). In order to show that this posterior is also sharply peaked, we consider the differential entropy,  $\mathrm{H}[\mathbf{\Lambda}_k]$  given by (B.82), and show that, as  $N_k \to \infty$ ,  $\mathrm{H}[\mathbf{\Lambda}_k] \to 0$ , corresponding to the density collapsing to a spike. First consider the normalizing constant  $B(\mathbf{W}_k, \nu_k)$  given by (B.79). Since  $\mathbf{W}_k \to N_k^{-1} \mathbf{S}_k^{-1}$  and  $\nu_k \to N_k$ ,

$$-\ln B(\mathbf{W}_k, \nu_k) \to -\frac{N_k}{2} \left( D \ln N_k + \ln |\mathbf{S}_k| - D \ln 2 \right) + \sum_{i=1}^{D} \ln \Gamma \left( \frac{N_k + 1 - i}{2} \right).$$

We then make use of Stirling's approximation (1.146) to obtain

$$\ln\Gamma\left(\frac{N_k+1-i}{2}\right) \simeq \frac{N_k}{2} \left(\ln N_k - \ln 2 - 1\right)$$

which leads to the approximate limit

$$-\ln B(\mathbf{W}_{k}, \nu_{k}) \rightarrow -\frac{N_{k}D}{2} (\ln N_{k} - \ln 2 - \ln N_{k} + \ln 2 + 1) - \frac{N_{k}}{2} \ln |\mathbf{S}_{k}|$$

$$= -\frac{N_{k}}{2} (\ln |\mathbf{S}_{k}| + D). \tag{280}$$

Next, we use (10.241) and (B.81) in combination with  $\mathbf{W}_k \to N_k^{-1} \mathbf{S}_k^{-1}$  and  $\nu_k \to N_k$  to obtain the limit

$$\mathbb{E} \left[ \ln |\mathbf{\Lambda}| \right] \quad \to \quad D \ln \frac{N_k}{2} + D \ln 2 - D \ln N_k - \ln |\mathbf{S}_k|$$
$$= \quad -\ln |\mathbf{S}_k|,$$

where we approximated the argument to the digamma function by  $N_k/2$ . Substituting this and (280) into (B.82), we get

$$H[\mathbf{\Lambda}] \to 0$$

when  $N_k \to \infty$ .

Next consider the posterior distribution over the mean  $\mu_k$  of the  $k^{\mathrm{th}}$  component given by

From (10.61) we see that for larger N the mean  $m_k$  of this distribution reduces to  $\overline{\mathbf{x}}_k$  which is the corresponding maximum likelihood value. From (10.60) we see that

 $\beta_k \to N_k$  and Thus the precision  $\beta_k \Lambda_k \to \beta_k \nu_k \mathbf{W}_k \to N_k \mathbf{S}_k^{-1}$  which is large for large N and hence this distribution is sharply peaked around its mean.

Now consider the posterior distribution  $q^*(\pi)$  given by (10.57). For large N we have  $\alpha_k \to N_k$  and so from (B.17) and (B.19) we see that the posterior distribution becomes sharply peaked around its mean  $\mathrm{E}[\pi_k] = \alpha_k/\overline{\alpha} \to N_k/N$  which is the maximum likelihood solution.

For the distribution  $q^*(\mathbf{z})$  we consider the responsibilities given by (10.67). Using (10.65) and (10.66), together with the asymptotic result for the digamma function, we again obtain the maximum likelihood expression for the responsibilities for large N.

Finally, for the predictive distribution we first perform the integration over  $\pi$ , as in the solution to Exercise 10.19, to give

$$p(\widehat{\mathbf{x}}|D) = \sum_{k=1}^{K} \frac{\alpha_k}{\overline{\alpha}} \iint \mathcal{N}(\widehat{\mathbf{x}}|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) \, \mathrm{d}\boldsymbol{\mu}_k \, \mathrm{d}\boldsymbol{\Lambda}_k.$$

The integrations over  $\mu_k$  and  $\Lambda_k$  are then trivial for large N since these are sharply peaked and hence approximate delta functions. We therefore obtain

$$p(\widehat{\mathbf{x}}|D) = \sum_{k=1}^{K} \frac{N_k}{N} \mathcal{N}(\widehat{\mathbf{x}}|\overline{\mathbf{x}}_k, \mathbf{W}_k)$$

which is a mixture of Gaussians, with mixing coefficients given by  $N_k/N$ .

- **10.21** The number of equivalent parameter settings equals the number of possible assignments of K parameter sets to K mixture components: K for the first component, times K-1 for the second component, times K-2 for the third and so on, giving the result K!.
- **10.22** The mixture distribution over the parameter space takes the form

$$q(\mathbf{\Theta}) = \frac{1}{K!} \sum_{\kappa=1}^{K!} q_{\kappa}(\boldsymbol{\theta}_{\kappa})$$

where  $\theta_{\kappa} = \{\mu_k, \Sigma_k, \pi\}$ ,  $\kappa$  indexes the components of this mixture and  $\Theta = \{\theta_{\kappa}\}$ . With this model, (10.3) becomes

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

$$= \frac{1}{K!} \sum_{\kappa=1}^{K!} \int q_{\kappa}(\boldsymbol{\theta}_{\kappa}) \ln p(\mathbf{X}, \boldsymbol{\theta}_{\kappa}) d\boldsymbol{\theta}_{\kappa}$$

$$- \frac{1}{K!} \sum_{\kappa=1}^{K!} \int q_{\kappa}(\boldsymbol{\theta}_{\kappa}) \ln \left( \frac{1}{K!} \sum_{\kappa'=1}^{K!} q_{\kappa'}(\boldsymbol{\theta}_{\kappa'}) \right) d\boldsymbol{\theta}_{\kappa}$$

$$= \int q(\boldsymbol{\theta}) \ln p(\mathbf{X}, \boldsymbol{\theta}) d\boldsymbol{\theta} - \int q(\boldsymbol{\theta}) \ln q(\boldsymbol{\theta}) d\boldsymbol{\theta} + \ln K!$$

where  $q(\theta)$  corresponds to any one of the K! equivalent  $q_{\kappa}(\theta_{\kappa})$  distributions. Note that in the last step, we use the assumption that the overlap between these distributions is negligible and hence

$$\int q_{\kappa}(\boldsymbol{\theta}) \ln q_{\kappa'}(\boldsymbol{\theta}) \, \mathrm{d}\boldsymbol{\theta} \simeq 0$$

when  $\kappa \neq \kappa'$ .

10.23 When we are treating  $\pi$  as a parameter, there is neither a prior, nor a variational posterior distribution, over  $\pi$ . Therefore, the only term remaining from the lower bound, (10.70), that involves  $\pi$  is the second term, (10.72). Note however, that (10.72) involves the *expectations* of  $\ln \pi_k$  under  $q(\pi)$ , whereas here, we operate directly with  $\pi_k$ , yielding

$$\mathbb{E}_{q(\mathbf{Z})}[\ln p(\mathbf{Z}|\boldsymbol{\pi})] = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \ln \pi_k.$$

Adding a Langrange term, as in (9.20), taking the derivative w.r.t.  $\pi_k$  and setting the result to zero we get

$$\frac{N_k}{\pi_k} + \lambda = 0, (281)$$

where we have used (10.51). By re-arranging this to

$$N_k = -\lambda \pi_k$$

and summing both sides over k, we see that  $-\lambda = \sum_k N_k = N$ , which we can use to eliminate  $\lambda$  from (281) to get (10.83).

10.24 The singularities that may arise in maximum likelihood estimation are caused by a mixture component, k, collapsing on a data point,  $\mathbf{x}_n$ , i.e.,  $r_{kn}=1$ ,  $\mu_k=\mathbf{x}_n$  and  $|\mathbf{\Lambda}_k| \to \infty$ .

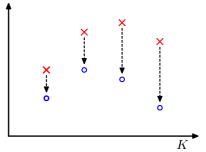
However, the prior distribution  $p(\mu, \Lambda)$  defined in (10.40) will prevent this from happening, also in the case of MAP estimation. Consider the product of the expected complete log-likelihood and  $p(\mu, \Lambda)$  as a function of  $\Lambda_k$ :

$$\begin{split} \mathbb{E}_{q(\mathbf{Z})} \left[ \ln p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \right] \\ &= \frac{1}{2} \sum_{n=1}^{N} r_{kn} \left( \ln |\boldsymbol{\Lambda}_{k}| - (\mathbf{x}_{n} - \boldsymbol{\mu}_{k})^{\mathrm{T}} \boldsymbol{\Lambda}_{k} (\mathbf{x}_{n} - \boldsymbol{\mu}_{k}) \right) \\ &+ \ln |\boldsymbol{\Lambda}_{k}| - \beta_{0} (\boldsymbol{\mu}_{k} - \mathbf{m}_{0})^{\mathrm{T}} \boldsymbol{\Lambda}_{k} (\boldsymbol{\mu}_{k} - \mathbf{m}_{0}) \\ &+ (\nu_{0} - D - 1) \ln |\boldsymbol{\Lambda}_{k}| - \mathrm{Tr} \left[ \mathbf{W}_{0}^{-1} \boldsymbol{\Lambda}_{k} \right] + \mathrm{const.} \end{split}$$

where we have used (10.38), (10.40) and (10.50), together with the definitions for the Gaussian and Wishart distributions; the last term summarizes terms independent of  $\Lambda_k$ . Using (10.51)–(10.53), we can rewrite this as

$$(\nu_0 + N_k - D) \ln |\mathbf{\Lambda}_k| - \operatorname{Tr} \left[ (\mathbf{W}_0^{-1} + \beta_0 (\boldsymbol{\mu}_k - \mathbf{m}_0) (\boldsymbol{\mu}_k - \mathbf{m}_0)^{\mathrm{T}} + N_k \mathbf{S}_k) \mathbf{\Lambda}_k \right],$$

Figure 9 Illustration of the true log marginal likelihood for a Gaussian mixture model  $(\times)$  and the corresponding variational bound obtained from a factorized approximation  $(\circ)$  as functions of the number of mixture components, K. The dashed arrows emphasize the typical increase in the difference between the true log marginal likelihood and the bound. As a consequence, the bound tends to have its peak at a lower value of K than the true log marginal likelihood.



where we have dropped the constant term. Using (C.24) and (C.28), we can compute the derivative of this w.r.t.  $\Lambda_k$  and setting the result equal to zero, we find the MAP estimate for  $\Lambda_k$  to be

$$\boldsymbol{\Lambda}_k^{-1} = \frac{1}{\nu_0 + N_k - D} (\mathbf{W}_0^{-1} + \beta_0 (\boldsymbol{\mu}_k - \mathbf{m}_0) (\boldsymbol{\mu}_k - \mathbf{m}_0)^{\mathrm{T}} + N_k \mathbf{S}_k).$$

From this we see that  $|\Lambda_k^{-1}|$  can never become 0, because of the presence of  $\mathbf{W}_0^{-1}$  (which we must chose to be positive definite) in the expression on the r.h.s.

- As the number of mixture components grows, so does the number of variables that may be correlated, but which are treated as independent under a variational approximation, as illustrated in Figure 10.2. As a result of this, the proportion of probability mass under the true distribution,  $p(\mathbf{Z}, \pi, \mu, \Sigma | \mathbf{X})$ , that the variational approximation  $q(\mathbf{Z}, \pi, \mu, \Sigma)$  does not capture, will grow. The consequence will be that the second term in (10.2), the KL divergence between  $q(\mathbf{Z}, \pi, \mu, \Sigma)$  and  $p(\mathbf{Z}, \pi, \mu, \Sigma | \mathbf{X})$ , will increase. Since this KL divergence is the difference between the true log marginal and the corresponding the lower bound, the latter must decrease compared to the former. Thus, as illustrated in Figure 9, chosing the number of components based on the lower bound will tend to underestimate the optimal number of components.
- **10.26** Extending the variational treatment of Section 10.3 to also include  $\beta$ , we specify the prior for  $\beta$

$$p(\beta) = \operatorname{Gam}(\beta|c_0, d_0) \tag{282}$$

and modify (10.90) as

$$p(\mathbf{t}, \mathbf{w}, \alpha, \beta) = p(\mathbf{t}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)p(\alpha)p(\beta)$$
(283)

where the first factor on the r.h.s. correspond to (10.87) with the dependence on  $\beta$  made explicit.

The formulae for  $q^*(\alpha)$ , (10.93)–(10.95), remain unchanged. For  $q(\mathbf{w})$ , we follow the path mapped out in Section 10.3, incorporating the modifications required by the

changed treatment of  $\beta$ ; (10.96)–(10.98) now become

$$\ln q^{\star}(\mathbf{w}) = \mathbb{E}_{\beta} \left[ \ln p(\mathbf{t}|\mathbf{w}, \beta) \right] + \mathbb{E}_{\alpha} \left[ \ln p(\mathbf{w}|\alpha) \right] + \text{const}$$

$$= -\frac{\mathbb{E}[\beta]}{2} \sum_{n=1}^{N} \left\{ \mathbf{w}^{T} \phi_{n} - t_{n} \right\}^{2} - \frac{\mathbb{E}[\alpha]}{2} \mathbf{w}^{T} \mathbf{w} + \text{const}$$

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

Accordingly, (10.100) and (10.101) become

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

For  $q(\beta)$ , we use (10.9), (282) and (283) to obtain

$$\ln q^{\star}(\beta) = \mathbb{E}_{\mathbf{w}} \left[ \ln p(\mathbf{t}|\mathbf{w}, \beta) \right] + \ln p(\beta) + \text{const}$$
$$= \frac{N}{2} \ln \beta - \frac{\beta}{2} \mathbb{E}_{\mathbf{w}} \left[ \sum_{n=1}^{N} \left\{ \mathbf{w}^{T} \phi_{n} - t_{n} \right\}^{2} \right] + (c_{0} - 1) \ln \beta - d_{0} \beta$$

which we recognize as the logarithm of a Gamma distribution with parameters

$$c_{N} = c_{0} + \frac{N}{2}$$

$$d_{N} = d_{0} + \frac{1}{2}\mathbb{E}\left[\sum_{n=1}^{N} (\mathbf{w}^{T}\phi_{n} - t_{n})^{2}\right]$$

$$= d_{0} + \frac{1}{2} \left(\operatorname{Tr}\left(\mathbf{\Phi}^{T}\mathbf{\Phi}\mathbb{E}\left[\mathbf{w}\mathbf{w}^{T}\right]\right) + \mathbf{t}^{T}\mathbf{t}\right) - \mathbf{t}^{T}\mathbf{\Phi}\mathbb{E}[\mathbf{w}]$$

$$= d_{0} + \frac{1}{2} \left(\|\mathbf{t} - \mathbf{\Phi}\mathbf{m}_{N}\|^{2} + \operatorname{Tr}\left(\mathbf{\Phi}^{T}\mathbf{\Phi}\mathbf{S}_{N}\right)\right)$$

where we have used (10.103) and, from (B.38),

$$\mathbb{E}[\mathbf{w}] = \mathbf{m}_N. \tag{284}$$

Thus, from (B.27),

$$\mathbb{E}[\beta] = \frac{c_N}{d_N}.\tag{285}$$

In the lower bound, (10.107), the first term will be modified and two new terms added on the r.h.s.We start with the modified log likelihood term:

$$\mathbb{E}_{\beta} \left[ \mathbb{E}_{\mathbf{w}} \left[ \ln p(\mathbf{t} | \mathbf{w}, \beta) \right] \right] = \frac{N}{2} \left( \mathbb{E}[\beta] - \ln(2\pi) \right) - \frac{\mathbb{E}[\beta]}{2} \mathbb{E} \left[ \|\mathbf{t} - \mathbf{\Phi} \mathbf{w}\|^2 \right]$$
$$= \frac{N}{2} \left( \psi(c_N) - \ln d_N - \ln(2\pi) \right)$$
$$- \frac{c_N}{2d_N} \left( |\mathbf{t} - \mathbf{\Phi} \mathbf{w}\|^2 + \operatorname{Tr} \left( \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \mathbf{S}_N \right) \right)$$

where we have used (284), (285), (10.103) and (B.30). Next we consider the term corresponding log prior over  $\beta$ :

$$\mathbb{E} \left[ \ln p(\beta) \right] = (c_0 - 1) \mathbb{E} [\ln \beta] - d_0 \mathbb{E} [\beta] + c_0 \ln d_0 - \ln \Gamma(c_0)$$

$$= (c_0 - 1) (\psi(c_N) - \ln d_N) - \frac{d_0 c_N}{d_N} + c_0 \ln d_0 - \ln \Gamma(c_0)$$

where we have used (285) and (B.30). Finally, from (B.31), we get the last term in the form of the negative entropy of the posterior over  $\beta$ :

$$-\mathbb{E}\left[\ln q^{\star}(\beta)\right] = (c_N - 1)\psi(c_N) + \ln d_N - c_N - \ln \Gamma(c_N).$$

Finally, the predictive distribution is given by (10.105) and (10.106), with  $1/\beta$  replaced by  $1/\mathbb{E}[\beta]$ .

**10.27** Consider each of the five terms in the lower bound (10.107) in turn. For the terms arising from the likelihood function we have

$$\mathbb{E}[\ln p(\mathbf{t}|\mathbf{w})] = -\frac{N}{2}\ln(2\pi) + \frac{N}{2}\ln\beta - \frac{\beta}{2}\mathbb{E}\left[\sum_{n=1}^{N}(t_n - \mathbf{w}^{\mathrm{T}}\boldsymbol{\phi}_n)^2\right]$$
$$= -\frac{N}{2}\ln(2\pi) + \frac{N}{2}\ln\beta$$
$$-\frac{\beta}{2}\left\{\mathbf{t}^{\mathrm{T}}\mathbf{t} - 2\mathbb{E}[\mathbf{w}^{\mathrm{T}}]\boldsymbol{\Phi}^{\mathrm{T}}\mathbf{t} + \mathrm{Tr}\left(\mathbb{E}[\mathbf{w}\mathbf{w}^{\mathrm{T}}]\boldsymbol{\Phi}^{\mathrm{T}}\boldsymbol{\Phi}\right)\right\}.$$

The prior over w gives rise to

$$\mathbb{E}[\ln p(\mathbf{w}|\alpha)] = -\frac{M}{2}\ln(2\pi) + \frac{M}{2}\mathbb{E}[\ln \alpha] - \frac{\mathbb{E}[\alpha]}{2}\mathbb{E}[\mathbf{w}^{\mathrm{T}}\mathbf{w}].$$

Similarly, the prior over  $\alpha$  gives

$$\mathbb{E}[\ln p(\alpha)] = a_0 \ln b_0 + (a_0 - 1)\mathbb{E}[\ln \alpha] - b_0 \mathbb{E}[\alpha] - \ln \Gamma(a_0).$$

The final two terms in  $\mathcal{L}$  represent the negative entropies of the Gaussian and gamma distributions, and are given by (B.41) and (B.31) respectively, so that

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

Similarly we have

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

Now we substitute the following expressions for the moments

$$\mathbb{E}[\mathbf{w}] = \mathbf{m}_{N}$$

$$\mathbb{E}[\mathbf{w}\mathbf{w}^{\mathrm{T}}] = \mathbf{m}_{N}\mathbf{m}_{N}^{\mathrm{T}} + \mathbf{S}_{N}$$

$$\mathbb{E}[\alpha] = \frac{a_{N}}{b_{N}}$$

$$\mathbb{E}[\ln \alpha] = \psi(a_{N}) - \ln b_{N}.$$

and combine the various terms together to obtain (10.107).

**10.28 NOTE**: In PRML, Equations (10.119)–(10.121) contain errors; please consult the PRML Errata for relevant corrections.

We start by writing the complete-data likelihood, given by (10.37) and (10.38) in a form corresponding to (10.113). From (10.37) and (10.38), we have

$$\begin{split} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) &= p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\mathbf{Z} | \boldsymbol{\pi}) \\ &= \prod_{n=1}^{N} \prod_{k=1}^{K} \left( \pi_{k} \mathcal{N} \left( \mathbf{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Lambda}_{k}^{-1} \right) \right)^{z_{nk}} \end{split}$$

which is a product over data points, just like (10.113). Focussing on the individual factors of this product, we have

$$p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{k=1}^K \left( \pi_k \mathcal{N} \left( \mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1} \right) \right)^{z_{nk}} = \exp \left\{ \sum_{k=1}^K z_{nk} \left( \ln \pi_k + \frac{1}{2} \ln |\boldsymbol{\Lambda}_k| - \frac{D}{2} \ln(2\pi) - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^{\mathrm{T}} \boldsymbol{\Lambda} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) \right\}.$$

Drawing on results from Solution 2.57, we can rewrite this in the form of (10.113), with

$$\boldsymbol{\eta} = \begin{bmatrix} \boldsymbol{\Lambda}_{k} \boldsymbol{\mu}_{k} \\ \boldsymbol{\Lambda}_{k} \\ \boldsymbol{\mu}_{k}^{\mathrm{T}} \boldsymbol{\Lambda}_{k} \boldsymbol{\mu}_{k} \\ \ln |\boldsymbol{\Lambda}_{k}| \\ \ln \pi_{k} \end{bmatrix}_{k=1,\dots,K}$$
(286)

$$\mathbf{u}(\mathbf{x}_{n}, \mathbf{z}_{n}) = \begin{bmatrix} \mathbf{z}_{nk} & \frac{\mathbf{x}_{n}}{\frac{1}{2}\mathbf{x}_{n}\mathbf{x}_{n}^{\mathrm{T}}} \\ -\frac{1}{2} \\ \frac{1}{2} \\ 1 \end{bmatrix} \bigg]_{k=1,\dots,K}$$
(287)

$$h(\mathbf{x}_n, \mathbf{z}_n) = \prod_{k=1}^K \left( (2\pi)^{-D/2} \right)^{z_{nk}}$$

$$g(\boldsymbol{\eta}) = 1$$
(288)

where we have introduce the notation

$$[\mathbf{v}_k]_{k=1,...,K} = \left[egin{array}{c} \mathbf{v}_1 \ \mathbf{v}_2 \ dots \ \mathbf{v}_K \end{array}
ight]$$

and the operator  $\overrightarrow{M}$  which returns a vector formed by stacking the columns of the argument matrix on top of each other.

Next we seek to rewrite the prior over the parameters, given by (10.39) and (10.40), in a form corresponding to (10.114) and which also matches (286). From, (10.39), (10.40) and Appendix B, we have

$$p(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Lambda}) = \operatorname{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}_0) \prod_{k=1}^K \mathcal{N} \left( \boldsymbol{\mu}_k | \mathbf{m}_0, (\beta_0 \boldsymbol{\Lambda}_k)^{-1} \right) \mathcal{W} \left( \boldsymbol{\Lambda}_k | \mathbf{W}_0, \nu_0 \right)$$

$$= C(\boldsymbol{\alpha}_0) \left( \frac{\beta_0}{2\pi} \right)^{KD/2} B(\mathbf{W}_0, \nu_0)^K \exp \left\{ \sum_{k=1}^K (\alpha_0 - 1) \ln \pi_k + \frac{\nu_0 - D}{2} \ln |\boldsymbol{\Lambda}_k| - \frac{1}{2} \operatorname{Tr} \left( \boldsymbol{\Lambda}_k \left[ \beta_0 (\boldsymbol{\mu}_k - \mathbf{m}_0) (\boldsymbol{\mu}_k - \mathbf{m}_0)^{\mathrm{T}} + \mathbf{W}_0 \right] \right) \right\}.$$

we can rewrite this in the form of (10.114) with  $\eta$  given by (286),

$$\chi_{0} = \begin{bmatrix} -\frac{1}{2} \left( \beta_{0} \overrightarrow{\mathbf{m}_{0}} \overrightarrow{\mathbf{m}_{0}^{T}} + \overrightarrow{\mathbf{W}_{0}^{-1}} \right) \\ -\beta_{0}/2 \\ (\nu_{0} - D)/2 \\ \alpha_{0} - 1 \end{bmatrix}_{h=1...K}$$

$$(289)$$

$$g(\eta) = 1$$

$$f(v_0, \boldsymbol{\chi}_0) = C(\boldsymbol{\alpha}_0) \left(\frac{\beta_0}{2\pi}\right)^{KD/2} B(\mathbf{W}_0, \nu_0)^K$$

and  $v_0$  replaces  $v_0$  in (10.114) to avoid confusion with  $v_0$  in (10.40).

Having rewritten the Bayesian mixture of Gaussians as a conjugate model from the exponential family, we now proceed to rederive (10.48), (10.57) and (10.59). By exponentiating both sides of (10.115) and making use of (286)–(288), we obtain (10.47), with  $\rho_{nk}$  given by (10.46), from which (10.48) follows.

Next we can use (10.50) to take the expectation w.r.t.  $\mathbf{Z}$  in (10.121), substituting  $r_{nk}$  for  $\mathbb{E}[z_{nk}]$  in (287). Combining this with (289), (10.120) and (10.121) become

$$\upsilon_N = \upsilon_0 + N = 1 + N$$

and

$$v_{N}\boldsymbol{\chi}_{N} = \begin{bmatrix} -\frac{1}{2} \left( \beta_{0} \overline{\mathbf{m}_{0}} \mathbf{m}_{0}^{\mathrm{T}} + \overline{\mathbf{W}_{0}^{-1}} \right) \\ -\beta_{0}/2 \\ (\nu_{0} - D)/2 \\ \alpha_{0} - 1 \end{bmatrix} + \sum_{n=1}^{N} \begin{bmatrix} r_{nk} \\ \frac{1}{2} \overline{\mathbf{x}_{n}} \mathbf{x}_{n}^{\mathrm{T}} \\ -\frac{1}{2} \\ \frac{1}{2} \\ 1 \end{bmatrix} \end{bmatrix}_{k=1,\dots,K}$$

$$= \begin{bmatrix} \beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k} \\ -\frac{1}{2} \left( \beta_{0} \overline{\mathbf{m}_{0}} \mathbf{m}_{0}^{\mathrm{T}} + \overline{\mathbf{W}_{0}^{-1}} + N_{k} \left( \overline{\mathbf{S}}_{k} + \overline{\mathbf{x}}_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}} \right) \\ -(\beta_{0} + N_{k})/2 \\ (\nu_{0} - D + N_{k})/2 \\ \alpha_{0} - 1 + N_{k} \end{bmatrix}_{k=1,\dots,K}$$

$$(290)$$

where we use  $v_N$  instead of  $v_N$  in (10.119)–(10.121), to avoid confusion with  $v_k$ , which appears in (10.59) and (10.63). From the bottom row of (287) and (290), we see that the inner product of  $\eta$  and  $v_N \chi_N$  gives us the r.h.s. of (10.56), from which (10.57) follows. The remaining terms of this inner product are

$$\begin{split} \sum_{k=1}^{K} & \Big\{ \boldsymbol{\mu}_{k}^{\mathrm{T}} \boldsymbol{\Lambda}_{k} \left( \beta_{0} \mathbf{m}_{0} + N_{k} \overline{\mathbf{x}}_{k} \right) \\ & - \frac{1}{2} \mathrm{Tr} \left( \boldsymbol{\Lambda}_{k} \left[ \beta_{0} \overrightarrow{\mathbf{m}_{0}} \overrightarrow{\mathbf{m}_{0}}^{\mathrm{T}} + \overrightarrow{\mathbf{W}_{0}}^{-1} + N_{k} \overline{\left( \mathbf{S}_{k} + \overline{\mathbf{x}}_{k} \overline{\mathbf{x}}_{k}^{\mathrm{T}} \right)} \right] \right) \\ & - \frac{1}{2} (\beta_{0} + N_{k}) \boldsymbol{\mu}_{k}^{\mathrm{T}} \boldsymbol{\Lambda}_{k} \boldsymbol{\mu}_{k} + \frac{1}{2} (\nu_{0} + N_{k} - D) \ln |\boldsymbol{\Lambda}| \Big\}. \end{split}$$

Restricting our attention to parameters corresponding to a single mixture component and making use of (10.60), (10.61) and (10.63), we can rewrite this as

$$-\frac{1}{2}\beta_{k}\boldsymbol{\mu}_{k}^{\mathrm{T}}\boldsymbol{\Lambda}_{k}\boldsymbol{\mu}_{k} + \beta_{k}\boldsymbol{\mu}_{k}^{\mathrm{T}}\boldsymbol{\Lambda}_{k}\mathbf{m}_{k} - \frac{1}{2}\beta_{k}\mathbf{m}_{k}^{\mathrm{T}}\boldsymbol{\Lambda}_{k}\mathbf{m}_{k} + \frac{1}{2}\ln|\boldsymbol{\Lambda}|$$

$$+\frac{1}{2}\beta_{k}\mathbf{m}_{k}^{\mathrm{T}}\boldsymbol{\Lambda}_{k}\mathbf{m}_{k} - \frac{1}{2}\mathrm{Tr}\left(\boldsymbol{\Lambda}_{k}\left[\beta_{0}\mathbf{m}_{0}\mathbf{m}_{0}^{\mathrm{T}} + \overrightarrow{\mathbf{W}}_{0}^{\mathrm{T}} + N_{k}(\overrightarrow{\mathbf{S}_{k}} + \overline{\mathbf{x}}_{k}\overline{\mathbf{x}}_{k}^{\mathrm{T}})\right]\right)$$

$$+\frac{1}{2}(\nu_{k} - D - 1)\ln|\boldsymbol{\Lambda}|.$$

The first four terms match the logarithm of  $\mathcal{N}\left(\boldsymbol{\mu}_{k}|\mathbf{m}_{k},\left(\beta_{k}\boldsymbol{\Lambda}_{k}\right)^{-1}\right)$  from the r.h.s. of (10.59); the missing  $D/2[\ln\beta_{k}-\ln(2\pi)]$  can be accounted for in  $f(\upsilon_{N},\boldsymbol{\chi}_{N})$ . To make the remaining terms match the logarithm of  $\mathcal{W}\left(\boldsymbol{\Lambda}_{k}|\mathbf{W}_{0},\nu_{0}\right)$  from the r.h.s. of (10.59), we need to show that

$$\beta_0 \mathbf{m}_0 \mathbf{m}_0^{\mathrm{T}} + N_k \overline{\mathbf{x}}_k \overline{\mathbf{x}}_k^{\mathrm{T}} - \beta_k \mathbf{m}_k \mathbf{m}_k^{\mathrm{T}}$$

equals the last term on the r.h.s. of (10.62). Using (10.60), (10.61) and (276), we get

$$\beta_{0}\mathbf{m}_{0}\mathbf{m}_{0}^{\mathrm{T}} + N_{k}\overline{\mathbf{x}}_{k}\overline{\mathbf{x}}_{k}^{\mathrm{T}} - \beta_{k}\mathbf{m}_{k}\mathbf{m}_{k}^{\mathrm{T}}$$

$$= \beta_{0}\mathbf{m}_{0}\mathbf{m}_{0}^{\mathrm{T}} + N_{k}\overline{\mathbf{x}}_{k}\overline{\mathbf{x}}_{k}^{\mathrm{T}} - \beta_{0}\mathbf{m}_{0}\mathbf{m}_{k}^{\mathrm{T}} - N_{k}\overline{\mathbf{x}}_{k}\mathbf{m}_{k}^{\mathrm{T}}$$

$$= \beta_{0}\mathbf{m}_{0}\mathbf{m}_{0}^{\mathrm{T}} - \beta_{0}\mathbf{m}_{0}\mathbf{m}_{k}^{\mathrm{T}} + N_{k}\overline{\mathbf{x}}_{k}\overline{\mathbf{x}}_{k}^{\mathrm{T}} - N_{k}\overline{\mathbf{x}}_{k}\mathbf{m}_{k}^{\mathrm{T}} + \beta_{k}\mathbf{m}_{k}\mathbf{m}_{k}^{\mathrm{T}} - \beta_{k}\mathbf{m}_{k}\mathbf{m}_{k}^{\mathrm{T}}$$

$$= \beta_{0}\mathbf{m}_{0}\mathbf{m}_{0}^{\mathrm{T}} - \beta_{0}\mathbf{m}_{0}\mathbf{m}_{k}^{\mathrm{T}} - \beta_{0}\mathbf{m}_{k}\mathbf{m}_{0}^{\mathrm{T}} + \beta_{0}\mathbf{m}_{k}\mathbf{m}_{k}^{\mathrm{T}}$$

$$+ N_{k}\overline{\mathbf{x}}_{k}\overline{\mathbf{x}}_{k}^{\mathrm{T}} - N_{k}\overline{\mathbf{x}}_{k}\mathbf{m}_{k}^{\mathrm{T}} - N_{k}\mathbf{m}_{k}\overline{\mathbf{x}}_{k}^{\mathrm{T}} + N_{k}\mathbf{m}_{k}\mathbf{m}_{k}^{\mathrm{T}}$$

$$= \beta_{0}(\mathbf{m}_{k} - \mathbf{m}_{0})(\mathbf{m}_{k} - \mathbf{m}_{0})^{\mathrm{T}} + N_{k}(\overline{\mathbf{x}}_{k} - \mathbf{m}_{k})(\overline{\mathbf{x}}_{k} - \mathbf{m}_{k})^{\mathrm{T}}$$

$$= \frac{\beta_{0}N_{k}}{\beta_{0} + N_{k}}(\overline{\mathbf{x}}_{k} - \mathbf{m}_{0})(\overline{\mathbf{x}}_{k} - \mathbf{m}_{0})^{\mathrm{T}}.$$

Thus we have recovered  $\ln \mathcal{W}(\mathbf{\Lambda}_k | \mathbf{W}_0, \nu_0)$  (missing terms are again accounted for by  $f(v_N, \chi_N)$ ) and thereby (10.59).

**10.29 NOTE**: In the 1<sup>st</sup> printing of PRML, the use of  $\lambda$  to denote the varitional parameter leads to inconsistencies w.r.t. exisiting literature. To remedy this  $\lambda$  should be replaced by  $\eta$  from the beginning of Section 10.5 up to and including the last line before equation (10.141). For further details, please consult the PRML Errata.

Standard rules of differentiation give

$$\frac{d\ln(x)}{dx} = \frac{1}{x}$$
$$\frac{d^2\ln(x)}{dx^2} = -\frac{1}{x^2}.$$

Since its second derivative is negative for all value of x,  $\ln(x)$  is concave for  $0 < x < \infty$ .

From (10.133) we have

$$g(\eta) = \min_{x} \{ \eta x - f(x) \}$$
$$= \min_{x} \{ \eta x - \ln(x) \}.$$

We can minimize this w.r.t. x by setting the corresponding derivative to zero and solving for x:

$$\frac{dg}{dx} = \eta - \frac{1}{x} = 0 \quad \Longrightarrow \quad x = \frac{1}{\eta}.$$

Substituting this in (10.133), we see that

$$g(\eta) = 1 - \ln\left(\frac{1}{\eta}\right).$$

If we substitute this into (10.132), we get

$$f(x) = \min_{\eta} \left\{ \eta x - 1 + \ln \left( \frac{1}{\eta} \right) \right\}.$$

Again, we can minimize this w.r.t.  $\eta$  by setting the corresponding derivative to zero and solving for  $\eta$ :

$$\frac{df}{d\eta} = x - \frac{1}{\eta} = 0 \quad \Longrightarrow \quad \eta = \frac{1}{x},$$

and substituting this into (10.132), we find that

$$f(x) = \frac{1}{x}x - 1 + \ln\left(\frac{1}{1/x}\right) = \ln(x).$$

**10.30 NOTE**: Please consult note preceding Solution 10.29 for relevant corrections.

Differentiating the log logistic function, we get

$$\frac{d}{dx}\ln\sigma = (1 + e^{-x})^{-1}e^{-x} = \sigma(x)e^{-x}$$
 (291)

and, using (4.88),

$$\frac{d^2}{dx^2} \ln \sigma = \sigma(x) (1 - \sigma(x)) e^{-x} - \sigma(x) e^{-x} = -\sigma(x)^2 e^{-x}$$

which will always be negative and hence  $\ln \sigma(x)$  is concave.

From (291), we see that the first order Taylor expansion of  $\ln \sigma(x)$  around  $\xi$  becomes

$$\ln \sigma(x) = \ln \sigma(\xi) + (x - \xi)\sigma(\xi)e^{-\xi} + O\left((x - xi)^2\right).$$

Since  $\ln \sigma(x)$  is concave, its tangent line will be an upper bound and hence

$$\ln \sigma(x) \leqslant \ln \sigma(\xi) + (x - \xi)\sigma(\xi)e^{-\xi}.$$
 (292)

Following the presentation in Section 10.5, we define

$$\eta = \sigma(\xi)e^{-\xi}. (293)$$

Using (4.60), we have

$$\eta = \sigma(\xi)e^{-\xi} = \frac{e^{-\xi}}{1 + e^{-\xi}}$$
$$= \frac{1}{1 + e^{\xi}} = \sigma(-\xi)$$
$$= 1 - \sigma(\xi)$$

and hence

$$\sigma(\xi) = 1 - \eta.$$

From this and (293)

## 欢迎关注公众号@枫器学织与算法之道

Using these results in (292), we have

$$\ln \sigma(x) \leq \ln(1-\eta) + x\eta - \eta \left[\ln(1-\eta) - \ln \eta\right].$$

By exponentiating both sides and making use of (10.135), we obtain (10.137).

**10.31 NOTE**: Please consult note preceding Solution 10.29 for relevant corrections. Taking the derivative of f(x) w.r.t. x we get

$$\frac{\mathrm{d}f}{\mathrm{d}x} = -\frac{1}{e^{x/2} + e^{-x/2}} \frac{1}{2} \left( e^{x/2} - e^{-x/2} \right) = -\frac{1}{2} \tanh\left(\frac{x}{2}\right)$$

where we have used (5.59). From (5.60), we get

$$f''(x) = \frac{d^2f}{dx^2} = -\frac{1}{4}\left(1 - \tanh\left(\frac{x}{2}\right)^2\right).$$

Since  $\tanh(x/2)^2 < 1$  for finite values of x, f''(x) will always be negative and so f(x) is concave.

Next we define  $y=x^2$ , noting that y will always be non-negative, and express f as a function of y:

$$f(y) = -\ln\left\{\exp\left(\frac{\sqrt{y}}{2}\right) + \exp\left(-\frac{\sqrt{y}}{2}\right)\right\}.$$

We then differentiate f w.r.t. y, yielding

$$\frac{\mathrm{d}f}{\mathrm{d}y} = -\left\{ \exp\left(\frac{\sqrt{y}}{2}\right) + \exp\left(-\frac{\sqrt{y}}{2}\right) \right\}^{-1}$$

$$\frac{1}{4\sqrt{y}} \left\{ \exp\left(\frac{\sqrt{y}}{2}\right) - \exp\left(-\frac{\sqrt{y}}{2}\right) \right\}$$

$$= -\frac{1}{4\sqrt{y}} \tanh\left(\frac{\sqrt{y}}{2}\right).$$
(294)

and, using (5.60),

$$\frac{\mathrm{d}^2 f}{\mathrm{d}y^2} = \frac{1}{8y^{3/2}} \tanh\left(\frac{\sqrt{y}}{2}\right) - \frac{1}{16y} \left\{ 1 - \tanh\left(\frac{\sqrt{y}}{2}\right)^2 \right\}$$

$$= \frac{1}{8y} \left( \tanh\left(\frac{\sqrt{y}}{2}\right) \left\{ \frac{1}{\sqrt{y}} + \frac{1}{2} \tanh\left(\frac{\sqrt{y}}{2}\right) \right\} - \frac{1}{2} \right). \quad (296)$$

We see that this will be positive if the factor inside the outermost parenthesis is positive, which is equivalent to

$$\frac{1}{\sqrt{y}}\tanh\left(\frac{\sqrt{y}}{2}\right) > \frac{1}{2}\left\{1-\tanh^2\left(\frac{\sqrt{y}}{2}\right)\right\}.$$

If we divide both sides by  $\tanh (\sqrt{y}/2)$ , substitute a for  $\sqrt{y}/2$  and then make use of (5.59), we can write this as

$$\frac{1}{a} > \frac{e^{a} + e^{-a}}{e^{a} - e^{-a}} - \frac{e^{a} - e^{-a}}{e^{a} + e^{-a}}$$

$$= \frac{(e^{a} + e^{-a})^{2} - (e^{a} - e^{-a})^{2}}{(e^{a} - e^{-a})(e^{a} + e^{-a})}$$

$$= \frac{4}{e^{2a} - e^{-2a}}.$$

Taking the inverse of both sides of this inequality we get

$$a < \frac{1}{4} \left( e^{2a} - e^{-2a} \right).$$

If differentiate both sides w.r.t. a we see that the derivatives are equal at a=0 and for a>0, the derivative of the r.h.s. will be greater than that of the l.h.s. Thus, the r.h.s. will grow faster and the inequality will hold for a>0. Consequently (296) will be positive for y>0 and approach  $+\infty$  as y approaches 0.

Now we use (295) to make a Taylor expansion of  $f(x^2)$  around  $\xi^2$ , which gives

$$\begin{split} f(x^2) &= f(\xi^2) + (x^2 - \xi^2) f'(\xi^2) + O\left((x^2 - \xi^2)^2\right) \\ &\geqslant -\ln\left\{\exp\left(\frac{\xi}{2}\right) + \exp\left(-\frac{\xi}{2}\right)\right\} - (x^2 - \xi^2) \frac{1}{4\xi} \tanh\left(\frac{\xi}{2}\right). \end{split}$$

where we have used the fact that f is convex function of  $x^2$  and hence its tangent will be a lower bound. Defining

$$\lambda(\xi) = \frac{1}{4\xi} \tanh\left(\frac{\xi}{2}\right)$$

we recover (10.143), from which (10.144) follows.

10.32 We can see this from the lower bound (10.154), which is simply a sum of the prior and indepedent contributions from the data points, all of which are quadratic in w. A new data point would simply add another term to this sum and we can regard terms from the previously arrived data points and the original prior collectively as a revised prior, which should be combined with the contributions from the new data point.

The corresponding sufficient statistics, (10.157) and (10.158), can be rewritten di-

rectly in the corresponding sequential form,

$$\begin{split} \mathbf{m}_N &= \mathbf{S}_N \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \sum_{n=1}^N (t_n - 1/2) \phi_n \right) \\ &= \mathbf{S}_N \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \sum_{n=1}^{N-1} (t_n - 1/2) \phi_n + (t_N - 1/2) \phi_N \right) \\ &= \mathbf{S}_N \left( \mathbf{S}_{N-1}^{-1} \mathbf{S}_{N-1} \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \sum_{n=1}^{N-1} (t_n - 1/2) \phi_n \right) + (t_N - 1/2) \phi_N \right) \\ &= \mathbf{S}_N \left( \mathbf{S}_{N-1}^{-1} \mathbf{m}_{N-1} + (t_N - 1/2) \phi_N \right) \end{split}$$

and

$$\mathbf{S}_{N}^{-1} = \mathbf{S}_{0}^{-1} + 2\sum_{n=1}^{N} \lambda(\xi_{n}) \phi_{n} \phi_{n}^{\mathrm{T}}$$

$$= \mathbf{S}_{0}^{-1} + 2\sum_{n=1}^{N-1} \lambda(\xi_{n}) \phi_{n} \phi_{n}^{\mathrm{T}} + 2\lambda(\xi_{N}) \phi_{N} \phi_{N}^{\mathrm{T}}$$

$$= \mathbf{S}_{N-1}^{-1} + 2\lambda(\xi_{N}) \phi_{N} \phi_{N}^{\mathrm{T}}.$$

The update formula for the variational parameters, (10.163), remain the same, but each parameter is updated only once, although this update will be part of an iterative scheme, alternating between updating  $\mathbf{m}_N$  and  $\mathbf{S}_N$  with  $\xi_N$  kept fixed, and updating  $\xi_N$  with  $\mathbf{m}_N$  and  $\mathbf{S}_N$  kept fixed. Note that updating  $\xi_N$  will not affect  $\mathbf{m}_{N-1}$  and  $\mathbf{S}_{N-1}$ . Note also that this updating policy differs from that of the batch learning scheme, where all variational parameters are updated using statistics based on all data points.

**10.33** Taking the derivative of (10.161) w.r.t.  $\xi_n$ , we get

$$\frac{\partial Q}{\partial \xi_n} = \frac{1}{\sigma(\xi_n)} \sigma'(\xi_n) - \frac{1}{2} - \lambda'(\xi_n) \left( \phi_n^{\mathrm{T}} \mathbb{E} \left[ \mathbf{w} \mathbf{w}^{\mathrm{T}} \right] \phi - \xi_n^2 \right) + \lambda(\xi_n) 2 \xi_n 
= \frac{1}{\sigma(\xi_n)} \sigma(\xi_n) (1 - \sigma(x)) - \frac{1}{2} - \lambda'(\xi_n) \left( \phi_n^{\mathrm{T}} \mathbb{E} \left[ \mathbf{w} \mathbf{w}^{\mathrm{T}} \right] \phi - \xi_n^2 \right) 
+ \frac{1}{2\xi_n} \left[ \sigma(x) - \frac{1}{2} \right] 2 \xi_n 
= -\lambda'(\xi_n) \left( \phi_n^{\mathrm{T}} \mathbb{E} \left[ \mathbf{w} \mathbf{w}^{\mathrm{T}} \right] \phi - \xi_n^2 \right)$$

where we have used (4.88) and (10.141). Setting this equal to zero, we obtain (10.162), from which (10.163) follows.

**10.34 NOTE**: In the 1<sup>st</sup> printing of PRML, there are a number of sign errors in Equation (10.164); the correct form is

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

We can differentiate  $\mathcal{L}$  w.r.t.  $\xi_n$  using (3.117) and results from Solution 10.33, to obtain

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = \frac{1}{2} \text{Tr} \left( \mathbf{S}_N^{-1} \frac{\partial \mathbf{S}_N}{\partial \xi_n} \right) + \frac{1}{2} \text{Tr} \left( \mathbf{a}_N \mathbf{a}_N^{\mathrm{T}} \frac{\partial \mathbf{S}_N}{\partial \xi_n} \right) + \lambda'(\xi_n) \xi_n^2$$
(297)

where we have defined

$$\mathbf{a}_N = \mathbf{S}_N^{-1} \mathbf{m}_N. \tag{298}$$

From (10.158) and (C.21), we get

$$\frac{\partial \mathbf{S}_{N}}{\partial \xi_{n}} = \frac{\partial \left(\mathbf{S}_{N}^{-1}\right)^{-1}}{\partial \xi_{n}} = -\mathbf{S}_{N} \frac{\partial \mathbf{S}_{N}^{-1}}{\partial \xi_{n}} \mathbf{S}_{N}$$
$$= -\mathbf{S}_{N} 2\lambda'(\xi_{n}) \phi_{n} \phi_{n}^{\mathrm{T}} \mathbf{S}_{N}.$$

Substituting this into (297) and setting the result equal to zero, we get

$$-\frac{1}{2}\operatorname{Tr}\left(\left(\mathbf{S}_{N}^{-1}+\mathbf{a}_{N}\mathbf{a}_{N}^{\mathrm{T}}\right)\mathbf{S}_{N}2\lambda'(\xi_{n})\phi_{n}\phi_{n}^{\mathrm{T}}\mathbf{S}_{N}\right)+\lambda'(\xi_{n})\xi_{n}^{2}=0.$$

Rearranging this and making use of (298) we get

$$\xi_n^2 = \phi_n^{\mathrm{T}} \mathbf{S}_N \left( \mathbf{S}_N^{-1} + \mathbf{a}_N \mathbf{a}_N^{\mathrm{T}} \right) \mathbf{S}_N \phi_n$$
$$= \phi_n^{\mathrm{T}} \left( \mathbf{S}_N + \mathbf{m}_N \mathbf{m}_N^{\mathrm{T}} \right) \phi_n$$

where we have also used the symmetry of  $S_N$ .

**10.35 NOTE**: See note in Solution 10.34.

From (2.43), (4.140) and (10.153), we see that

$$p(\mathbf{w})h(\mathbf{w}, \boldsymbol{\xi}) = (2\pi)^{-W/2} |\mathbf{S}_0|^{-1/2}$$

$$\exp \left\{ -\frac{1}{2} \mathbf{w}^{\mathrm{T}} \left( \mathbf{S}_0^{-1} + 2 \sum_{n=1}^N \lambda(\xi_n) \phi_n \phi_n^{\mathrm{T}} \right) \mathbf{w} + \mathbf{w}^{\mathrm{T}} \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \sum_{n=1}^N \phi_n \left[ t_n - \frac{1}{2} \right] \right) \right\}$$

$$\exp \left\{ -\frac{1}{2} \mathbf{m}_0^{\mathrm{T}} \mathbf{S}_0^{-1} \mathbf{m}_0 + \sum_{n=1}^N \frac{\xi_n}{2} + \lambda(\xi_n) \xi_n^2 \right\} \prod_{n=1}^N \sigma(\xi_n).$$

Using (10.157) and (10.158), we can complete the square over w, yielding

$$p(\mathbf{w})h(\mathbf{w}, \boldsymbol{\xi}) = (2\pi)^{-W/2} |\mathbf{S}_0|^{-1/2} \prod_{n=1}^N \sigma(\xi_n)$$

$$\exp\left\{-\frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^{\mathrm{T}} \mathbf{S}_N^{-1} (\mathbf{w} - \mathbf{m}_N)\right\}$$

$$\exp\left\{\frac{1}{2} \mathbf{m}_N^{\mathrm{T}} \mathbf{S}_N^{-1} \mathbf{m}_N - \frac{1}{2} \mathbf{m}_0^{\mathrm{T}} \mathbf{S}_0^{-1} \mathbf{m}_0 \sum_{n=1}^N \frac{\xi_n}{2} + \lambda(\xi_n) \xi_n^2\right\}.$$

Now we can do the integral over  ${\bf w}$  in (10.159), in effect replacing the first exponential factor with  $(2\pi)^{W/2} |{\bf S}_N|^{1/2}$ . Taking logarithm, we then obtain (10.164).

**10.36** If we denote the joint distribution corresponding to the first j factors by  $p_j(\theta, \mathcal{D})$ , with corresponding evidence  $p_j(\mathcal{D})$ , then we have

$$p_{j}(\mathcal{D}) = \int p_{j}(\boldsymbol{\theta}, \mathcal{D}) d\boldsymbol{\theta} = \int p_{j-1}(\boldsymbol{\theta}, \mathcal{D}) f_{j}(\boldsymbol{\theta}) d\boldsymbol{\theta}$$
$$= p_{j-1}(\mathcal{D}) \int p_{j-1}(\boldsymbol{\theta}|\mathcal{D}) f_{j}(\boldsymbol{\theta}) d\boldsymbol{\theta}$$
$$\simeq p_{j-1}(\mathcal{D}) \int q_{j-1}(\boldsymbol{\theta}) f_{j}(\boldsymbol{\theta}) d\boldsymbol{\theta} = p_{j-1}(\mathcal{D}) Z_{j}.$$

By applying this result recursively we see that the evidence is given by the product of the normalization constants

$$p(\mathcal{D}) = \prod_{j} Z_{j}.$$

**10.37** Here we use the general expectation-propagation equations (10.204)–(10.207). The initial  $q(\theta)$  takes the form

$$q_{\text{init}}(\boldsymbol{\theta}) = \widetilde{f}_0(\boldsymbol{\theta}) \prod_{i \neq 0} \widetilde{f}_i(\boldsymbol{\theta})$$

where  $\widetilde{f}_0(\boldsymbol{\theta}) = f_0(\boldsymbol{\theta})$ . Thus

$$q^{\setminus 0}(\boldsymbol{\theta}) \propto \prod_{i \neq 0} \widetilde{f}_i(\boldsymbol{\theta})$$

and  $q^{\text{new}}(\theta)$  is determined by matching moments (sufficient statistics) against

$$q^{\setminus 0}(\boldsymbol{\theta}) f_0(\boldsymbol{\theta}) = q_{\text{init}}(\boldsymbol{\theta}).$$

Since by definition this belongs to the same exponential family form as  $q^{\text{new}}(\theta)$  it follows that

$$q^{\mathrm{new}}(\boldsymbol{\theta}) = q_{\mathrm{init}}(\boldsymbol{\theta}) = q^{\setminus 0}(\boldsymbol{\theta}) f_0(\boldsymbol{\theta}).$$

Thus

$$\widetilde{f}_0(oldsymbol{ heta}) = rac{Z_0 \, q^{
m new}(oldsymbol{ heta})}{q^{ackslash 0}(oldsymbol{ heta})} = Z_0 f_0(oldsymbol{ heta})$$

where

$$Z_0 = \int q^{\setminus 0}(\boldsymbol{\theta}) f_0(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int q^{\text{new}}(\boldsymbol{\theta}) d\boldsymbol{\theta} = 1.$$

**10.38** The ratio is given by

$$q^{\setminus n}(\boldsymbol{\theta}) \propto \exp\left\{-\frac{1}{2v}\|\boldsymbol{\theta} - \mathbf{m}\|^2 + \frac{1}{2v_n}\|\boldsymbol{\theta} - \mathbf{m}_n\|^2\right\}$$
$$\propto \exp\left\{-\frac{1}{2}\boldsymbol{\theta}^{\mathrm{T}}\boldsymbol{\theta}\left(\frac{1}{v} - \frac{1}{v_n}\right) + \boldsymbol{\theta}^{\mathrm{T}}\left(\frac{1}{v}\mathbf{m} - \frac{1}{v_n}\mathbf{m}_n\right)\right\}$$

from which we obtain the variance given by (10.215). The mean is then obtained by completing the square and is therefore given by

$$\mathbf{m}^{\backslash n} = v^{\backslash n} \left( v^{-1} \mathbf{m} - v_n^{-1} \mathbf{m}_n \right)$$

$$= v^{\backslash n} \left( v^{-1} \mathbf{m} - v_n^{-1} \mathbf{m}_n \right) + v^{\backslash n} v_n^{-1} \mathbf{m} - v^{\backslash n} v_n^{-1} \mathbf{m}$$

$$= v^{\backslash n} \left( v^{-1} - v_n^{-1} \right) \mathbf{m} + v^{\backslash n} v_n^{-1} \left( \mathbf{m} - \mathbf{m}_n \right)$$

Hence we obtain (10.214).

The normalization constant is given by

$$Z_n = \int \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}^{n}, v^{n}\mathbf{I}) \left\{ (1 - w)\mathcal{N}(\mathbf{x}_n|\boldsymbol{\theta}, \mathbf{I}) + w\mathcal{N}(\mathbf{x}_n|\mathbf{0}, a\mathbf{I}) \right\} d\boldsymbol{\theta}.$$

The first term can be integrated by using the result (2.115) while the second term is trivial since the background distribution does not depend on  $\theta$  and hence can be taken outside the integration. We therefore obtain (10.216).

**10.39 NOTE**: In PRML, a term  $v^{n}D$  should be added to the r.h.s. of (10.245). We derive (10.244) by noting

$$\nabla_{\mathbf{m}^{\backslash n}} \ln Z_n = \frac{1}{Z_n} \nabla_{\mathbf{m}^{\backslash n}} \int q^{\backslash n}(\boldsymbol{\theta}) f_n(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

$$= \frac{1}{Z_n} \int q^{\backslash n}(\boldsymbol{\theta}) f_n(\boldsymbol{\theta}) \left\{ -\frac{1}{v^{\backslash n}} (\mathbf{m}^{\backslash n} - \boldsymbol{\theta}) \right\} d\boldsymbol{\theta}$$

$$= -\frac{\mathbf{m}^{\backslash n}}{v^{\backslash n}} + \frac{\mathbb{E}[\boldsymbol{\theta}]}{v^{\backslash n}}.$$

We now use this to derive (10.217) by substituting for  $Z_n$  using (10.216) to give

$$\nabla_{\mathbf{m}^{\backslash n}} \ln Z_n = \frac{1}{Z_n} (1 - w) \mathcal{N}(\mathbf{x}_n | \mathbf{m}^{\backslash n}, (v^{\backslash n} + 1) \mathbf{I}) \frac{1}{v^{\backslash n} + 1} (\mathbf{x}_n - \mathbf{m}^{\backslash n})$$
$$= \rho_n \frac{1}{v^{\backslash n} + 1} (\mathbf{x}_n - \mathbf{m}^{\backslash n})$$

where we have defined

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

Similarly for (10.245) we have

$$\nabla_{v^{\setminus n}} \ln Z_n = \frac{1}{Z_n} \nabla_{v^{\setminus n}} \int q^{\setminus n}(\boldsymbol{\theta}) f_n(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

$$= \frac{1}{Z_n} \int q^{\setminus n}(\boldsymbol{\theta}) f_n(\boldsymbol{\theta}) \left\{ \frac{1}{2(v^{\setminus n})^2} \left( \mathbf{m}^{\setminus n} - \boldsymbol{\theta} \right)^{\mathrm{T}} \left( \mathbf{m}^{\setminus n} - \boldsymbol{\theta} \right) - \frac{D}{2v^{\setminus n}} \right\} d\boldsymbol{\theta}$$

$$= \frac{1}{2(v^{\setminus n})^2} \left\{ \mathbb{E}[\boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\theta}] - 2\mathbb{E}[\boldsymbol{\theta}^{\mathrm{T}}] \mathbf{m}^{\setminus n} + \|\mathbf{m}^{\setminus n}\|^2 \right\} - \frac{D}{2v^{\setminus n}}.$$

Re-arranging we obtain (10.245). Now we substitute for  $Z_n$  using (10.216) to give

$$\nabla_{v^{\setminus n}} \ln Z_n = \frac{1}{Z_n} (1 - w) \mathcal{N}(\mathbf{x}_n | \mathbf{m}^{\setminus n}, (v^{\setminus n} + 1) \mathbf{I})$$
$$\left[ \frac{1}{2(v^{\setminus n} + 1)^2} ||\mathbf{x}_n - \mathbf{m}^{\setminus n}||^2 - \frac{D}{2(v^{\setminus n} + 1)} \right].$$

Next we note that the variance is given by

$$v\mathbf{I} = \mathbb{E}[\boldsymbol{\theta}\boldsymbol{\theta}^{\mathrm{T}}] - \mathbb{E}[\boldsymbol{\theta}]\mathbb{E}[\boldsymbol{\theta}^{\mathrm{T}}]$$

and so, taking the trace, we have

$$Dv = \mathbb{E}[\boldsymbol{\theta}^{\mathrm{T}}\boldsymbol{\theta}] - \mathbb{E}[\boldsymbol{\theta}^{\mathrm{T}}]\mathbb{E}[\boldsymbol{\theta}]$$

where D is the dimensionality of  $\theta$ . Combining the above results we obtain (10.218).

## **Chapter 11 Sampling Methods**

**11.1** Since the samples are independent, for the mean, we have

$$\mathbb{E}\left[\hat{f}\right] = \frac{1}{L} \sum_{l=1}^{L} \int f(z^{(l)}) p(z^{(l)}) \, dz^{(l)} = \frac{1}{L} \sum_{l=1}^{L} \mathbb{E}\left[f\right] = \mathbb{E}\left[f\right].$$

Using this together with (1.38) and (1.39), for the variance, we have

$$\operatorname{var}\left[\widehat{f}\right] = \mathbb{E}\left[\left(\widehat{f} - \mathbb{E}\left[\widehat{f}\right]\right)^{2}\right]$$
$$= \mathbb{E}\left[\widehat{f}^{2}\right] - \mathbb{E}\left[f\right]^{2}.$$

$$\mathbb{E}\left[f(z^{(k)}), f(z^{(m)})\right] = \begin{cases} \operatorname{var}[f] + \mathbb{E}[f^2] & \text{if } n = k, \\ \mathbb{E}[f^2] & \text{otherwise,} \end{cases}$$
$$= \mathbb{E}[f^2] + \delta_{mk} \operatorname{var}[f],$$

where we again exploited the fact that the samples are independent.

Hence

$$\operatorname{var}\left[\widehat{f}\right] = \mathbb{E}\left[\frac{1}{L}\sum_{m=1}^{L}f(z^{(m)})\frac{1}{L}\sum_{k=1}^{L}f(z^{(k)})\right] - \mathbb{E}[f]^{2}$$

$$= \frac{1}{L^{2}}\sum_{m=1}^{L}\sum_{k=1}^{L}\left\{\mathbb{E}[f^{2}] + \delta_{mk}\operatorname{var}[f]\right\} - \mathbb{E}[f]^{2}$$

$$= \frac{1}{L}\operatorname{var}[f]$$

$$= \frac{1}{L}\mathbb{E}\left[\left(f - \mathbb{E}[f]\right)^{2}\right].$$

**11.2** From (1.27) we have,

$$p_y(y) = p_z(h(y)) |h'(y)|.$$

Differentiating (11.6) w.r.t. y and using the fact that  $p_z(h(y)) = 1$ , we see that

$$p_y(y) = p(y).$$

**11.3** Using the standard integral

$$\int \frac{1}{a^2 + u^2} du = \frac{1}{a} \tan^{-1} \left( \frac{u}{a} \right) + C$$

where C is a constant, we can integrate the r.h.s. of (11.8) to obtain

$$z = h(y) = \int_{-\infty}^{y} p(\hat{y}) d\hat{y} = \frac{1}{\pi} \tan^{-1}(y) + \frac{1}{2}$$

where we have chosen the constant C=1/2 to ensure that the range of the cumulative distribution function is [0,1].

Thus the required transformation function becomes

$$y = h^{-1}(z) = \tan\left(\pi\left(z - \frac{1}{2}\right)\right).$$

**11.4** We need to calculate the determinant of the Jacobian

$$\frac{\partial(z_1,z_2)}{\partial(y_1,y_2)}.$$

In doing so, we will find it helpful to make use of intermediary variables in polar coordinates

$$\theta = \tan^{-1} \frac{z_2}{z_1} \tag{299}$$

$$r^2 = z_1^2 + z_2^2 (300)$$

from which it follows that

$$z_1 = r\cos\theta \tag{301}$$

$$z_2 = r \sin \theta. \tag{302}$$

From (301) and (302) we have

$$\frac{\partial(z_1, z_2)}{\partial(r, \theta)} = \begin{pmatrix} \cos \theta & \sin \theta \\ -r \sin \theta & r \cos \theta \end{pmatrix}$$

and thus

$$\left| \frac{\partial(z_1, z_2)}{\partial(r, \theta)} \right| = r(\cos^2 \theta + \sin^2 \theta) = r.$$
 (303)

From (11.10), (11.11) and (300)–(302) we have

$$y_1 = z_1 \left(\frac{-2\ln r^2}{r^2}\right)^{1/2} = \left(-2\ln r^2\right)^{1/2} \cos\theta$$
 (304)

$$y_2 = z_2 \left(\frac{-2\ln r^2}{r^2}\right)^{1/2} = \left(-2\ln r^2\right)^{1/2} \sin\theta$$
 (305)

which give

$$\frac{\partial(y_1, y_2)}{\partial(r, \theta)} = \begin{pmatrix} -2\cos\theta \left(-2\ln r^2\right)^{-1/2} r^{-1} & -2\sin\theta \left(-2\ln r^2\right)^{-1/2} r^{-1} \\ -\sin\theta \left(-2\ln r^2\right)^{1/2} & \cos\theta \left(-2\ln r^2\right)^{1/2} \end{pmatrix}$$

and thus

$$\left| \frac{\partial(r,\theta)}{\partial(y_1,y_2)} \right| = \left| \frac{\partial(y_1,y_2)}{\partial(r,\theta)} \right|^{-1} = \left( -2r^{-1}(\cos^2\theta + \sin^2\theta) \right)^{-1} = -\frac{r}{2}.$$

Combining this with (303), we get

$$\left| \frac{\partial(z_1, z_2)}{\partial(y_1, y_2)} \right| = \left| \frac{\partial(z_1, z_2)}{\partial(r, \theta)} \frac{\partial(r, \theta)}{\partial(y_1, y_2)} \right| 
= \left| \frac{\partial(z_1, z_2)}{\partial(r, \theta)} \right| \left| \frac{\partial(r, \theta)}{\partial(y_1, y_2)} \right| = -\frac{r^2}{2}$$
(306)

However, we only retain the absolute value of this, since both sides of (11.12) must be non-negative. Combining this with

$$p(z_1, z_2) = \frac{1}{\pi}$$

which follows from the fact that  $z_1$  and  $z_2$  are uniformly distributed on the unit circle, we can rewrite (11.12) as

$$p(y_1, y_2) = \frac{1}{2\pi}r^2. (307)$$

By squaring the left- and rightmost sides of (304) and (305), adding up the results and rearranging, we see that

$$r^2 = \exp\left(-\frac{1}{2}\left(y_1^2 + y_2^2\right)\right)$$

which toghether with (307) give (11.12).

**11.5** Since  $\mathbb{E}[\mathbf{z}] = \mathbf{0}$ ,

$$\mathbb{E}\left[\mathbf{y}\right] = \mathbb{E}\left[\boldsymbol{\mu} + \mathbf{L}\mathbf{z}\right] = \boldsymbol{\mu}.$$

Similarly, since  $\mathbb{E}\left[\mathbf{z}\mathbf{z}^{\mathrm{T}}\right] = \mathbf{I}$ ,

$$\begin{aligned} \operatorname{cov}\left[\mathbf{y}\right] &= & \mathbb{E}\left[\mathbf{y}\mathbf{y}^{\mathrm{T}}\right] - \mathbb{E}\left[\mathbf{y}\right]\mathbb{E}\left[\mathbf{y}^{\mathrm{T}}\right] \\ &= & \mathbb{E}\left[\left(\boldsymbol{\mu} + \mathbf{L}\mathbf{z}\right)\left(\boldsymbol{\mu} + \mathbf{L}\mathbf{z}\right)^{\mathrm{T}}\right] - \boldsymbol{\mu}\boldsymbol{\mu}^{\mathrm{T}} \\ &= & \mathbf{L}\mathbf{L}^{\mathrm{T}} \\ &= & \boldsymbol{\Sigma}. \end{aligned}$$

**11.6** The probability of acceptance follows directly from the mechanism used to accept or reject the sample. The probability of a sample  $\mathbf{z}$  being accepted equals the probability of a sample u, drawn uniformly from the interval  $[0, kq(\mathbf{z})]$ , being less than or equal to a value  $\widetilde{p}(\mathbf{z}) \leq kq(\mathbf{z})$ , and is given by is given by

$$p(\text{acceptance}|\mathbf{z}) = \int_0^{\widetilde{p}(\mathbf{z})} \frac{1}{kq(\mathbf{z})} du = \frac{\widetilde{p}(\mathbf{z})}{kq(\mathbf{z})}.$$

Therefore, the probability of drawing a sample, z, is

$$q(\mathbf{z})p(\text{acceptance}|\mathbf{z}) = q(\mathbf{z})\frac{\widetilde{p}(\mathbf{z})}{kq(\mathbf{z})} = \frac{\widetilde{p}(\mathbf{z})}{k}.$$
 (308)

Integrating both sides w.r.t. **z**, we see that  $kp(\text{acceptance}) = Z_p$ , where

$$Z_p = \int \widetilde{p}(\mathbf{z}) \, \mathrm{d}\mathbf{z}.$$

Combining this with (308) and (11.13), we obtain

欢迎关注公众号 $\overline{\mathscr{Q}}$  如常学  $\overline{\mathscr{Z}}$   $\overline{\mathscr{Z} }$   $\overline{\mathscr{Z}}$   $\overline{\mathscr{Z} }$   $\overline{\mathscr{Z}}$   $\overline$ 

**11.7 NOTE**: In PRML, the roles of y and z in the text of the exercise should be swapped in order to be consistent with the notation used in Section 11.1.2, including (11.16); this is opposite to the notation used in Section 11.1.1.

We will suppose that y has a uniform distribution on [0,1] and we seek the distribution of z, which is derived from

$$z = b \tan y + c$$
.

From (11.5), we have

$$q(z) = p(y) \left| \frac{\mathrm{d}y}{\mathrm{d}z} \right| \tag{309}$$

From the inverse transformation

$$y = \tan^{-1}(u(z)) = \tan^{-1}\left(\frac{z-c}{b}\right)$$

where we have implicitly defined u(z), we see that

$$\frac{\mathrm{d}y}{\mathrm{d}z} = \frac{\mathrm{d}}{\mathrm{d}u} \tan^{-1}(u) \frac{\mathrm{d}u}{\mathrm{d}z}$$
$$= \frac{1}{1+u^2} \frac{\mathrm{d}u}{\mathrm{d}z}$$
$$= \frac{1}{1+(z-c)^2/b^2} \frac{1}{b}.$$

Substituting this into (309), using the fact that p(y)=1 and finally absorbing the factor 1/b into k, we obtain (11.16).

**11.8 NOTE**: In PRML, equation (11.17) and the following end of the sentence need to modified as follows:

$$q(z) = k_i \lambda_i \exp \left\{ -\lambda_i (z - z_i) \right\}$$
  $\widehat{z}_{i-1,i} < z \leqslant \widehat{z}_{i,i+1}$ 

where  $\widehat{z}_{i-1,i}$  is the point of intersection of the tangent lines at  $z_{i-1}$  and  $z_i$ ,  $\lambda_i$  is the slope of the tangent at  $z_i$  and  $k_i$  accounts for the corresponding offset.

We start by determining  $\widetilde{q}(z)$  with coefficients  $\widetilde{k}_i$ , such that  $\widetilde{q}(z) \geqslant p(z)$  everywhere. From Figure 11.6, we see that

$$\widetilde{q}(z_i) = p(z_i)$$

and thus, from (11.17),

$$\widetilde{q}(z_i) = \widetilde{k}_i \lambda_i \exp(-\lambda_i (z_i - z_i))$$

$$= \widetilde{k}_i \lambda_i = p(z_i). \tag{310}$$

Next we compute the normalization constant for q in terms of  $k_i$ ,

$$Z_{q} = \int \widetilde{q}(z) dz$$

$$= \sum_{i=1}^{K} \widetilde{k}_{i} \lambda_{i} \int_{\widehat{z}_{i-1,i}}^{\widehat{z}_{i,i+1}} \exp(-\lambda_{i}(z-z_{i})) dz$$

$$= \sum_{i=1}^{K} k_{i}$$
(311)

where K denotes the number of grid points and

$$k_{i} = \widetilde{k}_{i} \lambda_{i} \int_{\widehat{z}_{i-1,i}}^{\widehat{z}_{i,i+1}} \exp\left(-\lambda_{i}(z-z_{i})\right) dz$$

$$= \widetilde{k}_{i} \left(\exp\left\{-\lambda_{i} \left(\widehat{z}_{i-1,i}-z_{i}\right)\right\} - \exp\left\{-\lambda_{i} \left(\widehat{z}_{i,i+1}-z_{i}\right)\right\}\right). \quad (312)$$

Note that  $\widehat{z}_{0,1}$  and  $\widehat{z}_{K,K+1}$  equal the lower and upper limits on z, respectively, or  $-\infty/+\infty$  where no such limits exist.

## **11.9 NOTE**: See correction detailed in Solution 11.8

To generate a sample from q(z), we first determine the segment of the envelope function from which the sample will be generated. The probability mass in segment i is given from (311) as  $k_i/Z_q$ . Hence we draw v from U(v|0,1) and obtain

$$i = \begin{cases} 1 & \text{if } v \leqslant k_1/Z_q \\ m & \text{if } \sum_{j=1}^{m-1} k_j/Z_q < v \leqslant \sum_{j=1}^m k_j/Z_q, \quad 1 < m < K \\ K & \text{otherwise.} \end{cases}$$

Next, we can use the techniques of Section 11.1.1 to sample from the exponential distribution corresponding to the chosen segment i. We must, however, take into account that we now only want to sample from a finite interval of the exponential distribution and so the lower limit in (11.6) will be  $\widehat{z}_{i-1,i}$ . If we consider the uniformly distributed variable w, U(w|0,1), (11.6), (310) and (311) give

$$w = h(z) = \int_{\widehat{z}_{i-1,i}}^{z} q(\widetilde{z}) d\widetilde{z}$$

$$= \frac{\widetilde{k}_{i}}{k_{i}} \lambda_{i} \exp(\lambda_{i} z_{i}) \int_{\widehat{z}_{i-1,i}}^{z} \exp(-\lambda_{i} \widetilde{z}) d\widetilde{z}$$

$$= \frac{\widetilde{k}_{i}}{k_{i}} \exp(\lambda_{i} z_{i}) \left[ \exp(-\lambda_{i} \widehat{z}_{i-1,i}) - \exp(-\lambda_{i} z) \right].$$

Thus, by drawing a sample  $w^{\star}$  and transforming it according to

$$z^{\star} = \frac{1}{\lambda_{i}} \ln \left[ w^{\star} \frac{k_{i}}{\widetilde{k}_{i} \exp(\lambda_{i} z_{i})} - \exp(-\lambda_{i} \widehat{z}_{i-1,i}) \right]$$
$$= \frac{1}{\lambda_{i}} \ln \left[ w^{\star} \left( \exp\left\{ -\lambda_{i} \widehat{z}_{i-1,i} \right\} - \exp\left\{ -\lambda_{i} \widehat{z}_{i,i+1} \right\} \right) - \exp\left( -\lambda_{i} \widehat{z}_{i-1,i} \right) \right]$$

where we have used (312), we obtain a sample from q(z).

**11.10 NOTE**: In PRML, " $z^{(1)} = 0$ " should be " $z^{(0)} = 0$ " on the first line following Equation (11.36)

From (11.34)–(11.36) and the fact that  $\mathbb{E}\left[z^{(\tau)}\right]=0$  for all values of  $\tau$ , we have

$$\mathbb{E}_{z^{(\tau)}} \left[ \left( z^{(\tau)} \right)^{2} \right] = 0.5 \, \mathbb{E}_{z^{(\tau-1)}} \left[ \left( z^{(\tau-1)} \right)^{2} \right] + 0.25 \, \mathbb{E}_{z^{(\tau-1)}} \left[ \left( z^{(\tau-1)} + 1 \right)^{2} \right] \\
+ 0.25 \, \mathbb{E}_{z^{(\tau-1)}} \left[ \left( z^{(\tau-1)} - 1 \right)^{2} \right] \\
= \mathbb{E}_{z^{(\tau-1)}} \left[ \left( z^{(\tau-1)} \right)^{2} \right] + \frac{1}{2}. \tag{313}$$

With  $z^{(0)} = 0$  specified in the text following (11.36), (313) gives

$$\mathbb{E}_{z^{(1)}} \left[ \left( z^{(1)} \right)^2 \right] = \mathbb{E}_{z^{(0)}} \left[ \left( z^{(0)} \right)^2 \right] + \frac{1}{2} = \frac{1}{2}.$$

Assuming that

$$\mathbb{E}_{z^{(k)}}\left[\left(z^{(k)}\right)^2\right] = \frac{k}{2}$$

(313) immediatly gives

$$\mathbb{E}_{z^{(k+1)}}\left[\left(z^{(k+1)}\right)^2\right] = \frac{k}{2} + \frac{1}{2} = \frac{k+1}{2}$$

and thus

$$\mathbb{E}_{z^{(\tau)}}\left[\left(z^{(\tau)}\right)^2\right] = \frac{\tau}{2}.$$

**11.11** This follows from the fact that in Gibbs sampling, we sample a single variable,  $z_k$ , at the time, while all other variables,  $\{z_i\}_{i\neq k}$ , remain unchanged. Thus,  $\{z_i'\}_{i\neq k}=\{z_i\}_{i\neq k}$  and we get

$$p^{\star}(\mathbf{z})T(\mathbf{z}, \mathbf{z}') = p^{\star}(z_{k}, \{z_{i}\}_{i\neq k})p^{\star}(z'_{k}|\{z_{i}\}_{i\neq k})$$

$$= p^{\star}(z_{k}|\{z_{i}\}_{i\neq k})p^{\star}(\{z_{i}\}_{i\neq k})p^{\star}(z'_{k}|\{z_{i}\}_{i\neq k})$$

$$= p^{\star}(z_{k}|\{z'_{i}\}_{i\neq k})p^{\star}(\{z'_{i}\}_{i\neq k})p^{\star}(z'_{k}|\{z'_{i}\}_{i\neq k})$$

$$= p^{\star}(z_{k}|\{z'_{i}\}_{i\neq k})p^{\star}(z'_{k}, \{z'_{i}\}_{i\neq k})$$

$$= p^{\star}(\mathbf{z}')T(\mathbf{z}', \mathbf{z}),$$

where we have used the product rule together with  $T(\mathbf{z}, \mathbf{z}') = p^*(z_k' | \{z_i\}_{i \neq k})$ .

- 11.12 Gibbs sampling is *not* ergodic w.r.t. the distribution shown in Figure 11.15, since the two regions of non-zero probability do not overlap when projected onto either the  $z_1$ -ot the  $z_2$ -axis. Thus, as the initial sample will fall into one and only one of the two regions, all subsequent samples will also come from that region. However, had the initial sample fallen into the other region, the Gibbs sampler would have remained in that region. Thus, the corresponding Markov chain would have two stationary distributions, which is counter to the definition of the equilibrium distribution.
- **11.13** The joint distribution over x,  $\mu$  and  $\tau$  can be written

$$p(x, \mu, \tau | \mu_0, s_0, a, b) = \mathcal{N}(x | \mu, \tau^{-1}) \mathcal{N}(\mu | \mu_0, s_0) \operatorname{Gam}(\tau | a, b).$$

From Bayes' theorem we see that

$$p(\mu|x,\tau,\mu_0,s_0) \propto \mathcal{N}\left(x|\mu,\tau^{-1}\right) \mathcal{N}\left(\mu|\mu_0,s_0\right)$$

which, from (2.113)–(2.117), is also a Gaussian,

$$\mathcal{N}\left(\mu|\widehat{\mu},\widehat{s}\right)$$

with parameters

$$\widehat{s}^{-1} = s_0^{-1} + \tau \widehat{\mu} = \widehat{s} (\tau x + s_0^{-1} \mu_0) .$$

Similarly,

$$p(\tau|x, \mu, a, b) \propto \mathcal{N}\left(x|\mu, \tau^{-1}\right) \operatorname{Gam}\left(\tau|a, b\right)$$

which, we know from Section 2.3.6, is a gamma distribution

$$\operatorname{Gam}\left(\tau|\widehat{a},\widehat{b}\right)$$

with parameters

$$\widehat{a} = a + \frac{1}{2}$$

$$\widehat{b} = b + \frac{1}{2} (x - \mu)^2.$$

**11.14 NOTE**: In PRML,  $\alpha_i^2$  should be  $\alpha^2$  in the last term on the r.h.s. of (11.50). If we take the expectation of (11.50), we obtain

$$\mathbb{E}\left[z_{i}^{\prime}\right] = \mathbb{E}\left[\mu_{i} + \alpha\left(z_{i} - \mu_{i}\right) + \sigma_{i}\left(1 - \alpha^{2}\right)^{1/2}\nu\right]$$

$$= \mu_{i} + \alpha\left(\mathbb{E}\left[z_{i}\right] - \mu_{i}\right) + \sigma_{i}\left(1 - \alpha^{2}\right)^{1/2}\mathbb{E}\left[\nu\right]$$

$$= \mu_{i}.$$

Now we can use this together with (1.40) to compute the variance of  $z'_i$ ,

$$\operatorname{var}\left[z_{i}^{\prime}\right] = \mathbb{E}\left[\left(z_{i}^{\prime}\right)^{2}\right] - \mathbb{E}\left[z_{i}^{\prime}\right]^{2}$$

$$= \mathbb{E}\left[\left(\mu_{i} + \alpha\left(z_{i} - \mu_{i}\right) + \sigma_{i}\left(1 - \alpha^{2}\right)^{1/2}\nu\right)^{2}\right] - \mu_{i}^{2}$$

$$= \alpha^{2}\mathbb{E}\left[\left(z_{i} - \mu_{i}\right)^{2}\right] + \sigma_{i}^{2}\left(1 - \alpha^{2}\right)\mathbb{E}\left[\nu^{2}\right]$$

$$= \sigma_{i}^{2}$$

where we have used the first and second order moments of  $z_i$  and  $\nu$ .

**11.15** Using (11.56), we can differentiate (11.57), yielding

$$\frac{\partial H}{\partial r_i} = \frac{\partial K}{\partial r_i} = r_i$$

and thus (11.53) and (11.58) are equivalent.

Similarly, differentiating (11.57) w.r.t.  $z_i$  we get

$$\frac{\partial H}{\partial z_i} = \frac{\partial E}{\partial z_i},$$

and from this, it is immediately clear that (11.55) and (11.59) are equivalent.

**11.16** From the product rule we know that

$$p(\mathbf{r}|\mathbf{z}) \propto p(\mathbf{r},\mathbf{z}).$$

Using (11.56) and (11.57) to rewrite (11.63) as

$$p(\mathbf{z}, \mathbf{r}) = \frac{1}{Z_H} \exp(-H(\mathbf{z}, \mathbf{r}))$$

$$= \frac{1}{Z_H} \exp(-E(\mathbf{z}) - K(\mathbf{r}))$$

$$= \frac{1}{Z_H} \exp\left(-\frac{1}{2} ||\mathbf{r}||^2\right) \exp(-E(\mathbf{z})).$$

Thus we see that  $p(\mathbf{z}, \mathbf{r})$  is Gaussian w.r.t.  $\mathbf{r}$  and hence  $p(\mathbf{r}|\mathbf{z})$  will be Gaussian too.

**11.17 NOTE**: In the 1<sup>st</sup> printing of PRML, there are sign errors in equations (11.68) and (11.69). In both cases, the sign of the argument to the exponential forming the second argument to the min-function should be changed.

First we note that, if  $H(\mathcal{R}) = H(\mathcal{R}')$ , then the detailed balance clearly holds, since in this case, (11.68) and (11.69) are identical.

Otherwise, we either have  $H(\mathcal{R}) > H(\mathcal{R}')$  or  $H(\mathcal{R}) < H(\mathcal{R}')$ . We consider the former case, for which (11.68) becomes

$$\frac{1}{Z_H} \exp(-H(\mathcal{R})) \delta V \frac{1}{2},$$

since the min-function will return 1. (11.69) in this case becomes

$$\frac{1}{Z_H} \exp(-H(\mathcal{R}')) \delta V \frac{1}{2} \exp(H(\mathcal{R}') - H(\mathcal{R})) = \frac{1}{Z_H} \exp(-H(\mathcal{R})) \delta V \frac{1}{2}.$$

In the same way it can be shown that both (11.68) and (11.69) equal

$$\frac{1}{Z_H} \exp(-H(\mathcal{R}')) \delta V \frac{1}{2}$$

when  $H(\mathcal{R}) < H(\mathcal{R}')$ .

## **Chapter 12 Continuous Latent Variables**

12.1 Suppose that the result holds for projection spaces of dimensionality M. The M+1 dimensional principal subspace will be defined by the M principal eigenvectors  $\mathbf{u}_1,\ldots,\mathbf{u}_M$  together with an additional direction vector  $\mathbf{u}_{M+1}$  whose value we wish to determine. We must constrain  $\mathbf{u}_{M+1}$  such that it cannot be linearly related to  $\mathbf{u}_1,\ldots,\mathbf{u}_M$  (otherwise it will lie in the M-dimensional projection space instead of defining an M+1 independent direction). This can easily be achieved by requiring that  $\mathbf{u}_{M+1}$  be orthogonal to  $\mathbf{u}_1,\ldots,\mathbf{u}_M$ , and these constraints can be enforced using Lagrange multipliers  $\eta_1,\ldots,\eta_M$ .

Following the argument given in section 12.1.1 for  $\mathbf{u}_1$  we see that the variance in the direction  $\mathbf{u}_{M+1}$  is given by  $\mathbf{u}_{M+1}^{\mathrm{T}}\mathbf{S}\mathbf{u}_{M+1}$ . We now maximize this using a Lagrange multiplier  $\lambda_{M+1}$  to enforce the normalization constraint  $\mathbf{u}_{M+1}^{\mathrm{T}}\mathbf{u}_{M+1}=1$ . Thus we seek a maximum of the function

$$\mathbf{u}_{M+1}^{\mathrm{T}}\mathbf{S}\mathbf{u}_{M+1} + \lambda_{M+1}\left(1 - \mathbf{u}_{M+1}^{\mathrm{T}}\mathbf{u}_{M+1}\right) + \sum_{i=1}^{M} \eta_{i}\mathbf{u}_{M+1}^{\mathrm{T}}\mathbf{u}_{i}.$$

with respect to  $\mathbf{u}_{M+1}$ . The stationary points occur when

$$0 = 2\mathbf{S}\mathbf{u}_{M+1} - 2\lambda_{M+1}\mathbf{u}_{M+1} + \sum_{i=1}^{M} \eta_i \mathbf{u}_i.$$

Left multiplying with  $\mathbf{u}_j^{\mathrm{T}}$ , and using the orthogonality constraints, we see that  $\eta_j=0$  for  $j=1,\ldots,M$ . We therefore obtain

$$\mathbf{S}\mathbf{u}_{M+1} = \lambda_{M+1}\mathbf{u}_{M+1}$$

and so  $\mathbf{u}_{M+1}$  must be an eigenvector of  $\mathbf{S}$  with eigenvalue  $\mathbf{u}_{M+1}$ . The variance in the direction  $\mathbf{u}_{M+1}$  is given by  $\mathbf{u}_{M+1}^{\mathrm{T}}\mathbf{S}\mathbf{u}_{M+1} = \lambda_{M+1}$  and so is maximized by choosing  $\mathbf{u}_{M+1}$  to be the eigenvector having the largest eigenvalue amongst those not previously selected. Thus the result holds also for projection spaces of dimensionality M+1, which completes the inductive step. Since we have already shown this result explicitly for M=1 if follows that the result must hold for any  $M \leq D$ .

**12.2** Using the result (C.24) we can set the derivative of  $\widetilde{J}$  with respect to  $\widehat{\mathbf{U}}$  to zero to obtain

$$0 = (\mathbf{S}^{\mathrm{T}} + \mathbf{S})\widehat{\mathbf{U}} - \widehat{\mathbf{U}}(\mathbf{H}^{\mathrm{T}} + \mathbf{H}).$$

We note that S is symmetric so that  $S^T = S$ . Similarly we can choose H to be symmetric without loss of generality since any non-symmetric component would cancel from the expression for  $\widetilde{J}$  since the latter involves a trace of H times a symmetric matrix. (See Exercise 1.14 and its solution.) Thus we have

$$\widehat{SU} = \widehat{U}H$$
.

Clearly one solution is take the columns of  $\widehat{\mathbf{U}}$  to be eigenvectors of  $\mathbf{S}$ . To discuss the general solution, consider the eigenvector equation for  $\mathbf{H}$  given by

$$H\Psi=\Psi L.$$

Since  $\mathbf{H}$  is a symmetric matrix its eigenvectors can be chosen to be a complete orthonormal set in the (D-M)-dimensional space, and  $\mathbf{L}$  will be a diagonal matrix containing the corresponding eigenvalues, with  $\mathbf{\Psi}$  a  $(D-M)\times(D-M)$ -dimensional orthogonal matrix satisfying  $\mathbf{\Psi}^{\mathrm{T}}\mathbf{\Psi}=\mathbf{I}$ .

If we right multiply the eigenvector equation for  ${\bf S}$  by  ${f \Psi}$  we obtain

$$\mathbf{S}\widehat{\mathbf{U}}\mathbf{\Psi} = \widehat{\mathbf{U}}\mathbf{H}\mathbf{\Psi} = \widehat{\mathbf{U}}\mathbf{\Psi}\mathbf{L}$$

and defining  $\widetilde{\mathbf{U}} = \widehat{\mathbf{U}} \mathbf{\Psi}$  we obtain

$$\widetilde{SU} = \widetilde{UL}$$

so that the columns of  $\widetilde{\mathbf{U}}$  are the eigenvectors of  $\mathbf{S}$ , and the elements of the diagonal matrix  $\mathbf{L}$  are the corresponding eigenvalues.

Using the cyclic property of the trace, together with the orthogonality property  $\Psi^T\Psi$ , the distortion function can be written

$$J = \mathrm{Tr}(\widehat{\mathbf{U}}^{\mathrm{T}}\mathbf{S}\widehat{\mathbf{U}}) = \mathrm{Tr}(\mathbf{\Psi}^{\mathrm{T}}\widehat{\mathbf{U}}^{\mathrm{T}}\mathbf{S}\widehat{\mathbf{U}}\mathbf{\Psi}) = \mathrm{Tr}(\widetilde{\mathbf{U}}\mathbf{S}\widetilde{\mathbf{U}}) = \mathrm{Tr}(\mathbf{L}).$$

Thus the distortion measure can be expressed in terms of the sum of the eigenvalues of **S** corresponding to the (D-M) eigenvectors orthogonal to the principal subspace.

**12.3** By left-multiplying both sides of (12.28) by  $\mathbf{v}_i^{\mathrm{T}}$ , we obtain

$$\frac{1}{N} \mathbf{v}_i^{\mathrm{T}} \mathbf{X} \mathbf{X}^{\mathrm{T}} \mathbf{v}_i = \lambda_i \mathbf{v}_i^{\mathrm{T}} \mathbf{v}_i = \lambda_i$$

where we have used the fact that  $v_i$  is orthonormal. From this we see that

$$\left\|\mathbf{X}^{\mathrm{T}}\mathbf{v}_{i}\right\|^{2} = N\lambda_{i}$$

from which we in turn see that  $u_i$  defined by (12.30) will have unit length.

**12.4** Using the results of Section 8.1.4, the marginal distribution for this modified probabilistic PCA model can be written

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{m} + \boldsymbol{\mu}, \sigma^2 \mathbf{I} + \mathbf{W}^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \mathbf{W}).$$

If we now define new parameters

$$\widetilde{\mathbf{W}} = \mathbf{\Sigma}^{1/2} \mathbf{W}$$
 $\widetilde{\boldsymbol{\mu}} = \mathbf{W} \mathbf{m} + \boldsymbol{\mu}$ 

then we obtain a marginal distribution having the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\widetilde{\boldsymbol{\mu}}, \sigma^2 \mathbf{I} + \widetilde{\mathbf{W}}^{\mathrm{T}} \widetilde{\mathbf{W}}).$$

Thus any Gaussian form for the latent distribution therefore gives rise to a predictive distribution having the same functional form, and so for convenience we choose the simplest form, namely one with zero mean and unit covariance.

12.5 Since y = Ax + b,

$$p(\mathbf{y}|\mathbf{x}) = \delta(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})$$

i.e. a delta function at Ax + b. From the sum and product rules, we have

$$p(\mathbf{y}) = \int p(\mathbf{y}, \mathbf{x}) d\mathbf{x} = \int p(\mathbf{y}|\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$
$$= \int \delta(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) p(\mathbf{x}) d\mathbf{x}.$$

When M = D and **A** is assumed to have full rank, we have

$$\mathbf{x} = \mathbf{A}^{-1}(\mathbf{y} - \mathbf{b})$$

and thus

$$p(\mathbf{y}) = \mathcal{N} \left( \mathbf{A}^{-1} (\mathbf{y} - \mathbf{b}) | \boldsymbol{\mu}, \boldsymbol{\Sigma} \right)$$
$$= \mathcal{N} \left( \mathbf{y} | \mathbf{A} \boldsymbol{\mu} + \mathbf{b}, \mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^{\mathrm{T}} \right).$$

When M > D, y will be strictly confined to a D-dimensional subspace and hence p(y) will be singular. In this case we have

$$\mathbf{x} = \mathbf{A}^{-L}(\mathbf{y} - \mathbf{b})$$

where  $A^{-L}$  is the left inverse of A and thus

$$p(\mathbf{y}) = \mathcal{N}\left(\mathbf{A}^{-L}(\mathbf{y} - \mathbf{b})|\boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$$
$$= \mathcal{N}\left(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \left(\left(\mathbf{A}^{-L}\right)^{T} \boldsymbol{\Sigma}^{-1} \mathbf{A}^{-L}\right)^{-1}\right).$$

The covariance matrix on the last line cannot be computed, but we can still compute p(y), by using the corresponding precision matrix and constraining the density to be zero outside the column space of A:

$$p(\mathbf{y}) = \mathcal{N}\left(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \left(\left(\mathbf{A}^{-L}\right)^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}\mathbf{A}^{-L}\right)^{-1}\right)\delta\left(\mathbf{y} - \mathbf{A}\mathbf{A}^{-L}\left(\mathbf{y} - \mathbf{b}\right) - \mathbf{b}\right).$$

Finally, when M < D, we can make use of (2.113)–(2.115) and set  $\mathbf{L}^{-1} = \mathbf{0}$  in (2.114). While this means that  $p(\mathbf{y}|\mathbf{x})$  is singular, the marginal distribution  $p(\mathbf{y})$ , given by (2.115), is non-singular as long as  $\mathbf{A}$  and  $\mathbf{\Sigma}$  are assumed to be of full rank.

- **12.6** Omitting the parameters, W,  $\mu$  and  $\sigma$ , leaving only the stochastic variables z and x, the graphical model for probabilistic PCA is identical with the 'naive Bayes' model shown in Figure 8.24 in Section 8.2.2. Hence these two models exhibit the same independence structure.
- **12.7** From (2.59), the multivariate form of (2.270), (12.31) and (12.32), we get

$$\begin{split} \mathbb{E}[\mathbf{x}] &= \mathbb{E}_{\mathbf{z}} \left[ \mathbb{E}_{\mathbf{x}} \left[ \mathbf{x} | \mathbf{z} \right] \right] \\ &= \mathbb{E}_{\mathbf{z}} \left[ \mathbf{W} \mathbf{z} + \boldsymbol{\mu} \right] \\ &= \boldsymbol{\mu}. \end{split}$$

Combining this with (2.63), the covariance formula corresponding to (2.271), (12.31) and (12.32), we get

$$\begin{aligned} & \operatorname{cov}[\mathbf{x}] &= & \mathbb{E}_{\mathbf{z}} \left[ \operatorname{cov}_{\mathbf{x}}[\mathbf{x}|\mathbf{z}] \right] + \operatorname{cov}_{\mathbf{x}} \left[ \mathbb{E}_{\mathbf{x}}[\mathbf{x}|\mathbf{z}] \right] \\ &= & \mathbb{E}_{\mathbf{z}} \left[ \sigma^{2} \mathbf{I} \right] + \operatorname{cov}_{\mathbf{z}} \left[ \mathbf{W}\mathbf{z} + \boldsymbol{\mu} \right] \\ &= & \sigma^{2} \mathbf{I} + \mathbb{E}_{\mathbf{z}} \left[ \left( \mathbf{W}\mathbf{z} + \boldsymbol{\mu} - \mathbb{E}_{\mathbf{z}} \left[ \mathbf{W}\mathbf{z} + \boldsymbol{\mu} \right] \right) \left( \mathbf{W}\mathbf{z} + \boldsymbol{\mu} - \mathbb{E}_{\mathbf{z}} \left[ \mathbf{W}\mathbf{z} + \boldsymbol{\mu} \right] \right)^{\mathrm{T}} \right] \\ &= & \sigma^{2} \mathbf{I} + \mathbb{E}_{\mathbf{z}} \left[ \mathbf{W}\mathbf{z}\mathbf{z}^{\mathrm{T}}\mathbf{W}^{\mathrm{T}} \right] \\ &= & \sigma^{2} \mathbf{I} + \mathbf{W}\mathbf{W}^{\mathrm{T}}. \end{aligned}$$

**12.8 NOTE**: In the 1<sup>st</sup> printing of PRML, equation (12.42) contains a mistake; the covariance on the r.h.s. should be  $\sigma^2 \mathbf{M}^{-1}$ .

By matching (12.31) with (2.113) and (12.32) with (2.114), we have from (2.116) and (2.117) that

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}\left(\mathbf{z}|(\mathbf{I} + \sigma^{-2}\mathbf{W}^{\mathrm{T}}\mathbf{W})^{-1}\mathbf{W}^{\mathrm{T}}\sigma^{-2}\mathbf{I}(\mathbf{x} - \boldsymbol{\mu}), (\mathbf{I} + \sigma^{-2}\mathbf{W}^{\mathrm{T}}\mathbf{W})^{-1}\right)$$
$$= \mathcal{N}\left(\mathbf{z}|\mathbf{M}^{-1}\mathbf{W}^{\mathrm{T}}(\mathbf{x} - \boldsymbol{\mu}), \sigma^{2}\mathbf{M}^{-1}\right),$$

where we have also used (12.41).

**12.9** By expanding the square in the last term of (12.43) and then making use of results from Appendix C, we can calculate the derivative w.r.t.  $\mu$  and set this equal to zero, yielding

$$-N\mathbf{C}^{-1}\boldsymbol{\mu} + \mathbf{C}^{-1} \sum_{n=1}^{N} \mathbf{x}_n = \mathbf{0}.$$
 (314)

Rearranging this and making use of (12.1), we get

$$\mu = \overline{\mathbf{x}}$$
.

**12.10** Using results from Appendix C, we can differentiate the r.h.s. of (314) w.r.t.  $\mu$ , giving

$$-N\mathbf{C}^{-1}$$
.

If  $\sigma^2 > 0$ , C will be positive definite, in which case  $C^{-1}$  is also positive definite, and hence the log likelihood function will be concave with a unique maximum at  $\mu$ .

**12.11** Taking  $\sigma^2 \to 0$  in (12.41) and substituting into (12.48) we obtain the posterior mean for probabilistic PCA in the form

$$(\mathbf{W}_{\mathrm{ML}}^{\mathrm{T}}\mathbf{W}_{\mathrm{ML}})^{-1}\mathbf{W}_{\mathrm{ML}}^{\mathrm{T}}(\mathbf{x}-\overline{\mathbf{x}}).$$

Now substitute for  $\mathbf{W}_{\mathrm{ML}}$  using (12.45) in which we take  $\mathbf{R} = \mathbf{I}$  for compatibility with conventional PCA. Using the orthogonality property  $\mathbf{U}_{M}^{\mathrm{T}}\mathbf{U}_{M} = \mathbf{I}$  and setting  $\sigma^{2} = 0$ , this reduces to

$$\mathbf{L}^{-1/2}\mathbf{U}_{M}^{\mathrm{T}}(\mathbf{x}-\overline{\mathbf{x}})$$

which is the orthogonal projection is given by the conventional PCA result (12.24).

**12.12** For  $\sigma^2 > 0$  we can show that the projection is shifted towards the origin of latent space by showing that the magnitude of the latent space vector is reduced compared to the  $\sigma^2 = 0$  case. The orthogonal projection is given by

$$\mathbf{z}_{\mathrm{orth}} = \mathbf{L}_{M}^{-1/2} \mathbf{U}_{M}^{\mathrm{T}} (\mathbf{x} - \overline{\mathbf{x}})$$

where  $L_M$  and  $U_M$  are defined as in (12.45). The posterior mean projection is given by (12.48) and so the difference between the squared magnitudes of each of these is given by

$$\begin{aligned} & \left\| \mathbf{z}_{\text{orth}} \right\|^2 - \left\| \mathbb{E}[\mathbf{z}|\mathbf{x}] \right\|^2 \\ &= \left( \mathbf{x} - \overline{\mathbf{x}} \right)^{\text{T}} \left( \mathbf{U}_M \mathbf{L}_M^{-1/2} \mathbf{U}_M - \mathbf{W}_{\text{ML}} \mathbf{M}^{-1} \mathbf{M}^{-1} \mathbf{W}_{\text{ML}}^{\text{T}} \right) (\mathbf{x} - \overline{\mathbf{x}}) \\ &= \left( \mathbf{x} - \overline{\mathbf{x}} \right)^{\text{T}} \mathbf{U}_M \left\{ \mathbf{L}^{-1} - (\mathbf{L} + \sigma^2 \mathbf{I})^{-1} \right\} \mathbf{U}_M^{\text{T}} (\mathbf{x} - \overline{\mathbf{x}}) \end{aligned}$$

where we have use (12.41), (12.45) and the fact that  ${\bf L}$  and  ${\bf M}$  are symmetric. The term in curly braces on the last line is diagonal and has elements  $\sigma^2/\lambda_i(\lambda_i+\sigma^2)$  which are all positive. Thus the matrix is positive definite and so the contraction with the vector  ${\bf U}_M({\bf x}-\overline{\bf x})$  must be positive, and so there is a positive (non-zero) shift towards the origin.



From (12.45) we see that

$$\left(\mathbf{W}_{\mathrm{ML}}^{\mathrm{T}}\mathbf{W}_{\mathrm{ML}}\right)^{-1} = \left(\mathbf{L}_{M} - \sigma^{2}\mathbf{I}\right)^{-1}$$

and so

$$\widetilde{\mathbf{x}}_n = \mathbf{U}_M \mathbf{U}_M^{\mathrm{T}} (\mathbf{x}_n - \overline{\mathbf{x}}).$$

This is the reconstruction of  $\mathbf{x}_n - \overline{\mathbf{x}}$  using the M eigenvectors corresponding to the M largest eigenvalues, which we know from Section 12.1.2 minimizes the least squares projection cost (12.11).

**12.14** If we substitute D - 1 for M in (12.51), we get

$$D(D-1) + 1 - \frac{(D-1)((D-1)-1)}{2} = \frac{2D^2 - 2D + 2 - D^2 + 3D - 2}{2}$$
$$= \frac{D^2 + D}{2} = \frac{D(D+1)}{2}$$

as required. Setting M=0 in (12.51) given the value 1 for the number of parameters in  $\mathbb{C}$ , corresponding to the scalar variance parameter,  $\sigma^2$ .

**12.15 NOTE**: In PRML, a term  $M/2 \ln(2\pi)$  is missing from the summand on the r.h.s. of (12.53). However, this is only stated here for completeness as it actually does not affect this solution.

Using standard derivatives together with the rules for matrix differentiation from Appendix C, we can compute the derivatives of (12.53) w.r.t.  $\mathbf{W}$  and  $\sigma^2$ :

$$\frac{\partial}{\partial \mathbf{W}} \mathbb{E}[\ln p\left(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2\right)] = \sum_{n=1}^{N} \left\{ \frac{1}{\sigma^2} (\mathbf{x}_n - \overline{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^{\mathrm{T}} - \frac{1}{\sigma^2} \mathbf{W} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\mathrm{T}}] \right\}$$

and

$$\frac{\partial}{\partial \sigma^2} \mathbb{E}[\ln p\left(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2\right)] = \sum_{n=1}^{N} \left\{ \frac{1}{2\sigma^4} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\mathrm{T}}] \mathbf{W}^{\mathrm{T}} \mathbf{W} + \frac{1}{2\sigma^4} \|\mathbf{x}_n - \overline{\mathbf{x}}\|^2 - \frac{1}{\sigma^4} \mathbb{E}[\mathbf{z}_n]^{\mathrm{T}} \mathbf{W}^{\mathrm{T}} (\mathbf{x}_n - \overline{\mathbf{x}}) - \frac{D}{2\sigma^2} \right\}$$

Setting these equal to zero and re-arranging we obtain (12.56) and (12.57), respectively.

**12.16** We start by noting that the marginal likelihood factorizes over data points as well as the individual elements of the data points,

$$p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^{2}) = \int p(\mathbf{Z})p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\mu}, \mathbf{W}, \sigma^{2}) d\mathbf{Z}$$

$$= \prod_{n=1}^{N} \int p(\mathbf{z}_{n})p(\mathbf{x}_{n}|\mathbf{z}_{n}, \boldsymbol{\mu}, \mathbf{W}, \sigma^{2}) d\mathbf{z}_{n}$$

$$= \prod_{n=1}^{N} \int p(\mathbf{z}_{n}) \prod_{i=1}^{D} \mathcal{N}(x_{ni}|\mathbf{w}_{i}\mathbf{z}_{n} + \mu_{i}, \sigma^{2}) d\mathbf{z}_{n}$$
(315)

where  $x_{ni}$  denotes the  $i^{\rm th}$  element of  $\mathbf{x}_n$ ,  $\mu_i$  denotes the  $i^{\rm th}$  element of  $\boldsymbol{\mu}$  and  $\mathbf{w}_i$  the  $i^{\rm th}$  row of  $\mathbf{W}$ . If we assume that any missing values are missing at random (see page 441 of PRML), we can deal with these by integrating them out of (315). Let  $\mathbf{x}_n^{\rm o}$  and  $\mathbf{x}_n^{\rm m}$  denote the observed and missing parts of  $\mathbf{x}_n$ , respectively. Using this notation, we can rewrite (315) as

$$p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^{2}) = \prod_{n=1}^{N} \int p(\mathbf{z}_{n}) \prod_{x_{ni} \in \mathbf{x}_{n}^{\circ}} \mathcal{N}(x_{ni}|\mathbf{w}_{i}\mathbf{z}_{n} + \mu_{i}, \sigma^{2})$$

$$= \prod_{x_{nj} \in \mathbf{x}_{n}^{m}} \mathcal{N}(x_{nj}|\mathbf{w}_{j}\mathbf{z}_{n} + \mu_{j}, \sigma^{2}) d\mathbf{x}_{n}^{m} d\mathbf{z}_{n}$$

$$= \prod_{n=1}^{N} \int p(\mathbf{z}_{n}) \prod_{x_{ni} \in \mathbf{x}_{n}^{\circ}} \mathcal{N}(x_{ni}|\mathbf{w}_{i}\mathbf{z}_{n} + \mu_{i}, \sigma^{2}) d\mathbf{z}_{n}$$

$$= \prod_{n=1}^{N} p(\mathbf{x}_{n}^{\circ}|\boldsymbol{\mu}, \mathbf{W}, \sigma^{2}).$$

Thus we are left with a 'reduced' marginal likelihood, where for each data point,  $\mathbf{x}_n$ , we only need to consider the observed elements,  $\mathbf{x}_n^{\circ}$ .

Now we can derive an EM algorithm for finding the parameter values that maximizes this 'reduced' marginal likelihood. In doing so, we shall find it convenient to introduce indicator variables,  $\iota_{ni}$ , such that  $\iota_{ni}=1$  if  $x_{ni}$  is observed and  $\iota_{ni}=0$  otherwise. This allows us to rewrite (12.32) as

$$p(\mathbf{x}|\mathbf{z}) = \prod_{i=1}^{D} \mathcal{N}(x_i|\mathbf{w}_i\mathbf{z} + \mu_i, \sigma^2)^{\iota_{ni}}$$

and the complete-data log likelihood as

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma) = \sum_{n=1}^{N} \left\{ \ln p(\mathbf{z}_n) + \sum_{i=1}^{D} \iota_{ni} \ln \mathcal{N}(x_{ni} | \mathbf{w}_i \mathbf{z}_n + \mu_i, \sigma^2) \right\}.$$

Following the path taken in Section 12.2.2, making use of (12.31) and taking the expectation w.r.t. the latent variables, we obtain

$$\mathbb{E}\left[\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma)\right] = -\sum_{n=1}^{N} \left\{ \frac{M}{2} \ln(2\pi) + \frac{1}{2} \operatorname{Tr}\left(\mathbb{E}\left[\mathbf{z}_{n} \mathbf{z}_{n}^{\mathrm{T}}\right]\right) + \sum_{i=1}^{D} \iota_{ni} \left\{ \ln(2\pi\sigma^{2}) + \frac{1}{2\sigma^{2}} (x_{ni} - \mu_{ni})^{2} - \frac{1}{\sigma^{2}} \mathbb{E}[\mathbf{z}_{n}]^{\mathrm{T}} \mathbf{w}_{i}^{\mathrm{T}}(x_{ni} - \mu_{ni}) + \frac{1}{2\sigma^{2}} \operatorname{Tr}\left(\mathbb{E}\left[\mathbf{z}_{n} \mathbf{z}_{n}^{\mathrm{T}}\right] \mathbf{w}_{i}^{\mathrm{T}} \mathbf{w}_{i}\right) \right\} \right\}.$$

Taking the derivative w.r.t.  $\mu_i$  and setting the result equal to zero, we obtain

$$\mu_i^{\text{new}} = \frac{1}{\sum_{m=1}^N \iota_{mi}} \sum_{n=1}^N \iota_{ni} x_{ni}.$$

In the E step, we compute the sufficient statistics, which due to the altered form of  $p(\mathbf{x}|\mathbf{z})$  now take slightly different shapes. Equation (12.54) becomes

$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}_n^{-1} \mathbf{W}_n^{\mathrm{T}} \mathbf{y}_n$$

where  $\mathbf{y}_n$  is a vector containing the observed elements of  $\mathbf{x}_n$  minus the corresponding elements of  $\mu^{\text{new}}$ ,  $\mathbf{W}_n$  is a matrix formed by the rows of  $\mathbf{W}$  corresponding to the observed elements of  $\mathbf{x}_n$  and, accordingly, from (12.41)

$$\mathbf{M}_n = \mathbf{W}_n^{\mathrm{T}} \mathbf{W}_n + \sigma^2 \mathbf{I}.$$

Similarly,

$$\mathbb{E}\left[\mathbf{z}_{n}\mathbf{z}_{n}^{\mathrm{T}}\right] = \sigma^{2}\mathbf{M}_{n}^{-1} + \mathbb{E}[\mathbf{z}_{n}]\mathbb{E}[\mathbf{z}_{n}]^{\mathrm{T}}.$$

The M step is similar to the fully observed case, with

$$\mathbf{W}_{\text{new}} = \left[ \sum_{n=1}^{N} \mathbf{y}_{n} \mathbb{E}[\mathbf{z}_{n}]^{\text{T}} \right] \left[ \sum_{n=1}^{N} \mathbb{E}\left[\mathbf{z}_{n} \mathbf{z}_{n}^{\text{T}}\right]^{-1} \right]$$

$$\sigma_{\text{new}}^{2} = \frac{1}{\sum_{n=1}^{N} \sum_{i=1}^{D} \iota_{ni}} \sum_{n=1}^{N} \sum_{i=1}^{D} \iota_{ni} \left\{ (x_{ni} - \mu_{i}^{\text{new}})^{2} - 2\mathbb{E}[\mathbf{z}_{n}]^{\text{T}} \left(\mathbf{w}_{i}^{\text{new}}\right)^{\text{T}} \left(x_{ni} - \mu_{i}^{\text{new}}\right) + \text{Tr} \left( \mathbb{E}\left[\mathbf{z}_{n} \mathbf{z}_{n}^{\text{T}}\right] \left(\mathbf{w}_{i}^{\text{new}}\right)^{\text{T}} \mathbf{w}_{i}^{\text{new}} \right) \right\}$$

where  $\mathbf{w}_i^{\mathrm{new}}$  equals the  $i^{\mathrm{th}}$  row  $\mathbf{W}_{\mathrm{new}}.$ 

In the fully observed case, all  $\iota_{ni}=1$ ,  $\mathbf{y}_n=\mathbf{x}_n$ ,  $\mathbf{W}_n=\mathbf{W}$  and  $\boldsymbol{\mu}^{\text{new}}=\overline{\mathbf{x}}$ , and hence we recover (12.54)–(12.57).

**12.17 NOTE**: In PRML, there are errors in equation (12.58) and the preceding text. In (12.58),  $\widetilde{\mathbf{X}}$  should be  $\widetilde{\mathbf{X}}^{\mathrm{T}}$  and in the preceding text we define  $\Omega$  to be a matrix of size  $\underline{M} \times \underline{N}$  whose  $n^{\mathrm{th}}$  column is given by the vector  $\mathbb{E}[\mathbf{z}_n]$ .

Setting the derivative of J with respect to  $\mu$  to zero gives

$$0 = -\sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu} - \mathbf{W}\mathbf{z}_n)$$

from which we obtain

$$\mu = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n - \frac{1}{N} \sum_{n=1}^{N} \mathbf{W} \mathbf{z}_n = \overline{\mathbf{x}} - \mathbf{W} \overline{\mathbf{z}}.$$

Back-substituting into J we obtain

$$J = \sum_{n=1}^{N} \|(\mathbf{x}_n - \overline{\mathbf{x}} - \mathbf{W}(\mathbf{z}_n - \overline{\mathbf{z}})\|^2.$$

We now define  $\mathbf{X}$  to be a matrix of size  $N \times D$  whose  $n^{\text{th}}$  row is given by the vector  $\mathbf{x}_n - \overline{\mathbf{x}}$  and similarly we define  $\mathbf{Z}$  to be a matrix of size  $D \times M$  whose  $n^{\text{th}}$  row is given by the vector  $\mathbf{z}_n - \overline{\mathbf{z}}$ . We can then write J in the form

$$J = \operatorname{Tr}\left\{ (\mathbf{X} - \mathbf{Z}\mathbf{W}^{\mathrm{T}})(\mathbf{X} - \mathbf{Z}\mathbf{W}^{\mathrm{T}})^{\mathrm{T}} \right\}.$$

Differentiating with respect to  $\mathbf{Z}$  keeping  $\mathbf{W}$  fixed gives rise to the PCA E-step (12.58). Similarly setting the derivative of J with respect to  $\mathbf{W}$  to zero with  $\{\mathbf{z}_n\}$  fixed gives rise to the PCA M-step (12.59).

**12.18** Analysis of the number of independent parameters follows the same lines as for probabilistic PCA except that the one parameter noise covariance  $\sigma^2 \mathbf{I}$  is replaced by a D parameter diagonal covariance  $\Psi$ . Thus the number of parameters is increased by D-1 compared to the probabilistic PCA result (12.51) giving a total number of independent parameters of

$$D(M+1) - M(M-1)/2$$
.

12.19 To see this we define a rotated latent space vector  $\widetilde{\mathbf{z}} = \mathbf{R}\mathbf{z}$  where  $\mathbf{R}$  is an  $M \times M$  orthogonal matrix, and similarly defining a modified factor loading matrix  $\widetilde{\mathbf{W}} = \mathbf{W}\mathbf{R}$ . Then we note that the latent space distribution  $p(\mathbf{z})$  depends only on  $\mathbf{z}^T\mathbf{z} = \widetilde{\mathbf{z}}^T\widetilde{\mathbf{z}}$ , where we have used  $\mathbf{R}^T\mathbf{R} = \mathbf{I}$ . Similarly, the conditional distribution of the observed variable  $p(\mathbf{x}|\mathbf{z})$  depends only on  $\mathbf{W}\mathbf{z} = \widetilde{\mathbf{W}}\widetilde{\mathbf{z}}$ . Thus the joint distribution takes the same form for any choice of  $\mathbf{R}$ . This is reflected in the predictive distribution  $p(\mathbf{x})$  which depends on  $\mathbf{W}$  only through the quantity  $\mathbf{W}\mathbf{W}^T = \widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^T$  and hence is also invariant to different choices of  $\mathbf{R}$ .

**12.20** The log likelihood function is given by

$$\ln L(\boldsymbol{\mu}, \mathbf{W}, \boldsymbol{\Psi}) = \sum_{n=1}^{N} \ln p(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{C})$$
$$= \sum_{n=1}^{N} \left\{ -\ln |\mathbf{C}| - (\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}} \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\}$$

where C is defined by (12.65). Differentiating with respect to  $\mu^{T}$  and setting the derivative to zero we obtain

$$0 = \sum_{n=1}^{N} \mathbf{C}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}).$$

Pre-multiplying by C and re-arranging shows that  $\mu$  is given by the sample mean defined by (12.1). Taking the second derivative of the log likelihood we obtain

$$\frac{\partial^2 \ln L}{\partial \boldsymbol{\mu}^{\mathrm{T}} \partial \boldsymbol{\mu}} = -N \mathbf{C}^{-1}.$$

Since C is a positive definite matrix, its inverse will also be positive definite (see Appendix C) and hence the stationary point will be a unique maximum of the log likelihood.

**12.21** By making use of (2.113)–(2.117) together with (12.31) and (12.64), we obtain the posterior distribution of the latent variable **z**, for a given value of the observed variable **x**, in the form

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mathbf{G}\mathbf{W}^{\mathrm{T}}\mathbf{\Psi}^{-1}(\mathbf{x} - \overline{\mathbf{x}}).$$

where G is defined by (12.68). Since the data points are drawn independently from the distribution, the posterior distribution for  $\mathbf{z}_n$  depends only on the observation  $\mathbf{x}_n$  (for given values of the parameters). Thus (12.66) follows directly. For the second order statistic we use the general result

$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^{\mathrm{T}}] = \text{cov}[\mathbf{z}_n] + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^{\mathrm{T}}$$

from which we obtain (12.67).

**12.22 NOTE**: In PRML, Equations (12.69) and (12.70) contain minor typographical errors. On the l.h.s.  $\mathbf{W}^{\text{new}}$  and  $\mathbf{\Psi}^{\text{new}}$  should be  $\mathbf{W}_{\text{new}}$  and  $\mathbf{\Psi}_{\text{new}}$ , respectively.

For the M step we first write down the complete-data log likelihood function, which takes the form

$$\ln L_{\mathbf{C}} = \sum_{n=1}^{N} \{ \ln p(\mathbf{z}_n) + \ln p(\mathbf{x}_n | \mathbf{z}_n) \}$$

$$= \frac{1}{2} \sum_{n=1}^{N} \{ -M \ln(2\pi) - \mathbf{z}_n^{\mathrm{T}} \mathbf{z}_n - D \ln(2\pi) - \ln |\mathbf{\Psi}| - (\mathbf{x}_n - \overline{\mathbf{x}} - \mathbf{W} \mathbf{z}_n)^{\mathrm{T}} \mathbf{\Psi}^{-1} (\mathbf{x}_n - \overline{\mathbf{x}} - \mathbf{W} \mathbf{z}_n) \}.$$

Now take the expectation with respect to  $\{\mathbf{z}_n\}$  to give

$$\mathbb{E}_{\mathbf{z}} \left[ \ln L_{\mathbf{C}} \right] = \frac{1}{2} \sum_{n=1}^{N} \left\{ -\ln |\mathbf{\Psi}| - \operatorname{Tr} \left( \mathbb{E}[\mathbf{z}_{n} \mathbf{z}_{n}^{\mathsf{T}}] \mathbf{W}^{\mathsf{T}} \mathbf{\Psi}^{-1} \mathbf{W} \right) + 2 \mathbb{E}[\mathbf{z}_{n}]^{\mathsf{T}} \mathbf{W}^{\mathsf{T}} \mathbf{\Psi}^{-1} (\mathbf{x}_{n} - \overline{\mathbf{x}}) \right\} - N \operatorname{Tr} \left( \mathbf{S} \mathbf{\Psi}^{-1} \right) + \operatorname{const.}$$

where S is the sample covariance matrix defined by (12.3), and the constant terms are those which are independent of W and  $\Psi$ . Recall that we are making a joint optimization with respect to W and  $\Psi$ . Setting the derivative with respect to  $W^{T}$  equal to zero, making use of the result (C.24), we obtain

$$0 = -2\mathbf{\Psi}^{-1}\mathbf{W}\sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_{n}\mathbf{z}_{n}^{\mathrm{T}}] + 2\mathbf{\Psi}^{-1}\sum_{n=1}^{N} \left[ (\mathbf{x}_{n} - \overline{\mathbf{x}})\mathbb{E}[\mathbf{z}_{n}]^{\mathrm{T}} \right].$$

Pre-multiplying by  $\Psi$  and re-arranging we then obtain (12.69). Note that this result is independent of  $\Psi$ .

Next we maximize the expected complete-data log likelihood with respect to  $\Psi$ . For convenience we set the derivative with respect to  $\Psi^{-1}$  equal to zero, and make use of (C.28) to give

$$0 = N\mathbf{\Psi} - \mathbf{W} \left[ \sum_{n=1}^{N} \mathbb{E}[\mathbf{z}_{n} \mathbf{z}_{n}^{\mathrm{T}}] \right] \mathbf{W}^{\mathrm{T}} + 2 \left[ \sum_{n=1}^{N} (\mathbf{x}_{n} - \overline{\mathbf{x}}) \mathbb{E}[\mathbf{z}_{n}]^{\mathrm{T}} \right] \mathbf{W} - N\mathbf{S}.$$

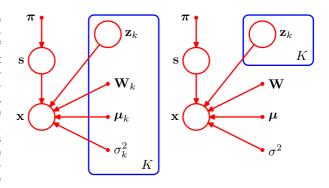
This depends on  $\mathbf{W}$ , and so we can substitute for  $\mathbf{W}_{\rm new}$  in the second term, using the result (12.69), which simplifies the expression. Finally, since  $\Psi$  is constrained to be diagonal, we take set all of the off-diagonal components to zero giving (12.70) as required.

- **12.23** The solution is given in figure 10. The model in which all parameters are shared (left) is not particularly useful, since all mixture components will have identical parameters and the resulting density model will not be any different to one offered by a single PPCA model. Different models would have arisen if only some of the parameters, e.g. the mean  $\mu$ , would have been shared.
- We can derive an EM algorithm by treating  $\eta$  in (2.160) as a latent variable. Thus given a set of i.i.d. data points,  $\mathbf{X} = \{\mathbf{x}_n\}$ , we define the complete-data log likelihood as

$$\ln p(\mathbf{X}, \boldsymbol{\eta} | \boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) = \sum_{n=1}^{N} \left\{ \ln \mathcal{N} \left( \mathbf{x}_{n} | \boldsymbol{\mu}, (\eta_{n} \boldsymbol{\Lambda})^{-1} \right) + \ln \operatorname{Gam} \left( \eta_{n} | \nu / 2, \nu / 2 \right) \right\}$$

where  $\eta$  is an N-dimensional vector with elements  $\eta_n$ . The corresponding expected

Figure 10 The left plot shows the graphical model corresponding to the general mixture of probabilistic PCA. The right plot shows the corresponding model were the parameter of all probabilist PCA models  $(\mu, \mathbf{W} \text{ and } \sigma^2)$  are shared across components. In both plots,  $\mathbf{s}$  denotes the K-nomial latent variable that selects mixture components; it is governed by the parameter,  $\pi$ .



complete-data log likelihood is then given by

$$\mathbb{E}_{\boldsymbol{\eta}} \left[ \ln p(\mathbf{X}, \boldsymbol{\eta} | \boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\nu}) \right] = -\frac{1}{2} \sum_{n=1}^{N} \left\{ D(\ln(2\pi) - \mathbb{E}[\ln \eta_n]) - \ln |\boldsymbol{\Lambda}| + \mathbb{E}[\eta_n] \left( \mathbf{x}^{\mathrm{T}} \boldsymbol{\Lambda} \mathbf{x} - 2 \mathbf{x}^{\mathrm{T}} \boldsymbol{\Lambda} \boldsymbol{\mu} + \boldsymbol{\mu}^{\mathrm{T}} \boldsymbol{\Lambda} \boldsymbol{\mu} \right) + 2 \ln \Gamma(\nu/2) - \nu (\ln \nu - \ln 2) - (\nu - 2) \mathbb{E}[\ln \eta_n] + \mathbb{E}[\eta_n] \right\}$$
(316)

where we have used results from Appendix B. In order to compute the necessary expectations, we need the distribution over  $\eta$ , given by

$$p(\boldsymbol{\eta}|\mathbf{X}, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) = \prod_{n=1}^{N} p(\eta_n|\mathbf{x}_n, \boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu)$$

$$\propto \prod_{n=1}^{N} \mathcal{N}\left(\mathbf{x}_n|\boldsymbol{\mu}, (\eta_n \boldsymbol{\Lambda})^{-1}\right) \operatorname{Gam}\left(\eta_n|\nu/2, \nu/2\right).$$

From Section 2.3.6, we know that the factors in this product are independent Gamma distributions with parameters

$$a_n = \frac{\nu + D}{2}$$

$$b_n = \frac{\nu + (\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Lambda} (\mathbf{x}_n - \boldsymbol{\mu})}{2}$$

and the necessary expectations are given by

$$\mathbb{E}[\eta_n] = \frac{a_n}{b_n}$$

$$\mathbb{E}[\ln \eta_n] = \psi(a_n) - \ln b_n.$$

In the M step, we calculate the derivatives of (316) w.r.t.  $\mu$  and  $\Sigma$ , set these equal to

zero and solve for the respective parameter, to obtain

$$oldsymbol{\mu}_{ ext{ML}} = rac{\sum_{n=1}^{N} \mathbb{E}[\eta_{n}] \mathbf{x}_{n}}{\sum_{n=1}^{N} \mathbb{E}[\eta_{n}]} 
otag 
ot$$

Also for  $\nu$ , we calculate the derivative of (316) and set the result equal to zero, to get

$$1 + \ln\left(\frac{\nu}{2}\right) - \psi\left(\frac{\nu}{2}\right) + \frac{1}{2}\sum_{n=1}^{N} \left\{ \mathbb{E}[\ln \eta_n] - \mathbb{E}[\eta_n] \right\} = 0.$$

Unfortunately, there is no closed form solution w.r.t.  $\nu$ , but since  $\nu$  is scalar, we can afford to solve this equation numerically.

**12.25** Following the discussion of section 12.2, the log likelihood function for this model can be written as

$$L(\boldsymbol{\mu}, \mathbf{W}, \boldsymbol{\Phi}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\mathbf{W}\mathbf{W}^{\mathrm{T}} + \boldsymbol{\Phi}|$$
$$-\frac{1}{2} \sum_{n=1}^{N} \left\{ (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}} (\mathbf{W}\mathbf{W}^{\mathrm{T}} + \boldsymbol{\Phi})^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu}) \right\},$$

where we have used (12.43).

If we consider the log likelihood function for the transformed data set we obtain

$$L_{\mathbf{A}}(\boldsymbol{\mu}, \mathbf{W}, \boldsymbol{\Phi}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\mathbf{W}\mathbf{W}^{\mathrm{T}} + \boldsymbol{\Phi}|$$
$$-\frac{1}{2} \sum_{n=1}^{N} \left\{ (\mathbf{A}\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}} (\mathbf{W}\mathbf{W}^{\mathrm{T}} + \boldsymbol{\Phi})^{-1} (\mathbf{A}\mathbf{x}_{n} - \boldsymbol{\mu}) \right\}.$$

Solving for the maximum likelihood estimator for  $\mu$  in the usual way we obtain

$$\mu_{\mathbf{A}} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{A} \mathbf{x}_n = \mathbf{A} \overline{\mathbf{x}} = \mathbf{A} \mu_{\mathrm{ML}}.$$

Back-substituting into the log likelihood function, and using the definition of the sample covariance matrix (12.3), we obtain

$$L_{\mathbf{A}}(\boldsymbol{\mu}, \mathbf{W}, \boldsymbol{\Phi}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\mathbf{W}\mathbf{W}^{\mathrm{T}} + \boldsymbol{\Phi}|$$
$$-\frac{1}{2} \sum_{n=1}^{N} \operatorname{Tr} \left\{ (\mathbf{W}\mathbf{W}^{\mathrm{T}} + \boldsymbol{\Phi})^{-1} \mathbf{A} \mathbf{S} \mathbf{A}^{\mathrm{T}} \right\}.$$

We can cast the final term into the same form as the corresponding term in the original log likelihood function if we first define

$$\Phi_{\mathbf{A}} = \mathbf{A} \Phi^{-1} \mathbf{A}^{\mathrm{T}}, \qquad \mathbf{W}_{\mathbf{A}} = \mathbf{A} \mathbf{W}.$$

With these definitions the log likelihood function for the transformed data set takes the form

$$\begin{split} L_{\mathbf{A}}(\boldsymbol{\mu}_{\mathbf{A}}, \mathbf{W}_{\mathbf{A}}, \boldsymbol{\Phi}_{\mathbf{A}}) &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\mathbf{W}_{\mathbf{A}} \mathbf{W}_{\mathbf{A}}^{\mathrm{T}} + \boldsymbol{\Phi}_{\mathbf{A}}| \\ &- \frac{1}{2} \sum_{n=1}^{N} \left\{ (\mathbf{x}_{n} - \boldsymbol{\mu}_{\mathbf{A}})^{\mathrm{T}} (\mathbf{W}_{\mathbf{A}} \mathbf{W}_{\mathbf{A}}^{\mathrm{T}} + \boldsymbol{\Phi}_{\mathbf{A}})^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu}_{\mathbf{A}}) \right\} - N \ln|\mathbf{A}|. \end{split}$$

This takes the same form as the original log likelihood function apart from an additive constant  $-\ln |\mathbf{A}|$ . Thus the maximum likelihood solution in the new variables for the transformed data set will be identical to that in the old variables.

We now ask whether specific constraints on  $\Phi$  will be preserved by this re-scaling. In the case of probabilistic PCA the noise covariance  $\Phi$  is proportional to the unit matrix and takes the form  $\sigma^2 \mathbf{I}$ . For this constraint to be preserved we require  $\mathbf{A}\mathbf{A}^T = \mathbf{I}$  so that  $\mathbf{A}$  is an orthogonal matrix. This corresponds to a rotation of the coordinate system. For factor analysis  $\Phi$  is a diagonal matrix, and this property will be preserved if  $\mathbf{A}$  is also diagonal since the product of diagonal matrices is again diagonal. This corresponds to an independent re-scaling of the coordinate system. Note that in general probabilistic PCA is not invariant under component-wise re-scaling and factor analysis is not invariant under rotation. These results are illustrated in Figure 11.

12.26 If we multiply (12.80) by  $\mathbf{K}$  we obtain (12.79) so that any solution of the former will also be a solution of the latter. Let  $\widetilde{\mathbf{a}}_i$  be a solution of (12.79) with eigenvalue  $\lambda_i$  and let  $\mathbf{a}_i$  be a solution of (12.80) also having eigenvalue  $\lambda_i$ . If we write  $\widetilde{\mathbf{a}}_i = \mathbf{a}_i + \mathbf{b}_i$  we see that  $\mathbf{b}_i$  must satisfy  $\mathbf{K}\mathbf{b}_i = \mathbf{0}$  and hence is an eigenvector of  $\mathbf{K}$  with eigenvalue 0. It therefore satisfies

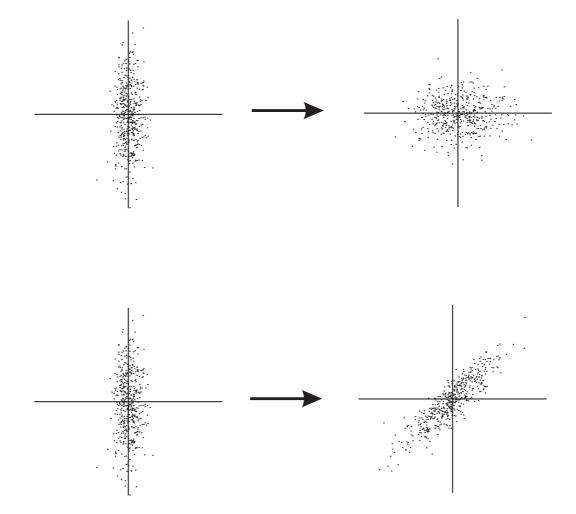
$$\sum_{n=1}^{N} b_{ni} k(\mathbf{x}_n, \mathbf{x}) = 0$$

for all values of x. Now consider the eigenvalue projection. We see that

$$\widetilde{\mathbf{v}}_{i}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}) = \sum_{n=1}^{N} \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \widetilde{a}_{ni} \boldsymbol{\phi}(\mathbf{x}_{n})$$

$$= \sum_{n=1}^{N} a_{ni} k(\mathbf{x}_{n}, \mathbf{x}) + \sum_{n=1}^{N} b_{ni} k(\mathbf{x}_{n}, \mathbf{x}) = \sum_{n=1}^{N} a_{ni} k(\mathbf{x}_{n}, \mathbf{x})$$

and so both solutions give the same projections. A slightly different treatment of the relationship between (12.79) and (12.80) is given by Schölkopf (1998).



**Figure 11** Factor analysis is covariant under a componentwise re-scaling of the data variables (top plots), while PCA and probabilistic PCA are covariant under rotations of the data space coordinates (lower plots).

# 欢迎关注公众号@机器学习与算法之道

**12.27** In the case of the linear kernel, we can rewrite the l.h.s. of (12.80) as

$$\mathbf{K}\mathbf{a}_{i} = \sum_{m=1}^{N} k(\mathbf{x}_{n}, \mathbf{x}_{m}) a_{im}$$
$$= \sum_{m=1}^{N} \mathbf{x}_{n}^{\mathrm{T}} \mathbf{x}_{m} a_{mi}$$

and substituting this in (12.80), we get

$$\sum_{m=1}^{N} \mathbf{x}_{m}^{\mathrm{T}} \mathbf{x}_{m} a_{mi} = \lambda_{i} N \mathbf{a}_{i}.$$

Next, we left-multiply both sides by  $x_n$  and sum over n to obtain

$$N\mathbf{S}\sum_{m=1}^{N}\mathbf{x}_{m}a_{im} = \lambda_{i}N\sum_{n=1}^{N}\mathbf{x}_{n}a_{in}.$$

Finally, we divide both sides by N and define

$$\mathbf{u}_i = \sum_{n=1}^{N} \mathbf{x}_n a_{in}$$

to recover (12.17).

**12.28** If we assume that the function y = f(x) is *strictly* monotonic, which is necessary to exclude the possibility for spikes of infinite density in p(y), we are guaranteed that the inverse function  $x = f^{-1}(y)$  exists. We can then use (1.27) to write

$$p(y) = q(f^{-1}(y)) \left| \frac{\mathrm{d}f^{-1}}{\mathrm{d}y} \right|.$$
 (317)

Since the only restriction on f is that it is monotonic, it can distribute the probability mass over x arbitrarily over y. This is illustrated in Figure 1 on page 9, as a part of Solution 1.4. From (317) we see directly that

$$|f'(x)| = \frac{q(x)}{p(f(x))}.$$

**12.29 NOTE**: In the 1<sup>st</sup> printing of PRML, this exercise contains two mistakes. In the second half of the exercise, we require that  $y_1$  is symmetrically distributed around 0, not just that  $-1 \le y_1 \le 1$ . Moreover,  $y_2 = y_1^2$  (not  $y_2 = y_2^2$ ).

If  $z_1$  and  $z_2$  are independent, then

$$cov[z_1, z_2] = \iint (z_1 - \bar{z}_1)(z_2 - \bar{z}_2)p(z_1, z_2) dz_1 dz_2 
= \iint (z_1 - \bar{z}_1)(z_2 - \bar{z}_2)p(z_1)p(z_2) dz_1 dz_2 
= \iint (z_1 - \bar{z}_1)p(z_1) dz_1 \int (z_2 - \bar{z}_2)p(z_2) dz_2 
= 0,$$

where

$$\bar{z}_i = \mathbb{E}[z_i] = \int z_i p(z_i) \, \mathrm{d}z_i.$$

For  $y_2$  we have

$$p(y_2|y_1) = \delta(y_2 - y_1^2),$$

i.e., a spike of probability mass one at  $y_1^2$ , which is clearly dependent on  $y_1$ . With  $\bar{y}_i$  defined analogously to  $\bar{z}_i$  above, we get

$$cov[y_1, y_2] = \iint (y_1 - \bar{y}_1)(y_2 - \bar{y}_2)p(y_1, y_2) dy_1 dy_2 
= \iint y_1(y_2 - \bar{y}_2)p(y_2|y_1)p(y_1) dy_1 dy_2 
= \iint (y_1^3 - y_1\bar{y}_2)p(y_1) dy_1 
= 0,$$

where we have used the fact that all odd moments of  $y_1$  will be zero, since it is symmetric around zero.

### **Chapter 13** Sequential Data

**13.1** Since the arrows on the path from  $x_m$  to  $x_n$ , with m < n - 1, will meet head-to-tail at  $x_{n-1}$ , which is in the conditioning set, all such paths are blocked by  $x_{n-1}$  and hence (13.3) holds.

The same argument applies in the case depicted in Figure 13.4, with the modification that m < n - 2 and that paths are blocked by  $x_{n-1}$  or  $x_{n-2}$ .

**13.2** We first of all find the joint distribution  $p(\mathbf{x}_1, \dots, \mathbf{x}_n)$  by marginalizing over the variables  $\mathbf{x}_{n+1}, \dots, \mathbf{x}_N$ , to give

$$p(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_{\mathbf{x}_{n+1}} \dots \sum_{\mathbf{x}_N} p(\mathbf{x}_1, \dots, \mathbf{x}_N)$$

$$= \sum_{\mathbf{x}_{n+1}} \dots \sum_{\mathbf{x}_N} p(\mathbf{x}_1) \prod_{m=2}^N p(\mathbf{x}_m | \mathbf{x}_{m-1})$$

$$= p(\mathbf{x}_1) \prod_{m=2}^n p(\mathbf{x}_m | \mathbf{x}_{m-1}).$$

Now we evaluate the required conditional distribution

$$p(\mathbf{x}_n|\mathbf{x}_1,\dots,\mathbf{x}_{n-1}) = \frac{p(\mathbf{x}_1,\dots,\mathbf{x}_n)}{\sum_{\mathbf{x}_n} p(\mathbf{x}_1,\dots,\mathbf{x}_n)}$$
$$= \frac{p(\mathbf{x}_1) \prod_{m=2}^n p(\mathbf{x}_m|\mathbf{x}_{m-1})}{\sum_{\mathbf{x}_n} p(\mathbf{x}_1) \prod_{m=2}^n p(\mathbf{x}_m|\mathbf{x}_{m-1})}.$$

We now note that any factors which do not depend on  $\mathbf{x}_n$  will cancel between numerator and denominator, giving

$$p(\mathbf{x}_n|\mathbf{x}_1,\dots,\mathbf{x}_{n-1}) = \frac{p(\mathbf{x}_n|\mathbf{x}_{n-1})}{\sum_{\mathbf{x}_n} p(\mathbf{x}_n|\mathbf{x}_{n-1})}$$
$$= p(\mathbf{x}_n|\mathbf{x}_{n-1})$$

as required.

For the second order Markov model, the joint distribution is given by (13.4). The marginal distribution over the variables  $\mathbf{x}_1, \dots, \mathbf{x}_n$  is given by

$$p(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_{\mathbf{x}_{n+1}} \dots \sum_{\mathbf{x}_N} p(\mathbf{x}_1, \dots, \mathbf{x}_N)$$

$$= \sum_{\mathbf{x}_{n+1}} \dots \sum_{\mathbf{x}_N} p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{m=3}^N p(\mathbf{x}_m | \mathbf{x}_{m-1}, \mathbf{x}_{m-2})$$

$$= p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{m=3}^n p(\mathbf{x}_m | \mathbf{x}_{m-1}, \mathbf{x}_{m-2}).$$

The required conditional distribution is then given by

$$p(\mathbf{x}_n|\mathbf{x}_1,\dots,\mathbf{x}_{n-1}) = \frac{p(\mathbf{x}_1,\dots,\mathbf{x}_n)}{\sum_{\mathbf{x}_n} p(\mathbf{x}_1,\dots,\mathbf{x}_n)}$$

$$= \frac{p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{x}_1) \prod_{m=3}^n p(\mathbf{x}_m|\mathbf{x}_{m-1},\mathbf{x}_{m-2})}{\sum_{\mathbf{x}_n} p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{x}_1) \prod_{m=3}^n p(\mathbf{x}_m|\mathbf{x}_{m-1},\mathbf{x}_{m-2})}.$$

Again, cancelling factors independent of  $\mathbf{x}_n$  between numerator and denominator we obtain

$$p(\mathbf{x}_n|\mathbf{x}_1,\dots,\mathbf{x}_{n-1}) = \frac{p(\mathbf{x}_n|\mathbf{x}_{n-1},\mathbf{x}_{n-2})}{\sum_{\mathbf{x}_n} p(\mathbf{x}_n|\mathbf{x}_{n-1},\mathbf{x}_{n-2})}$$
$$= p(\mathbf{x}_n|\mathbf{x}_{n-1},\mathbf{x}_{n-2}).$$

Thus the prediction at step n depends only on the observations at the two previous steps  $\mathbf{x}_{n-1}$  and  $\mathbf{x}_{n-2}$  as expected.

- **13.3** From Figure 13.5 we see that for any two variables  $\mathbf{x}_n$  and  $\mathbf{x}_m$ ,  $m \neq n$ , there is a path between the corresponding nodes that will only pass through one or more nodes corresponding to  $\mathbf{z}$  variables. None of these nodes will be in the conditioning set and the edges on the path meet head-to-tail. Thus, there will be an unblocked path between  $\mathbf{x}_n$  and  $\mathbf{x}_m$  and the model will not satisfy any conditional independence or finite order Markov properties.
- 13.4 The learning of w would follow the scheme for maximum learning described in Section 13.2.1, with w replacing  $\phi$ . As discussed towards the end of Section 13.2.1, the precise update formulae would depend on the form of regression model used and how it is being used.

The most obvious situation where this would occur is in a HMM similar to that depicted in Figure 13.18, where the emmission densities not only depends on the latent variable **z**, but also on some input variable **u**. The regression model could then be used to map **u** to **x**, depending on the state of the latent variable **z**.

Note that when a nonlinear regression model, such as a neural network, is used, the M-step for w may not have closed form.

**13.5** Consider first the maximization with respects to the components  $\pi_k$  of  $\pi$ . To do this we must take account of the summation constraint

$$\sum_{k=1}^{K} \pi_k = 1.$$

We therefore first omit terms from  $Q(\theta, \theta_{\rm old})$  which are independent of  $\pi$ , and then add a Lagrange multiplier term to enforce the constraint, giving the following function to be maximized

$$\widetilde{Q} = \sum_{k=1}^{K} \gamma(z_{1k}) \ln \pi_k + \lambda \left( \sum_{k=1}^{K} \pi_k - 1 \right).$$

Setting the derivative with respect to  $\pi_k$  equal to zero we obtain

$$0 = \gamma(z_{1k}) \frac{1}{\pi_k} + \lambda. \tag{318}$$

We now multiply through by  $\pi_k$  and then sum over k and make use of the summation constraint to give

$$\lambda = -\sum_{k=1}^{K} \gamma(z_{1k}).$$

Substituting back into (318) and solving for  $\lambda$  we obtain (13.18).

For the maximization with respect to  $\mathbf{A}$  we follow the same steps and first omit terms from  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{\text{old}})$  which are independent of  $\mathbf{A}$ , and then add appropriate Lagrange multiplier terms to enforce the summation constraints. In this case there are K constraints to be satisfied since we must have

$$\sum_{k=1}^{K} A_{jk} = 1$$

for  $j=1,\ldots,K$ . We introduce K Lagrange multipliers  $\lambda_j$  for  $j=1,\ldots,K$ , and maximize the following function

$$\widehat{Q} = \sum_{n=2}^{N} \sum_{j=1}^{K} \sum_{k=1}^{K} \xi(z_{n-1,j}, z_{nk}) \ln A_{jk} + \sum_{j=1}^{K} \lambda_j \left( \sum_{k=1}^{K} A_{jk} - 1 \right).$$

Setting the derivative of  $\widehat{Q}$  with respect to  $A_{jk}$  to zero we obtain

$$0 = \sum_{n=2}^{N} \xi(z_{n-1,j}, z_{nk}) \frac{1}{A_{jk}} + \lambda_j.$$
 (319)

Again we multiply through by  $A_{jk}$  and then sum over k and make use of the summation constraint to give

$$\lambda_j = -\sum_{n=2}^{N} \sum_{k=1}^{K} \xi(z_{n-1,j}, z_{nk}).$$

Substituting for  $\lambda_j$  in (319) and solving for  $A_{jk}$  we obtain (13.19).

**13.6** Suppose that a particular element  $\pi_k$  of  $\pi$  has been initialized to zero. In the first E-step the quantity  $\alpha(z_{1k})$  is given from (13.37) by

$$\alpha(z_{1k}) = \pi_k p(\mathbf{x}_1 | \boldsymbol{\phi}_k)$$

and so will be zero. From (13.33) we see that  $\gamma(z_{1k})$  will also be zero, and hence in the next M-step the new value of  $\pi_k$ , given by (13.18) will again be zero. Since this is true for any subsequent EM cycle, this quantity will remain zero throughout.

Similarly, suppose that an element  $A_{jk}$  of **A** has been set initially to zero. From (13.43) we see that  $\xi(z_{n-1,j},z_{nk})$  will be zero since  $p(z_{nk}|z_{n-1,j})\equiv A_{jk}$  equals zero. In the subsequent M-step, the new value of  $A_{jk}$  is given by (13.19) and hence will also be zero.

13.7 Using the expression (13.17) for  $Q(\theta, \theta_{\text{old}})$  we see that the parameters of the Gaussian emission densities appear only in the last term, which takes the form

$$\begin{split} & \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \ln p(\mathbf{x}_n | \boldsymbol{\phi}_k) = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_j) \\ & = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \left\{ -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^{\mathrm{T}} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right\}. \end{split}$$

We now maximize this quantity with respect to  $\mu_k$  and  $\Sigma_k$ . Setting the derivative with respect to  $\mu_k$  to zero and re-arranging we obtain (13.20). Next if we define

$$N_k = \sum_{n=1}^{N} \gamma(z_{nk})$$

$$\widehat{\mathbf{S}}_k = \sum_{n=1}^{N} \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^{\mathrm{T}}$$

then we can rewrite the final term from  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{\mathrm{old}})$  in the form

$$-\frac{N_k D}{2} \ln(2\pi) - \frac{N_k}{2} \ln|\mathbf{\Sigma}_k| - \frac{1}{2} \text{Tr}\left(\mathbf{\Sigma}_k^{-1} \widehat{\mathbf{S}}_k\right).$$

Differentiating this w.r.t.  $\Sigma_k^{-1}$ , using results from Appendix C, we obtain (13.21).

**13.8** Only the final term of  $Q(\theta, \theta^{\text{old}})$  given by (13.17) depends on the parameters of the emission model. For the multinomial variable  $\mathbf{x}$ , whose D components are all zero except for a single entry of 1,

$$\sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \ln p(\mathbf{x}_n | \phi_k) = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \sum_{i=1}^{D} x_{ni} \ln \mu_{ki}.$$

Now when we maximize with respect to  $\mu_{ki}$  we have to take account of the constraints that, for each value of k the components of  $\mu_{ki}$  must sum to one. We therefore introduce Lagrange multipliers  $\{\lambda_k\}$  and maximize the modified function given by

$$\sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \sum_{i=1}^{D} x_{ni} \ln \mu_{ki} + \sum_{k=1}^{K} \lambda_k \left( \sum_{i=1}^{D} \mu_{ki} - 1 \right).$$

Setting the derivative with respect to  $\mu_{ki}$  to zero we obtain

$$0 = \sum_{n=1}^{N} \gamma(z_{nk}) \frac{x_{ni}}{\mu_{ki}} + \lambda_k.$$

Multiplying through by  $\mu_{ki}$ , summing over i, and making use of the constraint on  $\mu_{ki}$  together with the result  $\sum_i x_{ni} = 1$  we have

$$\lambda_k = -\sum_{n=1}^N \gamma(z_{nk}).$$

Finally, back-substituting for  $\lambda_k$  and solving for  $\mu_{ki}$  we again obtain (13.23).

Similarly, for the case of a multivariate Bernoulli observed variable  ${\bf x}$  whose D components independently take the value 0 or 1, using the standard expression for the multivariate Bernoulli distribution we have

$$\sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \ln p(\mathbf{x}_n | \boldsymbol{\phi}_k)$$

$$= \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}) \sum_{i=1}^{D} \left\{ x_{ni} \ln \mu_{ki} + (1 - x_{ni}) \ln(1 - \mu_{ki}) \right\}.$$

Maximizing with respect to  $\mu_{ki}$  we obtain

$$\mu_{ki} = \frac{\sum_{n=1}^{N} \gamma(z_{nk}) x_{ni}}{\sum_{n=1}^{N} \gamma(z_{nk})}$$

which is equivalent to (13.23).

**13.9** We can verify all these independence properties using d-separation by refering to Figure 13.5.

(13.24) follows from the fact that arrows on paths from any of  $\mathbf{x}_1, \dots, \mathbf{x}_n$  to any of  $\mathbf{x}_{n+1}, \dots, \mathbf{x}_N$  meet head-to-tail or tail-to-tail at  $\mathbf{z}_n$ , which is in the conditioning set.

- (13.25) follows from the fact that arrows on paths from any of  $\mathbf{x}_1, \dots, \mathbf{x}_{n-1}$  to  $\mathbf{x}_n$  meet head-to-tail at  $\mathbf{z}_n$ , which is in the conditioning set.
- (13.26) follows from the fact that arrows on paths from any of  $\mathbf{x}_1, \dots, \mathbf{x}_{n-1}$  to  $\mathbf{z}_n$  meet head-to-tail or tail-to-tail at  $\mathbf{z}_{n-1}$ , which is in the conditioning set.
- (13.27) follows from the fact that arrows on paths from  $\mathbf{z}_n$  to any of  $\mathbf{x}_{n+1}, \dots, \mathbf{x}_N$  meet head-to-tail at  $\mathbf{z}_{n+1}$ , which is in the conditioning set.
- (13.28) follows from the fact that arrows on paths from  $\mathbf{x}_{n+1}$  to any of  $\mathbf{x}_{n+2}, \dots, \mathbf{x}_N$  to meet tail-to-tail at  $\mathbf{z}_{n+1}$ , which is in the conditioning set.
- (13.29) follows from (13.24) and the fact that arrows on paths from any of  $\mathbf{x}_1, \ldots, \mathbf{x}_{n-1}$  to  $\mathbf{x}_n$  meet head-to-tail or tail-to-tail at  $\mathbf{z}_{n-1}$ , which is in the conditioning set.
- (13.30) follows from the fact that arrows on paths from any of  $\mathbf{x}_1, \dots, \mathbf{x}_N$  to  $\mathbf{x}_{N+1}$  meet head-to-tail at  $\mathbf{z}_{N+1}$ , which is in the conditioning set.
- (13.31) follows from the fact that arrows on paths from any of  $\mathbf{x}_1, \dots, \mathbf{x}_N$  to  $\mathbf{z}_{N+1}$  meet head-to-tail or tail-to-tail at  $\mathbf{z}_N$ , which is in the conditioning set.
- **13.10** We begin with the expression (13.10) for the joint distribution of observed and latent variables in the hidden Markov model, reproduced here for convenience

$$p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = p(\mathbf{z}_1) \left[ \prod_{n=2}^{N} p(\mathbf{z}_n | \mathbf{z}_{n-1}) \right] \prod_{m=1}^{N} p(\mathbf{x}_m | \mathbf{z}_m)$$

where we have omitted the parameters in order to keep the notation uncluttered. By marginalizing over all of the latent variables except  $\mathbf{z}_n$  we obtain the joint distribution of the remaining variables, which can be factorized in the form

$$p(\mathbf{X}, \mathbf{z}_n) = \sum_{\mathbf{z}_1} \cdots \sum_{\mathbf{z}_{n-1}} \sum_{\mathbf{z}_{n+1}} \cdots \sum_{\mathbf{z}_N} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$$

$$= \left[ \sum_{\mathbf{z}_1} \cdots \sum_{\mathbf{z}_{n-1}} p(\mathbf{z}_1) \prod_{m=2}^n p(\mathbf{z}_m | \mathbf{z}_{m-1}) \prod_{l=1}^n p(\mathbf{x}_l | \mathbf{z}_l) \right]$$

$$\times \left[ \sum_{\mathbf{z}_{n+1}} \cdots \sum_{\mathbf{z}_N} \prod_{m=n+1}^N p(\mathbf{z}_m | \mathbf{z}_{m-1}) \prod_{l=n+1}^N p(\mathbf{x}_l | \mathbf{z}_l) \right].$$

The first factor in square brackets on the r.h.s. we recognize as  $p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n)$ . Next we note from the product rule that

$$p(\mathbf{x}_{n+1},\ldots,\mathbf{x}_N|\mathbf{x}_1,\ldots,\mathbf{x}_n,\mathbf{z}_n) = \frac{p(\mathbf{x}_1,\ldots,\mathbf{x}_N,\mathbf{z}_n)}{p(\mathbf{x}_1,\ldots,\mathbf{x}_n,\mathbf{z}_n)}.$$

Thus we can identify the second term in square brackets with the conditional distribution  $p(\mathbf{x}_{n+1},\ldots,\mathbf{x}_N|\mathbf{x}_1,\ldots,\mathbf{x}_n,\mathbf{z}_n)$ . However, we note that the second term in square brackets does not depend on  $\mathbf{x}_1,\ldots,\mathbf{x}_n$ . Thus we have the result

$$p(\mathbf{x}_1,\ldots,\mathbf{x}_N,\mathbf{z}_n)=p(\mathbf{x}_1,\ldots,\mathbf{x}_n,\mathbf{z}_n)p(\mathbf{x}_{n+1},\ldots,\mathbf{x}_N|\mathbf{z}_n).$$

Dividing both sides by  $p(\mathbf{z}_n)$  we obtain

$$p(\mathbf{x}_1,\ldots,\mathbf{x}_N|\mathbf{z}_n) = p(\mathbf{x}_1,\ldots,\mathbf{x}_n|\mathbf{z}_n)p(\mathbf{x}_{n+1},\ldots,\mathbf{x}_N|\mathbf{z}_n).$$

which is the required result (13.24).

Similarly, from (13.10) we have

$$p(\mathbf{x}_1,\ldots,\mathbf{x}_n,\mathbf{z}_1,\ldots,\mathbf{z}_n) = p(\mathbf{x}_1,\ldots,\mathbf{x}_{n-1},\mathbf{z}_1,\ldots,\mathbf{z}_{n-1})p(\mathbf{z}_n|\mathbf{z}_{n-1})p(\mathbf{x}_n|\mathbf{z}_n)$$

It follows that

$$p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_1, \dots, \mathbf{z}_{n-1} | \mathbf{x}_n, \mathbf{z}_n) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_1, \dots, \mathbf{z}_n)}{p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n)}$$
$$= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_1, \dots, \mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})}{p(\mathbf{z}_n)}$$

where we see that the right hand side is independent of  $\mathbf{x}_n$ , and hence the left hand side must be also. We therefore have the following conditional independence property

$$p(\mathbf{x}_1,\ldots,\mathbf{x}_{n-1},\mathbf{z}_1,\ldots,\mathbf{z}_{n-1}|\mathbf{x}_n,\mathbf{z}_n)=p(\mathbf{x}_1,\ldots,\mathbf{x}_{n-1},\mathbf{z}_1,\ldots,\mathbf{z}_{n-1}|\mathbf{z}_n).$$

Marginalizing both sides of this result over  $\mathbf{z}_1, \dots, \mathbf{z}_{n-1}$  then gives the required result (13.25).

Again, from the joint distribution we can write

$$p(\mathbf{x}_1,\ldots,\mathbf{x}_{n-1},\mathbf{z}_1,\ldots,\mathbf{z}_n)=p(\mathbf{x}_1,\ldots,\mathbf{x}_{n-1},\mathbf{z}_1,\ldots,\mathbf{z}_{n-1})p(\mathbf{z}_n|\mathbf{z}_{n-1}).$$

We therefore have

$$p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_1, \dots, \mathbf{z}_{n-2} | \mathbf{z}_{n-1}, \mathbf{z}_n) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_1, \dots, \mathbf{z}_n)}{p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{z}_{n-1})}$$
$$= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_1, \dots, \mathbf{z}_{n-1})}{p(\mathbf{z}_{n-1})}$$

where we see that the right hand side is independent of  $\mathbf{z}_n$  and hence the left hand side must be also. This implies the conditional independence property

$$p(\mathbf{x}_1,\ldots,\mathbf{x}_{n-1},\mathbf{z}_1,\ldots,\mathbf{z}_{n-2}|\mathbf{z}_{n-1},\mathbf{z}_n) = p(\mathbf{x}_1,\ldots,\mathbf{x}_{n-1},\mathbf{z}_1,\ldots,\mathbf{z}_{n-2}|\mathbf{z}_{n-1}).$$

Marginalizing both sides with respect to  $\mathbf{z}_1, \dots, \mathbf{z}_{n-2}$  then gives the required result (13.26).

To prove (13.27) we marginalize the both sides of the expression (13.10) for the joint distribution with respect to the variables  $\mathbf{x}_1, \dots, \mathbf{x}_n$  to give

$$p(\mathbf{x}_{n+1},\dots,\mathbf{x}_N,\mathbf{z}_n,\mathbf{z}_{n+1}) = \left[\sum_{\mathbf{z}_1}\dots\sum_{\mathbf{z}_{n-1}}p(\mathbf{z}_1)\prod_{m=2}^n p(\mathbf{z}_m|\mathbf{z}_{m-1})\prod_{l=1}^n p(\mathbf{x}_l|\mathbf{z}_l)\right]$$
$$p(\mathbf{z}_{n+1}|\mathbf{z}_n)p(\mathbf{x}_{n+1}|\mathbf{z}_{n+1})\left[\sum_{\mathbf{z}_{n+1}}\dots\sum_{\mathbf{z}_N}\prod_{m=n+2}^N p(\mathbf{z}_m|\mathbf{z}_{m-1})\prod_{l=n+1}^N p(\mathbf{x}_l|\mathbf{z}_l)\right].$$

The first factor in square brackets is just

$$\sum_{\mathbf{z}_1} \cdots \sum_{\mathbf{z}_{n-1}} p(\mathbf{z}_1, \dots, \mathbf{z}_n) = p(\mathbf{z}_n)$$

and so by the product rule the final factor in square brackets must be

$$p(\mathbf{x}_{n+1},\ldots,\mathbf{x}_N|\mathbf{z}_n,\mathbf{z}_{n+1}).$$

However, the second factor in square brackets is itself independent of  $\mathbf{z}_n$  which proves (13.27).

To prove (13.28) we first note that the decomposition of the joint distribution  $p(\mathbf{X}, \mathbf{Z})$  implies the following factorization

$$p(\mathbf{X}, \mathbf{Z}) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_1, \dots, \mathbf{z}_n) p(\mathbf{z}_{n+1} | \mathbf{z}_n) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1})$$
$$p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N, \mathbf{z}_{n+1}, \dots, \mathbf{z}_N | \mathbf{z}_{n+1}).$$

Next we make use of

$$p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, \mathbf{z}_{n+1}) = \sum_{\mathbf{x}_1} \dots \sum_{\mathbf{x}_n} \sum_{\mathbf{z}_1} \dots \sum_{\mathbf{z}_n} \sum_{\mathbf{z}_{n+2}} \dots \sum_{\mathbf{z}_N} p(\mathbf{X}, \mathbf{Z})$$
$$= p(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}).$$

If we now divide both sides by  $p(\mathbf{x}_{n+1}, \mathbf{z}_{n+1})$  we obtain

$$p(\mathbf{x}_{n+2},\ldots,\mathbf{x}_N|\mathbf{z}_{n+1},\mathbf{x}_{n+1})=p(\mathbf{x}_{n+2},\ldots,\mathbf{x}_N|\mathbf{z}_{n+1})$$

as required.

To prove (13.30) we first use the expression for the joint distribution of  $\mathbf{X}$  and  $\mathbf{Z}$  to give

$$p(\mathbf{x}_{N+1}, \mathbf{X}, \mathbf{z}_{N+1}) = \sum_{\mathbf{z}_1} \cdots \sum_{\mathbf{z}_N} p(\mathbf{X}, \mathbf{Z}) p(\mathbf{z}_{N+1} | \mathbf{z}_N) p(\mathbf{x}_{N+1} | \mathbf{z}_{N+1})$$
$$= p(\mathbf{X}, \mathbf{z}_{N+1}) p(\mathbf{x}_{N+1} | \mathbf{z}_{N+1})$$

from which it follows that

$$p(\mathbf{x}_{N+1}|\mathbf{X},\mathbf{z}_{N+1}) = p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1})$$

as required.

To prove (13.31) we first use the expression for the joint distribution of  ${\bf X}$  and  ${\bf Z}$  to give

欢迎关注公众等。机器等之事算法之道
$$p(\mathbf{X},\mathbf{Z}_N)p(\mathbf{z}_{N+1}|\mathbf{z}_N)$$

from which it follows that

$$p(\mathbf{z}_{N+1}|\mathbf{z}_N,\mathbf{X}) = p(\mathbf{z}_{N+1}|\mathbf{z}_N)$$

as required.

Finally, to prove (13.29) we first marginalize both sides of the joint distribution (13.10) with respect to  $\mathbf{z}_1, \dots, \mathbf{z}_{n-2}, \mathbf{z}_{n+1}, \dots \mathbf{z}_N$  to give

$$p(\mathbf{X}, \mathbf{z}_{n-1}, \mathbf{z}_n) = \left[ \sum_{\mathbf{z}_1} \cdots \sum_{\mathbf{z}_{n-2}} p(\mathbf{z}_1) \prod_{m=2}^{n-1} p(\mathbf{z}_m | \mathbf{z}_{m-1}) \prod_{l=1}^{n-1} p(\mathbf{x}_l | \mathbf{z}_l) \right]$$

$$p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n)$$

$$\left[ \sum_{\mathbf{z}_{n+1}} \cdots \sum_{\mathbf{z}_N} \prod_{m=n+1}^{N} p(\mathbf{z}_m | \mathbf{z}_{m-1}) \prod_{l=n+1}^{N} p(\mathbf{x}_l | \mathbf{z}_l) \right].$$

The first factor in square brackets is  $p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_{n-1})$ . The second factor in square brackets is

$$\sum_{\mathbf{z}_{n+1}} \cdots \sum_{\mathbf{z}_N} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, \mathbf{z}_n, \dots, \mathbf{z}_N) = p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, \mathbf{z}_n).$$

Thus we have

$$p(\mathbf{X}, \mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_{n-1})$$
$$p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, \mathbf{z}_n)$$

and dividing both sides by  $p(\mathbf{z}_n, \mathbf{z}_{n-1}) = p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{z}_{n-1})$  we obtain (13.28).

The final conclusion from all of this exhausting algebra is that it is much easier simply to draw a graph and apply the d-separation criterion!

#### **13.11** From the first line of (13.43), we have

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}).$$

This corresponds to the distribution over the variables associated with factor  $f_n$  in Figure 13.5, i.e.  $\mathbf{z}_{n-1}$  and  $\mathbf{z}_n$ .

From (8.69), (8.72), (13.50) and (13.52), we have

$$p(\mathbf{z}_{n-1}, \mathbf{z}_n) \propto f_n(\mathbf{z}_{n-1}, \mathbf{z}_n, \mathbf{x}_n) \mu_{\mathbf{z}_{n-1} \to f_n}(\mathbf{z}_{n-1}) \mu_{\mathbf{z}_n \to f_n}(\mathbf{z}_n)$$

$$= p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) \mu_{f_{n-1} \to \mathbf{z}_{n-1}}(\mathbf{z}_{n-1}) \mu_{f_{n+1} \to \mathbf{z}_n}(\mathbf{z}_n)$$

$$= p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) \alpha(\mathbf{z}_{n-1}) \beta(\mathbf{z}_n). \tag{320}$$

In order to normalize this, we use (13.36) and (13.41) to obtain

$$\sum_{\mathbf{z}_n} \sum_{\mathbf{z}_{n-1}} p(\mathbf{z}_{n-1}, \mathbf{z}_n) = \sum_{\mathbf{z}_n} \beta(\mathbf{z}_n) p(\mathbf{x}_n | \mathbf{z}_n \sum_{\mathbf{z}_{n-1}} p(\mathbf{z}_n | \mathbf{z}_{n-1}) \alpha(\mathbf{z}_n)$$
$$= \sum_{\mathbf{z}_n} \beta(\mathbf{z}_n) \alpha(\mathbf{z}_n) = p(\mathbf{X})$$

which together with (320) give (13.43).

**13.12** First of all, note that for every observed variable there is a corresponding latent variable, and so for every sequence  $\mathbf{X}^{(r)}$  of observed variables there is a corresponding sequence  $\mathbf{Z}^{(r)}$  of latent variables. The sequences are assumed to be independent given the model parameters, and so the joint distribution of all latent and observed variables will be given by

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) = \prod_{r=1}^{R} p(\mathbf{X}^{(r)}, \mathbf{Z}^{(r)} | \boldsymbol{\theta})$$

where X denotes  $\{X^{(r)}\}$  and Z denotes  $\{Z^{(r)}\}$ . Using the sum and product rules of probability we then see that posterior distribution for the latent sequences then factorizes with respect to those sequences, so that

$$p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) = \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{\sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}$$

$$= \frac{\prod_{r=1}^{R} p(\mathbf{X}^{(r)}, \mathbf{Z}^{(r)}|\boldsymbol{\theta})}{\sum_{\mathbf{Z}^{(1)}} \cdots \sum_{\mathbf{Z}^{(R)}} \prod_{r=1}^{R} p(\mathbf{X}^{(r)}, \mathbf{Z}^{(r)}|\boldsymbol{\theta})}$$

$$= \prod_{r=1}^{R} \left\{ \frac{p(\mathbf{X}^{(r)}, \mathbf{Z}^{(r)}|\boldsymbol{\theta})}{\sum_{\mathbf{Z}^{(r)}} p(\mathbf{X}^{(r)}, \mathbf{Z}^{(r)}|\boldsymbol{\theta})} \right\}$$

$$= \prod_{r=1}^{R} p(\mathbf{Z}^{(r)}|\mathbf{X}^{(r)}, \boldsymbol{\theta}).$$

Thus the evaluation of the posterior distribution of the latent variables, corresponding to the E-step of the EM algorithm, can be done independently for each of the sequences (using the standard alpha-beta recursions).

Now consider the M-step. We use the posterior distribution computed in the E-step using  $\theta_{\rm old}$  to evaluate the expectation of the complete-data log likelihood. From our

expression for the joint distribution we see that this is given by

$$Q(\theta, \theta_{\text{old}}) = \mathbb{E}_{\mathbf{Z}} \left[ \ln p(\mathbf{X}, \mathbf{Z} | \theta) \right]$$

$$= \mathbb{E}_{\mathbf{Z}} \left[ \sum_{r=1}^{R} \ln p(\mathbf{X}^{(r)}, \mathbf{Z}^{(r)} | \theta) \right]$$

$$= \sum_{r=1}^{R} p(\mathbf{Z}^{(r)} | \mathbf{X}^{(r)}, \theta_{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \theta)$$

$$= \sum_{r=1}^{R} \sum_{k=1}^{K} \gamma(z_{1k}^{(r)}) \ln \pi_{k} + \sum_{r=1}^{R} \sum_{n=2}^{N} \sum_{j=1}^{K} \sum_{k=1}^{K} \xi(z_{n-1,j}^{(r)}, z_{n,k}^{(r)}) \ln A_{jk}$$

$$+ \sum_{r=1}^{R} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma(z_{nk}^{(r)}) \ln p(\mathbf{x}_{n}^{(r)} | \phi_{k}).$$

We now maximize this quantity with respect to  $\pi$  and A in the usual way, with Lagrange multipliers to take account of the summation constraints (see Solution 13.5), yielding (13.124) and (13.125). The M-step results for the mean of the Gaussian follow in the usual way also (see Solution 13.7).

**13.13** Using (8.64), we can rewrite (13.50) as

$$\alpha(\mathbf{z}_n) = \sum_{\mathbf{z}_1, \dots, \mathbf{z}_{n-1}} F_n(\mathbf{z}_n, {\{\mathbf{z}_1, \dots, \mathbf{z}_{n-1}\}}), \tag{321}$$

where  $F_n(\cdot)$  is the product of all factors connected to  $\mathbf{z}_n$  via  $f_n$ , including  $f_n$  itself (see Figure 13.15), so that

$$F_n(\mathbf{z}_n, {\{\mathbf{z}_1, \dots, \mathbf{z}_{n-1}\}}) = h(\mathbf{z}_1) \prod_{i=2}^n f_i(\mathbf{z}_i, \mathbf{z}_{i-1}),$$
 (322)

where we have introduced  $h(\mathbf{z}_1)$  and  $f_i(\mathbf{z}_i, \mathbf{z}_{i-1})$  from (13.45) and (13.46), respectively. Using the corresponding r.h.s. definitions and repeatedly applying the product rule, we can rewrite (322) as

$$F_n(\mathbf{z}_n, {\mathbf{z}_1, \dots, \mathbf{z}_{n-1}}) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_1, \dots, \mathbf{z}_n).$$

Applying the sum rule, summing over  $\mathbf{z}_1, \dots, \mathbf{z}_{n-1}$  as on the r.h.s. of (321), we obtain (13.34).

**13.14 NOTE**: In PRML, the reference to (8.67) should refer to (8.64).

This solution largely follows Solution 13.13. Using (8.64), we can rewrite (13.52) as

$$\beta(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1},\dots,\mathbf{z}_N} F_{n+1}(\mathbf{z}_n, \{\mathbf{z}_{n+1},\dots,\mathbf{z}_N\}), \tag{323}$$

where  $F_{n+1}(\cdot)$  is the product of all factors connected to  $\mathbf{z}_n$  via  $f_{n+1}$ , including  $f_{n+1}$  itself so that

$$F_{n+1}(\mathbf{z}_n, {\{\mathbf{z}_{n+1}, \dots, \mathbf{z}_N\}}) = \prod_{i=n+1}^{N} f_i(\mathbf{z}_{i-1}, \mathbf{z}_i)$$

$$= p(\mathbf{z}_{n+1}|\mathbf{z}_n)p(\mathbf{x}_{n+1}|\mathbf{z}_{n+1}) \prod_{i=n+2}^{N} f_i(\mathbf{z}_{i-1}, \mathbf{z}_i)$$

$$= p(\mathbf{z}_{n+1}|\mathbf{z}_n)p(\mathbf{x}_{n+1}|\mathbf{z}_{n+1}) \cdots p(\mathbf{z}_N|\mathbf{z}_{N-1})p(\mathbf{x}_N|\mathbf{z}_N)$$
(324)

where we have used (13.46). Repeatedly applying the product rule, we can rewrite (324) as

$$F_{n+1}(\mathbf{z}_n, {\mathbf{z}_{n+1}, \dots, \mathbf{z}_N}) = p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, \mathbf{z}_{n+1}, \dots, \mathbf{z}_N | \mathbf{z}_n).$$

Substituting this into (323) and summing over  $\mathbf{z}_{n+1}, \dots, \mathbf{z}_N$ , we obtain (13.35).

**13.15** NOTE: In the 1<sup>st</sup> printing of PRML, there are typographic errors in (13.65);  $c_n$  should be  $c_n^{-1}$  and  $p(\mathbf{z}_n|\mathbf{z}_{-1})$  should be  $p(\mathbf{z}_n|\mathbf{z}_{n-1})$  on the r.h.s.

We can use (13.58), (13.60) and (13.63) to rewrite (13.33) as

$$\gamma(\mathbf{z}_n) = \frac{\alpha(\mathbf{z}_n)\beta(\mathbf{z}_n)}{p(\mathbf{X})} \\
= \frac{\widehat{\alpha}(\mathbf{z}_n) \left(\prod_{m=1}^n c_m\right) \left(\prod_{l=n+1}^N c_l\right) \widehat{\beta}(\mathbf{z}_n)}{p(\mathbf{X})} \\
= \frac{\widehat{\alpha}(\mathbf{z}_n) \left(\prod_{m=1}^N c_m\right) \widehat{\beta}(\mathbf{z}_n)}{p(\mathbf{X})} \\
= \widehat{\alpha}(\mathbf{z}_n) \widehat{\beta}(\mathbf{z}_n).$$

We can rewrite (13.43) in a similar fashion:

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = \frac{\alpha(\mathbf{z}_n)p(\mathbf{x}_n|\mathbf{z}_n)p(\mathbf{z}_n|\mathbf{z}_{n-1})\beta(\mathbf{z}_n)}{p(\mathbf{X})}$$

$$= \frac{\widehat{\alpha}(\mathbf{z}_n)\left(\prod_{m=1}^{n-1} c_m\right)\left(\prod_{l=n+1}^{N} c_l\right)\widehat{\beta}(\mathbf{z}_n)p(\mathbf{x}_n|\mathbf{z}_n)p(\mathbf{z}_n|\mathbf{z}_{n-1})}{p(\mathbf{X})}$$

$$= c_n^{-1}\widehat{\alpha}(\mathbf{z}_n)p(\mathbf{x}_n|\mathbf{z}_n)p(\mathbf{z}_n|\mathbf{z}_{n-1})\widehat{\beta}(\mathbf{z}_n).$$

**13.16** NOTE: In the 1<sup>st</sup> printing of PRML,  $\ln p(\mathbf{x}_{+1}|\mathbf{z}_n)$  should be  $\ln p(\mathbf{z}_{n+1}|\mathbf{z}_n)$  on the r.h.s. of (13.68) Moreover  $p(\ldots)$  should be  $\ln p(\ldots)$  on the r.h.s. of (13.70). We start by rewriting (13.6) as

$$p(\mathbf{x}_1,\ldots,\mathbf{x}_N,\mathbf{z}_1,\ldots,\mathbf{z}_N)=p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1)\prod_{n=2}^N p(\mathbf{x}_n|\mathbf{z}_n)p(\mathbf{z}_n|\mathbf{z}_{n-1}).$$

Taking the logarithm we get

$$\ln p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{z}_1, \dots, \mathbf{z}_N)$$

$$= \ln p(\mathbf{z}_1) + \ln p(\mathbf{x}_1 | \mathbf{z}_1) + \sum_{n=2}^{N} (\ln p(\mathbf{x}_n | \mathbf{z}_n) + \ln p(\mathbf{z}_n | \mathbf{z}_{n-1}))$$

where, with the first two terms we have recovered the r.h.s. of (13.69). We now use this to maximize over  $\mathbf{z}_1, \dots, \mathbf{z}_N$ ,

$$\max_{\mathbf{z}_{1},\dots,\mathbf{z}_{N}} \left\{ \omega(\mathbf{z}_{1}) + \sum_{n=2}^{N} \left[ \ln p(\mathbf{x}_{n}|\mathbf{z}_{n}) + \ln p(\mathbf{z}_{n}|\mathbf{z}_{n-1}) \right] \right\}$$

$$= \max_{\mathbf{z}_{2},\dots,\mathbf{z}_{N}} \left\{ \ln p(\mathbf{x}_{2}|\mathbf{z}_{2}) + \max_{\mathbf{z}_{1}} \left\{ \ln p(\mathbf{z}_{2}|\mathbf{z}_{1}) + \omega(\mathbf{z}_{1}) \right\} + \sum_{n=3}^{N} \left[ \ln p(\mathbf{x}_{n}|\mathbf{z}_{n}) + \ln p(\mathbf{z}_{n}|\mathbf{z}_{n-1}) \right] \right\}$$

$$= \max_{\mathbf{z}_{2},\dots,\mathbf{z}_{N}} \left\{ \omega(\mathbf{z}_{2}) + \sum_{n=3}^{N} \left[ \ln p(\mathbf{x}_{n}|\mathbf{z}_{n}) + \ln p(\mathbf{z}_{n}|\mathbf{z}_{n-1}) \right] \right\}$$
(325)

where we have exchanged the order of maximization and summation for  $\mathbf{z}_2$  to recover (13.68) for n=2, and since the first and the last line of (325) have identical forms, this extends recursively to all n>2.

**13.17** The emission probabilities over observed variables  $\mathbf{x}_n$  are absorbed into the corresponding factors,  $f_n$ , analogously to the way in which Figure 13.14 was transformed into Figure 13.15. The factors then take the form

$$h(\mathbf{z}_1) = p(\mathbf{z}_1|\mathbf{u}_1)p(\mathbf{x}_1|\mathbf{z}_1,\mathbf{u}_1) \tag{326}$$

$$f_n(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{u}_n) p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{u}_n). \tag{327}$$

**13.18** By combining the results from Solution 13.17 with those from Section 13.2.3, the desired outcome is easily obtained.

By combining (327) with (13.49) and (13.50), we see that

$$\alpha(\mathbf{z}_n) = \sum_{\mathbf{z}_{n-1}} p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{u}_n) p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{u}_n) \alpha(\mathbf{z}_{n-1})$$

corresponding to (13.36). The initial condition is given directly by (326) and corresponds to (13.37).

Similarly, from (327), (13.51) and (13.52), we see that

$$\beta(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} p(\mathbf{z}_{n+1}|\mathbf{z}_n, \mathbf{u}_{n+1}) p(\mathbf{x}_{n+1}|\mathbf{z}_{n+1}, \mathbf{u}_{n+1}) \beta(\mathbf{z}_{n+1})$$

which corresponds to (13.38). The presence of the input variables does not affect the initial condition  $\beta(\mathbf{z}_N) = 1$ .

- 13.19 Since the joint distribution over all variables, latent and observed, is Gaussian, we can maximize w.r.t. any chosen set of variables. In particular, we can maximize w.r.t. all the latent variables jointly or maximize each of the marginal distributions separately. However, from (2.98), we see that the resulting means will be the same in both cases and since the mean and the mode coincide for the Gaussian, maximizing w.r.t. to latent variables jointly and individually will yield the same result.
- **13.20** Making the following substitions from the l.h.s. of (13.87),

$$\mathbf{x} \Rightarrow \mathbf{z}_{n-1} \quad \boldsymbol{\mu} \Rightarrow \boldsymbol{\mu}_{n-1} \quad \boldsymbol{\Lambda}^{-1} \Rightarrow \mathbf{V}_{n-1}$$

$$\mathbf{y} \Rightarrow \mathbf{z}_n \quad \mathbf{A} \Rightarrow \mathbf{A} \quad \mathbf{b} \Rightarrow \mathbf{0} \quad \mathbf{L}^{-1} \Rightarrow \mathbf{\Gamma},$$

in (2.113) and (2.114), (2.115) becomes

$$p(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \mathbf{A}\boldsymbol{\mu}_{n-1}, \boldsymbol{\Gamma} + \mathbf{A}\mathbf{V}_{n-1}\mathbf{A}^{\mathrm{T}}),$$

as desired.

**13.21** If we substitute the r.h.s. of (13.87) for the integral on the r.h.s. of (13.86), we get

$$c_n \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_n, \mathbf{V}_n) = \mathcal{N}(\mathbf{x}_n | \mathbf{C}\mathbf{z}_n, \boldsymbol{\Sigma}) \mathcal{N}(\mathbf{z}_n | \mathbf{A}\boldsymbol{\mu}_{n-1}, \mathbf{P}_{n-1}).$$

The r.h.s. define the joint probability distribution over  $\mathbf{x}_n$  and  $\mathbf{z}_n$  in terms of a conditional distribution over  $\mathbf{x}_n$  given  $\mathbf{z}_n$  and a distribution over  $\mathbf{z}_n$ , corresponding to (2.114) and (2.113), respectively. What we need to do is to rewrite this into a conditional distribution over  $\mathbf{z}_n$  given  $\mathbf{x}_n$  and a distribution over  $\mathbf{x}_n$ , corresponding to (2.116) and (2.115), respectively.

If we make the substitutions

$$\mathbf{x} \Rightarrow \mathbf{z}_n \quad \boldsymbol{\mu} \Rightarrow \mathbf{A} \boldsymbol{\mu}_{n-1} \quad \boldsymbol{\Lambda}^{-1} \Rightarrow \mathbf{P}_{n-1}$$

$$\mathbf{y} \Rightarrow \mathbf{x}_n \quad \mathbf{A} \Rightarrow \mathbf{C} \quad \mathbf{b} \Rightarrow \mathbf{0} \quad \mathbf{L}^{-1} \Rightarrow \mathbf{\Sigma},$$

in (2.113) and (2.114), (2.115) directly gives us the r.h.s. of (13.91).

From (2.114), we have that

$$p(\mathbf{z}_n|\mathbf{x}_n) = \mathcal{N}(\mathbf{z}_n|\boldsymbol{\mu}_n, \mathbf{V}_n) = \mathcal{N}(\mathbf{z}_n|\mathbf{M}(\mathbf{C}^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}\mathbf{x}_n + \mathbf{P}_{n-1}^{-1}\mathbf{A}\boldsymbol{\mu}_{n-1}), \mathbf{M}), (328)$$

where we have used (2.117) to define

$$\mathbf{M} = (\mathbf{P}_{n-1} + \mathbf{C}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \mathbf{C})^{-1}. \tag{329}$$

Using (C.7) and (13.92), we can rewrite (329) as follows:

$$\begin{split} \mathbf{M} &= & (\mathbf{P}_{n-1} + \mathbf{C}^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \mathbf{C})^{-1} \\ &= & \mathbf{P}_{n-1} - \mathbf{P}_{n-1} \mathbf{C}^{\mathrm{T}} (\boldsymbol{\Sigma} + \mathbf{C} \mathbf{P}_{n-1} \mathbf{C}^{\mathrm{T}})^{-1} \mathbf{C} \mathbf{P}_{n-1} \\ &= & (\mathbf{I} - \mathbf{P}_{n-1} \mathbf{C}^{\mathrm{T}} (\boldsymbol{\Sigma} + \mathbf{C} \mathbf{P}_{n-1} \mathbf{C}^{\mathrm{T}})^{-1} \mathbf{C}) \mathbf{P}_{n-1} \\ &= & (\mathbf{I} - \mathbf{K}_n \mathbf{C}) \mathbf{P}_{n-1}, \end{split}$$

which equals the r.h.s. of (13.90).

Using (329), (C.5) and (13.92), we can derive the following equality:

$$\begin{split} \mathbf{M}\mathbf{C}^{\mathrm{T}}\mathbf{\Sigma}^{-1} &= & (\mathbf{P}_{n-1} + \mathbf{C}^{\mathrm{T}}\mathbf{\Sigma}^{-1}\mathbf{C})^{-1}\mathbf{C}^{\mathrm{T}}\mathbf{\Sigma}^{-1} \\ &= & \mathbf{P}_{n-1}\mathbf{C}^{\mathrm{T}}(\mathbf{C}\mathbf{P}_{n-1}\mathbf{C}^{\mathrm{T}} + \mathbf{\Sigma})^{-1} = \mathbf{K}_{n}. \end{split}$$

Using this and (13.90), we can rewrite the expression for the mean in (328) as follows:

$$\begin{split} \mathbf{M}(\mathbf{C}^{\mathrm{T}}\mathbf{\Sigma}^{-1}\mathbf{x}_n + \mathbf{P}_{n-1}^{-1}\mathbf{A}\boldsymbol{\mu}_{n-1}) &= \mathbf{M}\mathbf{C}^{\mathrm{T}}\mathbf{\Sigma}^{-1}\mathbf{x}_n + (\mathbf{I} - \mathbf{K}_n\mathbf{C})\mathbf{A}\boldsymbol{\mu}_{n-1} \\ &= \mathbf{K}_n\mathbf{x}_n + \mathbf{A}\boldsymbol{\mu}_{n-1} - \mathbf{K}_n\mathbf{C}\mathbf{A}\boldsymbol{\mu}_{n-1} \\ &= \mathbf{A}\boldsymbol{\mu}_{n-1} + \mathbf{K}_n(\mathbf{x}_n - \mathbf{C}\mathbf{A}\boldsymbol{\mu}_{n-1}), \end{split}$$

which equals the r.h.s. of (13.89).

**13.22** Using (13.76), (13.77) and (13.84), we can write (13.93), for the case n = 1, as

$$c_1 \mathcal{N}(\mathbf{z}_1 | \boldsymbol{\mu}_1, \mathbf{V}_1) = \mathcal{N}(\mathbf{z}_1 | \boldsymbol{\mu}_0, \mathbf{V}_0) \mathcal{N}(\mathbf{x}_1 | \mathbf{C} \mathbf{z}_1, \boldsymbol{\Sigma}).$$

The r.h.s. define the joint probability distribution over  $\mathbf{x}_1$  and  $\mathbf{z}_1$  in terms of a conditional distribution over  $\mathbf{x}_1$  given  $\mathbf{z}_1$  and a distribution over  $\mathbf{z}_1$ , corresponding to (2.114) and (2.113), respectively. What we need to do is to rewrite this into a conditional distribution over  $\mathbf{z}_1$  given  $\mathbf{x}_1$  and a distribution over  $\mathbf{x}_1$ , corresponding to (2.116) and (2.115), respectively.

If we make the substitutions

$$\begin{split} \mathbf{x} \Rightarrow \mathbf{z}_1 \quad \boldsymbol{\mu} \Rightarrow \boldsymbol{\mu}_0 \quad \boldsymbol{\Lambda}^{-1} \Rightarrow \mathbf{V}_0 \\ \mathbf{y} \Rightarrow \mathbf{x}_1 \quad \mathbf{A} \Rightarrow \mathbf{C} \quad \mathbf{b} \Rightarrow \mathbf{0} \quad \mathbf{L}^{-1} \Rightarrow \boldsymbol{\Sigma}, \end{split}$$

in (2.113) and (2.114), (2.115) directly gives us the r.h.s. of (13.96).

**13.23** Using (13.76) and (13.77) we can rewrite (13.93) as

$$c_1\widehat{\alpha}(\mathbf{z}_1) = \mathcal{N}(\mathbf{z}_1|\boldsymbol{\mu}_0, \mathbf{V}_0) \mathcal{N}(\mathbf{x}_1|\mathbf{C}\mathbf{z}_1, \boldsymbol{\Sigma}).$$

Making the same substitutions as in Solution 13.22, (2.115) and (13.96) give

$$p(\mathbf{x}_1) = \mathcal{N}\left(\mathbf{x}_1 | \mathbf{C}\boldsymbol{\mu}_0, \boldsymbol{\Sigma} + \mathbf{C}\mathbf{V}_0\mathbf{C}^{\mathrm{T}}\right) = c_1.$$

Hence, from the product rule and (2.116),

$$\widehat{\alpha}(\mathbf{z}_1) = p(\mathbf{z}_1|\mathbf{x}_1) = \mathcal{N}(\mathbf{z}_1|\boldsymbol{\mu}_1, \mathbf{V}_1)$$

where, from (13.97) and (C.7),

$$\begin{aligned} \mathbf{V}_1 &= & \left( \mathbf{V}_0^{-1} + \mathbf{C}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \mathbf{C} \right)^{-1} \\ &= & \mathbf{V}_0 - \mathbf{V}_0 \mathbf{C}^{\mathrm{T}} \left( \mathbf{\Sigma} + \mathbf{C} \mathbf{V}_0 \mathbf{C}^{\mathrm{T}} \right)^{-1} \mathbf{C} \mathbf{V}_0 \\ &= & \left( \mathbf{I} - \mathbf{K}_1 \mathbf{C} \right) \mathbf{V}_0 \end{aligned}$$

and

$$egin{array}{lll} oldsymbol{\mu}_1 &=& \mathbf{V}_1 \left( \mathbf{C}^{\mathrm{T}} oldsymbol{\Sigma}^{-1} \mathbf{x}_1 + \mathbf{V}_0^{-1} oldsymbol{\mu}_0 
ight) \ &=& oldsymbol{\mu}_0 + \mathbf{K}_1 \left( \mathbf{x}_1 - \mathbf{C} oldsymbol{\mu}_0 
ight) \end{array}$$

where we have used

$$\begin{split} \mathbf{V}_1 \mathbf{C}^T \boldsymbol{\Sigma}^{-1} &=& \mathbf{V}_0 \mathbf{C}^T \boldsymbol{\Sigma}^{-1} - \mathbf{K}_1 \mathbf{C} \mathbf{V}_0 \mathbf{C}^T \boldsymbol{\Sigma}^{-1} \\ &=& \mathbf{V}_0 \mathbf{C}^T \left( \mathbf{I} - \left( \boldsymbol{\Sigma} + \mathbf{C} \mathbf{V}_0 \mathbf{C}^T \right)^{-1} \mathbf{C} \mathbf{V}_0 \mathbf{C}^T \right) \boldsymbol{\Sigma}^{-1} \\ &=& \mathbf{V}_0 \mathbf{C}^T \left( \mathbf{I} - \left( \boldsymbol{\Sigma} + \mathbf{C} \mathbf{V}_0 \mathbf{C}^T \right)^{-1} \mathbf{C} \mathbf{V}_0 \mathbf{C}^T \right. \\ &+& \left. \left( \boldsymbol{\Sigma} + \mathbf{C} \mathbf{V}_0 \mathbf{C}^T \right)^{-1} \boldsymbol{\Sigma} - \left( \boldsymbol{\Sigma} + \mathbf{C} \mathbf{V}_0 \mathbf{C}^T \right)^{-1} \boldsymbol{\Sigma} \right) \boldsymbol{\Sigma}^{-1} \\ &=& \mathbf{V}_0 \mathbf{C}^T \left( \boldsymbol{\Sigma} + \mathbf{C} \mathbf{V}_0 \mathbf{C}^T \right)^{-1} = \mathbf{K}_1. \end{split}$$

**13.24** This extension can be embedded in the existing framework by adopting the following modifications:

$$\boldsymbol{\mu}_0' = \left[ \begin{array}{c} \boldsymbol{\mu}_0 \\ 1 \end{array} \right] \quad \mathbf{V}_0' = \left[ \begin{array}{cc} \mathbf{V}_0 & \mathbf{0} \\ \mathbf{0} & 0 \end{array} \right] \quad \boldsymbol{\Gamma}' = \left[ \begin{array}{cc} \boldsymbol{\Gamma} & \mathbf{0} \\ \mathbf{0} & 0 \end{array} \right]$$
$$\mathbf{A}' = \left[ \begin{array}{cc} \mathbf{A} & \mathbf{a} \\ \mathbf{0} & 1 \end{array} \right] \quad \mathbf{C}' = \left[ \begin{array}{cc} \mathbf{C} & \mathbf{c} \end{array} \right].$$

This will ensure that the constant terms a and c are included in the corresponding Gaussian means for  $z_n$  and  $x_n$  for n = 1, ..., N.

Note that the resulting covariances for  $\mathbf{z}_n$ ,  $\mathbf{V}_n$ , will be singular, as will the corresponding prior covariances,  $\mathbf{P}_{n-1}$ . This will, however, only be a problem where these matrices need to be inverted, such as in (13.102). These cases must be handled separately, using the 'inversion' formula

$$(\mathbf{P}_{n-1}')^{-1} = \left[ \begin{array}{cc} \mathbf{P}_{n-1}^{-1} & \mathbf{0} \\ \mathbf{0} & 0 \end{array} \right],$$

nullifying the contribution from the (non-existent) variance of the element in  $\mathbf{z}_n$  that accounts for the constant terms  $\mathbf{a}$  and  $\mathbf{c}$ .

**13.25** NOTE: In PRML, the second half on the third sentence in the exercise should read: "... in which C = 1, A = 1 and  $\Gamma = 0$ ." Moreover, in the following sentence  $m_0$  and  $V_0$  should be replaced by  $\mu_0$  and  $P_0$ ; see also the PRML errata.

Since  ${f C}=1,$   ${f P}_0=\sigma_0^2$  and  ${f \Sigma}=\sigma^2,$  (13.97) gives

$$\mathbf{K}_1 = \frac{\sigma_0^2}{\sigma_0^2 + \sigma}.$$

Substituting this into (13.94) and (13.95), we get

$$\mu_{1} = \mu_{0} + \frac{\sigma_{0}^{2}}{\sigma_{0}^{2} + \sigma} (x_{1} - \mu_{0})$$

$$= \frac{1}{\sigma_{0}^{2} + \sigma} (\sigma_{0}^{2} x_{1} + \sigma \mu_{0})$$

$$\sigma_{1}^{2} = \left(1 - \frac{\sigma_{0}^{2}}{\sigma_{0}^{2} + \sigma}\right) \sigma_{0}^{2}$$

$$= \frac{\sigma_{0}^{2} \sigma_{0}^{2}}{\sigma_{0}^{2} + \sigma}$$

where  $\sigma_1^2$  replaces  $V_1$ . We note that these agree with (2.141) and (2.142), respectively.

We now assume that (2.141) and (2.142) hold for N, and we rewrite them as

$$\mu_N = \sigma_N^2 \left( \frac{1}{\sigma_0^2} \mu_0 + \frac{N}{\sigma^2} \mu_{\rm ML}^{(N)} \right)$$
 (330)

$$\sigma_N^2 = \frac{\sigma_0^2 \sigma^2}{N \sigma_0^2 + \sigma} \tag{331}$$

where, analogous to (2.143),

$$\mu_{\rm ML}^{(N)} = \frac{1}{N} \sum_{n=1}^{N} x_n. \tag{332}$$

Since A = 1 and  $\Gamma = 0$ , (13.88) gives

$$\mathbf{P}_N = \sigma_N^2 \tag{333}$$

substituting this into (13.92), we get

$$\mathbf{K}_{N+1} = \frac{\sigma_N^2}{N\sigma_N^2 + \sigma} \tag{334}$$

Using (331), (333), (334) and (13.90), we get

$$\sigma_{N+1}^{2} = \left(1 - \frac{\sigma_{N}^{2}}{\sigma_{N}^{2} + \sigma}\right) \sigma_{N}^{2} 
= \frac{\sigma_{N}^{2} \sigma^{2}}{\sigma_{N}^{2} + \sigma} 
= \frac{\sigma_{0}^{2} \sigma^{4} / (\sigma_{0}^{2} + \sigma)}{(\sigma^{2} \sigma_{0}^{2} \sigma^{4} + \sigma^{2} N \sigma_{0}^{2}) / (\sigma_{0}^{2} + \sigma)} 
= \frac{\sigma_{0}^{2} \sigma^{2}}{(N+1)\sigma_{0}^{2} + \sigma}.$$
(335)

Using (330), (332), (334), (335) and (13.89), we get

$$\mu_{N+1} = \mu_N + \frac{\sigma_N^2}{\sigma_N^2 + \sigma} (x_{N+1} - \mu_N)$$

$$= \frac{1}{\sigma_N^2 + \sigma} (\sigma_N^2 x_{N+1} + \sigma \mu_N)$$

$$= \frac{\sigma_N^2}{\sigma_0^2 + \sigma} x_{N+1} + \frac{\sigma_N^2 \sigma^2}{\sigma_N^2 + \sigma} \left( \frac{1}{\sigma_0^2} \mu_0 + \frac{1}{\sigma^2} \sum_{n=1}^N x_n \right)$$

$$= \sigma_{N+1}^2 \left( \frac{1}{\sigma_0^2} \mu_0 + \frac{N+1}{\sigma^2} \mu_{\text{ML}}^{(N+1)} \right).$$

Thus (330) and (331) must hold for all  $N \ge 1$ .

**NOTE**: In the 1<sup>st</sup> printing of PRML, equation (12.42) contains a mistake; the covariance on the r.h.s. should be  $\sigma^2 \mathbf{M}^{-1}$ . Furthermore, the exercise should make explicit the assumption that  $\mu = 0$  in (12.42).

From (13.84) and the assumption that A = 0, we have

$$p(\mathbf{z}_n|\mathbf{x}_1,\dots,\mathbf{x}_n) = p(\mathbf{z}_n|\mathbf{x}_n) = \mathcal{N}\left(\mathbf{z}_n|\boldsymbol{\mu}_n,\mathbf{V}_n\right)$$
(336)

where  $\mu_n$  and  $V_n$  are given by (13.89) and (13.90), respectively. Since A=0 and  $\Gamma = I$ ,  $P_{n-1} = I$  for all n and thus (13.92) becomes

$$\mathbf{K}_{n} = \mathbf{P}_{n-1}\mathbf{C}^{\mathrm{T}} \left(\mathbf{C}\mathbf{P}_{n-1}\mathbf{C}^{\mathrm{T}} + \mathbf{\Sigma}\right)^{-1}$$
$$= \mathbf{W}^{\mathrm{T}} \left(\mathbf{W}\mathbf{W}^{\mathrm{T}} + \sigma^{2}\mathbf{I}\right)^{-1}$$
(337)

where we have substituted W for C and  $\sigma^2 I$  for  $\Sigma$ . Using (337), (12.41), (C.6) and (C.7), (13.89) and (13.90) can be rewritten as

$$\mu_n = \mathbf{K}_n \mathbf{x}_n$$

$$= \mathbf{W}^{\mathrm{T}} (\mathbf{W} \mathbf{W}^{\mathrm{T}} + \sigma^2 \mathbf{I})^{-1} \mathbf{x}_n$$

$$= \mathbf{M}^{-1} \mathbf{W}^{\mathrm{T}} \mathbf{x}_n$$

$$\mathbf{V}_{n} = (\mathbf{I} - \mathbf{K}_{n} \mathbf{C}) \mathbf{P}_{n-1} = \mathbf{I} - \mathbf{K}_{n} \mathbf{W}$$

$$= \mathbf{I} - \mathbf{W}^{T} (\mathbf{W} \mathbf{W}^{T} + \sigma^{2} \mathbf{I})^{-1} \mathbf{W}$$

$$= (\sigma^{-2} \mathbf{W}^{T} \mathbf{W} + \mathbf{I})^{-1}$$

$$= \sigma^{2} (\mathbf{W}^{T} \mathbf{W} + \sigma^{2} \mathbf{I})^{-1} = \sigma^{2} \mathbf{M}^{-1}$$

which makes (336) equivalent with (12.42), assuming  $\mu=0$ .

NOTE in the 1<sup>st</sup> printing of PRML, this exercise small have made explicit the assumption that C=I in (13.86).

From (13.86), it is easily seen that if  $\Sigma$  goes to  $\mathbf{0}$ , the posterior over  $\mathbf{z}_n$  will become completely determined by  $\mathbf{x}_n$ , since the first factor on the r.h.s. of (13.86), and hence also the l.h.s., will collapse to a spike at  $\mathbf{x}_n = \mathbf{C}\mathbf{z}_n$ .

**13.28** NOTE: In PRML, this exercise should also assume that  ${\bf C}={\bf I}$ ., Moreover,  ${\bf V}_0$  should be replaced by  ${\bf P}_0$  in the text of the exercise; see also the PRML errata.

Starting from (13.75) and (13.77), we can use (2.113)–(2.117) to obtain

$$p(\mathbf{z}_1|\mathbf{x}_1) = \mathcal{N}(\mathbf{z}_1|\boldsymbol{\mu}_1, \mathbf{V}_1)$$

where

$$\boldsymbol{\mu}_{1} = \mathbf{V}_{1} \left( \mathbf{C}^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \mathbf{x}_{1} + \mathbf{P}_{0}^{-1} \boldsymbol{\mu}_{0} \right) = \mathbf{x}_{1}$$
 (338)

$$\mathbf{V}_{1} = \left(\mathbf{P}_{0}^{-1} + \mathbf{C}^{\mathrm{T}} \mathbf{\Sigma}^{-1} \mathbf{C}\right)^{-1} = \mathbf{\Sigma}$$
 (339)

since  $P_0 \to \infty$  and C = I; note that these results can equally well be obtained from (13.94), (13.95) and (13.97).

Now we assume that for N

$$\boldsymbol{\mu}_N = \overline{\mathbf{x}}_N = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \tag{340}$$

$$\mathbf{V}_{N} = \frac{1}{N}\mathbf{\Sigma} \tag{341}$$

and we note that these assumptions are met by (338) and (339), respectively. From (13.88) and (341), we then have

$$\mathbf{P}_N = \mathbf{V}_N = \frac{1}{N} \mathbf{\Sigma} \tag{342}$$

since C = I and  $\Gamma = 0$ . Using this together with (13.92), we obtain

$$\mathbf{K}_{N+1} = \mathbf{P}_{N} \mathbf{C}^{\mathrm{T}} \left( \mathbf{C} \mathbf{P}_{N} \mathbf{C}^{\mathrm{T}} + \mathbf{\Sigma} \right)^{-1}$$

$$= \mathbf{P}_{N} \left( \mathbf{P}_{N} + \mathbf{\Sigma} \right)^{-1}$$

$$= \frac{1}{N} \mathbf{\Sigma} \left( \frac{1}{N} \mathbf{\Sigma} + \mathbf{\Sigma} \right)^{-1}$$

$$= \frac{1}{N} \mathbf{\Sigma} \left( \frac{N+1}{N} \mathbf{\Sigma} \right)^{-1}$$

$$= \frac{1}{N+1} \mathbf{I}.$$

Substituting this into (13.89) and (13.90), making use of (340) and (342), we have

$$\mu_{N+1} = \mu_N + \frac{1}{N+1} (\mathbf{x}_{N+1} - \mu_N)$$

$$= \overline{\mathbf{x}}_N + \frac{1}{N+1} (\mathbf{x}_{N+1} - \overline{\mathbf{x}}_N)$$

$$= \frac{1}{N+1} \mathbf{x}_{N+1} + \left(1 - \frac{1}{N+1}\right) \frac{1}{N} \sum_{n=1}^{N}$$

$$= \frac{1}{N+1} \sum_{n=1}^{N+1} \mathbf{x}_n = \overline{\mathbf{x}}_{N+1}$$

$$\mathbf{V}_{N+1} = \left(\mathbf{I} - \frac{1}{N+1} \mathbf{I}\right) \frac{1}{N} \mathbf{\Sigma}$$

$$= \frac{1}{N+1} \mathbf{\Sigma}.$$

Thus, (340) and (341) holds for all  $N \geqslant 1$ .

**13.29 NOTE**: In the 1<sup>st</sup> printing of PRML,  $\mu_N$  should be  $\mu_n$  on the r.h.s. of (13.100) Multiplying both sides of (13.99) by  $\widehat{\alpha}(\mathbf{z}_n)$ , and then making use of (13.98), we get

$$c_{n+1}\mathcal{N}\left(\mathbf{z}_{n}|\widehat{\boldsymbol{\mu}}_{n},\widehat{\mathbf{V}}_{n}\right) = \widehat{\alpha}(\mathbf{z}_{n}) \int \widehat{\beta}(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1}|\mathbf{z}_{n+1}) p(\mathbf{z}_{n+1}|\mathbf{z}_{n}) \, d\mathbf{z}_{n+1}$$
(343)

Using (2.113)–(2.117), (13.75) and (13.84), we have

$$\widehat{\alpha}(\mathbf{z}_n)p(\mathbf{z}_{n+1}|\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n|\boldsymbol{\mu}_n, \mathbf{V}_n)\mathcal{N}(\mathbf{z}_{n+1}|\mathbf{A}\mathbf{z}_n, \boldsymbol{\Gamma})$$

$$= \mathcal{N}(\mathbf{z}_{n+1}|\mathbf{A}\boldsymbol{\mu}_n, \mathbf{A}\mathbf{V}_n\mathbf{A} + \boldsymbol{\Gamma})\mathcal{N}(\mathbf{z}_n|\mathbf{m}_n, \mathbf{M}_n) \quad (344)$$

where

$$\mathbf{m}_n = \mathbf{M}_n \left( \mathbf{A}^{\mathrm{T}} \mathbf{\Gamma}^{-1} \mathbf{z}_{n+1} + \mathbf{V}_n^{-1} \boldsymbol{\mu}_n \right)$$
 (345)

and, using (C.7) and (13.102),

$$\mathbf{M}_{n} = (\mathbf{A}^{\mathrm{T}} \mathbf{\Gamma}^{-1} \mathbf{A} + \mathbf{V}_{n})^{-1}$$

$$= \mathbf{V}_{n} - \mathbf{V}_{n} \mathbf{A}^{\mathrm{T}} (\mathbf{\Gamma} + \mathbf{A} \mathbf{V}_{n} \mathbf{A}^{\mathrm{T}})^{-1} \mathbf{A} \mathbf{V}_{n}$$

$$= \mathbf{V}_{n} - \mathbf{V}_{n} \mathbf{A}^{\mathrm{T}} \mathbf{P}_{n}^{-1} \mathbf{A} \mathbf{V}_{n}$$

$$= (\mathbf{I} - \mathbf{V}_{n} \mathbf{A}^{\mathrm{T}} \mathbf{P}_{n}^{-1} \mathbf{A}) \mathbf{V}_{n}$$

$$= (\mathbf{I} - \mathbf{J}_{n} \mathbf{A}) \mathbf{V}_{n}$$

$$(348)$$

Substituting the r.h.s. of (344) into (343) and then making use of (13.85)–(13.88) and

(13.98), we have

$$c_{n+1}\mathcal{N}\left(\mathbf{z}_{n}|\widehat{\boldsymbol{\mu}}_{n},\widehat{\mathbf{V}}_{n}\right) = \int \widehat{\boldsymbol{\beta}}(\mathbf{z}_{n+1})p(\mathbf{x}_{n+1}|\mathbf{z}_{n+1})\mathcal{N}\left(\mathbf{z}_{n+1}|\mathbf{A}\boldsymbol{\mu}_{n},\mathbf{P}_{n}\right)$$

$$\mathcal{N}\left(\mathbf{z}_{n}|\mathbf{m}_{n},\mathbf{M}_{n}\right) d\mathbf{z}_{n+1}$$

$$= \int \widehat{\boldsymbol{\beta}}(\mathbf{z}_{n+1})c_{n+1}\widehat{\boldsymbol{\alpha}}(\mathbf{z}_{n+1})\mathcal{N}\left(\mathbf{z}_{n}|\mathbf{m}_{n},\mathbf{M}_{n}\right) d\mathbf{z}_{n+1}$$

$$= c_{n+1} \int \gamma(\mathbf{z}_{n+1})\mathcal{N}\left(\mathbf{z}_{n}|\mathbf{m}_{n},\mathbf{M}_{n}\right) d\mathbf{z}_{n+1}$$

$$= c_{n+1} \int \mathcal{N}\left(\mathbf{z}_{n+1}|\widehat{\boldsymbol{\mu}}_{n},\widehat{\mathbf{V}}_{n}\right)\mathcal{N}\left(\mathbf{z}_{n}|\mathbf{m}_{n},\mathbf{M}_{n}\right) d\mathbf{z}_{n+1}.$$

Thus, from this, (345) and (2.113)–(2.115), we see that

$$\widehat{\boldsymbol{\mu}}_n = \mathbf{M}_n \left( \mathbf{A}^{\mathrm{T}} \boldsymbol{\Gamma}^{-1} \widehat{\boldsymbol{\mu}}_{n+1} + \mathbf{V}_n^{-1} \boldsymbol{\mu}_n \right)$$
 (349)

$$\widehat{\mathbf{V}}_n = \mathbf{M}_n \mathbf{A}^{\mathrm{T}} \mathbf{\Gamma}^{-1} \widehat{\mathbf{V}}_{n+1} \mathbf{\Gamma}^{-1} \mathbf{A} \mathbf{M}_n + \mathbf{M}_n.$$
 (350)

Using (347) and (13.102), we see that

$$\mathbf{M}_{n}\mathbf{A}^{\mathrm{T}}\mathbf{\Gamma}^{-1} = (\mathbf{V} - \mathbf{V}_{n}\mathbf{A}^{\mathrm{T}}\mathbf{P}_{n}^{-1}\mathbf{A}\mathbf{V}_{n})\mathbf{A}^{\mathrm{T}}\mathbf{\Gamma}^{-1}$$

$$= \mathbf{V}_{n}\mathbf{A}^{\mathrm{T}}(\mathbf{I} - \mathbf{P}_{n}^{-1}\mathbf{A}\mathbf{V}_{n}\mathbf{A}^{\mathrm{T}})\mathbf{\Gamma}^{-1}$$

$$= \mathbf{V}_{n}\mathbf{A}^{\mathrm{T}}(\mathbf{I} - \mathbf{P}_{n}^{-1}\mathbf{A}\mathbf{V}_{n}\mathbf{A}^{\mathrm{T}} - \mathbf{P}_{n}^{-1}\mathbf{\Gamma} + \mathbf{P}_{n}^{-1}\mathbf{\Gamma})\mathbf{\Gamma}^{-1}$$

$$= \mathbf{V}_{n}\mathbf{A}^{\mathrm{T}}(\mathbf{I} - \mathbf{P}_{n}^{-1}\mathbf{P}_{n} + \mathbf{P}_{n}^{-1}\mathbf{\Gamma})\mathbf{\Gamma}^{-1}$$

$$= \mathbf{V}_{n}\mathbf{A}^{\mathrm{T}}\mathbf{P}_{n}^{-1} = \mathbf{J}_{n}$$
(351)

and using (348) together with (351), we can rewrite (349) as (13.100). Similarly, using (13.102), (347) and (351), we rewrite (350) as

$$\widehat{\mathbf{V}}_{n} = \mathbf{M}_{n} \mathbf{A}^{\mathrm{T}} \mathbf{\Gamma}^{-1} \widehat{\mathbf{V}}_{n+1} \mathbf{\Gamma}^{-1} \mathbf{A} \mathbf{M}_{n} + \mathbf{M}_{n} 
= \mathbf{J}_{n} \widehat{\mathbf{V}}_{n+1} \mathbf{J}_{n}^{\mathrm{T}} + \mathbf{V}_{n} - \mathbf{V}_{n} \mathbf{A}^{\mathrm{T}} \mathbf{P}_{n}^{-1} \mathbf{A} \mathbf{V}_{n} 
= \mathbf{V}_{n} + \mathbf{J}_{n} \left( \widehat{\mathbf{V}}_{n+1} - \mathbf{P}_{n} \right) \mathbf{J}_{n}^{\mathrm{T}}$$

#### **13.30 NOTE**: See note in Solution 13.15.

The first line of (13.103) corresponds exactly to (13.65). We then use (13.75), (13.76), (13.84) and (13.98) to rewrite  $p(\mathbf{z}_n|\mathbf{z}_{n-1})$ ,  $p(\mathbf{x}_n|\mathbf{z}_n)$ ,  $\widehat{\alpha}(\mathbf{z}_{n-1})$  and  $\widehat{\beta}(\mathbf{z}_n)$ , respectively, yielding the second line of (13.103).

**13.31** Substituting the r.h.s. of (13.84) for  $\widehat{\alpha}(\mathbf{z}_n)$  in (13.103) and then using (2.113)–(2.117) and (13.86), we get

$$\begin{aligned}
& \left\{ \left( \mathbf{z}_{n-1}, \mathbf{z}_{n} \right) \right. \\
& = \frac{\mathcal{N} \left( \mathbf{z}_{n-1} | \boldsymbol{\mu}_{n-1}, \mathbf{V}_{n-1} \right) \mathcal{N} \left( \mathbf{z}_{n} | \mathbf{A} \mathbf{z}_{n-1}, \boldsymbol{\Gamma} \right) \mathcal{N} \left( \mathbf{x}_{n} | \mathbf{C} \mathbf{z}_{n}, \boldsymbol{\Sigma} \right) \mathcal{N} \left( \mathbf{z}_{n} | \widehat{\boldsymbol{\mu}}_{n}, \widehat{\mathbf{V}}_{n} \right)}{\mathcal{N} \left( \mathbf{z}_{n} | \boldsymbol{\mu}_{n}, \mathbf{V}_{n} \right) c_{n}} \\
& = \frac{\mathcal{N} \left( \mathbf{z}_{n} | \mathbf{A} \boldsymbol{\mu}_{n-1}, \mathbf{P}_{n-1} \right) \mathcal{N} \left( \mathbf{z}_{n-1} | \mathbf{m}_{n-1}, \mathbf{M}_{n-1} \right) \mathcal{N} \left( \mathbf{z}_{n} | \widehat{\boldsymbol{\mu}}_{n}, \widehat{\mathbf{V}}_{n} \right)}{\mathcal{N} \left( \mathbf{z}_{n} | \mathbf{A} \boldsymbol{\mu}_{n-1}, \mathbf{P}_{n-1} \right)} \\
& = \mathcal{N} \left( \mathbf{z}_{n-1} | \mathbf{m}_{n-1}, \mathbf{M}_{n-1} \right) \mathcal{N} \left( \mathbf{z}_{n} | \widehat{\boldsymbol{\mu}}_{n}, \widehat{\mathbf{V}}_{n} \right) \end{aligned} \tag{352}$$

where  $m_{n-1}$  and  $M_{n-1}$  are given by (345) and (346). Equation (13.104) then follows from (345), (351) and (352).

**13.32 NOTE**: In PRML,  $V_0$  should be replaced by  $P_0$  in the text of the exercise; see also the PRML errata.

We can write the expected complete log-likelihood, given by the equation after (13.109), as a function of  $\mu_0$  and  $P_0$ , as follows:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = -\frac{1}{2} \ln |\mathbf{P}_{0}|$$

$$-\frac{1}{2} \mathbb{E}_{\mathbf{Z}|\boldsymbol{\theta}^{\text{old}}} \left[ \mathbf{z}_{1}^{\text{T}} \mathbf{P}_{0}^{-1} \mathbf{z}_{1} - \mathbf{z}_{1}^{\text{T}} \mathbf{P}_{0}^{-1} \boldsymbol{\mu}_{0} - \boldsymbol{\mu}_{0}^{\text{T}} \mathbf{P}_{0}^{-1} \mathbf{z}_{1} + \boldsymbol{\mu}_{0}^{\text{T}} \mathbf{P}_{0}^{-1} \boldsymbol{\mu}_{0} \right]$$
(353)
$$= \frac{1}{2} \left( \ln |\mathbf{P}_{0}^{-1}| - \text{Tr} \left[ \mathbf{P}_{0}^{-1} \mathbb{E}_{\mathbf{Z}|\boldsymbol{\theta}^{\text{old}}} \left[ \mathbf{z}_{1} \mathbf{z}_{1}^{\text{T}} - \mathbf{z}_{1} \boldsymbol{\mu}_{0}^{\text{T}} - \boldsymbol{\mu}_{0} \mathbf{z}_{1}^{\text{T}} + \boldsymbol{\mu}_{0} \boldsymbol{\mu}_{0}^{\text{T}} \right] \right),$$
(354)

where we have used (C.13) and omitted terms independent of  $\mu_0$  and  $\mathbf{P}_0$ .

From (353), we can calculate the derivative w.r.t.  $\mu_0$  using (C.19), to get

$$\frac{\partial Q}{\partial \boldsymbol{\mu}_0} = 2\mathbf{P}_0^{-1}\boldsymbol{\mu}_0 - 2\mathbf{P}_0^{-1}\mathbb{E}[\mathbf{z}_1].$$

Setting this to zero and rearranging, we immediately obtain (13.110).

Using (354), (C.24) and (C.28), we can evaluate the derivatives w.r.t.  $P_0^{-1}$ ,

$$\frac{\partial Q}{\partial \mathbf{P}_0^{-1}} = \frac{1}{2} \left( \mathbf{P}_0 - \mathbb{E}[\mathbf{z}_1 \mathbf{z}_1^T] - \mathbb{E}[\mathbf{z}_1] \boldsymbol{\mu}_0^T - \boldsymbol{\mu}_0 \mathbb{E}[\mathbf{z}_1^T] + \boldsymbol{\mu}_0 \boldsymbol{\mu}_0^T \right).$$

Setting this to zero, rearranging and making use of (13.110), we get (13.111).

**13.33 NOTE**: In PRML, the first instance of  $A^{new}$  on the second line of equation (13.114) should be transposed.

Expanding the square in the second term of (13.112) and making use of the trace operator, we obtain

$$\mathbb{E}_{\mathbf{Z}|\boldsymbol{\theta}} \left[ \frac{1}{2} \sum_{n=2}^{N} (\mathbf{z}_{n} - \mathbf{A} \mathbf{z}_{n-1})^{\mathrm{T}} \boldsymbol{\Gamma}^{-1} (\mathbf{z}_{n} - \mathbf{A} \mathbf{z}_{n-1}) \right]$$

$$= \mathbb{E}_{\mathbf{Z}|\boldsymbol{\theta}} \left[ \frac{1}{2} \sum_{n=2}^{N} \mathrm{Tr} \left( \boldsymbol{\Gamma}^{-1} \left\{ \mathbf{A} \mathbf{z}_{n-1} \mathbf{z}_{n-1}^{\mathrm{T}} \mathbf{A}^{\mathrm{T}} - \mathbf{z}_{n} \mathbf{z}_{n-1}^{\mathrm{T}} \mathbf{A}^{\mathrm{T}} - \mathbf{z}_{n} \mathbf{z}_{n}^{\mathrm{T}} + \mathbf{z}_{n} \mathbf{z}_{n}^{\mathrm{T}} \right\} \right) \right]. \quad (355)$$

Using results from Appendix C, we can calculate the derivative of this w.r.t. A, yielding

$$\frac{\partial Q}{\partial \mathbf{A}} = \mathbf{\Gamma}^{-1} \mathbf{A} \mathbb{E} \left[ \mathbf{z}_{n-1} \mathbf{z}_{n-1}^{\mathrm{T}} \right] - \mathbf{\Gamma}^{-1} \mathbb{E} \left[ \mathbf{z}_{n} \mathbf{z}_{n-1}^{\mathrm{T}} \right].$$

Setting this equal to zero and solving for A, we obtain (13.113).

Using (355) and results from Appendix C, we can calculate the derivative of (13.112) w.r.t.  $\Gamma^{-1}$ , to obtain

$$\frac{\partial Q}{\partial \mathbf{\Gamma}} = \frac{N-1}{2} \mathbf{\Gamma} - \frac{1}{2} \sum_{n=2}^{N} \left( \mathbf{A} \mathbb{E} \left[ \mathbf{z}_{n-1} \mathbf{z}_{n-1}^{\mathrm{T}} \right] \mathbf{A}^{\mathrm{T}} - \mathbb{E} \left[ \mathbf{z}_{n} \mathbf{z}_{n-1}^{\mathrm{T}} \right] \mathbf{A}^{\mathrm{T}} - \mathbf{A} \mathbb{E} \left[ \mathbf{z}_{n-1} \mathbf{z}_{n}^{\mathrm{T}} \right] + \mathbb{E} \left[ \mathbf{z}_{n} \mathbf{z}_{n}^{\mathrm{T}} \right] \right).$$

Setting this equal to zero, substituting  $A^{\rm new}$  for A and solving for  $\Gamma$ , we obtain (13.114).

**13.34 NOTE**: In PRML, the first and third instances of  $C^{\text{new}}$  on the second line of equation (13.116) should be transposed.

By making use of (C.28), equations (13.115) and (13.116) are obtained in an identical manner to (13.113) and (13.114), respectively, in Solution 13.33.

## **Chapter 14** Combining Models

**14.1** The required predictive distribution is given by

$$p(\mathbf{t}|\mathbf{x}, \mathbf{X}, \mathbf{T}) = \sum_{h} p(h) \sum_{\mathbf{z}_{h}} p(\mathbf{z}_{h}) \int p(\mathbf{t}|\mathbf{x}, \boldsymbol{\theta}_{h}, \mathbf{z}_{h}, h) p(\boldsymbol{\theta}_{h}|\mathbf{X}, \mathbf{T}, h) d\boldsymbol{\theta}_{h}, \quad (356)$$

where

$$p(\boldsymbol{\theta}_{h}|\mathbf{X}, \mathbf{T}, h) = \frac{p(\mathbf{T}|\mathbf{X}, \boldsymbol{\theta}_{h}, h)p(\boldsymbol{\theta}_{h}|h)}{p(\mathbf{T}|\mathbf{X}, h)}$$

$$\propto p(\boldsymbol{\theta}|h) \prod_{n=1}^{N} p(\mathbf{t}_{n}|\mathbf{x}_{n}, \boldsymbol{\theta}, h)$$

$$= p(\boldsymbol{\theta}|h) \prod_{n=1}^{N} \left( \sum_{\mathbf{z}_{nh}} p(\mathbf{t}_{n}, \mathbf{z}_{nh}|\mathbf{x}_{n}, \boldsymbol{\theta}, h) \right)$$
(357)

The integrals and summations in (356) are examples of Bayesian averaging, accounting for the uncertainty about which model, h, is the correct one, the value of the corresponding parameters,  $\theta_h$ , and the state of the latent variable,  $\mathbf{z}_h$ . The summation in (357), on the other hand, is an example of the use of latent variables, where different data points correspond to different latent variable states, although all the data are assumed to have been generated by a single model, h.

**14.2** Using (14.13), we can rewrite (14.11) as

$$E_{\text{COM}} = \mathbb{E}_{\mathbf{x}} \left[ \left\{ \frac{1}{M} \sum_{m=1}^{M} \epsilon_{m}(\mathbf{x}) \right\}^{2} \right]$$

$$= \frac{1}{M^{2}} \mathbb{E}_{\mathbf{x}} \left[ \left\{ \sum_{m=1}^{M} \epsilon_{m}(\mathbf{x}) \right\}^{2} \right]$$

$$= \frac{1}{M^{2}} \sum_{m=1}^{M} \sum_{l=1}^{M} \mathbb{E}_{\mathbf{x}} \left[ \epsilon_{m}(\mathbf{x}) \epsilon_{l}(\mathbf{x}) \right]$$

$$= \frac{1}{M^{2}} \sum_{m=1}^{M} \mathbb{E}_{\mathbf{x}} \left[ \epsilon_{m}(\mathbf{x})^{2} \right] = \frac{1}{M} E_{\text{AV}}$$

where we have used (14.10) in the last step.

**14.3** We start by rearranging the r.h.s. of (14.10), by moving the factor 1/M inside the sum and the expectation operator outside the sum, yielding

$$\mathbb{E}_{\mathbf{x}} \left[ \sum_{m=1}^{M} \frac{1}{M} \epsilon_m(\mathbf{x})^2 \right].$$

If we then identify  $\epsilon_m(\mathbf{x})$  and 1/M with  $x_i$  and  $\lambda_i$  in (1.115), respectively, and take  $f(x) = x^2$ , we see from (1.115) that

$$\left(\sum_{m=1}^{M} \frac{1}{M} \epsilon_m(\mathbf{x})\right)^2 \leqslant \sum_{m=1}^{M} \frac{1}{M} \epsilon_m(\mathbf{x})^2.$$

Since this holds for all values of x, it must also hold for the expectation over x, proving (14.54).

**14.4** If  $E(y(\mathbf{x}))$  is convex, we can apply (1.115) as follows:

$$E_{\text{AV}} = \frac{1}{M} \sum_{m=1}^{M} \mathbb{E}_{\mathbf{x}} \left[ E(y(\mathbf{x})) \right]$$
$$= \mathbb{E}_{\mathbf{x}} \left[ \sum_{m=1}^{M} \frac{1}{M} E(y(\mathbf{x})) \right]$$
$$\geqslant \mathbb{E}_{\mathbf{x}} \left[ E \left( \sum_{m=1}^{M} \frac{1}{M} y(\mathbf{x}) \right) \right]$$
$$= E_{\text{COM}}$$

where  $\lambda_i = 1/M$  for i = 1, ..., M in (1.115) and we have implicitly defined versions of  $E_{\text{AV}}$  and  $E_{\text{COM}}$  corresponding to  $E(y(\mathbf{x}))$ .

14.5 To prove that (14.57) is a sufficient condition for (14.56) we have to show that (14.56) follows from (14.57). To do this, consider a fixed set of  $y_m(\mathbf{x})$  and imagine varying the  $\alpha_m$  over all possible values allowed by (14.57) and consider the values taken by  $y_{\text{COM}}(\mathbf{x})$  as a result. The maximum value of  $y_{\text{COM}}(\mathbf{x})$  occurs when  $\alpha_k = 1$  where  $y_k(\mathbf{x}) \geqslant y_m(\mathbf{x})$  for  $m \neq k$ , and hence all  $\alpha_m = 0$  for  $m \neq k$ . An analogous result holds for the minimum value. For other settings of  $\alpha$ ,

$$y_{\min}(\mathbf{x}) < y_{\text{COM}}(\mathbf{x}) < y_{\max}(\mathbf{x}),$$

since  $y_{\text{COM}}(\mathbf{x})$  is a convex combination of points,  $y_m(\mathbf{x})$ , such that

$$\forall m: y_{\min}(\mathbf{x}) \leqslant y_m(\mathbf{x}) \leqslant y_{\max}(\mathbf{x}).$$

Thus, (14.57) is a sufficient condition for (14.56).

Showing that (14.57) is a necessary condition for (14.56) is equivalent to showing that (14.56) is a sufficient condition for (14.57). The implication here is that if (14.56) holds for any choice of values of the committee members  $\{y_m(\mathbf{x})\}$  then (14.57) will be satisfied. Suppose, without loss of generality, that  $\alpha_k$  is the smallest of the  $\alpha$  values, i.e.  $\alpha_k \leqslant \alpha_m$  for  $k \neq m$ . Then consider  $y_k(\mathbf{x}) = 1$ , together with  $y_m(\mathbf{x}) = 0$  for all  $m \neq k$ . Then  $y_{\min}(\mathbf{x}) = 0$  while  $y_{\text{COM}}(\mathbf{x}) = \alpha_k$  and hence from (14.56) we obtain  $\alpha_k \geqslant 0$ . Since  $\alpha_k$  is the smallest of the  $\alpha$  values it follows that all of the coefficients must satisfy  $\alpha_k \geqslant 0$ . Similarly, consider the case in which  $y_m(\mathbf{x}) = 1$  for all m. Then  $y_{\min}(\mathbf{x}) = y_{\max}(\mathbf{x}) = 1$ , while  $y_{\text{COM}}(\mathbf{x}) = \sum_m \alpha_m$ . From (14.56) it then follows that  $\sum_m \alpha_m = 1$ , as required.

**14.6** If we differentiate (14.23) w.r.t.  $\alpha_m$  we obtain

$$\frac{\partial E}{\partial \alpha_m} = \frac{1}{2} \left( (e^{\alpha_m/2} + e^{-\alpha_m/2}) \sum_{n=1}^{N} w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) - e^{-\alpha_m/2} \sum_{n=1}^{N} w_n^{(m)} \right).$$

Setting this equal to zero and rearranging, we get

$$\frac{\sum_{n} w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n} w_n^{(m)}} = \frac{e^{-\alpha_m/2}}{e^{\alpha_m/2} + e^{-\alpha_m/2}} = \frac{1}{e^{\alpha_m} + 1}.$$

Using (14.16), we can rewrite this as

$$\frac{1}{e^{\alpha_m} + 1} = \epsilon_m,$$

which can be further rewritten as

$$e^{\alpha_m} = \frac{1 - \epsilon_m}{\epsilon_m},$$

from which (14.17) follows directly.

**14.7** Taking the functional derivative of (14.27) w.r.t.  $y(\mathbf{x})$ , we get

$$\frac{\delta}{\delta y(\mathbf{x})} \mathbb{E}_{\mathbf{x},t} \left[ \exp\left\{ -ty(\mathbf{x}) \right\} \right] = -\sum_{t} t \exp\left\{ -ty(\mathbf{x}) \right\} p(t|\mathbf{x}) p(\mathbf{x})$$
$$= \left\{ \exp\left\{ y(\mathbf{x}) \right\} p(t=-1|\mathbf{x}) - \exp\left\{ -y(\mathbf{x}) \right\} p(t=+1|\mathbf{x}) \right\} p(\mathbf{x}).$$

Setting this equal to zero and rearranging, we obtain (14.28).

**14.8** Assume that (14.20) is a negative log likelihood function. Then the corresponding likelihood function is given by

$$\exp(-E) = \prod_{n=1}^{N} \exp\left(-\exp\left\{-t_n f_m(\mathbf{x}_n)\right\}\right)$$

and thus

$$p(t_n|\mathbf{x}_n) \propto \exp\left(-\exp\left\{-t_n f_m(\mathbf{x}_n)\right\}\right).$$

We can normalize this probability distribution by computing the normalization constant

$$Z = \exp\left(-\exp\left\{f_m(\mathbf{x}_n)\right\}\right) + \exp\left(-\exp\left\{-f_m(\mathbf{x}_n)\right\}\right)$$

but since Z involves  $f_m(\mathbf{x})$ , the log of the resulting normalized probability distribution no longer corresponds to (14.20) as a function of  $f_m(\mathbf{x})$ .

**14.9** The sum-of-squares error for the additive model of (14.21) is defined as

$$E = \frac{1}{2} \sum_{n=1}^{N} (t_n - f_m(\mathbf{x}_n))^2.$$

Using (14.21), we can rewrite this as

$$\frac{1}{2} \sum_{n=1}^{N} (t_n - f_{m-1}(\mathbf{x}_n) - \frac{1}{2} \alpha_m y_m(\mathbf{x}))^2,$$

where we recognize the two first terms inside the square as the residual from the (m-1)-th model. Minimizing this error w.r.t.  $y_m(\mathbf{x})$  will be equivalent to fitting  $y_m(\mathbf{x})$  to the (scaled) residuals.

**14.10** The error function that we need to minimize is

$$E(t) = \frac{1}{2} \sum_{\{t_n\}} (t_n - t)^2.$$

Taking the derivative of this w.r.t. t and setting it equal to zero we get

$$\frac{\mathrm{d}E}{\mathrm{d}t} = -\sum_{\{t_n\}} (t_n - t) = 0.$$

Solving for t yields

$$t = \frac{1}{N'} \sum_{\{t_n\}} t_n$$

where  $N' = |\{t_n\}|$ , i.e. the number of values in  $\{t_n\}$ .

**14.11 NOTE**: In PRML, the text of this exercise contains mistakes; please refer to the PRML Errata for relevant corrections.

The misclassification rates for the two tree models are given by

$$R_{\rm A} = \frac{100 + 100}{400 + 400} = \frac{1}{4}$$

$$R_{\rm B} = \frac{0 + 200}{400 + 400} = \frac{1}{4}$$

From (14.31) and (14.32) we see that the pruning criterion for the cross-entropy case evaluates to

$$C_{\text{Xent}}(T_{\text{A}}) = -2\left(\frac{100}{400}\ln\frac{100}{400} + \frac{300}{400}\ln\frac{300}{400}\right) + 2\lambda \simeq 1.12 + 2\lambda$$

$$C_{\text{Xent}}(T_{\text{B}}) = -\frac{400}{400}\ln\frac{400}{400} - \frac{200}{400}\ln\frac{200}{400} - \frac{0}{400}\ln\frac{0}{400} - \frac{200}{400}\ln\frac{200}{400} + 2\lambda$$

$$\simeq 0.69 + 2\lambda$$

Finally, from (14.31) and (14.33) we see that the pruning criterion for the Gini index case become

$$C_{\text{Gini}}(T_{\text{A}}) = 2 \left[ \frac{300}{400} \left( 1 - \frac{300}{400} \right) + \frac{100}{400} \left( 1 - \frac{100}{400} \right) \right] + 2\lambda = \frac{3}{4} + 2\lambda$$

$$C_{\text{Gini}}(T_{\text{B}}) = \frac{400}{400} \left( 1 - \frac{400}{400} \right) + \frac{200}{400} \left( 1 - \frac{200}{400} \right)$$

$$+ \frac{0}{400} \left( 1 - \frac{0}{400} \right) + \frac{200}{400} \left( 1 - \frac{200}{400} \right) + 2\lambda = \frac{1}{2} + 2\lambda.$$

Thus we see that, while both trees have the same misclassification rate, B performs better in terms of cross-entropy as well as Gini index.

**14.12** Drawing on (3.32), we redefine (14.34) as

$$p(\mathbf{t}|\boldsymbol{\theta}) = \sum_{k=1}^{K} \pi_k \mathcal{N}\left(\mathbf{t}|\mathbf{W}^{\mathrm{T}}\boldsymbol{\phi}, \beta^{-1}\mathbf{I}\right)$$

and then make the corresponding changes to (14.35)–(14.37) and  $Q(\theta, \theta^{\text{old}})$ ; also **t** will be replaced by **T** to align with the notation used in Section 3.1.5. Equation (14.39) will now take the form

$$Q\left(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}\right) = \sum_{n=1}^{N} \gamma_{nk} \left\{ -\frac{\beta}{2} \left\| \mathbf{t}_{n} - \mathbf{W}^{\text{T}} \boldsymbol{\phi}_{n} \right\|^{2} \right\} + \text{const.}$$

Following the same steps as in the single target case, we arrive at a corresponding version of (14.42):

$$\mathbf{W}_k = \left(\mathbf{\Phi}^{\mathrm{T}} \mathbf{R}_k \mathbf{\Phi}\right)^{-1} \mathbf{\Phi}^{\mathrm{T}} \mathbf{R}_k \mathbf{T}.$$

For  $\beta$ , (14.43) becomes

$$Q\left(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}\right) = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \left\{ \frac{D}{2} \ln \beta - \frac{\beta}{2} \left\| \mathbf{t}_{n} - \mathbf{W}^{\text{T}} \boldsymbol{\phi}_{n} \right\|^{2} \right\}$$

and consequently (14.44) becomes

$$\frac{1}{\beta} = \frac{1}{ND} \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \left\| \mathbf{t}_{n} - \mathbf{W}^{\mathrm{T}} \boldsymbol{\phi}_{n} \right\|^{2}.$$

**14.13** Starting from the mixture distribution in (14.34), we follow the same steps as for mixtures of Gaussians, presented in Section 9.2. We introduce a *K*-nomial latent variable, **z**, such that the joint distribution over **z** and *t* equals

$$p(t, \mathbf{z}) = p(t|\mathbf{z})p(\mathbf{z}) = \prod_{k=1}^{K} \left( \mathcal{N}\left(t|\mathbf{w}_{k}^{\mathrm{T}}\boldsymbol{\phi}, \beta^{-1}\right) \pi_{k} \right)^{z_{k}}.$$

Given a set of observations,  $\{(t_n, \phi_n)\}_{n=1}^N$ , we can write the complete likelihood over these observations and the corresponding  $\mathbf{z}_1, \dots, \mathbf{z}_N$ , as



- **14.14** Since  $Q(\theta, \theta^{\text{old}})$  (defined by the unnumbered equation preceding (14.38)) has exactly the same dependency on  $\pi$  as (9.40), (14.38) can be derived just like the corresponding result in Solution 9.9.
- **14.15** The predictive distribution from the mixture of linear regression models for a new input feature vector,  $\hat{\phi}$ , is obtained from (14.34), with  $\phi$  replaced by  $\hat{\phi}$ . Calculating the expectation of t under this distribution, we obtain

$$\mathbb{E}[t|\widehat{\boldsymbol{\phi}}, \boldsymbol{\theta}] = \sum_{k=1}^{K} \pi_k \mathbb{E}[t|\widehat{\boldsymbol{\phi}}, \mathbf{w}_k, \beta].$$

Depending on the parameters, this expectation is potentially K-modal, with one mode for each mixture component. However, the weighted combination of these modes output by the mixture model may not be close to any single mode. For example, the combination of the two modes in the left panel of Figure 14.9 will end up in between the two modes, a region with no signicant probability mass.

**14.16** This solution is analogous to Solution 14.12. It makes use of results from Section 4.3.4 in the same way that Solution 14.12 made use of results from Section 3.1.5. Note, however, that Section 4.3.4 uses k as class index and K to denote the number of classes, whereas here we will use k and k for mixture component indexing and number of mixture components, as used elsewhere in Chapter 14.

Using 1-of-C coding for the targets, we can look to (4.107) to rewrite (14.45) as

$$p(\mathbf{t}|\boldsymbol{\phi}, \boldsymbol{\theta}) = \sum_{k=1}^{K} \pi_k \prod_{c=1}^{C} y_{kc}^{t_c}$$

and making the corresponding changes to (14.46)–(14.48), which lead to an expected complete-data log likelihood function,

$$Q\left(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}\right) = \sum_{n=1}^{N} \sum_{k=1}^{K} \gamma_{nk} \left\{ \ln \pi_k + \sum_{c=1}^{C} t_{nc} \ln y_{nkc} \right\}$$

corresponding to (14.49).

As in the case of the mixture of logistic regression models, the M step for  $\pi$  is the same as for other mixture models, given by (14.50). In the M step for  $\mathbf{W}_1, \dots, \mathbf{W}_K$ , where

$$\mathbf{W}_k = [\mathbf{w}_{k1}, \dots, \mathbf{w}_{kC}]$$

we can again deal with each mixture component separately, using an iterative method such as IRLS, to solve

$$\nabla_{\mathbf{w}_{kc}} Q = \sum_{n=1}^{N} \gamma_{nk} \left( y_{nkc} - t_{nc} \right) \boldsymbol{\phi}_n = 0$$

where we have used (4.109) and (14.51). We obtain the corresponding Hessian from (4.110) (**NOTE**: In the 1<sup>st</sup> printing of PRML, the leading minus sign on the r.h.s. should be removed.) and (14.52) as

$$\mathbf{H}_{k} = \nabla_{\mathbf{w}_{kc}} \nabla_{\mathbf{w}_{k\hat{c}}} Q = \sum_{n=1}^{N} \gamma_{nk} y_{nkc} \left( I_{c\hat{c}} - y_{nk\hat{c}} \right) \boldsymbol{\phi}_{n} \boldsymbol{\phi}_{n}^{\mathrm{T}}.$$

**14.17** If we define  $\psi_k(t|\mathbf{x})$  in (14.58) as

$$\psi_k(t|\mathbf{x}) = \sum_{m=1}^{M} \lambda_{mk} \phi_{mk}(t|\mathbf{x}),$$

we can rewrite (14.58) as

$$p(t|\mathbf{x}) = \sum_{k=1}^{K} \pi_k \sum_{m=1}^{M} \lambda_{mk} \phi_{mk}(t|\mathbf{x})$$
$$= \sum_{k=1}^{K} \sum_{m=1}^{M} \pi_k \lambda_{mk} \phi_{mk}(t|\mathbf{x}).$$

By changing the indexation, we can write this as

$$p(t|\mathbf{x}) = \sum_{l=1}^{L} \eta_l \phi_l(t|\mathbf{x}),$$

where L=KM, l=(k-1)M+m,  $\eta_l=\pi_k\lambda_{mk}$  and  $\phi_l(\cdot)=\phi_{mk}(\cdot)$ . By construction,  $\eta_l\geqslant 0$  and  $\sum_{l=1}^L\eta_l=1$ .

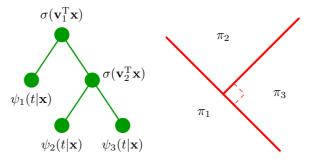
Note that this would work just as well if  $\pi_k$  and  $\lambda_{mk}$  were to be dependent on  $\mathbf{x}$ , as long as they both respect the constraints of being non-negative and summing to 1 for every possible value of  $\mathbf{x}$ .

Finally, consider a tree-structured, hierarchical mixture model, as illustrated in the left panel of Figure 12. On the top (root) level, this is a mixture with two components. The mixing coefficients are given by a linear logistic regression model and hence are input dependent. The left sub-tree correspond to a local conditional density model,  $\psi_1(t|\mathbf{x})$ . In the right sub-tree, the structure from the root is replicated, with the difference that both sub-trees contain local conditional density models,  $\psi_2(t|\mathbf{x})$  and  $\psi_3(t|\mathbf{x})$ .

We can write the resulting mixture model on the form (14.58) with mixing coefficients

$$\pi_1(\mathbf{x}) = \sigma(\mathbf{v}_1^{\mathrm{T}}\mathbf{x}) 
\pi_2(\mathbf{x}) = (1 - \sigma(\mathbf{v}_1^{\mathrm{T}}\mathbf{x}))\sigma(\mathbf{v}_2^{\mathrm{T}}\mathbf{x}) 
\pi_3(\mathbf{x}) = (1 - \sigma(\mathbf{v}_1^{\mathrm{T}}\mathbf{x}))(1 - \sigma(\mathbf{v}_2^{\mathrm{T}}\mathbf{x})),$$

Figure 12 Left: an illustration of a hierarchical mixture model, where the input dependent mixing coefficients are determined by linear logistic models associated with interior nodes; leaf nodes correspond to local (conditional) density models. Right: a possible division of the input space into regions where different mixing coefficients dominate, under the model illustrated left.



where  $\sigma(\cdot)$  is defined in (4.59) and  $\mathbf{v}_1$  and  $\mathbf{v}_2$  are the parameter vectors of the logistic regression models. Note that  $\pi_1(\mathbf{x})$  is independent of the value of  $\mathbf{v}_2$ . This would not be the case if the mixing coefficients were modelled using a single level softmax model,

$$\pi_k(\mathbf{x}) = \frac{e^{\mathbf{u}_k^{\mathrm{T}} \mathbf{x}}}{\sum_{j=1}^{3} e^{\mathbf{u}_j^{\mathrm{T}} \mathbf{x}}},$$

where the parameters  $\mathbf{u}_k$ , corresponding to  $\pi_k(\mathbf{x})$ , will also affect the other mixing coefficients,  $\pi_{j\neq k}(\mathbf{x})$ , through the denominator. This gives the hierarchical model different properties in the modelling of the mixture coefficients over the input space, as compared to a linear softmax model. An example is shown in the right panel of Figure 12, where the red lines represent borders of equal mixing coefficients in the input space. These borders are formed from two straight lines, corresponding to the two logistic units in the left panel of 12. A corresponding division of the input space by a softmax model would involve three straight lines joined at a single point, looking, e.g., something like the red lines in Figure 4.3 in PRML; note that a linear three-class softmax model could not implement the borders show in right panel of Figure 12.