

SOLUTIONS MANUAL FOR
A First Course In
Machine Learning
(2Nd Edition)
Exercise Solutions

_____ by _____

*Simon Rogers and
Mark Girolami*



SOLUTIONS MANUAL FOR
A First Course In
Machine Learning
(2Nd Edition)
Exercise Solutions

by

*Simon Rogers and
Mark Girolami*



CRC Press

Taylor & Francis Group
Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business

CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2017 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works

Printed on acid-free paper
Version Date: 20160404

International Standard Book Number-13: 978-1-4987-3859-0 (Ancillary)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

Chapter 1

EX 1.1. A high positive value of w_0 and a small negative value for w_1 . These reflect the high intercept on the t axis (corresponding to the theoretical time winning time at $x = 0$ and the small decrease in winning time over the years.

EX 1.2. The following would do the job:

```

1      % Attributes are stored in Nx1 vector x
2      % Targets are stored in Nx1 vector t
3      xb = mean(x);
4      tb = mean(t);
5      x2b = mean(x.*x);
6      xtb = mean(x.*t);
7      w1 = (xtb - tb*xb) / (x2b - xb^2);
8      w0 = tb - w1*xb;
9      % Plot the data
10     plot(x,t,'b.','markersize',25);
11     % Plot the model
12     hold on;
13     plot(x,w0+w1*x,'r','linewidth',2);

```

EX 1.3. We need to find $\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}$. We'll start with $\mathbf{X}^T \mathbf{X}$. Multiplying \mathbf{X}^T by \mathbf{X} gives:

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} \sum_{n=1}^N x_{n1}^2 & \sum_{n=1}^N x_{n1} x_{n2} \\ \sum_{n=1}^N x_{n2} x_{n1} & \sum_{n=1}^N x_{n2}^2 \end{bmatrix}$$

Multiplying this by \mathbf{w} gives:

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \begin{bmatrix} w_0 \sum_{n=1}^N x_{n1}^2 + w_1 \sum_{n=1}^N x_{n1} x_{n2} \\ w_0 \sum_{n=1}^N x_{n2} x_{n1} + w_1 \sum_{n=1}^N x_{n2}^2 \end{bmatrix}$$

Finally, pre-multiplying this by \mathbf{w}^\top gives:

$$\begin{aligned}\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} &= w_0 \left(w_0 \sum_{n=1}^N x_{n1}^2 + w_1 \sum_{n=1}^N x_{n1} x_{n2} \right) + \\ &\quad w_1 \left(w_0 \sum_{n=1}^N x_{n2} x_{n1} + w_1 \sum_{n=1}^N x_{n2}^2 \right) \\ &= w_0^2 \sum_{n=1}^N x_{n1}^2 + 2w_0 w_1 \sum_{n=1}^N x_{n1} x_{n2} + w_1^2 \sum_{n=1}^N x_{n2}^2\end{aligned}$$

as required.

EX 1.4. Let's first work out $\mathbf{X}\mathbf{w}$:

$$\mathbf{X}\mathbf{w} = \begin{bmatrix} w_0 x_{11} + w_1 x_{12} \\ w_0 x_{21} + w_1 x_{22} \\ \vdots \\ w_0 x_{N1} + w_1 x_{N2} \end{bmatrix}$$

Therefore

$$(\mathbf{X}\mathbf{w})^\top = [w_0 x_{11} + w_1 x_{12}, w_0 x_{21} + w_1 x_{22}, \dots, w_0 x_{N1} + w_1 x_{N2}]$$

Finally, work out $\mathbf{w}^\top \mathbf{X}^\top$:

$$\mathbf{w}^\top \mathbf{X}^\top = [w_0 x_{11} + w_1 x_{12}, w_0 x_{21} + w_1 x_{22}, \dots, w_0 x_{N1} + w_1 x_{N2}]$$

as required.

EX 1.5. Starting with $\sum_n \mathbf{x}_n t_n$. The result of this is a column vector of the same size as \mathbf{x} (2×1). Now, using the definition of \mathbf{X} ,

$$\mathbf{X}^\top = \begin{bmatrix} x_{11}, x_{21}, \dots, x_{N1} \\ x_{12}, x_{22}, \dots, x_{N2} \end{bmatrix}$$

(which is a $2 \times N$ vector). Multiplying this by \mathbf{t} gives a 2×1 vector that looks like this:

$$\mathbf{X}^\top \mathbf{t} = \begin{bmatrix} \sum_{n=1}^N x_{n1} t_n \\ \sum_{n=1}^N x_{n2} t_n \end{bmatrix}$$

which is $\sum_n \mathbf{x}_n t_n$ as required. The second example, $\mathbf{X}^\top \mathbf{X} \mathbf{w}$. We already know what $\mathbf{X}^\top \mathbf{X} \mathbf{w}$ is (Exercise 1.3)

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \begin{bmatrix} w_0 \sum_{n=1}^N x_{n1}^2 + w_1 \sum_{n=1}^N x_{n1} x_{n2} \\ w_0 \sum_{n=1}^N x_{n2} x_{n1} + w_1 \sum_{n=1}^N x_{n2}^2 \end{bmatrix}$$

Now, $\mathbf{x}_n \mathbf{x}_n^\top$ is the following matrix:

$$\mathbf{x}_n \mathbf{x}_n^\top = \begin{bmatrix} x_{n1}^2 & x_{n1} x_{n2} \\ x_{n2} x_{n1} & x_{n2}^2 \end{bmatrix}$$

Multiplying this by \mathbf{w} gives:

$$\mathbf{x}_n \mathbf{x}_n^T \mathbf{w} = \begin{bmatrix} w_0 x_{n1}^2 + w_1 x_{n1} x_{n2} \\ w_0 x_{n2} x_{n1} + w_1 x_{n2}^2 \end{bmatrix}$$

Summing over the N terms leads us to the matrix we derived previously.

EX 1.6. Code below:

```

1 %% Women's 100m data
2 % Load all Olympic data
3 load olympics;
4 % Copy the necessary variables
5 x = female100(:,1); % Olympic year
6 t = female100(:,2); % Winning time
7 % Augment x
8 X = [repmat(1,size(x)) x];
9 % Get solution
10 w = inv(X'*X)*X'*t;
```

The fitted model is:

$$t = 40.9242 - 0.0151x$$

EX 1.7. Plugging 2012 and 2016 into the above expression yields winning times of 10.5997 and 10.5394 respectively.

EX 1.8. The men's model is:

$$t = 36.4165 - 0.0133x$$

The women's model is:

$$t = 40.9242 - 0.0151x$$

The women's time is decreasing faster than the men's. Therefore, the women will be faster at the first Olympics after the x that gives identical winning times:

$$\begin{aligned} 40.9242 - 0.0151x &= 36.4165 - 0.0133x \\ x &= 2589 \end{aligned}$$

The next Olympic year after 2589 is (assuming they continue to be held every four years) is the year 2592. The winning times are the unrealistically fast 1.8580 seconds and 1.8628 seconds for women and men respectively.

EX 1.9. Code below (`synthdata_cv.m`):

```

1 clear all;close all;
2 load synthdata
3
4 % Augment x
5 X = repmat(1,size(x));
6 for k = 1:4
7     X = [X x.^k];
8 end
9
```

```

10 % Fit the model
11 w = inv(X'*X)*X'*t;
12
13 % Randomise the data order
14 N = size(X,1);
15 order = randperm(N);
16 sizes = repmat(floor(N/10),1,10);
17 sizes(end) = sizes(end) + N-sum(sizes);
18 sizes = [0 cumsum(sizes)];
19
20 X = repmat(1,size(x));
21
22 loss = zeros(4,10);
23 for poly_order = 1:4
24     % Augment x
25     X = [X x.^poly_order];
26     for k = 1:10 % 10-fold CV
27         % Extract the train and test data
28         traindata = X(order,:);
29         traint = t(order);
30         testdata = X(order(sizes(k)+1:sizes(k+1)),:);
31         testt = t(order(sizes(k)+1:sizes(k+1)));
32         traindata(sizes(k)+1:sizes(k+1),:) = [];
33         traint(sizes(k)+1:sizes(k+1)) = [];
34
35         % Fit the model
36         w = inv(traindata'*traindata)*traindata'*traint;
37
38         % Compute loss on test data
39         predictions = testdata*w;
40         loss(poly_order,k) = sum((predictions - testt).^2);
41     end
42 end
43
44 % Plot the loss
45 plot([1:4],mean(loss,2));

```

EX 1.10. The total loss is

$$\mathcal{L} = \sum_{n=1}^N (t_n - \mathbf{w}^T \mathbf{x}_n)^2.$$

Writing this in matrix form, differentiating and solving gives us:

$$\begin{aligned}
 \mathcal{L} &= (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w}) \\
 &= \mathbf{t}^T \mathbf{t} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{t} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \\
 \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= -2\mathbf{X}^T \mathbf{t} + 2\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{0} \\
 \mathbf{w} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}.
 \end{aligned}$$

This is identical to the value obtained for the average loss. It is not surprising as all we are doing is multiplying the loss by a constant and this will not change the value of \mathbf{w} at the minimum.

EX 1.11. The loss is given by

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N \alpha_n (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

If we define the matrix:

$$\mathbf{A} = \begin{bmatrix} \alpha_1 & 0 & \dots & 0 \\ 0 & \alpha_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \alpha_N \end{bmatrix}$$

we can write the loss in vector/matrix form as:

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N (\mathbf{t} - \mathbf{X}\mathbf{w})^\top \mathbf{A} (\mathbf{t} - \mathbf{X}\mathbf{w})$$

Multiplying out, differentiating, equating to zero and solving:

$$\begin{aligned} \mathcal{L} &= \frac{1}{N} (\mathbf{t}^\top \mathbf{A} \mathbf{t} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{A} \mathbf{t} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{A} \mathbf{X} \mathbf{w}) \\ \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= -\frac{2}{N} \mathbf{X}^\top \mathbf{A} \mathbf{t} + \frac{2}{N} \mathbf{X}^\top \mathbf{A} \mathbf{X} \mathbf{w} = \mathbf{0} \\ \mathbf{w} &= (\mathbf{X}^\top \mathbf{A} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{A} \mathbf{t}. \end{aligned}$$

Try this out in Matlab (set some α_n very low and some very high) to see the effect on the solution.

EX 1.12. Code below (**regls100m.m**):

```

1 clear all;close all;
2 load olympics;
3 % Extract men's 100m data
4 x = male100(:,1);
5 t = male100(:,2);
6
7 % Choose number of folds
8 K = 5;
9
10 % Randomise the data order
11 N = size(x,1);
12 order = randperm(N);
13 sizes = repmat(floor(N/K),1,K);
14 sizes(end) = sizes(end) + N-sum(sizes);
15 sizes = [0 cumsum(sizes)];
16
17 % Rescale x
18 x = x - x(1);
19 x = x./4;
20
21 X = [repmat(1,size(x)) x];
22 % Comment out the following line for linear
23 X = [X x.^2 x.^3 x.^4];

```

```
24
25 % Scan a wide range of values of the regularisation parameter
26 regvals = 10.^[-12:1:12];
27
28 for r = 1:length(regvals)
29     for k = 1:K
30         % Extract the train and test data
31         traindata = X(order,:);
32         traint = t(order);
33         testdata = X(order(sizes(k)+1:sizes(k+1)),:);
34         testt = t(order(sizes(k)+1:sizes(k+1)));
35         traindata(sizes(k)+1:sizes(k+1),:) = [];
36         traint(sizes(k)+1:sizes(k+1)) = [];
37
38         % Fit the model
39         w = inv(traindata'*traindata + regvals(r)*eye(size(X,2)))*...
40             traindata'*traint;
41
42         % Compute loss on test data
43         predictions = testdata*w;
44         loss(r,k) = sum((predictions - testt).^2);
45     end
46 end
```

Chapter 2

EX 2.1. The errors are real valued and hence a continuous random variable would be more appropriate.

EX 2.2. If all outcomes are equally likely, they have the same probability of occurring. Defining Y to be the random variable taking the value shown on a die, we can state the following:

$$P(Y = y) = r,$$

where r is a constant. From the definition of probabilities, we know that:

$$\sum_{y=1}^6 P(Y = y) = 1.$$

Substituting r into this gives us the following:

$$\sum_{y=1}^6 r = 1, \quad 6r = 1, \quad r = 1/6.$$

EX 2.3. (a) Y is a discrete random variable that can take any value from 0 to inf. The probability that $Y \leq 4$ is equal to the sum of all of the probabilities that satisfy $Y \leq 4$, $Y = 0, Y = 1, Y = 2, Y = 3, Y = 4$:

$$P(Y \leq 4) = \sum_{y=0}^4 P(Y = y).$$

When $\lambda = 5$, we can compute these probabilities as:

$$P(Y \leq 4) = 0.0067379 + 0.0336897 + 0.0842243 + 0.1403739 + 0.1754674 = 0.44049.$$

(b) Because Y has to satisfy either $P(Y \leq 4)$ or $P(Y > 4)$, we know that $P(Y > 4) = 1 - P(Y \leq 4)$:

$$P(Y > 4) = 0.5591.$$

EX 2.4. We require $\mathbf{E}_{p(y)} \{ \sin(y) \}$ where $p(y) = \mathcal{U}(a, b)$. The uniform density is given by:

$$p(y) = \begin{cases} \frac{1}{b-a} & a \leq y \leq b \\ 0 & \text{otherwise} \end{cases}$$

The required expectation is given by:

$$\begin{aligned}
 \mathbf{E}_{p(y)} \{ \sin(y) \} &= \int \sin(y) p(y) dy \\
 &= \int_{y=a}^b \sin(y) \frac{1}{b-a} dy \\
 &= \frac{1}{b-a} [-\cos(y)]_a^b \\
 &= \frac{\cos(a) - \cos(b)}{b-a}.
 \end{aligned}$$

When $a = 0$, $b = 1$, this is equal to

$$\mathbf{E}_{p(y)} \{ \sin(y) \} = \frac{\cos(0) - \cos(1)}{1} = 0.45970.$$

Code to compute a sample-based approximation below (`sampleexpect.m`):

```

1 clear all;
2 close all;
3 % Compute a sample based approximation to the required expectation
4 u = rand(10000,1); % Take 10000 samples
5 su = sin(u);
6 % Plot how the approximation changes as more samples are used
7 ns = 10:100:10000;
8 stages = zeros(size(ns));
9 for i = 1:length(ns)
10     stages(i) = mean(su(1:ns(i)));
11 end
12 plot(ns,stages)
13 % Plot the true value
14 hold on
15 plot([0 ns(end)], [0.4597 0.4597], 'k—')
```

EX 2.5. The multivariate Gaussian pdf is given by:

$$p(\mathbf{w}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right\}.$$

Setting $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$ gives:

$$p(\mathbf{w}) = \frac{1}{(2\pi)^{D/2} |\sigma^2 \mathbf{I}|^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{w} - \boldsymbol{\mu})^T \mathbf{I}^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right\}.$$

Because it only has elements on the diagonal, the determinant of $\sigma^2 \mathbf{I}$ is given by the product of these diagonal elements. As they are all the same, $|\sigma^2 \mathbf{I}|^{1/2} = \left(\prod_{d=1}^D \sigma^2 \right)^{1/2} = (\sigma^2)^{D/2}$. $\mathbf{I}^{-1} = \mathbf{I}$ and multiplying a vector/matrix by \mathbf{I} leaves the matrix/vector unchanged. Therefore, the argument within the expectation can be written as $-\frac{1}{2\sigma^2} (\mathbf{w} - \boldsymbol{\mu})^T (\mathbf{w} - \boldsymbol{\mu})$ and recalling that $\mathbf{b}^T \mathbf{b} = \sum_i b_i^2$, we can rewrite the pdf as:

$$p(\mathbf{w}) = \frac{1}{(2\pi)^{D/2} (\sigma^2)^{D/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{d=1}^D (w_d - \mu_d)^2 \right\}.$$

Where w_d and μ_d are the d th elements of \mathbf{w} and $\boldsymbol{\mu}$ respectively. The exponential of a sum is the same as a product of exponentials. Hence,

$$\begin{aligned} p(\mathbf{w}) &= \frac{1}{(2\pi)^{D/2}(\sigma^2)^{D/2}} \prod_{d=1}^D \exp \left\{ -\frac{1}{2\sigma^2} (w_d - \mu_d)^2 \right\} \\ &= \prod_{d=1}^D \frac{1}{(2\pi)^{1/2}\sigma} \exp \left\{ -\frac{1}{2\sigma^2} (w_d - \mu_d)^2 \right\} \\ &= \prod_{d=1}^D p(w_d | \mu_d, \sigma^2), \end{aligned}$$

where $p(w_d | \mu_d, \sigma^2) = \mathcal{N}(\mu_d, \sigma^2)$. Hence, the diagonal covariance is equivalent to assuming that the elements of \mathbf{w} are distributed as independent, univariate Gaussians with mean μ_d and variance σ^2 .

EX 2.6. Using the same methods as in the previous exercise, we see that the determinant of the covariance matrix is given by $\prod_{d=1}^D \sigma_d^2$ and we have the following:

$$p(\mathbf{w}) = \frac{1}{(2\pi)^{D/2} \left(\prod_{d=1}^D \sigma_d^2 \right)^{1/2}} \exp \left\{ -\frac{1}{2} \sum_{d=1}^D \frac{(w_d - \mu_d)^2}{\sigma_d^2} \right\}$$

Changing the sum to a product leaves us with

$$\begin{aligned} p(\mathbf{w}) &= \frac{1}{(2\pi)^{D/2} \left(\prod_{d=1}^D \sigma_d^2 \right)^{1/2}} \prod_{d=1}^D \exp \left\{ -\frac{1}{2\sigma_d^2} (w_d - \mu_d)^2 \right\} \\ &= \prod_{d=1}^D \frac{1}{(2\pi)^{1/2}\sigma_d} \exp \left\{ -\frac{1}{2\sigma_d^2} (w_d - \mu_d)^2 \right\}. \end{aligned}$$

This is the product of D independent univariate Gaussian densities.

EX 2.7. The Hessian for a general model of our form is given by:

$$-\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

For the linear model, \mathbf{X} is defined as:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$$

Therefore $-\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$ is:

$$\mathbf{X}^\top \mathbf{X} = \begin{bmatrix} N & \sum_{n=1}^N x_n \\ \sum_{n=1}^N x_n & \sum_{n=1}^N x_n^2 \end{bmatrix}$$

The diagonal elements are $-N/\sigma^2$ and $-(1/\sigma^2) \sum_{n=1}^N x_n^2$ which are equivalent (they differ only by multiplication with a negative constant) the expressions obtained in Chapter 1.

EX 2.8. We have N values, x_1, \dots, x_N . Assuming that these values came from a Gaussian, we want to find the maximum likelihood estimate of the μ and want to find the maximum likelihood estimates of the mean and variance of the Gaussian. The Gaussian pdf is:

$$\frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2\sigma^2} (x_n - \mu)^2 \right\}$$

Assuming the IID assumption, the likelihood of all N points is given by a product over the N objects:

$$\prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2\sigma^2} (x_n - \mu)^2 \right\}.$$

We'll work with the log of the likelihood:

$$\log L = \sum_{n=1}^N \left(-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (x_n - \mu)^2 \right)$$

To find the maximum likelihood estimate for μ , we differentiate with respect to μ , equate to zero and solve:

$$\begin{aligned} \frac{\partial \log L}{\partial \mu} &= \sum_{n=1}^N \frac{1}{\sigma^2} (x_n - \mu) \\ 0 &= \frac{1}{\sigma^2} \sum_{n=1}^N (x_n - \mu) \\ 0 &= \sum_{n=1}^N x_n - \sum_{n=1}^N \mu \\ &= \sum_{n=1}^N x_n - N\mu \\ \mu &= \frac{1}{N} \sum_{n=1}^N x_n \end{aligned}$$

Similarly, for σ^2 ,

$$\begin{aligned} \frac{\partial \log L}{\partial \sigma^2} &= \sum_{n=1}^N \left(-\frac{1}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} (x_n - \mu)^2 \right) = 0 \\ N\sigma^2 &= \sum_{n=1}^N (x_n - \mu)^2 \end{aligned} \tag{2.1}$$

$$\sigma^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2 \tag{2.2}$$

EX 2.9. The Bernoulli distribution is defined as:

$$P(X_n = x|r) = r^x(1-r)^{1-x}$$

where x is either 0 or 1. Using the IID assumption, we have:

$$L = \prod_{n=1}^N r^{x_n}(1-r)^{1-x_n}$$

and the log likelihood is:

$$\log L = \sum_{n=1}^N x_n \log r + (1-x_n) \log(1-r)$$

Differentiating with respect to r gives us:

$$\begin{aligned} \frac{\partial \log L}{\partial r} &= \sum_{n=1}^N \left(\frac{x_n}{r} - \frac{1-x_n}{1-r} \right) = 0 \\ \sum_{n=1}^N \frac{x_n}{r} &= \sum_{n=1}^N \frac{1-x_n}{1-r} \\ \sum_{n=1}^N x_n - r \sum_{n=1}^N x_n &= rN - r \sum_{n=1}^N x_n \\ r &= \frac{1}{N} \sum_{n=1}^N x_n. \end{aligned}$$

EX 2.10. The Fisher information is defined as the expectation of the negative second derivative. From the above expression, we can see that the second derivative of the Gaussian likelihood (assuming N observations, x_1, \dots, x_N is:

$$\frac{\partial^2 \log L}{\partial \mu^2} = -\frac{N}{\sigma^2}.$$

Hence the Fisher information is equal to N/σ^2 .

EX 2.11. Starting from the second expression, we have

$$\widehat{\sigma^2} = \frac{1}{N} \left[\sum_{n=1}^N t_n^2 - 2 \sum_{n=1}^N t_n \mathbf{x}_n^T \widehat{\mathbf{w}} + \sum_{n=1}^N (\mathbf{x}_n \widehat{\mathbf{w}})^2 \right].$$

Concentrating on the final term,

$$\begin{aligned}
 \sum_{n=1}^N (\mathbf{x}^\top \hat{\mathbf{w}})^2 &= \sum_{n=1}^N \mathbf{x}_n^\top \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{x}_n \\
 &= \text{Tr}(\mathbf{X} \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{X}^\top) \\
 &= \text{Tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} \mathbf{t}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\
 &= \text{Tr}(\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} \mathbf{t}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}) \\
 &= \text{Tr}(\mathbf{X}^\top \mathbf{t} \mathbf{t}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}) \\
 &= \text{Tr}(\mathbf{t} \mathbf{t}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\
 &= \text{Tr}(\mathbf{t}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}) \\
 &= \mathbf{t}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} \\
 &= \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}} \\
 &= \sum_{n=1}^N t_n \mathbf{x}_n^\top \hat{\mathbf{w}}.
 \end{aligned}$$

Therefore,

$$\hat{\sigma}^2 = \frac{1}{N} \left[\sum_{n=1}^N t_n^2 - \sum_{n=1}^N t_n \mathbf{x}_n^\top \hat{\mathbf{w}} \right].$$

Now, $\sum_{n=1}^N t_n^2 = \mathbf{t}^\top \mathbf{t}$ and we already know that $\sum_{n=1}^N t_n \mathbf{x}_n^\top \hat{\mathbf{w}} = \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}}$. So,

$$\hat{\sigma}^2 = \frac{1}{N} [\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}}],$$

as required.

EX 2.12. Code below (`predvar.m`):

```

1 clear all;close all;
2 % Relevant code extracted from predictive_variance_example.m
3 x = rand(50,1)*10-5;
4 x = sort(x);
5 % Compute true function values
6 f = 5*x.^3 - x.^2 + x;
7 % Generate some test locations
8 testx = [min(x):0.2:max(x)]';
9 % Add some noise
10 t = f+randn(50,1)*sqrt(1000);
11 % Remove all training data between -1.5 and 1.5
12 pos = find(x>-1.5 & x<1.5);
13 x(pos) = [];
14 f(pos) = [];
15 t(pos) = [];
16
17 % Choose model order
18 K = 5;
19
20 X = repmat(1,size(x));
21 testX = repmat(1,size(testx));

```



```

22 for k = 1:K
23     X = [X x.^k];
24     testX = [testX testx.^k];
25 end
26
27
28 w_hat = inv(X'*X)*X'*t;
29 ss_hat = mean((t - X*w_hat).^2);
30 pred_va = ss_hat*diag(testX*inv(X'*X)*testX');
31 % Make a plot
32 figure(1);hold off
33 plot(x,t,'b. ');
34 hold on
35 errorbar(testx,testX*w_hat,pred_va,'r');

```

EX 2.13. The Bernoulli distribution for a binary random variable x is:

$$p(x|\theta) = \theta^x (1 - \theta)^{1-x}$$

The Fisher information is defined as the negative expected value of the second derivative of the log density evaluated at some parameter value:

$$\mathcal{F} = -\mathbf{E}_{p(x|\theta)} \left\{ \frac{\partial^2 \log p(x|\theta)}{\partial \theta^2} \bigg|_{\theta} \right\}$$

Differentiating $\log p(x|\theta)$ twice gives:

$$\begin{aligned} \frac{\partial \log p(x|\theta)}{\partial \theta} &= \frac{x}{\theta} - \frac{1-x}{1-\theta} \\ \frac{\partial^2 \log p(x|\theta)}{\partial \theta^2} &= -\frac{x}{\theta^2} - \frac{1-x}{(1-\theta)^2}. \end{aligned}$$

The Fisher information is therefore:

$$\mathcal{F} = \frac{1}{\theta^2} \mathbf{E}_{p(x|\theta)} \{x\} + \frac{1}{(1-\theta)^2} \mathbf{E}_{p(x|\theta)} \{1-x\}.$$

Substituting in the expectations (θ and $1-\theta$ respectively) gives:

$$\mathcal{F} = \frac{\theta}{\theta^2} + \frac{1-\theta}{(1-\theta)^2} = \frac{1}{\theta(1-\theta)}$$

EX 2.14. The multivariate Gaussian pdf is given by:

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}.$$

Logging and removing terms not including $\boldsymbol{\mu}$:

$$\log p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}.$$

First and second derivatives are:

$$\begin{aligned}\frac{\partial \log p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}} &= \boldsymbol{\Sigma}^{-1} \mathbf{x} - \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \\ \frac{\partial^2 \log p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu} \partial \boldsymbol{\mu}^\top} &= -\boldsymbol{\Sigma}^{-1}.\end{aligned}$$

Therefore, the Fisher information is:

$$\mathcal{F} = \boldsymbol{\Sigma}^{-1}.$$

Chapter 3

EX 3.1. Our Beta prior (with $\alpha = \beta = 1$) is defined as:

$$p(r) = 1 \quad (0 \leq r \leq 1)$$

The binomial likelihood is given by:

$$P(Y = y|r, N) = \binom{N}{y} r^y (1-r)^{N-y}$$

We know that the posterior density for r is proportional to the likelihood multiplied by the prior, and we also know, because our prior is a particular Beta density, and the Beta prior is conjugate to the binomial likelihood that the posterior must also be a Beta density:

$$\begin{aligned} p(r|Y, N) &\propto p(Y = y|r, N)p(r) \\ &\propto r^y (1-r)^{N-y} \times 1 \\ &= r^{\alpha'-1} (1-r)^{\beta'-1} \end{aligned}$$

suggesting that the posterior is a Beta density with parameters $\alpha' = y + 1$ and $\beta' = N - y + 1$.

EX 3.2. Using the same steps as the previous exercise:

$$\begin{aligned} p(r|Y, N) &\propto p(Y = y|r, N)p(r) \\ &\propto r^y (1-r)^{N-y} \times 2r \\ &\propto r^{y+1} (1-r)^{N-y} \\ &= r^{\alpha'-1} (1-r)^{\beta'-1} \end{aligned}$$

which suggests a Beta density with parameters $\alpha' = y + 2$ and $\beta' = N - y + 1$. To find the prior parameters corresponding to the prior $p(r) = 2r$, consider the form of the Beta density (ignoring the constant term):

$$\begin{aligned} p(r) \propto r^{\alpha-1} (1-r)^{\beta-1} &\propto 2r \\ &\propto r^1 (1-r)^0 \end{aligned}$$

From this, it is clear that the parameters are $\alpha = 2, \beta = 1$.

EX 3.3. Using the same steps as the previous exercise:

$$\begin{aligned}
 p(r|Y, N) &\propto p(Y = y|r, N)p(r) \\
 &\propto r^y(1-r)^{N-y} \times 3r^2 \\
 &\propto r^{y+2}(1-r)^{N-y} \\
 &= r^{\alpha'-1}(1-r)^{\beta'-1}
 \end{aligned}$$

which suggests a Beta density with parameters $\alpha' = y + 3$ and $\beta' = N - y + 1$. To find the prior parameters corresponding to the prior $p(r) = 3r^2$, consider the form of the Beta density (ignoring the constant term):

$$\begin{aligned}
 p(r) \propto r^{\alpha-1}(1-r)^{\beta-1} &\propto 3r^2 \\
 &\propto r^2(1-r)^0
 \end{aligned}$$

From this, it is clear that the parameters are $\alpha = 3, \beta = 1$.

EX 3.4. The effective sample size is $\alpha - 1$ heads and $\beta - 1$ tails. In the first example, this is 0 of each. In the second, we have 1 head and 0 tails and in the third, 2 heads and 0 tails.

EX 3.5. From the definition of expectations,

$$\begin{aligned}
 \mathbf{E}_{p(r)} \{r\} &= \int_{r=0}^{r=1} rp(r) dr \\
 &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_{r=0}^{r=1} r \times r^{\alpha-1}(1-r)^{\beta-1} dr \\
 &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_{r=0}^{r=1} r^\alpha(1-r)^{\beta-1} dr \\
 &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_{r=0}^{r=1} r^{\alpha'-1}(1-r)^{\beta-1} dr \\
 \text{where } \alpha' &= \alpha + 1.
 \end{aligned}$$

Now, the integrand is an unnormalised Beta density so its integral must be the inverse of the Beta normalisation constant. Therefore,

$$\begin{aligned}
 \mathbf{E}_{p(r)} \{r\} &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha')\Gamma(\beta)}{\Gamma(\alpha')\Gamma(\alpha + \beta)} \\
 &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha + 1)\Gamma(\beta)}{\Gamma(\alpha + \beta + 1)}.
 \end{aligned}$$

Now we require the following Gamma identity, $\Gamma(n + 1) = n\Gamma(n)$.

$$\begin{aligned}
 \mathbf{E}_{p(r)} \{r\} &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)} \frac{\alpha}{\alpha + \beta} \\
 &= \frac{\alpha}{\alpha + \beta}.
 \end{aligned}$$

EX 3.6. We require

$$\text{var}\{r\} = \mathbf{E}_{p(r)}\{r^2\} - (\mathbf{E}_{p(r)}\{r\})^2.$$

From the previous Exercise, the second term is

$$(\mathbf{E}_{p(r)}\{r\})^2 = \left(\frac{\alpha}{\alpha + \beta}\right)^2.$$

The first term is computed as follows:

$$\begin{aligned} \mathbf{E}_{p(r)}\{r^2\} &= \int_{r=0}^{r=1} r^2 p(r) \, dr \\ &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_{r=0}^{r=1} r^2 \times r^{\alpha-1} (1-r)^{\beta-1} \, dr \\ &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_{r=0}^{r=1} r^{\alpha+1} (1-r)^{\beta-1} \, dr \\ &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_{r=0}^{r=1} r^{\alpha'-1} (1-r)^{\beta-1} \, dr \\ &\quad \text{where } \alpha' = \alpha + 2. \end{aligned}$$

As in the previous Exercise, the integrand is just an unnormalised Beta density. Therefore,

$$\begin{aligned} \mathbf{E}_{p(r)}\{r^2\} &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha')\Gamma(\beta)}{\Gamma(\alpha')\Gamma(\beta)} \\ &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha + 1 + 1)\Gamma(\beta)}{\Gamma(\alpha + 1 + 1)\Gamma(\beta)} \\ &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha + 1)\Gamma(\beta)}{\Gamma(\alpha + 1)\Gamma(\beta)} \frac{\alpha + 1}{\alpha + \beta + 1} \\ &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)} \frac{\alpha + 1}{\alpha + \beta + 1} \frac{\alpha}{\alpha + \beta} \\ &= \frac{\alpha(\alpha + 1)}{(\alpha + \beta)(\alpha + \beta + 1)}, \end{aligned}$$

where we used the Gamma identity twice. Combining this with the expression for $(\mathbf{E}_{p(r)}\{r\})^2$ gives

$$\begin{aligned} \text{var}\{r\} &= \frac{\alpha(\alpha + 1)}{(\alpha + \beta)(\alpha + \beta + 1)} - \left(\frac{\alpha}{\alpha + \beta}\right)^2 \\ &= \frac{\alpha(\alpha + 1)}{(\alpha + \beta)(\alpha + \beta + 1)} - \frac{\alpha^2}{(\alpha + \beta)^2} \\ &= \frac{\alpha(\alpha + 1)(\alpha + \beta)}{(\alpha + \beta)^2(\alpha + \beta + 1)} - \frac{\alpha^2(\alpha + \beta + 1)}{(\alpha + \beta)^2(\alpha + \beta + 1)} \\ &= \frac{\alpha(\alpha + 1)(\alpha + \beta) - \alpha^2(\alpha + \beta + 1)}{(\alpha + \beta)^2(\alpha + \beta + 1)} \\ &= \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}. \end{aligned}$$

EX 3.7. Assuming that the probability of heads is given by r , we observe y heads in N tosses, and r has a Beta prior with parameters α and β , the posterior density is a Beta density with parameters $\delta = \alpha + y$ and $\gamma = \beta + N - y$. The marginal likelihood is given by:

$$p(y_N|\alpha, \beta) = \binom{N}{y_N} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha + y_N)\Gamma(\beta + N - y_N)}{\Gamma(\alpha + \beta + N)}.$$

The probability of winning is given by:

$$\mathbf{E}_{p(r|y_N)} \{P(Y_{\text{new}} \leq 6|r)\} = 1 - \sum_{y_{\text{new}}=7}^{y_{\text{new}}=10} \mathbf{E}_{p(r|y_N)} \{P(Y_{\text{new}} = y_{\text{new}}|r)\}$$

where

$$\mathbf{E}_{p(r|y_N)} \{P(Y_{\text{new}} = y_{\text{new}}|r)\} = \binom{N_{\text{new}}}{y_{\text{new}}} \frac{\Gamma(\delta + \gamma)}{\Gamma(\delta)\Gamma(\gamma)} \frac{\Gamma(\delta + y_{\text{new}})\Gamma(\gamma + N_{\text{new}} - y_{\text{new}})}{\Gamma(\delta + \gamma + N_{\text{new}})}$$

From the question, $y_N = 9$, $N = 20$.

Scenario 1: $\alpha = \beta = 1$. The posterior density has parameters $\delta = \alpha + y_N = 10$, $\gamma = \beta + N - y_N = 1 + 20 - 9 = 12$. The marginal likelihood comes out as: 0.0476. The probability of winning as: 0.84812.

Scenario 2: $\alpha = \beta = 50$. The posterior density has parameters $\delta = \alpha + y_N = 59$, $\gamma = \beta + N - y_N = 50 + 20 - 9 = 61$. The marginal likelihood comes out as: 0.1486. The probability of winning as: 0.83162.

Scenario 3: $\alpha = 5, \beta = 1$. The posterior density has parameters $\delta = \alpha + y_N = 14$, $\gamma = \beta + N - y_N = 1 + 20 - 9 = 12$. The marginal likelihood comes out as: 0.0135. The probability of winning as: 0.7211.

EX 3.8. Code below (`scenarios.m`):

```

1 clear all;
2 close all;
3 % Define an array to hold the three scenarios
4 % Each row is alpha,beta for one scenario
5 hypers = [1,1;50,50;5,1];
6
7 scenario = 1; % Change this to look at different ones
8
9 % Generate the 100 tosses
10 N = 100;
11 toss = rand(N,1)<0.7;
12 Nheads = sum(toss);
13 Ntails = N - Nheads;
14
15 % Compute the posterior
16 alpha = hypers(scenario,1);
17 beta = hypers(scenario,2);
18 postalpha = alpha + Nheads;
19 postbeta = beta + Ntails;
20
```

```

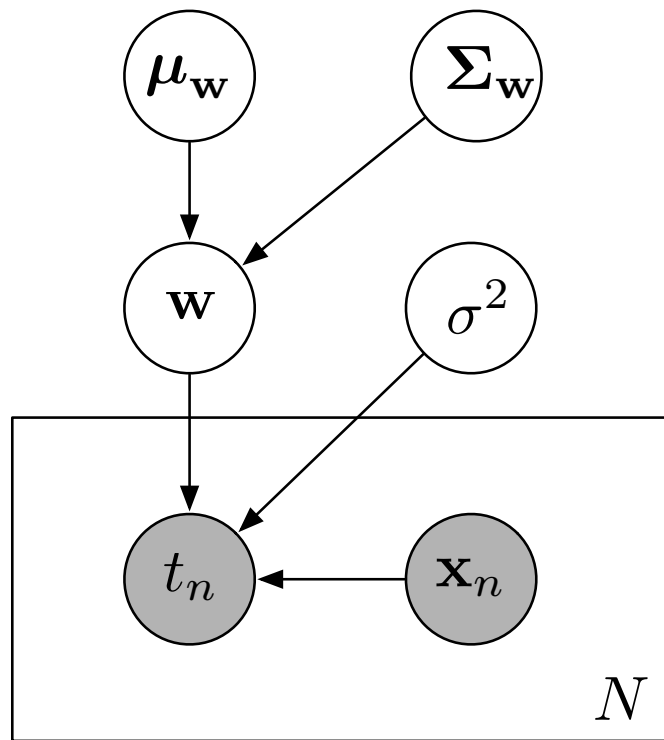
21 % Plot the posterior and prior
22 figure(1);hold off
23 x = [0:0.01:1];
24 plot(x,betapdf(x,hypers(scenario,1),hypers(scenario,2)));
25 hold on
26 plot(x,betapdf(x,postalpha,postbeta),'r')
27 legend('Prior','Posterior');
28
29 % Compute the marginal likelihood (Equation 3.14 in book)
30 % Note we do it in log space first, because of numerical issues
31 ml = log(factorial(N)/(factorial(Nheads)*factorial(Ntails)))+...
32     gammaln(alpha+beta)-gammaln(alpha)-gammaln(beta)+...
33     gammaln(postalpha)+gammaln(postbeta)-gammaln(postalpha+postbeta+N);
34 ml = exp(ml);
35
36 % Probability of winning
37 % Compute prob of y=7,8,9,10
38 tot = 0;
39 for i = 7:10
40     tot = tot + ...
41         exp(...
42             log(factorial(10)/(factorial(i)*factorial(10-i))) + ...
43             gammaln(postalpha+postbeta) - gammaln(postalpha) - gammaln(postbeta) + ...
44             gammaln(postalpha+i) + gammaln(postbeta + 10 - i) - ...
45             gammaln(postalpha+postbeta+10)...
46         )
47 end
48
49 % Compute probability of winning
50 prob = 1 - tot;

```

EX 3.9. Firstly equate the two expressions and then re-arrange to find Σ_0 :

$$\begin{aligned}
 \frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right)^{-1} \mathbf{X}^T \mathbf{t} &= (\mathbf{X}^T \mathbf{X} + N\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t} \\
 \frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right)^{-1} &= (\mathbf{X}^T \mathbf{X} + N\lambda \mathbf{I})^{-1} \\
 (\mathbf{X}^T \mathbf{X} + N\lambda \mathbf{I}) \frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right)^{-1} &= \mathbf{I} \\
 (\mathbf{X}^T \mathbf{X} + N\lambda \mathbf{I}) \frac{1}{\sigma^2} &= \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right) \\
 \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \frac{N\lambda}{\sigma^2} \mathbf{I} &= \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \\
 \Sigma_0 &= \frac{\sigma^2}{N\lambda} \mathbf{I}
 \end{aligned} \tag{3.1}$$

EX 3.10. See figure below:



EX 3.11. Code below (`marglike.m`):

```

1 clear all;close all;
2 % Generate the data
3 N = 100;
4 x = rand(N,1)*10-5;
5 x = sort(x);
6 ss = 100;
7 f = 5*x.^3 - x.^2 + x;
8 t = f + randn(size(f))*sqrt(ss);
9
10 % Choose \sigma_0^2
11 ss0 = 0.3; % Try increasing this
12
13 % Loop over model orders
14 order = [1:7];
15 X = repmat(1,size(x));
16 for i = 1:length(order)
17     mu0 = zeros(order(i)+1,1);
18     s0 = ss0*eye(order(i)+1);
19     X = [X x.^order(i)];
20
21     mlmean = X*mu0;
22     mlcov = ss*eye(N) + X*s0*X';
23
24     logml(i) = -(N/2)*log(2*pi) - (N/2)*log(det(mlcov)) - ...
25         0.5*(t-mlmean)'*inv(mlcov)*(t-mlmean);

```



```

26
27     % Note that the non log value may be numerically instable.
28     ml(i) = exp(logml(i));
29 end

```

EX 3.12. Assuming \mathbf{w} is fixed, and the following inverse gamma prior on σ^2 :

$$p(\sigma^2|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} \exp \left\{ -\frac{\beta}{\sigma^2} \right\}.$$

The likelihood is the standard Gaussian likelihood:

$$p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2) = \frac{1}{(2\pi)^{N/2} (\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}) \right\}.$$

The posterior density over σ^2 is therefore proportional to the product of these two densities. Collecting similar terms:

$$p(\sigma^2|\mathbf{w}, \mathbf{X}, \alpha, \beta) \propto (\sigma^2)^{-\alpha-N/2-1} \exp \left\{ -\frac{1}{\sigma^2} \left(\beta + \frac{1}{2} (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}) \right) \right\}.$$

This is another inverse gamma with parameters α^* and β^* given by:

$$\begin{aligned} \alpha^* &= \alpha + \frac{N}{2} \\ \beta^* &= \beta + \frac{1}{2} (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}). \end{aligned}$$

Chapter 4

EX 4.1. Assuming the following prior density for \mathbf{w} :

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_D),$$

the posterior is:

$$p(\mathbf{w}|t_1, \dots, t_N, \mathbf{x}_1, \dots, \mathbf{x}_N) \propto \mathcal{N}(\mathbf{0}, \mathbf{I}_D) \prod_{n=1}^N p(t_n|\mathbf{w}, \mathbf{x}_n)$$

Combining all of the D -dimensional data objects into an $N \times D$ matrix \mathbf{X} and the labels into a vector \mathbf{t} allows us to write this as the product of two multi-variate Gaussians:

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}) \propto \exp \left\{ -\frac{1}{2} \mathbf{w}^\top \mathbf{w} \right\} \exp \left\{ -\frac{1}{2} (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}) \right\}.$$

Re-arranging and equating coefficients gives the posterior:

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}, \mathbf{X}) &\propto \exp \left\{ -\frac{1}{2} (\mathbf{w}^\top \mathbf{w} + \mathbf{t}^\top \mathbf{t} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{t}) \right\} \\ &= \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \boldsymbol{\Sigma} &= (\mathbf{X}^\top \mathbf{X} + \mathbf{I}_D)^{-1} \\ \boldsymbol{\mu} &= \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{t}. \end{aligned}$$

The Laplace approximation approximates the posterior with a Gaussian with mean equal to the posterior mode and a covariance matrix equal to the negative inverse of the second derivate of the log posterior. Removing constants, the log posterior is equal to:

$$g = -\frac{1}{2} \mathbf{w}^\top \mathbf{w} - \frac{1}{2} (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}).$$

To find the maximum value, we take partial derivatives with respect to \mathbf{w} :

$$\frac{\partial g}{\partial \mathbf{w}} = -\frac{1}{2} (2(\mathbf{I} + \mathbf{X}^\top \mathbf{X})\mathbf{w} - 2\mathbf{X}^\top \mathbf{t}).$$

Setting to zero and solving for \mathbf{w} gives:

$$\mathbf{w} = (\mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}.$$

This will be the mean of the Gaussian approximation and we can already see that it is equal to $\boldsymbol{\mu}$, the posterior mean.

To find the covariance matrix, we take second derivatives of g :

$$\frac{\partial^2 g}{\partial \mathbf{w} \partial \mathbf{w}^\top} = -(\mathbf{I}_D + \mathbf{X}^\top \mathbf{X}).$$

The covariance of the Laplace approximation is therefore:

$$(\mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1}.$$

This is identical to the posterior covariance matrix thus confirming that the approximation is, in this case, exact.

EX 4.2. We are interested in the Laplace approximation to the posterior $p(r|N, y, \alpha, \beta)$. The first step is to find the MAP value of r . The log posterior is (removing constants):

$$g = (\alpha - 1) \log r + (\beta - 1) \log(1 - r) + y \log r + (N - y) \log(1 - r).$$

Taking partial derivatives with respect to r leaves:

$$\frac{\partial g}{\partial r} = \frac{\alpha + y - 1}{r} - \frac{\beta + N - y - 1}{1 - r}.$$

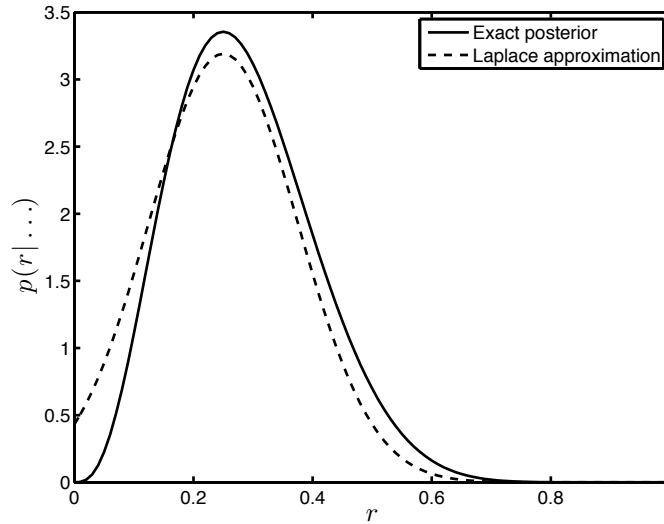
Re-arranging this expression gives the MAP solution for r :

$$\hat{r} = \frac{\alpha + y - 1}{\alpha + \beta + N - 2}.$$

This will be the mean of the Gaussian approximation. The variance is given by:

$$\begin{aligned} \sigma^2 &= - \left(\frac{\partial^2 g}{\partial r^2} \right)^{-1} \\ &= - \left(-\frac{\alpha + y - 1}{r^2} - \frac{\beta + N - y - 1}{(1 - r)^2} \right)^{-1} \end{aligned}$$

EX 4.3. The plot below shows an example where $\alpha = \beta = 2$, $y = 2$, and $N = 10$.



EX 4.4. This is a nice example to illustrate the power of sampling. Imagine sampling pairs of values from uniform distributions between -0.5 and 0.5. These values would be uniformly distributed in the square shown below:

The probability that one of these points lies within the circle must be equal to the ratio of the area of the circle to the area of the square. We can obtain an empirical estimate of this value by say generating N pairs (i.e. points) and then computing the number of pairs that lie within the circle (s). Then:

$$\frac{s}{N} \approx \frac{\pi r^2}{1 \times 1} = 0.5^2 \pi.$$

Re-arranging this expression will give an approximate value of π :

$$\pi \approx \frac{s}{0.5^2 N}.$$

For example, one random sample of 100 points yields a value of:

$$\pi \approx 3.240.$$

As we would expect, this approximation will, on average, get better as we take more samples. An example after 10000 is:

$$\pi \approx 3.129.$$

EX 4.5. Re-arrange as follows:

$$\begin{aligned} P(T_{\text{new}} = 1 | \mathbf{x}_{\text{new}}, \hat{\mathbf{w}}) &= \frac{1}{1 + \exp(-\hat{\mathbf{w}}^T \mathbf{x}_{\text{new}})} \\ 0.5 &= \frac{1}{2} = \frac{1}{1 + \exp(-\hat{\mathbf{w}}^T \mathbf{x}_{\text{new}})} \\ 2 &= 1 + \exp(-\hat{\mathbf{w}}^T \mathbf{x}_{\text{new}}) \\ 1 &= \exp(0) = \exp(-\hat{\mathbf{w}}^T \mathbf{x}_{\text{new}}) \\ \text{therefore} &\quad \hat{\mathbf{w}}^T \mathbf{x}_{\text{new}} = 0. \end{aligned}$$

EX 4.6. Assuming a Gaussian prior:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_D),$$

and:

$$v_n = \exp\{\mathbf{w}^T \mathbf{x}_n\}$$

the log posterior is proportional to:

$$g = -\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_n (t_n \log v_n - v_n).$$

The gradient is obtained by taking partial derivatives with respect to \mathbf{w} :

$$\frac{\partial g}{\partial \mathbf{w}} = -\mathbf{w} + \sum_n \frac{t_n}{v_n} \frac{\partial v_n}{\partial \mathbf{w}} - \frac{\partial v_n}{\partial \mathbf{w}}.$$

Now,

$$\frac{\partial v_n}{\partial \mathbf{w}} = \mathbf{x}_n v_n$$

and therefore:

$$\frac{\partial g}{\partial \mathbf{w}} = -\mathbf{w} + \sum_n \mathbf{x}_n (t_n - v_n).$$

The Hessian is obtained by differentiating a second time. You might find it worthwhile to work out each element of the matrix individually for (say) $D = 2$ if this step appears difficult:

$$\frac{\partial^2 g}{\partial \mathbf{w} \partial \mathbf{w}^\top} = -\mathbf{I} - \sum_n v_n \mathbf{x}_n \mathbf{x}_n^\top$$

Given some initial value of \mathbf{w} , the optimisation scheme repeatedly updates $\hat{\mathbf{w}}$ as follows until it converges to a maximum:

$$\mathbf{w} \leftarrow \mathbf{w} - \left(-\mathbf{I} - \sum_n v_n \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left(-\mathbf{w} + \sum_n \mathbf{x}_n (t_n - v_n) \right)$$

EX 4.7. Given the result of the previous exercise, this is straightforward. The Laplace approximation is a Gaussian with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ where $\boldsymbol{\mu} = \hat{\mathbf{w}}$, the maximum obtained using the Newton-Raphson scheme derived in the previous exercise and the covariance is:

$$\boldsymbol{\Sigma} = - \left(-\mathbf{I} - \sum_n \hat{v}_n \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1},$$

where:

$$\hat{v}_n = \exp\{\hat{\mathbf{w}}^\top \mathbf{x}_n\}.$$

EX 4.8. Example code below (`poisscount.m`):

```

1  %% Laplace and Metropolis-Hastings for a Poisson Counts model
2  % Simon Rogers, August 2011
3  clear all;
4  close all;
5
6  %% Generate the 'true' data.
7  N = 200; % Number of data objects
8  x = [repmat(1,N,1) randn(N,1)];
9  w = randn(2,1);
10
11 temp = exp(x*w);
12
13 t = poissrnd(temp);
14
15 truew = w;
16 %% Newton Raphson
17 % Initialise w
18 w = rand(2,1)*5;
19 wall = [];
```

```

20 change = inf;
21 temp = exp(x*w);
22 ll = -0.5*w'*w + sum(t.*log(temp) - temp);
23 while change > 1e-3
24     % Compute gradient
25     temp = exp(x*w);
26     gra = -w + sum(repmat(t,1,2).*x - repmat(temp,1,2).*x,1)';
27     % Compute Hessian
28     hess = -eye(2) - x'*(x.*repmat(temp,1,2));
29     % Update w
30     neww = w - inv(hess)*gra;
31     % Compute change
32     change = sum((neww-w).^2);
33     w = neww;
34
35     % Store all of the w
36     wall = [wall;w'];
37     % Compute the likelihood (should always increase)
38     temp = exp(x*w);
39     ll = [ll;-0.5*w'*w + sum(t.*log(temp) - temp)];
40 end
41
42 %% Metropolis-Hastings
43 % Initialise w
44 wm = randn(2,1);
45 temp = exp(x*wm);
46 wall_m = [];
47 % Compute old log posterior
48 oldll = -0.5*wm'*wm + sum(t.*log(temp) - temp);
49 for s = 1:10000 % Generate 10000 samples
50     % Propose a new candidate
51     wnew = wm + randn(2,1)*0.1; % Gaussian proposal, std = 0.1
52     % Compute new log likelihood
53     temp = exp(x*wnew);
54     newll = -0.5*wnew'*wnew + sum(t.*log(temp) - temp);
55     % Accept or reject
56     if rand<exp(newll-oldll)
57         wm = wnew;
58         oldll = newll;
59     end
60     % Store all w
61     wall_m = [wall_m;wm'];
62 end
63
64 % Plot samples
65 % No thinning or ignoring of burn-in
66 figure(1);
67 hold off
68 plot(wall_m(:,1),wall_m(:,2),'k.','color',[0.6 0.6 0.6])
69
70
71 %% Plot the approximate posterior
72 % Compute posterior (negative inverse of Hessian)
73 gcov = -inv(-eye(2) - x'*(x.*repmat(temp,1,2)));
74 % Evaluate over a set of points
75 [W1,W2] = meshgrid([-1:0.01:1],[-1:0.01:1]);
76 gli = 1./(sqrt(2*pi)*det(gcov)).*repmat(1,size(W1));
77 for i = 1:size(W1,1)
78     for j = 1:size(W1,2)

```

```
79         wv = [W1(i,j);W2(i,j)];
80         gli(i,j) = gli(i,j).*exp(-0.5*(wv-w)'*inv(gcov)*(wv-w));
81     end
82 end
83 % Plot the contours
84 hold on
85 contour(W1,W2,gli,'k')
86
87
88 % Plot the true value
89 plot(truew(1),truew(2),'ko','markersize',20,'linewidth',2)
```


Chapter 5

EX 5.1. The class conditional density is given by:

$$p(\mathbf{x}_n | t_n = c, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) = \mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c).$$

Assuming a Gaussian prior on $\boldsymbol{\mu}_c$:

$$p(\boldsymbol{\mu}_c | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0),$$

the posterior density for $\boldsymbol{\mu}$ is given by:

$$p(\boldsymbol{\mu}_c | \mathbf{x}_1, \dots, \mathbf{x}_{N_c}, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0, \boldsymbol{\Sigma}_c) \propto \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \prod_{n=1}^{N_c} p(\mathbf{x}_n | t_n = c, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c).$$

Removing constant terms and re-arranging:

$$p(\boldsymbol{\mu}_c | \dots) \propto \exp \left\{ -\frac{1}{2} (\boldsymbol{\mu}_c - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\mu}_c - \boldsymbol{\mu}_0) \right\} \exp \left\{ -\frac{1}{2} \sum_{n=1}^{N_c} (\mathbf{x}_n - \boldsymbol{\mu}_c)^\top (\mathbf{x}_n - \boldsymbol{\mu}_c) \right\}.$$

Keeping only terms dependent on $\boldsymbol{\mu}_c$:

$$p(\boldsymbol{\mu}_c | \dots) \propto \exp \left\{ \frac{1}{2} \left(\boldsymbol{\mu}_c^\top \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_c + N_c \boldsymbol{\mu}_c^\top \boldsymbol{\mu}_c - 2 \boldsymbol{\mu}_c^\top \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 - 2 \boldsymbol{\mu}_c^\top \sum_{n=1}^{N_c} \mathbf{x}_n \right) \right\}.$$

Therefore, the posterior is a Gaussian:

$$p(\boldsymbol{\mu}_c | \dots) = \mathcal{N}(\mathbf{m}, \mathbf{S}),$$

where:

$$\begin{aligned} \mathbf{S} &= (\boldsymbol{\Sigma}_0^{-1} + N_c \mathbf{I})^{-1} \\ \mathbf{m} &= \mathbf{S} \left(\boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \sum_{n=1}^{N_c} \mathbf{x}_n \right). \end{aligned}$$

EX 5.2. This requires the expectation of a Gaussian with respect to another Gaussian:

$$\begin{aligned} p(\mathbf{x}_{\text{new}} | T_{\text{new}} = c, \mathbf{X}, \mathbf{t}) &= \mathbf{E}_{p(\boldsymbol{\mu}_c | \mathbf{m}, \mathbf{S})} \{ p(\mathbf{x}_{\text{new}} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \} \\ &= \int \mathcal{N}_{\boldsymbol{\mu}_c}(\mathbf{m}, \mathbf{S}) \mathcal{N}_{\mathbf{x}_{\text{new}}}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) d\boldsymbol{\mu}_c \\ &= \int \mathcal{N}_{\boldsymbol{\mu}_c}(\mathbf{m}, \mathbf{S}) \mathcal{N}_{\boldsymbol{\mu}_c}(\mathbf{x}_n, \boldsymbol{\Sigma}_c) d\boldsymbol{\mu}_c. \end{aligned}$$

We now make use of the following standard result:

$$\mathcal{N}_{\mathbf{x}}(\mathbf{a}, \mathbf{A})\mathcal{N}_{\mathbf{x}}(\mathbf{b}, \mathbf{B}) = K\mathcal{N}_{\mathbf{x}}(\mathbf{d}, \mathbf{D}),$$

where:

$$K = \mathcal{N}_{\mathbf{a}}(\mathbf{b}, \mathbf{A} + \mathbf{B}).$$

The quantities \mathbf{d} and \mathbf{D} are not important as this density will disappear in our integral. Therefore, our expectation becomes:

$$p(\mathbf{x}_{\text{new}} | T_{\text{new}} = c, \mathbf{X}, \mathbf{t}) = \mathcal{N}_{\mathbf{x}_{\text{new}}}(\mathbf{m}, \mathbf{S} + \mathbf{\Sigma}_c).$$

EX 5.3. The total likelihood for class c is:

$$\mathcal{L} = \prod_{n=1}^{N_c} \frac{1}{(2\pi)^{D/2} |\mathbf{\Sigma}_c|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_c)^\top \mathbf{\Sigma}_c^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_c) \right\}.$$

As normal, it is easiest to maximise the log likelihood. Logging and multiplying out gives:

$$\log \mathcal{L} = -\frac{DN_c}{2} \log(2\pi) - \frac{N_c}{2} \log |\mathbf{\Sigma}_c| - \frac{1}{2} \sum_{n=1}^{N_c} (\mathbf{x}_n^\top \mathbf{\Sigma}_c^{-1} \mathbf{x}_n - 2\boldsymbol{\mu}_c^\top \mathbf{\Sigma}_c^{-1} \mathbf{x}_n + \boldsymbol{\mu}_c^\top \mathbf{\Sigma}_c^{-1} \boldsymbol{\mu}_c).$$

Taking the partial derivative with respect to $\boldsymbol{\mu}_c$, equating to zero and solving gives a solution for $\boldsymbol{\mu}_c$:

$$\begin{aligned} \frac{\partial \log \mathcal{L}}{\partial \boldsymbol{\mu}_c} = 0 &= \mathbf{\Sigma}_c^{-1} \sum_{n=1}^{N_c} \mathbf{x}_n - N_c \mathbf{\Sigma}_c^{-1} \boldsymbol{\mu}_c \\ \mathbf{\Sigma}_c^{-1} \sum_{n=1}^{N_c} \mathbf{x}_n &= N_c \mathbf{\Sigma}_c^{-1} \boldsymbol{\mu}_c \\ \boldsymbol{\mu}_c &= \frac{1}{N_c} \sum_{n=1}^{N_c} \mathbf{x}_n. \end{aligned}$$

Similarly, for $\mathbf{\Sigma}_c$:

$$\frac{\partial \log \mathcal{L}}{\partial \mathbf{\Sigma}_c} = 0 = -\frac{N_c}{2} \mathbf{\Sigma}_c^{-1} + \frac{1}{2} \mathbf{\Sigma}_c^{-1} \sum_{n=1}^{N_c} (\mathbf{x}_n - \boldsymbol{\mu}_c)(\mathbf{x}_n - \boldsymbol{\mu}_c)^\top \mathbf{\Sigma}_c^{-1},$$

where we have used the following two results (for symmetric \mathbf{A}):

$$\begin{aligned} \frac{\partial \log |\mathbf{A}|}{\partial \mathbf{A}} &= \mathbf{A}^{-1} \\ \frac{\partial \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a}}{\partial \mathbf{A}} &= -\mathbf{A}^{-1} \mathbf{a} \mathbf{a}^\top \mathbf{A}^{-1}. \end{aligned}$$

Re-arranging gives the expression for $\mathbf{\Sigma}_c$:

$$\mathbf{\Sigma}_c = \frac{1}{N_c} \sum_{n=1}^{N_c} (\mathbf{x}_n - \boldsymbol{\mu}_c)(\mathbf{x}_n - \boldsymbol{\mu}_c)^\top$$

where $\boldsymbol{\mu}_c$ is the maximum likelihood estimate derived above.

EX 5.4. The data consists of N_c M -dimensional vectors of integer counts. The multinomial likelihood for the N_c data objects is:

$$p(\mathbf{x}_1, \dots, \mathbf{x}_{N_c} | q_{c1}, \dots, q_{cM}) = \mathcal{L} \propto \prod_{n=1}^{N_c} \prod_{m=1}^M q_{cm}^{x_{nm}},$$

(where the constant that does not depend on q_{cm} has been omitted for brevity). Once again, we work with the log likelihood and take partial derivatives with respect to q_{cm} . We must also add a Lagrangian term due to the constraining that $\sum_m q_{cm} = 1$. Calling this new objective function g gives:

$$\begin{aligned} \log \mathcal{L} &= \sum_{n=1}^{N_c} \sum_{m=1}^M x_{nm} \log q_{cm} \\ g &= \sum_{n=1}^{N_c} \sum_{m=1}^M x_{nm} \log q_{cm} - \lambda \left(\sum_m q_{cm} - 1 \right) \\ \frac{\partial \log \mathcal{L}}{\partial q_{cm}} = 0 &= \sum_{n=1}^{N_c} \frac{x_{nm}}{q_{cm}} - \lambda \\ q_{cm} &= \frac{\sum_{n=1}^{N_c} x_{nm}}{\lambda}. \end{aligned}$$

To compute λ we sum both sides over m :

$$\begin{aligned} \sum_m q_{cm} &= \frac{1}{\lambda} \sum_{m,n} x_{nm} \\ 1 &= \frac{1}{\lambda} \sum_{m,n} x_{nm} \\ \lambda &= \frac{1}{\sum_{n,m} x_{nm}}. \end{aligned}$$

Therefore, the maximum likelihood estimate of q_{cm} is:

$$q_{cm} = \frac{\sum_n x_{nm}}{\sum_{n,m'} x_{nm'}}.$$

EX 5.5. The likelihood (ignoring terms not involving q_{cm}) is:

$$p(\mathbf{x}_1, \dots, \mathbf{x}_{N_c} | \mathbf{q}_c) \propto \prod_{n=1}^{N_c} \prod_{m=1}^M q_{cm}^{x_{nm}}.$$

The prior over \mathbf{q}_c is a Dirichlet with constant parameter α (again, ignoring terms not involving q_{cm}):

$$p(\mathbf{q}_c | \alpha) \propto \prod_{m=1}^M q_{cm}^{\alpha-1}.$$

The Dirichlet is the conjugate prior to the multinomial likelihood. Therefore, we know the posterior will be another Dirichlet. Multiplying the prior and likelihood together and combining q_{cm} terms gives:

$$p(\mathbf{q}_c | \alpha, \mathbf{x}_1, \dots, \mathbf{x}_{N_c}) \propto \prod_{m=1}^M q_{cm}^{\alpha-1+\sum_n x_{nm}}.$$

This is the form of a Dirichlet, where the m th parameter is $\beta_m = \alpha + \sum_n x_{nm}$.

EX 5.6. The required expectation is:

$$p(\mathbf{x}_{\text{new}} | T_{\text{new}} = c, \mathbf{X}, \mathbf{t}) = \mathbf{E}_{p(\mathbf{q}_c | \mathbf{X}^c, \alpha)} \{p(\mathbf{x}_{\text{new}} | \mathbf{q}_c)\}.$$

Writing out the expectation in full gives:

$$\begin{aligned} p(\mathbf{x}_{\text{new}} | T_{\text{new}} = c, \mathbf{X}, \mathbf{t}) &= \int \left[\frac{\Gamma(\sum_m \beta_m)}{\prod_m \Gamma(\beta_m)} \prod_m q_{cm}^{\beta_m-1} \right] \left[\frac{(\sum_m x_{\text{new},m})!}{\prod_m (x_{\text{new},m}!)} \prod_m q_{cm}^{x_{\text{new},m}} \right] d\mathbf{q}_c \\ &= \frac{(\sum_m x_{\text{new},m})!}{\prod_m (x_{\text{new},m}!)} \frac{\Gamma(\sum_m \beta_m)}{\prod_m \Gamma(\beta_m)} \int \prod_m q_{cm}^{\beta_m+x_{\text{new},m}-1} d\mathbf{q}_c \\ &= \frac{(\sum_m x_{\text{new},m})!}{\prod_m (x_{\text{new},m}!)} \frac{\Gamma(\sum_m \beta_m)}{\prod_m \Gamma(\beta_m)} \frac{\prod_m \Gamma(\beta_m + x_{\text{new},m})}{\Gamma(\sum_m \beta_m + x_{\text{new},m})}. \end{aligned}$$

EX 5.7. The MAP estimate is obtained by maximising the likelihood multiplied by the prior. Extracting only the terms that include q_{cm} , and logging give us:

$$\log \mathcal{L} \propto \sum_m (\alpha - 1) \log q_{cm} + \sum_m \left(\sum_n x_{nm} \right) \log q_{cm}.$$

Introducing the necessary Lagrangian (to ensure $\sum_m q_{cm} = 1$) we end up with the following objective function:

$$g = \sum_m \left(\sum_n x_{nm} + \alpha - 1 \right) \log q_{cm} - \lambda \left(\sum_m q_{cm} - 1 \right).$$

Taking partial derivatives, equating to zero and solving gives:

$$\begin{aligned} \frac{\partial g}{\partial q_{cm}} = 0 &= \frac{\alpha - 1 + \sum_n x_{nm}}{q_{cm}} - \lambda \\ q_{cm} &= \frac{\alpha - 1 + \sum_n x_{nm}}{M(\alpha - 1) + \sum_{n,m'} x_{nm'}}, \end{aligned}$$

where we summed both sides over m to obtain the value for λ .

EX 5.8. The starting point is:

$$\underset{\mathbf{w}}{\operatorname{argmin}} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_n \xi_n$$

subject to:

$$\xi_n \geq 0, \quad t_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi_n.$$

Adding the constraints as Lagrangian terms gives:

$$\frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_n \xi_n - \sum_n \alpha_n (t_n(\mathbf{w}^\top \mathbf{x}_n + b) - 1 + \xi_n) - \sum_n \xi_n \gamma_n$$

subject to $\alpha_n \geq 0$, $\gamma_n \geq 0$, which must be minimised with respect to \mathbf{w} , b and ξ_n and maximised with respect to α_n and γ_n . At the solution:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} &= \mathbf{w} - \sum_n \alpha_n t_n = 0 \\ \frac{\partial}{\partial b} &= - \sum_n \alpha_n t_n = 0 \\ \frac{\partial}{\partial \xi_n} &= C - \alpha_n - \gamma_n = 0. \end{aligned}$$

Substituting the first and third of these expressions back into the objective function and re-arranging gives:

$$-\frac{1}{2} \sum_{n,m} \alpha_n \alpha_m t_n t_m \mathbf{x}_n^\top \mathbf{x}_m + \sum_n \alpha_n$$

which must be maximised with respect to α subject to:

$$\sum_n \alpha_n t_n = 0, \quad 0 \leq \alpha_n \leq C.$$

Where the first constraint comes from the partial derivative with respect to b and the second from the partial derivative with respect to ξ_n and the fact that $\gamma_n \geq 0$.

Chapter 6

EX 6.1. The log likelihood for all N data objects is:

$$\mathcal{L} = \sum_n \log \sum_k \pi_k \prod_d p(x_{nd} | \mu_{kd}, \sigma_{kd}^2).$$

Inserting a set of variational parameters, q_{nk} such that $\sum_k q_{nk} = 1$ and $q_{nk} \geq 0$ and subsequently applying Jensen's inequality results in:

$$\begin{aligned} \mathcal{L} &= \sum_n \log \sum_k \pi_k \prod_d \frac{q_{nk}}{q_{nk}} p(x_{nd} | \mu_{kd}, \sigma_{kd}^2) \\ &\geq \sum_n \sum_k q_{nk} \log \pi_k + q_{nk} \sum_d \log p(x_{nd} | \mu_{kd}, \sigma_{kd}^2) - q_{nk} \log q_{nk}. \end{aligned}$$

Only the second term involves σ_{kd}^2 . Writing this term out in full leaves:

$$\sum_n \sum_k q_{nk} \sum_d \left(-\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma_{kd}^2 - \frac{1}{2\sigma_{kd}^2} (x_{nd} - \mu_{kd})^2 \right).$$

Taking partial derivatives with respect to σ_{kd}^2 results in:

$$\frac{\partial}{\partial \sigma_{kd}^2} = -\frac{\sum_n q_{nk}}{2\sigma_{kd}^2} + \sum_n \frac{q_{nk}}{2(\sigma_{kd}^2)^2} (x_{nd} - \mu_{kd})^2.$$

Equating to zero and re-arranging results in:

$$\sigma_{kd}^2 = \frac{\sum_n q_{nk} (x_{nd} - \mu_{kd})^2}{\sum_n q_{nk}}.$$

EX 6.2. Starting from the log likelihood above and applying Jensen's inequality:

$$\begin{aligned} \mathcal{L} &= \sum_n \log \sum_k \pi_k \prod_d \frac{q_{nk}}{q_{nk}} p(x_{nd} | \mu_{kd}, \sigma_k^2) \\ &\geq \sum_n \sum_k q_{nk} \log \pi_k + q_{nk} \sum_d \log p(x_{nd} | \mu_{kd}, \sigma_k^2) - q_{nk} \log q_{nk}. \end{aligned}$$

σ_k^2 only appears in the second term. Expanding:

$$\sum_n \sum_k q_{nk} \left(-\frac{D}{2} \log 2\pi - \frac{D}{2} \log \sigma_k^2 - \frac{1}{2\sigma_k^2} \sum_d (x_{nd} - \mu_{kd})^2 \right).$$

Taking partial derivatives with respect to σ_k^2 gives:

$$\frac{\partial}{\partial \sigma_k^2} = -\frac{D}{2\sigma_k^2} \sum_n q_{nk} + \frac{1}{2(\sigma_k^2)^2} \sum_{d,n} (x_{nd} - \mu_{kd})^2.$$

Equating to zero and re-arranging gives:

$$\sigma_k^2 = \frac{\sum_n q_{nk} \sum_d (x_{nd} - \mu_{kd})^2}{D \sum_n q_{nk}}$$

EX 6.3. The log likelihood is given by:

$$\mathcal{L} = \sum_n \log \sum_k \pi_k \prod_d p_{kd}^{x_{nd}} (1 - p_{kd})^{1-x_{nd}}.$$

Applying Jensen's inequality gives:

$$\mathcal{L} \geq \sum_n \sum_k \left[q_{nk} \log \pi_k + q_{nk} \sum_d (x_{nd} \log p_{kd} + (1 - x_{nd}) \log(1 - p_{kd})) - q_{nk} \log q_{nk} \right]$$

Taking partial derivatives with respect to p_{kd} setting to zero and solving gives:

$$\begin{aligned} \frac{\partial}{\partial p_{kd}} = 0 &= \sum_n q_{nk} \frac{x_{nd}}{p_{kd}} - \sum_n q_{nk} \frac{1 - x_{nd}}{1 - p_{kd}} \\ \sum_n q_{nk} \frac{x_{nd}}{p_{kd}} &= \sum_n q_{nk} \frac{1 - x_{nd}}{1 - p_{kd}} \\ \sum_n q_{nk} x_{nd} - p_{kd} \sum_n q_{nk} x_{nd} &= p_{kd} \sum_n q_{nk} - p_{kd} \sum_n q_{nk} x_{nd} \\ p_{kd} &= \frac{\sum_n q_{nk} x_{nd}}{\sum_n q_{nk}} \end{aligned}$$

EX 6.4. This proceeds in exactly the same manner as the previous exercise except that after applying Jensen's inequality, the bound will also include a term from the Beta prior. Including only terms that involve q_{nk} (and removing summations over k and d):

$$(\alpha - 1) \log p_{kd} + (\beta - 1) \log(1 - p_{kd}) + \sum_n q_{nk} (x_{nd} \log p_{kd} + (1 - x_{nd}) \log(1 - p_{kd})).$$

Taking partial derivatives with respect to p_{kd} , equating to zero and solving gives:

$$\begin{aligned} \frac{\partial}{\partial p_{kd}} = 0 &= \frac{\alpha - 1}{p_{kd}} - \frac{\beta - 1}{1 - p_{kd}} + \sum_n q_{nk} \frac{x_{nd}}{p_{kd}} - \sum_n q_{nk} \frac{1 - x_{nd}}{1 - p_{kd}} \\ \sum_n q_{nk} \frac{x_{nd}}{p_{kd}} + \frac{\alpha - 1}{p_{kd}} &= \sum_n q_{nk} \frac{1 - x_{nd}}{1 - p_{kd}} + \frac{\beta - 1}{1 - p_{kd}} \\ (\alpha - 1) - p_{kd}(\alpha - 1) + \sum_n q_{nk} x_{nd} &= p_{kd}(\beta - 1) + p_{kd} \sum_n q_{nk} \\ p_{kd} &= \frac{\alpha - 1 + \sum_n q_{nk} x_{nd}}{\alpha + \beta - 2 + \sum_n q_{nk}}. \end{aligned}$$

EX 6.5. The log likelihood (after applying Jensen's inequality and omitting terms not featuring μ_{kd}) is:

$$\mathcal{L} \geq -\frac{1}{2s^2}(\mu_{kd} - m)^2 - \sum_n \sum_k q_{nk} \sum_d \frac{1}{2\sigma_{kd}^2}(x_{nd} - \mu_{kd})^2 + \dots$$

Taking partial derivatives with respect to μ_{kd} , equating to zero and re-arranging gives:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mu_{kd}} &= -\frac{1}{s^2}(\mu_{kd} - m) + \frac{1}{\sigma_{kd}^2} \sum_n q_{nk}(x_{nd} - \mu_{kd}) = 0 \\ \mu_{kd} &= \frac{\frac{1}{s^2}m + \frac{1}{\sigma_{kd}^2} \sum_n q_{nk}x_{nd}}{\frac{1}{s^2} + \frac{1}{\sigma_{kd}^2} \sum_n q_{nk}}. \end{aligned}$$

EX 6.6. The log likelihood is given by:

$$\mathcal{L} = \sum_n \log \sum_k \pi_k \frac{\lambda_k^{x_n} \exp\{-\lambda_k\}}{x_n!}$$

Incorporating variational parameters q_{nk} and applying Jensen's inequality gives:

$$\mathcal{L} \geq \sum_n \sum_k q_{nk} \log \pi_k + \sum_n \sum_k q_{nk} \log \frac{\lambda_k^{x_n} \exp\{-\lambda_k\}}{x_n!} - \sum_n \sum_k q_{nk} \log q_{nk}.$$

We shall first derive an update for π_k . The relevant terms in the bound, incorporating a Lagrangian term to ensure $\sum_k \pi_k = 1$ are:

$$\sum_n \sum_k q_{nk} \log \pi_k - \lambda \left(\sum_k \pi_k - 1 \right).$$

Taking partial derivatives with respect to π_k gives:

$$\frac{\partial}{\partial \pi_k} = \frac{1}{\pi_k} \sum_n q_{nk} - \lambda.$$

Equating to zero:

$$\lambda \pi_k = \sum_n q_{nk}.$$

To compute λ we sum both sides over k :

$$\begin{aligned} \lambda \sum_k \pi_k &= \sum_n \sum_k q_{nk} \\ \lambda &= N. \end{aligned}$$

Therefore:

$$\pi_k = \frac{1}{N} \sum_n q_{nk}.$$

Now q_{nk} . Each term is relevant, and adding the Lagrangian:

$$\sum_n \sum_k q_{nk} \log \pi_k + \sum_n \sum_k q_{nk} \log p(x_n | \lambda_k) - \sum_n \sum_k q_{nk} \log q_{nk} - \lambda \left(\sum_k q_{nk} \right).$$

Taking partial derivatives with respect to q_{nk} :

$$\frac{\partial}{\partial q_{nk}} = \log \pi_k + \log p(x_n | \lambda_k) - (\log q_{nk} + 1) - \lambda.$$

Equating to zero, taking exponentials and grouping all constants into γ :

$$\gamma q_{nk} = \pi_k p(x_n | \lambda_k).$$

To find γ , we sum both sides over k :

$$\begin{aligned} \gamma \sum_k q_{nk} &= \sum_k \pi_k p(x_n | \lambda_k) \\ \gamma &= \sum_k \pi_k p(x_n | \lambda_k). \end{aligned}$$

Therefore:

$$q_{nk} = \frac{\pi_k p(x_n | \lambda_k)}{\sum_j \pi_j p(x_n | \lambda_j)}.$$

Finally, the update for λ_k . Expanding the relevant component of the bound:

$$\sum_n \sum_k q_{nk} x_n \log \lambda_k - \sum_n \sum_k q_{nk} \lambda_k.$$

Taking partial derivatives with respect to λ_k :

$$\frac{\partial}{\partial \lambda_k} = \frac{1}{\lambda_k} \sum_n q_{nk} x_n - \sum_n q_{nk}.$$

Equating to zero and solving for λ_k gives:

$$\lambda_k = \frac{\sum_n q_{nk} x_n}{\sum_n q_{nk}}.$$

Example code to run this mixture is given below (`poissmix.m`):

```

1 clear all;close all;
2
3 % Generate data from 3 poissons
4 N = 20;
5 x = poissrnd(10,[N 1]);
6 x = [x;poissrnd(2,[N,1])];
7 x = [x;poissrnd(20,[N,1])];
8 % Fit a two component mixture
9 K = 3;
10 % Initialise

```

```

11 pr = repmat(1/K,1,K);
12 mu = rand(K,1);
13 N = size(x,1);
14 allmu = [];
15 % Do 100 iterations
16 for it = 1:100
17     % Update q
18     temp = poisspdf(repmat(x,1,K), repmat(mu',N,1));
19     q = temp.*repmat(pr,N,1);
20     q = q./repmat(sum(q,2),1,K);
21
22     % Update pr
23     pr = mean(q,1);
24
25     % Update means
26     mu = (sum(repmat(x,1,K).*q,1)./sum(q,1))';
27     allmu = [allmu;mu'];
28
29     % Plot the evolution of the means, and show the q
30     subplot(121);plot(allmu);
31     subplot(122);imagesc(q);drawnow
32 end

```


Chapter 7

EX 7.1. The bound is given by:

$$\log p(Y) \geq \int Q(\theta) \log \frac{p(Y, \theta)}{Q(\theta)} d\theta.$$

To show that the bound is maximised if $Q(\theta)$ is the true posterior, we will substitute the posterior $p(\theta|Y)$ for $Q(\theta)$ and re-arrange the right hand side such that the inequality becomes an equality. Firstly, substituting and expanding:

$$\log p(Y) \geq \int p(\theta|Y) \log p(Y, \theta) d\theta - \int p(\theta|Y) \log p(\theta|Y) d\theta.$$

Next we expand the joint density: $p(Y, \theta) = p(\theta|Y)p(Y)$:

$$\log p(Y) \geq \int p(\theta|Y) \log p(\theta|Y) d\theta + \int p(\theta|Y) \log p(Y) d\theta - \int p(\theta|Y) \log p(\theta|Y) d\theta.$$

The first and third terms cancel. The second term is an expectation of a constant and hence the expectation disappears:

$$\int p(\theta|Y) \log p(Y) d\theta = \log p(Y) \int p(\theta|Y) d\theta = \log p(Y).$$

EX 7.2. We shall go through each term in the order they are given in Section 7.5.5 starting with $\mathbf{E}_{Q_\tau(\tau)} \{\log p(\tau|a, b)\}$:

$$\begin{aligned} \mathbf{E}_{Q_\tau(\tau)} \{\log p(\tau|a, b)\} &= \mathbf{E}_{Q_\tau(\tau)} \{a \log b - \log \Gamma(a) + (a-1) \log \tau - b\tau\} \\ &= a \log b - \log \Gamma(a) + (a-1) \langle \log \tau \rangle - b \langle \tau \rangle. \end{aligned}$$

Now $\sum_n \mathbf{E}_{Q_{\mathbf{x}_n}(\mathbf{x}_n)} \{\log p(\mathbf{x}_n)\}$. Recall that $p(\mathbf{x}_n) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $Q_{\mathbf{x}_n}(\mathbf{x}_n) = \mathcal{N}(\boldsymbol{\mu}_{\mathbf{x}_n}, \boldsymbol{\Sigma}_{\mathbf{x}_n})$:

$$\begin{aligned} \sum_n \mathbf{E}_{Q_{\mathbf{x}_n}(\mathbf{x}_n)} \{\log p(\mathbf{x}_n)\} &= \sum_n \mathbf{E}_{Q_{\mathbf{x}_n}(\mathbf{x}_n)} \left\{ -\frac{D}{2} \log 2\pi - \frac{1}{2} (\mathbf{x}_n^\top \mathbf{x}_n) \right\} \\ &= -\frac{ND}{2} \log 2\pi - \frac{1}{2} \sum_n (\text{Tr}(\boldsymbol{\Sigma}_{\mathbf{x}_n}) + \boldsymbol{\mu}_{\mathbf{x}_n}^\top \boldsymbol{\mu}_{\mathbf{x}_n}). \end{aligned}$$

(Note that $\mathbf{E}_{\mathcal{N}_a(\mathbf{m}, \mathbf{S})} \{\mathbf{a}^\top \mathbf{B} \mathbf{a}\} = \text{Tr}(\mathbf{S} \mathbf{B}) + \mathbf{m}^\top \mathbf{B} \mathbf{m}$).

The expression for $\sum_m \mathbf{E}_{Q_{\mathbf{w}_m}(\mathbf{w}_m)} \{\log p(\mathbf{w}_m)\}$ is arrived at in exactly the same manner as that for $\sum_n \mathbf{E}_{Q_{\mathbf{x}_n}(\mathbf{x}_n)} \{\log p(\mathbf{x}_n)\}$.

Next, $\sum_n \sum_m \mathbf{E}_{Q(\cdot)} \{\log p(y_{nm}|\mathbf{x}_n, \mathbf{w}_m, \tau)\}$ (where the expectation is over all parameters):

$$\begin{aligned} &= \sum_n \sum_m \mathbf{E}_{Q(\cdot)} \left\{ -\frac{1}{2} \log 2\pi + \frac{1}{2} \log \tau - \frac{1}{2} \tau (y_{nm} - \mathbf{w}_m^\top \mathbf{x}_n)^2 \right\} \\ &= -\frac{NM}{2} \log 2\pi + \frac{NM}{2} \langle \log \tau \rangle - \frac{1}{2} \langle \tau \rangle \sum_n \sum_m \langle (y_{nm} - \mathbf{w}_m^\top \mathbf{x}_n)^2 \rangle. \end{aligned}$$

Now, we shall look at $\mathbf{E}_{Q_\tau(\tau)} \{\log Q_\tau(\tau)\}$. $Q_\tau(\tau)$ is Gamma with parameters e and f (see Equation 7.12). Therefore it follows exactly the same form as $\mathbf{E}_{Q_\tau(\tau)} \{\log p(\tau|a, b)\}$:

$$\mathbf{E}_{Q_\tau(\tau)} \{\log Q_\tau(\tau)\} = e \log f - \log \Gamma(e) + (e-1) \langle \log \tau \rangle - f \langle \tau \rangle.$$

The final two expressions are derived in exactly the same way and we will only show the one for \mathbf{x}_n , $\sum_n \mathbf{E}_{Q_{\mathbf{x}_n}(\mathbf{x}_n)} \{\log Q_{\mathbf{x}_n}(\mathbf{x}_n)\}$. Recall that $Q_{\mathbf{x}_n}(\mathbf{x}_n) = \mathcal{N}(\boldsymbol{\mu}_{\mathbf{x}_n}, \boldsymbol{\Sigma}_{\mathbf{x}_n})$:

$$\begin{aligned} \sum_n \mathbf{E}_{Q_{\mathbf{x}_n}(\mathbf{x}_n)} \{\log Q_{\mathbf{x}_n}(\mathbf{x}_n)\} &= -\frac{ND}{2} \log 2\pi - \frac{1}{2} \sum_n \log |\boldsymbol{\Sigma}_{\mathbf{x}_n}| - \dots \\ &\quad \frac{1}{2} \sum_n \mathbf{E}_{Q_{\mathbf{x}_n}(\mathbf{x}_n)} \{(\mathbf{x}_n - \boldsymbol{\mu}_{\mathbf{x}_n})^\top \boldsymbol{\Sigma}_{\mathbf{x}_n}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_{\mathbf{x}_n})\}. \end{aligned}$$

To progress further, we need the following identity:

$$\mathbf{E}_{\mathcal{N}_{\mathbf{a}}(\mathbf{b}, \mathbf{C})} \{(\mathbf{a} - \mathbf{d})^\top \mathbf{E}(\mathbf{a} - \mathbf{d})\} = (\mathbf{b} - \mathbf{d})^\top \mathbf{E}(\mathbf{b} - \mathbf{d}) + \text{Tr}(\mathbf{E}\mathbf{C}).$$

Matching this to our expression, we can see that $\mathbf{b} = \mathbf{d} = \boldsymbol{\mu}_{\mathbf{x}_n}$ and so the first term on the RHS is zero. The second term is $\text{Tr}(\boldsymbol{\Sigma}_{\mathbf{x}_n}^{-1} \boldsymbol{\Sigma}_{\mathbf{x}_n}) = \text{Tr}(\mathbf{I}) = D$ (for D -dimensional \mathbf{x}_n). Therefore, our expectation becomes:

$$\sum_n \mathbf{E}_{Q_{\mathbf{x}_n}(\mathbf{x}_n)} \{\log Q_{\mathbf{x}_n}(\mathbf{x}_n)\} = -\frac{ND}{2} \log 2\pi - \frac{ND}{2} - \frac{1}{2} \sum_n \log |\boldsymbol{\Sigma}_{\mathbf{x}_n}|.$$

The expression for $\sum_m \mathbf{E}_{Q_{\mathbf{w}_m}(\mathbf{w}_m)} \{\log Q_{\mathbf{w}_m}(\mathbf{w}_m)\}$ follows in exactly the same way.

EX 7.3. Starting with $Q_{\mathbf{x}_n}(\mathbf{x}_n)$:

$$\begin{aligned} Q_{\mathbf{x}_n}(\mathbf{x}_n) &\propto \exp(\mathbf{E}_{Q_{\mathbf{W}}(\mathbf{W})Q_\tau(\tau)} \{\log p(\mathbf{Y}, \mathbf{W}, \mathbf{X}, \tau|\mathbf{Z})\}) \\ &\propto \exp\left(\mathbf{E}_{Q_{\mathbf{W}}(\mathbf{W})Q_\tau(\tau)} \left\{ \log p(\mathbf{x}_n) + \sum_m z_{nm} \log p(y_{nm}|\mathbf{x}_n, \mathbf{w}_m, \tau) \right\}\right) \\ &\propto \exp\left(-\frac{1}{2} \mathbf{x}_n^\top \mathbf{x}_n - \frac{1}{2} \sum_m z_{nm} \mathbf{E}_{Q_{\mathbf{W}}(\mathbf{W})Q_\tau(\tau)} \{\tau (y_{nm} - \mathbf{w}_m^\top \mathbf{x}_n)^2\}\right) \\ &\propto \exp\left(-\frac{1}{2} \mathbf{x}_n^\top \mathbf{x}_n - \frac{1}{2} \langle \tau \rangle \sum_m z_{nm} (\mathbf{x}_n^\top \langle \mathbf{w}_m \mathbf{w}_m^\top \rangle \mathbf{x}_n - 2y_{nm} \mathbf{x}_n^\top \langle \mathbf{w}_m \rangle)\right). \end{aligned}$$

This has the form of a Gaussian, $Q_{\mathbf{x}_n}(\mathbf{x}_n) = \mathcal{N}(\boldsymbol{\mu}_{\mathbf{x}_n}, \boldsymbol{\Sigma}_{\mathbf{x}_n})$. Equating coefficients, we get the following:

$$\begin{aligned}\boldsymbol{\Sigma}_{\mathbf{x}_n} &= \left(\mathbf{I} + \langle \tau \rangle \sum_m z_{nm} \langle \mathbf{w}_m \mathbf{w}_m^\top \rangle \right)^{-1} \\ \boldsymbol{\mu}_{\mathbf{x}_n} &= \langle \tau \rangle \boldsymbol{\Sigma}_{\mathbf{x}_n} \sum_m z_{nm} y_{nm} \langle \mathbf{w}_m \rangle.\end{aligned}$$

The expression for $Q_{\mathbf{w}_m}(\mathbf{w}_m)$ follows an almost identical argument so we omit it here. Finally, the expression for $Q_\tau(\tau)$:

$$\begin{aligned}Q_\tau(\tau) &\propto \exp(\mathbf{E}_{Q_{\mathbf{W}}(\mathbf{W})Q_{\mathbf{X}}(\mathbf{X})} \{\log p(\mathbf{Y}, \mathbf{W}, \mathbf{X}, \tau | \mathbf{Z})\}) \\ &\propto \exp\left(\mathbf{E}_{Q_{\mathbf{W}}(\mathbf{W})Q_{\mathbf{X}}(\mathbf{X})} \left\{ \log \Gamma(a, b) + \sum_{n,m} z_{nm} \log p(y_{nm} | \mathbf{x}_n, \mathbf{w}_m, \tau) \right\}\right) \\ &\propto \exp\left(-b\tau + (a-1) \log \tau - \frac{\sum_{n,m} z_{nm}}{2} \log \tau - \frac{1}{2} \tau \sum_{n,m} z_{nm} \mathbf{E}_{Q_{\mathbf{W}}(\mathbf{W})Q_{\mathbf{X}}(\mathbf{X})} \left\{ (y_{nm} - \mathbf{w}_m^\top \mathbf{x}_n)^2 \right\}\right) \\ &\propto \exp\left(-b\tau + (a-1) \log \tau - \frac{\sum_{n,m} z_{nm}}{2} \log \tau - \frac{1}{2} \tau \sum_{n,m} z_{nm} \langle (y_{nm} - \mathbf{w}_m^\top \mathbf{x}_n)^2 \rangle\right).\end{aligned}$$

This has the form of another Gamma distribution, with parameters:

$$\begin{aligned}e &= a + \frac{1}{2} \sum_{n,m} z_{nm} \\ f &= b + \frac{1}{2} \sum_{n,m} z_{nm} \langle (y_{nm} - \mathbf{w}_m^\top \mathbf{x}_n)^2 \rangle.\end{aligned}$$

Chapter 8

EX 8.1. The following code would do it:

```
1      x = [0:0.01:1];
2      gam = 10.0;
3      N = length(x);
4      K = zeros(N);
5      % Create the kernel matrix
6      for n = 1:N
7          for m = 1:N
8              K(n,m) = exp(-gam*(x(n)-x(m))^2);
9          end
10     end
11     % Add a small constant to the diagonal for numerical stability
12     K = K + 1e-6*eye(N);
13     % Generate 10 samples
14     f_samps = gausssamp(repmat(0,N,1),K,10)
15     plot(x,f_samps)
```

EX 8.2. Use the code above, adding a linear term to K .

EX 8.3. The following code gives a function to compute the marginal likelihood. Note that we optimise the log of the parameters to avoid having to use a constrained optimisation:

```
1      function ML = marg_like(hyp,y,x)
2      % Function to compute the marginal likelihood
3      % hyp contains log of the hyperparameters
4      N = length(x);
5      % Compute covariance function
6      CSS = exp(-exp(hyp(1))*(repmat(x,1,N) - repmat(x',N,1)).^2);
7      CSS = CSS + exp(hyp(2))*eye(N);
8      % Compute marginal likelihood
9      ML = -(N/2)*log(2*pi) - 0.5*log(det(CSS)) - 0.5*y'*inv(CSS)*y;
10     % Return negative as optimiser minimises
11     ML = -ML;
```

and the following code calls Matlab's `fminunc` function for some generated data:

```
1      x = [0:0.05:1];
```

```

2     gam = 10.0;
3     N = length(x);
4     K = zeros(N);
5
6     % Create the kernel matrix
7     for n = 1:N
8         for m = 1:N
9             K(n,m) = exp(-gam*(x(n)-x(m))^2);
10        end
11    end
12
13    K = K + 1e-6*eye(N);
14    % Generate a sample
15    f_samps = gausssamp(repmat(0,N,1),K,1)
16
17    % Add some noise
18    ss = 0.1;
19    true_f = f_samps(1,:) + randn(1,N)*sqrt(ss);
20
21
22    true_f = true_f';
23    x = x';
24
25    % Initial hyperparameters
26    inithyp = [0;0]
27
28    % Call optimisation code
29    hyp = fminunc(@ (h) marg_like(h,true_f,x),inithyp)

```

EX 8.4. The function value at the n th data point is f_n . This function is then exponentiated to give the rate of a Poisson $\lambda_n = \exp(f_n)$. The integer count sampled from the Poisson is z_n . The gradient of the log posterior is given by:

$$\frac{\delta}{\delta \mathbf{f}} = \mathbf{z} - \exp(\mathbf{f}) - \mathbf{C}^{-1} \mathbf{f}$$

$$\frac{\delta^2}{\delta \mathbf{f} \delta \mathbf{f}^T} = -\mathbf{C}^{-1} - \mathbf{E}$$

where \mathbf{E} is a matrix with zeros everywhere except the diagonal where the n th diagonal value is $\exp(f_n)$. This is demonstrated with the following code:

```

1     % Generate some data
2     x = [0:0.01:1];
3     gam = 10.0;
4     N = length(x);
5     C = zeros(N);
6
7     % Create the kernel matrix
8     for n = 1:N
9         for m = 1:N
10            C(n,m) = exp(-gam*(x(n)-x(m))^2);
11        end
12    end
13
14    C = C + 1e-6*eye(N);

```

```

15      % Generate a sample
16      f_samps = gausssamp(repmat(0,N,1),C,1)
17
18      % Plot the function
19      figure(1);
20      hold off
21      plot(x,f_samps);
22
23      % Generate the counts
24      rate = exp(f_samps);
25      n = poissrnd(rate)';
26      hold on
27      plot(x,n,'ro')
28
29      % Use Newton-Raphson to optimise f
30      f = zeros(N,1);
31      iK = inv(C);
32      oldf = f;
33      for it = 1:100
34          gr = n - exp(f) - iK*f;
35          he = -iK - diag(exp(f));
36          f = f - inv(he)*gr;
37          ch = sum((f - oldf).^2);
38          oldf = f;
39          if ch < 1e-6
40              break
41          end
42      end
43
44      plot(x,f,'g')

```

EX 8.5. The following code generates some data, uses M-H to sample from the posterior and uses the posterior samples to compute the predictions on a grid, in order to show the predictive contours.

```

1  % GP classification using M-H
2  x = [randn(20,2)-repmat(4,20,2);randn(20,2)+repmat(4,20,2)];
3  t = [repmat(0,20,1);repmat(1,20,1)];
4
5  N = size(x,1);
6
7  % Compute the covariance matrix
8  gam = 1;alp = 5;
9  C = zeros(N);
10 for n = 1:N
11     for m = 1:N
12         C(n,m) = alp*exp(-gam*(sum((x(n,:)-x(m,:)).^2)));
13     end
14 end
15
16 iC = inv(C);
17
18 nSamps = 10000;
19 pos0 = find(t==0);
20 pos1 = find(t==1);
21
22 % Initialise with f being -3 for the -ve class and +3 for the positive

```

```

23 %\ Could use anything, but this gets us to a good place a bit faster!
24 f = [repmat(-3,20,1);repmat(3,20,1)];
25 oldProb = -0.5*f'*iC*f - sum(t.*log(1+exp(-f)));
26 oldProb = oldProb + sum((1-t).*(-f - log(1+exp(-f))));
27
28 % Make a grid of prediction points and compute the
29 % test covariance matrix
30 [X,Y] = meshgrid(min(x(:,1))-0.5:0.1:max(x(:,1)),...
31                 min(x(:,2))-0.5:0.1:max(x(:,2)));
32
33 testData = [X(:),Y(:)];
34 n_test = size(testData,1);
35 testC = zeros(n_test,size(x,1));
36 for n = 1:n_test
37     for m = 1:size(x,1)
38         testC(n,m) = alp*exp(-gam*(sum((testData(n,:)-x(m,:)).^2)));
39     end
40 end
41
42 % This term is used for predictions and doesn't change so we can just
43 % compute it once
44 testCovTerm = diag(testC*iC*testC');
45
46 % Run the sampler
47 testProbs = zeros(n_test,1);
48 allf = zeros(N,nSamps/100);
49 for s = 1:nSamps
50     % Propose a new sample for each observation
51     % Note that we could do this for all samples at once, but
52     % it becomes very hard to get a sample accepted
53     order = randperm(N);
54     for n = 1:N
55         newf = f;
56         newf(order(n)) = f(order(n)) + randn(1,1)*0.1;
57         newProb = -0.5*newf'*iC*newf - sum(t.*log(1+exp(-newf)));
58         newProb = newProb + sum((1-t).*(-newf - log(1+exp(-newf))));
59         u = rand;
60         if u <= exp(newProb - oldProb)
61             f = newf;
62             oldProb = newProb;
63         end
64     end
65
66     % Do the predictions
67     testVar = alp - testCovTerm;
68     testMu = testC*iC*f;
69
70     f_samp = randn(n_test,1).*sqrt(testVar) + testMu;
71     testProbs = testProbs + 1./(1+exp(-f_samp));
72
73 end
74
75
76 % Plot the output
77 figure(4)
78 hold off
79 pos = find(t==0);
80 plot(x(pos,1),x(pos,2),'ko');
81 pos = find(t==1);

```

```
82 hold on
83 plot(x(pos,1),x(pos,2),'ko','markerfacecolor','k')
84 testContour = reshape(testProbs./nSamps,size(X));
85 contour(X,Y,testContour,5);
```


Chapter 9

EX 9.1. The following code does both methods of sampling and then plots the autocorrelation:

```
1  % Generate some points data in 1-D
2  x = [0:0.1:1];
3  alp = 5;gam = 1;
4
5  N = length(x);
6  C = zeros(N);
7  for n = 1:N
8      for m = 1:N
9          C(n,m) = alp*exp(-gam*(x(n)-x(m))^2);
10     end
11 end
12
13 C = C + 1e-6*eye(N);
14
15 nSamps = 10000;
16
17
18 % Sample 1000 values from the Gaussian
19 f_full = gausssamp(zeros(N,1),C,nSamps);
20
21 % Samples using Gibbs sampling
22 f = zeros(N,1);
23
24 all_f = zeros(nSamps,N);
25
26 for s = 1:nSamps
27     for n = 1:N
28         % Compute the conditional mean and covariance
29         subC = C;
30         subC(n,:) = [];
31         subC(:,n) = [];
32         subc = C(n,:);
33         subc(n) = [];
34         subf = f;
35         subf(n) = [];
36         % Note that the mean is zero so these
37         % expressions are a bit simpler than in
38         % the book
39     end
```

```

40         % Also — the following expression is
41         % much more accurate than subC*inv(subC)*subC'
42         co = C(n,n) - subC/subC*subC';
43         mu = subC*inv(subC)*subf;
44         f(n) = randn.*sqrt(co) + mu;
45     end
46     all_f(s,:) = f';
47 end
48
49
50 % Compute the autocorrelation of the nth point
51 n = 5;
52
53 k_vals = [0:1:200];
54 ac = zeros(length(k_vals),2);
55 N = nSamps;
56
57 full_vals = f_full(:,n);
58 full_vals = full_vals - mean(full_vals);
59 all_vals = all_f(:,n);
60 all_vals = all_vals - mean(all_vals);
61
62 ss = [var(full_vals), var(all_vals)];
63
64
65 for k = 1:length(k_vals)
66     ac(k,1) = sum(full_vals(1:end-k_vals(k)).*full_vals(1+k_vals(k):end));
67     ac(k,1) = ac(k,1) / (ss(1) * (N-k_vals(k)));
68
69     ac(k,2) = sum(all_vals(1:end-k_vals(k)).*all_vals(1+k_vals(k):end));
70     ac(k,2) = ac(k,2) / (ss(2) * (N-k_vals(k)));
71 end
72
73 figure(1)
74 plot(ac)
75 xlabel('Lag');
76 ylabel('Autocorrelation');

```

EX 9.2. See the solution to Exercise 8.5. The only difference is changing the sigmoid likelihood to the probit likelihood.

EX 9.3. Code available on accompanying website.

EX 9.4. Code available on accompanying website.

EX 9.5. The following code will sample a regression dataset and then obtain an approximate posterior via ABC:

```

1 % ABC for GP regression
2 x = [0:0.1:1];
3
4 alp = 5;
5 gam = 1;
6
7 N = length(x);
8 C = zeros(N);
9

```



```

10 for n = 1:N
11     for m = 1:N
12         C(n,m) = alp*exp(-gam*(x(n)-x(m))^2);
13     end
14 end
15
16 C = C + 1e-6*eye(N);
17 noise_var = 0.1;
18
19 t = gausssamp(repmat(0,N,1),C,1)' + randn(N,1).*sqrt(noise_var);
20
21 % ABC
22
23 nSamps = 5000;
24 accepted_samps = [];
25
26 err_thresh = 5.0;
27 figure(1);hold off
28 plot(x,t,'k','linewidth',2);hold on
29 for s = 1:nSamps
30     % Propose a value
31     proposal = gausssamp(repmat(0,N,1),C,1)' + randn(N,1).*sqrt(noise_var);
32     err = sum((proposal - t).^2);
33
34     if err < err_thresh
35         accepted_samps = [accepted_samps proposal];
36         plot(x,proposal,'k','color',[0.6 0.6 0.6]);drawnow
37     end
38
39 end
40
41 plot(x,t,'k','linewidth',2);
42
43 %% Compare with the true posterior
44
45 postCov = inv((eye(n)*(1/noise_var) + inv(C)));
46 postMu = (1/noise_var) * postCov * t;
47
48 % Generate samples
49 samp_f = gausssamp(postMu,postCov,1000);
50 figure(2);hold off
51 plot(x,samp_f,'k','color',[0.6 0.6 0.6]);
52 hold on
53 plot(x,t,'k','linewidth',2);

```


Chapter 10

EX 10.1. We have N observations, x_1, \dots, x_N . Our model is defined as a set of K Gaussian components, with means μ_k and unit variances ($\sigma_k^2 = 1$). At any point within the Gibbs sampling procedure, each of the N observations is assigned to one of these components ($z_{nk} = 1, z_{nj} = 0 \forall j \neq k$). Each mean value has a Gaussian prior density with its own mean and variance μ_0, σ_0^2 . The conditional density for μ_k can therefore be computed as:

$$\begin{aligned}
 p(\mu_k | \dots) &\propto p(\mu_k | \mu_0, \sigma_0^2) \prod_{n=1}^N [p(x_n | \mu_k)]^{z_{nk}} \\
 &\propto \exp \left\{ -\frac{1}{2\sigma_0^2} (\mu_k - \mu_0)^2 \right\} \exp \left\{ -\frac{1}{2} \sum_{n=1}^N z_{nk} (x_n - \mu_k)^2 \right\} \\
 &= \mathcal{N}(a, b^2) \\
 \frac{\mu_k^2}{b^2} &= \frac{\mu_k^2}{\sigma_0^2} + \mu_k^2 \sum_n z_{nk} \\
 b^2 &= \left(\sigma_0^{-2} + \sum_n z_{nk} \right)^{-1} \\
 \frac{a\mu_k}{b^2} &= \frac{\mu_k \mu_0}{\sigma_0^2} + \mu_k \sum_{n=1}^N z_{nk} x_n \\
 a &= b^2 \left(\frac{\mu_0}{\sigma_0^2} + \sum_{n=1}^N z_{nk} x_n \right)
 \end{aligned}$$

EX 10.2. The following code performs Gibbs sampling for the one-dimensional Gaussian mixture with fixed known variance:

```

1 % 1D Gaussian mixture
2 x = [randn(20,1); randn(20,1)+3];
3 x = sort(x);
4
5 N = length(x);
6
7 % Define the prior parameters

```

```

8 mu0 = 0;
9 ss0 = 1;
10
11 % Component variance
12 ss = 1;
13
14 % Number of components
15 K = 2;
16
17
18 nSamps = 1000;
19
20 % Parameters for the prior Dirichlet
21 alp = repmat(1,K,1);
22 prk = gamrnd(alp,1);
23 prk = prk./sum(prk);
24
25 Z = zeros(N,K);
26
27 % Structures to hold the output
28 allZZ = zeros(N);
29 allMu = zeros(nSamps,K);
30 allPrk = zeros(nSamps,K);
31
32 % Outer loop to run multiple chains for computing Rhat
33 nChain = 5;
34 samps = zeros(nSamps,nChain);
35 for chain = 1:nChain
36     % Randomly initialise the means
37     mu = randn(K,1).*sqrt(ss0) + mu0;
38
39
40     for s = 1:nSamps
41         % Update the Z
42         for n = 1:N
43             like = -0.5*(x(n) - mu).^2;
44             like = like + log(prk);
45             prob = exp(like - max(like));
46             prob = prob./sum(prob);
47             pos = find(rand<=cumsum(prob),1);
48             Z(n,:) = 0;
49             Z(n,pos) = 1;
50         end
51
52         % Update the mu
53
54         for k = 1:K
55             ssk = 1/(1/ss0 + sum(Z(:,k)));
56             muk = ssk * (mu0/ss0 + sum(Z(:,k).*x));
57             mu(k) = randn*sqrt(ssk) + muk;
58         end
59
60         % Update the prior
61
62         prk = gamrnd(alp + sum(Z,1)',1);
63         prk = prk./sum(prk);
64
65         allPrk(s,:) = prk';
66         allMu(s,:) = mu';

```

```

67     allZZ = allZZ + Z*Z';
68 end
69
70 figure(1);
71 plot(allMu);
72 figure(2);
73 plot(allPrk);
74
75
76 allZZ = allZZ./nSamps;
77 figure(3);
78 imagesc(allZZ);
79
80 % Compute the autocorrelation of the kth mu
81 lag_vals = [0:1:200];
82 ac = zeros(length(lag_vals),1);
83
84 k = 1;
85
86 full_vals = allMu(:,k);
87 full_vals = full_vals - mean(full_vals);
88
89 samp_ss = var(full_vals);
90
91
92 for l = 1:length(lag_vals)
93     ac(l) = sum(full_vals(1:end-lag_vals(l)).*full_vals(1+lag_vals(l):end));
94     ac(l) = ac(l)/(samp_ss*(nSamps-lag_vals(l)));
95 end
96
97 figure(4);
98 plot(lag_vals,ac,'ro');drawnow
99
100 % Store the values of the lower mean (they switch in different
101 % instances)
102 m = mean(allMu);
103 pos = find(m==min(m));
104 samps(:,chain) = allMu(:,pos);
105 end
106
107 % Compute Rhat
108 muc = zeros(nSamps,nChain);
109 vc = zeros(nSamps,nChain);
110 for s = 1:nSamps
111     muc(s,:) = mean(samps(1:s,:),1);
112     vc(s,:) = var(samps(1:s,:),1);
113 end
114 W = mean(vc,2);
115 count = [1:nSamps]';
116 B = (count./(nChain-1)).*sum((muc - repmat(mean(muc,2),1,nChain)).^2,2);
117 V = ((count-1)./count).*W + (1./count).*B;
118 Rhat = sqrt(V./W);
119
120 figure(5);
121 plot(count,Rhat);

```

EX 10.3. We observe N binary vectors, each of which has D dimensions. We model them as a mixture model with K components, each of which is parameterised by a vector

of D probabilities, $\mathbf{p}_k = [p_{k1}, \dots, p_{kD}]^T$. For each component of these vectors we assume a Beta prior with parameters α and β :

$$p(p_{kd}|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p_{kd}^{\alpha-1} (1 - p_{kd})^{\beta-1}.$$

The vector of prior probabilities for each component $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]^T$ is given a uniform Dirichlet prior with parameter γ . The likelihood of an observation $\mathbf{x}_n = [x_{n1}, \dots, x_{nD}]^T$ being in component k is:

$$P(\mathbf{x}_n|\mathbf{p}_k) = \prod_{d=1}^D p_{kd}^{x_{nd}} (1 - p_{kd})^{1-x_{nd}}.$$

The posterior over all parameters is given by:

$$p(\mathbf{Z}, \boldsymbol{\pi}, \mathbf{p}_1, \dots, \mathbf{p}_K | \mathbf{X}, \alpha, \beta, \gamma) \propto p(\boldsymbol{\pi}|\gamma) \left[\prod_{k=1}^K p(\mathbf{p}_k|\alpha, \beta) \right] \left[\prod_{n=1}^N \sum_{k=1}^K (\pi_k p(\mathbf{x}_n|\mathbf{p}_k))^{z_{nk}} \right].$$

To create a Gibbs sampler we require the conditional distributions for z_{nk} , p_{kd} and π_k . The conditional density for π_k is the same as in the Gaussian case – a Dirichlet with parameters $\gamma + \sum_{n=1}^N z_{nk}$:

$$p(\boldsymbol{\pi}|\dots) = \text{Dir} \left(\gamma + \sum_{n=1}^N z_{n1}, \dots, \gamma + \sum_{n=1}^N z_{nK} \right).$$

The conditional distribution for z_{nk} is:

$$P(z_{nk} = 1 | \dots) \propto \pi_k p(\mathbf{x}_n|\mathbf{p}_k),$$

where the normalisation constant will be the sum of these terms over $j = 1 \dots K$. Finally, the conditional distribution for p_{kd} is derived as follows:

$$\begin{aligned} p(p_{kd}|\dots) &\propto p_{kd}^{\alpha-1} (1 - p_{kd})^{\beta-1} \prod_{n=1}^N (p_{kd}^{x_{nd}} (1 - p_{kd})^{1-x_{nd}})^{z_{nk}} \\ &\propto p_{kd}^{\alpha-1} (1 - p_{kd})^{\beta-1} p_{kd}^{\sum_n z_{nk} x_{nd}} (1 - p_{kd})^{\sum_n z_{nk} (1-x_{nd})} \\ &\propto p_{kd}^{\alpha-1 + \sum_n z_{nk} x_{nd}} (1 - p_{kd})^{\beta-1 + \sum_n z_{nk} (1-x_{nd})} \end{aligned}$$

which is a Beta density with parameters $\alpha^* = \alpha + \sum_n z_{nk} x_{nd}$ and $\beta^* = \beta + \sum_n z_{nk} (1 - x_{nd})$.

EX 10.4. The conditional density of the mean parameter given in the book is:

$$p(\boldsymbol{\mu}_k | \dots) = \mathcal{N}(\mathbf{a}, \mathbf{B})$$

where:

$$\mathbf{B} = \left[\boldsymbol{\Sigma}_0^{-1} + \left(\sum_n z_{nk} \mathbf{I} \right) \right]^{-1},$$

and

$$\mathbf{a} = \mathbf{B} \left[\boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \sum_n z_{nk} \mathbf{x}_n \right].$$

We are interested in computing the conditional probability of $z_{nk} = 1$ with the mean parameter $\boldsymbol{\mu}_k$ collapsed. The posterior over $\boldsymbol{\mu}_k$ not including the n th observation is identical to that given above, with the summations omitting the n th component. From now on, assume that \mathbf{a} and \mathbf{B} do not include the contribution from the n th observation. The likelihood is given by:

$$p(\mathbf{x}|\boldsymbol{\mu}_k) = \mathcal{N}(\boldsymbol{\mu}_k, \mathbf{I}).$$

We therefore need to compute the following:

$$p(\mathbf{x}|\dots) = \int \mathcal{N}(\boldsymbol{\mu}_k, \mathbf{I}) p(\boldsymbol{\mu}_k|\mathbf{a}, \mathbf{B}) d\boldsymbol{\mu}_k.$$

As both components are Gaussian, this is another Gaussian:

$$p(\mathbf{x}|\dots) = \mathcal{N}(\mathbf{a}, \mathbf{B} + \mathbf{I}).$$

Therefore, the conditional probability for $z_{nk} = 1$ is proportional to this term multiplied by the term derived in the book from collapsing $\boldsymbol{\pi}$. It is normalised by dividing by the same term summed over $j = 1 \dots K$.

EX 10.5. The following code will perform collapsed Gibbs sampling for a Gaussian mixture with 1-dimensional data.

```

1 % 1D Gaussian mixture with collapsed Gibbs sampling
2 x = [randn(20,1);randn(20,1)+3];
3 x = sort(x);
4
5 N = length(x);
6
7 % Define the prior parameters
8 mu0 = 0;
9 ss0 = 1;
10
11 % Component variance
12 ss = 1;
13
14 % Number of components
15 K = 2;
16
17
18 nSamps = 1000;
19
20 % Parameters for the prior Dirichlet
21 alp = repmat(1,1,K);
22 prk = gamrnd(alp,1);
23 prk = prk./sum(prk);
24
25 % Randomly initialise Z
26 Z = rand(N,K);
```

```

27 Z = (Z == repmat(max(Z,[],2),1,K));
28
29 % Structures to hold the output
30 allZZ = zeros(N);
31 allMu = zeros(nSamps,K);
32 allPrk = zeros(nSamps,K);
33
34 % Outer loop to run multiple chains for computing Rhat
35 nChain = 5;
36 samps = zeros(nSamps,nChain);
37 for chain = 1:nChain
38
39
40     for s = 1:nSamps
41         % Update the Z
42         for n = 1:N
43             Z(n,:) = 0;
44             sumZ = sum(Z,1);
45             sumXZ = sum(Z.*repmat(x,1,K),1);
46             postvar = (1./ss0 + sumZ).^(-1);
47             postmean = postvar.*(mu0/ss0 + sumXZ);
48             predvar = postvar + ss;
49
50             like = -(1./(2*postvar)).*(x(n) - postmean).^2;
51             like = like - 0.5*log(2*pi*postvar);
52
53             % Add the collapsed prior term
54             like = like + log(sumZ + alp);
55
56             prob = exp(like - max(like));
57             prob = prob./sum(prob);
58             pos = find(rand<=cumsum(prob),1);
59             Z(n,pos) = 1;
60         end
61
62         % Sample the mu (even though we don't need it)
63
64         for k = 1:K
65             ssk = 1/(1/ss0 + sum(Z(:,k)));
66             muk = ssk * (mu0/ss0 + sum(Z(:,k).*x));
67             mu(k) = randn*sqrt(ssk) + muk;
68         end
69
70
71         allMu(s,:) = mu';
72         Z = 1.0*Z;
73         allZZ = allZZ + Z*Z';
74     end
75
76     figure(1);
77     plot(allMu);
78     figure(2);
79     plot(allPrk);
80
81
82     allZZ = allZZ./nSamps;
83     figure(3);
84     imagesc(allZZ);
85

```



```

86     % Compute the autocorrelation of the kth mu
87     lag_vals = [0:1:200];
88     ac = zeros(length(lag_vals),1);
89
90     k = 1;
91
92     full_vals = allMu(:,k);
93     full_vals = full_vals - mean(full_vals);
94
95     samp_ss = var(full_vals);
96
97
98     for l = 1:length(lag_vals)
99         ac(l) = sum(full_vals(1:end-lag_vals(l)).*full_vals(1+lag_vals(l):end))
100         ac(l) = ac(l)/(samp_ss*(nSamps-lag_vals(l)));
101     end
102
103     figure(4);
104     plot(lag_vals,ac,'ro');drawnow
105
106     % Store the values of the lower mean (they switch in different
107     % instances)
108     m = mean(allMu);
109     pos = find(m==min(m));
110     samps(:,chain) = allMu(:,pos);
111 end
112
113 % Compute Rhat
114 muc = zeros(nSamps,nChain);
115 vc = zeros(nSamps,nChain);
116 for s = 1:nSamps
117     muc(s,:) = mean(samps(1:s,:),1);
118     vc(s,:) = var(samps(1:s,:),1);
119 end
120 W = mean(vc,2);
121 count = [1:nSamps]';
122 B = (count./(nChain-1)).*sum((muc - repmat(mean(muc,2),1,nChain)).^2,2);
123 V = ((count-1)./count).*W + (1./count).*B;
124 Rhat = sqrt(V./W);
125
126 figure(5);
127 plot(count,Rhat);

```

EX 10.6. The following code samples from a CRP. Vary alpha to see the effect to the histogram:

```

1  % Sample from a CRP
2  N = 100; % 100 customers
3
4  Z = ones(N,1); % Initialise with all customers at one table
5
6  nSamps = 1000;
7
8  alp = 0.1; % Concentration parameter to vary
9  nK = zeros(nSamps,1);
10
11 for s = 1:nSamps

```

```

12   for n = 1:N
13       thistable = find(Z(n,:));
14       Z(n,:) = 0;
15       sumZ = sum(Z,1);
16       if sumZ(thistable) == 0
17           sumZ(thistable) = [];
18           Z(:,thistable) = [];
19       end
20       probs = [sumZ alp];
21       probs = probs./sum(probs);
22       pos = find(rand<=cumsum(probs),1);
23       Z(n,pos) = 1;
24   end
25   nK(s) = size(Z,2);
26 end
27
28 hist(nK,unique(nK));

```

EX 10.7. It will be most efficient to use a collapsed Gibbs sampling scheme for this by marginalising the probabilities p_{kd} . As shown above, the conditional posterior density over these parameters is a Beta distribution with parameters:

$$\alpha^* = \alpha + \sum_n z_{nk} x_{nd}, \quad \beta^* = \beta + \sum_n z_{nk} (1 - x_{nd}).$$

(remember that when considering the collapsed update for the n th observation we would omit the n th data point from the summation). The update for z_{nk} is:

$$P(z_{nk} = 1 | \dots) \propto \left(\sum_{m \neq n} z_{mk} \right) \prod_{d=1}^D \int P(x_{nd} | z_{nk} = 1, p_{kd}) p(p_{kd} | \dots) dp_{kd}.$$

Because of the conjugacy of the Beta prior and the binomial likelihood, we can analytically evaluate the integral:

$$\begin{aligned}
 &= \int p_{kd}^{x_{nd}} (1 - p_{kd})^{1-x_{nd}} \frac{\Gamma(\alpha^* + \beta^*)}{\Gamma(\alpha^*)\Gamma(\beta^*)} p_{kd}^{\alpha^*-1} (1 - p_{kd})^{\beta^*-1} dp_{kd} \\
 &= \frac{\Gamma(\alpha^* + \beta^*)}{\Gamma(\alpha^*)\Gamma(\beta^*)} \int p_{kd}^{\alpha^*+x_{nd}-1} (1 - p_{kd})^{\beta^*+(1-x_{nd})-1} dp_{kd} \\
 &= \frac{\Gamma(\alpha^* + \beta^*)}{\Gamma(\alpha^*)\Gamma(\beta^*)} \frac{\Gamma(\alpha^* + x_{nd})\Gamma(\beta^* + (1 - x_{nd}))}{\Gamma(\alpha^* + \beta^* + 1)} \\
 &= \frac{(\alpha^*)^{x_{nd}} (\beta^*)^{1-x_{nd}}}{\alpha^* + \beta^*}.
 \end{aligned}$$

Where in the final step we have used the identity $\Gamma(z+1) = z\Gamma(z)$. This completes the definition of the sampler. The update for z_{nk} is proportional to the product of these terms over the dimensions (D) multiplied by the number of objects in the cluster k . The probability of a new cluster is proportional to the concentration parameter multiplied by the product of these terms where $\alpha^* = \alpha$ and $\beta^* = \beta$.

EX 10.8. The following code samples from a DP with a Gaussian base distribution:

```

1 % DP sample
2 basemu = 0;
3 basess = 1;
4 alp = 10; % Concentration parameter
5
6 N = 100; % Number of samples to draw
7
8 samples = [];
9
10 for n = 1:N
11     if n==1
12         samples = randn.*sqrt(basess) + basemu;
13     else
14         if rand < alp/(alp + n-1)
15             % Sample a new value
16             samples(end+1) = randn.*sqrt(basess) + basemu;
17         else
18             % Copy a previous one
19             samples(end+1) = samples(randi(length(samples)));
20         end
21     end
22 end

```

EX 10.9. See script on accompanying webpage.

EX 10.10. We will derive the collapsed sampler. Assume that we observe $m = 1 \dots M$ datasets, each of which has $n = 1 \dots N_m$ observations, x_n^m . There are $k = 1 \dots K$ top-level components and the top-level concentration parameter is α . The concentration parameter in the m th file is γ_m . In the m th file there are $v = 1 \dots V_m$ components, each associated with one of the top-level components. Let $z_{nv}^m = 1$ if the n th observation in the m th file is associated with the v th component in that file. Let $w_{vk}^m = 1$ if the v th component in the m th file is associated with the k th top-level component.

For the collapsed sampler, we require the distributions for re-sampling z_{nv}^m . When we re-sample this assignment, there are three possibilities. 1: We put it into a current component in file m . 2: We put it into a new component that uses a pre-existing top-level component. 3: We put it into a new component that has a new top-level component. We will consider each of these in turn. Firstly, though we need the posterior density over μ_k (the mean of the k th top-level component). This is Gaussian and given by:

$$\begin{aligned}
 p(\mu_k | \dots) &\propto \mathcal{N}(\mu_0, \sigma_0^2) \prod_{m=1}^M \prod_{n=1}^{N_m} \prod_{v=1}^{V_m} \mathcal{N}(x_n^m | \mu_k, \sigma^2)^{z_{nv}^m w_{vk}^m} \\
 &= \mathcal{N}(a, b^2) \\
 b^2 &= \left[\frac{1}{\sigma_0^2 + \frac{1}{\sigma^2} \sum_{m,n,v} z_{nv}^m w_{vk}^m} \right]^{-1} \\
 a &= b^2 \left[\frac{\sum_{m,n,v} z_{nv}^m w_{vk}^m x_n^m}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right]
 \end{aligned} \tag{10.1}$$

The predictive density is therefore (where summations will not include the data point in question):

$$p(x_n^m | z_{nv}^m = 1, w_{vk}^m = 1) = \mathcal{N}(a, b^2 + \sigma^2) = L(x_n^m | k)$$

Note that the expression for a brand new top-level component is identical but all the summations will equate to zero.

The probability of placing the object in component v is therefore proportional to:

$$P(z_{nv}^m = 1 | w_{vk}^m = 1, \dots) \propto L(x_n^m | k) \left(\sum_n z_{nv}^m \right)$$

To compute the probability of placing the point in a new component in file m we have to marginalise over the top level components (admitting the possibility of a new one). Let s_k be the number of components across all files that are assigned to top-level component k :

$$P(z_{nv^*}^m = 1 | \dots) \propto \gamma_m \left(\frac{\alpha L(x_n^m | k^*) + \sum_l s_l L(x_n^m | l)}{\alpha + \sum_j s_j} \right)$$

We can therefore use these two probabilities to decide if the n th object ought to go into one of the current components in the m th file, or a new component. If we choose a new component, we then need to choose which top-level component it ought to be associated with (or a new one). The probability of it being a current one is proportional to $s_l L(x_n^m | l)$ and the probability of a new one is proportional to $\alpha L(x_n^m | l^*)$.

This completes the sampler, although it might not be particularly efficient. To improve its efficiency, one could re-sample the w_{vk}^m for each file-level component. It is fairly straightforward to compute the required distribution but, be careful: if, as in this case, the mean parameters have all been marginalised, the joint density of all of the objects assigned to this component cannot be assumed to be independent in the standard way (they are all dependent via the marginalisation of the mean parameters). This can be overcome by, for example, computing $p(x_1, x_2, x_3, \dots | \dots) = p(x_1 | \dots) p(x_2 | x_1, \dots) p(x_3 | x_1, x_2, \dots)$ etc.

EX 10.11. The equation we are trying to obtain is the probability that the i th word in the n th document comes from topic k . This will be the product of two components: the probability of the k th topic in document n multiplied by the probability of word w (assuming that the i th word is word w) in topic k . Based on the definitions in the text, in the un-collapsed world, this would be given by:

$$P(z_{ni} = k | x_{ni} = w, \beta_{kw}, \theta_{nk}) = \theta_{nk} \beta_{kw}.$$

To obtain the collapsed version, we need to marginalise the parameters θ_n and β_k . Both are multinomials and have Dirichlet priors with parameters α and γ respectively. We know from our work in mixture models that if we have one sample from a multinomial (say the k th value has been sampled) whose parameters are Dirichlet distributed, the probability of that sample is equal to the k th Dirichlet

parameter divided by the sum of the parameters. We also know that the posterior Dirichlet for a Dirichlet prior and multinomial likelihood has parameters equal to the prior parameters plus the observation counts. Therefore, defining c_{nk}^{-i} as the count of the number of words (excluding the i th one) in the n th document currently assigned to topic k , we know that the term associated with the marginalisation of θ_n is:

$$\frac{c_{nk}^{-i} + \alpha_k}{\sum_j c_{nj}^{-i} + \alpha_j}$$

and the term associated with the marginalisation of β_k is:

$$\frac{v_{kw} + \gamma_w}{\sum_{w'} v_{kw'} + \gamma_{w'}}$$

Multiplying the two together gives the update given in the text.