# Solutions Manual to Pattern Recognition and Machine Learning

Hiromichi Inawashiro October 11, 2025

## Contents

1	Introduction	1
2	Probability Distributions	43
3	Linear Models for Regression	106
4	Linear Models for Classification	131
5	Neural Networks	157
6	Kernel Methods	191
7	Sparse Kernel Machines	220

## 1 Introduction

#### 1.1

Let

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} (y(x_n, \mathbf{w}) - t_n)^2.$$
 (1.1)

Setting the derivative with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = \sum_{n=1}^{N} \frac{\partial y(x_n, \mathbf{w})}{\partial \mathbf{w}} \left( y(x_n, \mathbf{w}) - t_n \right). \tag{1.2}$$

If

$$y(x_n, \mathbf{w}) = \mathbf{w}^\mathsf{T} \boldsymbol{\phi}(x_n), \tag{1.3}$$

then

$$\mathbf{0} = \sum_{n=1}^{N} \boldsymbol{\phi}(x_n) \left( \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(x_n) - t_n \right). \tag{1.4}$$

Then,

$$\left(\sum_{n=1}^{N} \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^{\mathsf{T}}\right) \mathbf{w} = \sum_{n=1}^{N} t_n \boldsymbol{\phi}(x_n). \tag{1.5}$$

Therefore,

$$\underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w}) = \mathbf{A}^{-1} \mathbf{v}, \tag{1.6}$$

where

$$\mathbf{A} = \sum_{n=1}^{N} \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^{\mathsf{T}},$$

$$\mathbf{v} = \sum_{n=1}^{N} t_n \boldsymbol{\phi}(x_n).$$
(1.7)

If

$$\phi(x_n) = \begin{bmatrix} 1 \\ x_n \\ \vdots \\ x_n^M \end{bmatrix},$$

then

$$A_{mm'} = \sum_{n=1}^{N} x_n^{m+m'},$$

$$v_m = \sum_{n=1}^{N} t_n x_n^m.$$
(1.8)

#### 1.2

Let

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} (y(x_n, \mathbf{w}) - t_n)^2 + \frac{\lambda}{2} ||\mathbf{w}||^2.$$
 (1.9)

Setting the derivative with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = \sum_{n=1}^{N} \frac{\partial y(x_n, \mathbf{w})}{\partial \mathbf{w}} (y(x_n, \mathbf{w}) - t_n) + \lambda \mathbf{w}.$$
 (1.10)

If

$$y(x_n, \mathbf{w}) = \mathbf{w}^{\mathsf{T}} \phi(x_n), \tag{1.11}$$

then

$$\mathbf{0} = \sum_{n=1}^{N} \phi(x_n) \left( \mathbf{w}^{\mathsf{T}} \phi(x_n) - t_n \right) + \lambda \mathbf{w}. \tag{1.12}$$

Then,

$$\left(\sum_{n=1}^{N} \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^{\mathsf{T}} + \lambda \mathbf{I}\right) \mathbf{w} = \sum_{n=1}^{N} t_n \boldsymbol{\phi}(x_n). \tag{1.13}$$

Therefore,

$$\underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w}) = \mathbf{A}^{-1}\mathbf{v}, \tag{1.14}$$

where

$$\mathbf{A} = \sum_{n=1}^{N} \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^{\mathsf{T}} + \lambda \mathbf{I},$$

$$\mathbf{v} = \sum_{n=1}^{N} t_n \boldsymbol{\phi}(x_n).$$
(1.15)

If

$$\phi(x_n) = \begin{bmatrix} 1 \\ x_n \\ \vdots \\ x_n^M \end{bmatrix},$$

then

$$A_{mm'} = \sum_{n=1}^{N} x_n^{m+m'} + \lambda I_{mm'},$$

$$v_m = \sum_{n=1}^{N} t_n x_n^m.$$
(1.16)

#### 1.3

Let

$$p(a|r) = \frac{3}{10}, p(o|r) = \frac{2}{5}, p(l|r) = \frac{3}{10},$$

$$p(a|b) = \frac{1}{2}, p(o|b) = \frac{1}{2},$$

$$p(a|g) = \frac{3}{10}, p(o|g) = \frac{3}{10}, p(l|g) = \frac{2}{5}.$$

$$(1.17)$$

Let

$$p(r) = \frac{1}{5}, p(b) = \frac{1}{5}, p(g) = \frac{3}{5}.$$
 (1.18)

(a)

By the Bayes' theorem,

$$p(a) = p(a|r)p(r) + p(a|b)p(b) + p(a|g)p(g).$$
(1.19)

Therefore,

$$p(a) = \frac{17}{50}. (1.20)$$

(b)

By the Bayes' thorem,

$$p(g|o) = \frac{p(g,o)}{p(o)}.$$
 (1.21)

By the Bayes' throrem, the right hand side can be written as

$$\frac{p(o|g)p(g)}{p(o|r)p(r) + p(o|b)p(b) + p(o|g)p(g)}. (1.22)$$

Therefore,

$$p(g|o) = \frac{1}{2}. (1.23)$$

#### 1.4

Let x and y be variables such that

$$x = g(y). (1.24)$$

Let  $\hat{x}$  and  $\hat{y}$  be the locations of the maximum of  $p_x$  and  $p_y$  respectively. Let us assume that there exists a positive  $\epsilon$  such that if

$$|y - \hat{y}| < \epsilon, \tag{1.25}$$

then

$$g'(y) \neq 0. \tag{1.26}$$

(a)

Since

$$\int p_x(x)dx = \int p_x(g(y)) |g'(y)| dy,$$
(1.27)

we have

$$p_y(y) = p_x(g(y))|g'(y)|.$$
 (1.28)

Taking the derivative and substituting

$$y = \hat{y} \tag{1.29}$$

gives

$$0 = g'(\hat{y})p'_x(g(\hat{y})) + p_x(g(\hat{y}))g''(\hat{y}).$$
(1.30)

Therefore, in general,

$$\hat{x} \neq g\left(\hat{y}\right). \tag{1.31}$$

(b)

$$g(y) = ay + b, (1.32)$$

then

$$0 = ap'_{x}\left(g\left(\hat{y}\right)\right). \tag{1.33}$$

Therefore,

$$\hat{x} = g\left(\hat{y}\right). \tag{1.34}$$

#### 1.5

We have

$$\operatorname{var} f(x) = E(f(x) - Ef(x))^{2}.$$
 (1.35)

The right hand side can be written as

$$E((f(x))^{2} - 2f(x) E f(x) + (E f(x))^{2}) = E(f(x))^{2} - (E f(x))^{2}.$$
 (1.36)

Therefore,

$$var f(x) = E(f(x))^{2} - (E f(x))^{2}.$$
 (1.37)

#### 1.6

We have

$$cov(x, y) = E((x - Ex)(y - Ey)).$$
 (1.38)

The right hand side can be written as

$$Exy - E(x Ey) - E(y Ex) + E(Ex Ey) = Exy - Ex Ey.$$
 (1.39)

The right hand side can be written as

$$\int xyp(x,y)dxdy - \int xp(x)dx \int yp(y)dy.$$
 (1.40)

If x and y are independent, by the definition,

$$f(x,y) = f(x)f(y). (1.41)$$

Then,

$$\int xyp(x,y)dxdy = \int p(x)dx \int p(y)dy.$$
 (1.42)

$$cov(x,y) = 0. (1.43)$$

(a)

Let

$$I = \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2}x^2\right) dx. \tag{1.44}$$

Then,

$$I^{2} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^{2}}\left(x^{2} + y^{2}\right)\right) dx dy. \tag{1.45}$$

By the transformation from Cartesian coordinates (x, y) to polar coordinates  $(r, \theta)$ , the right hand side can be written as

$$\int_0^\infty \int_0^{2\pi} \exp\left(-\frac{1}{2\sigma^2}r^2\right) \begin{vmatrix} \cos\theta & -r\sin\theta \\ \sin\theta & r\cos\theta \end{vmatrix} dr d\theta = 2\pi \int_0^\infty \exp\left(-\frac{1}{2\sigma^2}r^2\right) r dr. \tag{1.46}$$

By the transformation  $s = \frac{r}{\sigma}$ , the right hand side can be written as

$$2\pi\sigma^2 \int_0^\infty \exp\left(-\frac{1}{2}s^2\right) s ds = 2\pi\sigma^2 \left[-\exp\left(-\frac{1}{2}s^2\right)\right]_0^\infty. \tag{1.47}$$

Therefore,

$$I = \left(2\pi\sigma^2\right)^{\frac{1}{2}}.\tag{1.48}$$

(b)

By the definition,

$$\mathcal{N}\left(x|\mu,\sigma^2\right) = \left(2\pi\sigma^2\right)^{-\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right). \tag{1.49}$$

Then,

$$\int_{-\infty}^{\infty} \mathcal{N}\left(x|\mu,\sigma^2\right) dx = \left(2\pi\sigma^2\right)^{-\frac{1}{2}} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) dx. \tag{1.50}$$

By the transformation  $t = x - \mu$ , the right hand side can be written as

$$(2\pi\sigma^2)^{-\frac{1}{2}} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2}t^2\right) dt = (2\pi\sigma^2)^{-\frac{1}{2}} I.$$
 (1.51)

$$\int_{-\infty}^{\infty} \mathcal{N}\left(x|\mu,\sigma^2\right) dx = 1. \tag{1.52}$$

(a)

Let x be a variable such that

$$p(x) = \mathcal{N}(x|\mu, \sigma^2). \tag{1.53}$$

Then,

$$E x = \int_{-\infty}^{\infty} x \mathcal{N}(x|\mu, \sigma^2) dx.$$
 (1.54)

The right hand side can be written as

$$\left(2\pi\sigma^2\right)^{-\frac{1}{2}} \int_{-\infty}^{\infty} x \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) dx. \tag{1.55}$$

By the transformation

$$y = x - \mu, \tag{1.56}$$

the integral can be written as

$$\int_{-\infty}^{\infty} (y+\mu) \exp\left(-\frac{1}{2\sigma^2}y^2\right) dy$$

$$= \int_{-\infty}^{\infty} y \exp\left(-\frac{1}{2\sigma^2}y^2\right) dy + \mu \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2}y^2\right) dy.$$
(1.57)

By 1.7(a), the right hand side can be written as

$$\mu \left(2\pi\sigma^2\right)^{\frac{1}{2}}.\tag{1.58}$$

Therefore,

$$\mathbf{E} x = \mu. \tag{1.59}$$

(b)

By 1.7(b),

$$\int_{-\infty}^{\infty} \mathcal{N}\left(x|\mu,\sigma^2\right) dx = 1,\tag{1.60}$$

so that

$$(2\pi\sigma^2)^{-\frac{1}{2}} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) dx = 1.$$
 (1.61)

Taking the derivative with respect to  $\sigma^2$  gives

$$(2\pi)^{-\frac{1}{2}} \left(-\frac{1}{2}\right) (\sigma^2)^{-\frac{3}{2}} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2} (x-\mu)^2\right) dx + (2\pi\sigma^2)^{-\frac{1}{2}} \int_{-\infty}^{\infty} \frac{1}{2} (\sigma^2)^{-2} (x-\mu)^2 \exp\left(-\frac{1}{2\sigma^2} (x-\mu)^2\right) dx = 0.$$
 (1.62)

The left hand side can be written as

$$-\frac{1}{2} (\sigma^{2})^{-1} \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^{2}) dx + \frac{1}{2} (\sigma^{2})^{-2} \int_{-\infty}^{\infty} (x-\mu)^{2} \mathcal{N}(x|\mu, \sigma^{2}) dx$$

$$= -\frac{1}{2} (\sigma^{2})^{-1} + \frac{1}{2} (\sigma^{2})^{-2} \operatorname{var} x.$$
(1.63)

Therefore,

$$var x = \sigma^2. (1.64)$$

#### 1.9

(a)

Let x be a variable such that

$$p(x) = \mathcal{N}\left(x|\mu, \sigma^2\right). \tag{1.65}$$

Setting the derivative of the right hand side with respect to x to zero gives

$$0 = (2\pi\sigma^2)^{-\frac{1}{2}} \left( -\frac{1}{\sigma^2} (x - \mu) \right) \exp\left( -\frac{1}{2\sigma^2} (x - \mu)^2 \right). \tag{1.66}$$

Therefore,

$$mode x = \mu. (1.67)$$

(b)

Let  $\mathbf{x}$  be a variable such that

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}).$$
 (1.68)

Setting the derivative of the right hand side with respect to  $\mathbf{x}$  to zero gives

$$\mathbf{0} = -(2\pi)^{-\frac{D}{2}} \left(\det \mathbf{\Sigma}\right)^{-\frac{1}{2}} \left(\mathbf{\Sigma}^{-1} + \left(\mathbf{\Sigma}^{-1}\right)^{\mathsf{T}}\right) (\mathbf{x} - \boldsymbol{\mu}) \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right).$$
(1.69)

Therefore,

$$mode \mathbf{x} = \boldsymbol{\mu}. \tag{1.70}$$

#### 1.10

(a)

We have

$$E(x+y) = \iint (x+y)p(x,y)dxdy.$$
 (1.71)

The right hand side can be written as

$$\int x \left( \int p(x,y) dy \right) dx + \int y \left( \int p(x,y) dx \right) dy = \int x p(x) dx + \int y p(y) dy.$$
(1.72)

The right hand side can be written as

$$\mathbf{E}\,x + \mathbf{E}\,y. \tag{1.73}$$

Therefore,

$$E(x+y) = E x + E y. (1.74)$$

(b)

We have

$$var(x+y) = E(x+y - E(x+y))^{2}$$
 (1.75)

The right hand side can be written as

$$E(x - Ex)^{2} + 2E((x - Ex)(y - Ey)) + E(y - Ey)^{2}$$

$$= var x + 2cov(x, y) + var y.$$
(1.76)

By 1.6, if x and y are independent, then

$$cov(x,y) = 0. (1.77)$$

$$var(x+y) = var x + var y. (1.78)$$

Let  $x_1, \dots, x_N$  be variables such that

$$p(x_n) = \mathcal{N}\left(x_n | \mu, \sigma^2\right). \tag{1.79}$$

Then,

$$\ln\left(\prod_{n=1}^{N} p(x_n)\right) = -\frac{N}{2}\ln\left(2\pi\sigma^2\right) - \frac{1}{2\sigma^2}\sum_{n=1}^{N} (x_n - \mu)^2.$$
 (1.80)

Setting the derivatives with respect to  $\mu$  and  $\sigma^2$  to zero gives

$$0 = \frac{1}{\sigma^2} \sum_{n=1}^{N} (x_n - \mu),$$

$$0 = -\frac{N}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{n=1}^{N} (x_n - \mu)^2.$$
(1.81)

Therefore, the maximum likelihood solutions for  $\mu$  and  $\sigma^2$  are given by

$$\mu_{\rm ML} = \frac{1}{N} \sum_{n=1}^{N} x_n,$$

$$\sigma_{\rm ML}^2 = \frac{1}{N} \sum_{n=1}^{N} (x_n - \mu_{\rm ML})^2.$$
(1.82)

#### 1.12

(a)

Let  $x_n$  and  $x_{n'}$  be independent variables such that

$$p(x_n) = \mathcal{N}\left(x_n|\mu, \sigma^2\right),$$
  

$$p(x_{n'}) = \mathcal{N}\left(x_{n'}|\mu, \sigma^2\right).$$
(1.83)

Then,

$$\operatorname{E} x_n x_{n'} = \mu^2. \tag{1.84}$$

By the property

$$\operatorname{E} x_n^2 = \operatorname{var} x_n + \left(\operatorname{E} x_n\right)^2, \tag{1.85}$$

we have

$$E x_n^2 = \sigma^2 + \mu^2. (1.86)$$

Therefore,

$$E x_n x_{n'} = \mu^2 + I_{nn'} \sigma^2. (1.87)$$

(b)

Let  $x_1, \dots, x_N$  be independent variables such that

$$p(x_n) = \mathcal{N}\left(x_n|\mu, \sigma^2\right). \tag{1.88}$$

By 1.11, the maximum likelihood solution for  $\mu$  is given by

$$\mu_{\rm ML} = \frac{1}{N} \sum_{n=1}^{N} x_n. \tag{1.89}$$

Then,

$$E \mu_{ML} = \frac{1}{N} \sum_{n=1}^{N} E x_n.$$
 (1.90)

Therefore,

$$E \mu_{ML} = \mu. \tag{1.91}$$

(c)

Let  $x_1, \dots, x_N$  be independent variables such that

$$p(x_n) = \mathcal{N}\left(x_n|\mu, \sigma^2\right). \tag{1.92}$$

By 1.11, the maximum likelihood solution for  $\sigma^2$  is given by

$$\sigma_{\rm ML}^2 = \frac{1}{N} \sum_{n=1}^{N} (x_n - \mu_{\rm ML})^2.$$
 (1.93)

Then,

$$E \sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^{N} E (x_n - \mu_{ML})^2.$$
 (1.94)

The right hand side can be writen as

$$\frac{1}{N} \sum_{n=1}^{N} \mathrm{E} \left( x_n^2 - 2\mu_{\mathrm{ML}} x_n + \mu_{\mathrm{ML}}^2 \right) 
= \frac{1}{N} \sum_{n=1}^{N} \mathrm{E} x_n^2 - \frac{2}{N} \mathrm{E} \left( \mu_{\mathrm{ML}} \left( \sum_{n=1}^{N} x_n \right) \right) + \mathrm{E} \mu_{\mathrm{ML}}^2.$$
(1.95)

The first term of the right hand side can be written as

$$\frac{1}{N} \sum_{n=1}^{N} (\mu^2 + \sigma^2) = \mu^2 + \sigma^2, \tag{1.96}$$

while, by 1.11, the second and third terms can be writen as

$$-\frac{2}{N} E\left(N\left(\frac{1}{N}\sum_{n=1}^{N} x_n\right)^2\right) + E\left(\frac{1}{N}\sum_{n=1}^{N} x_n\right)^2 = -E\left(\frac{1}{N}\sum_{n=1}^{N} x_n\right)^2.$$
 (1.97)

By (a), the right hand side can be written as

$$-\frac{1}{N^2} \sum_{n=1}^{N} \operatorname{E} x_n^2 - \frac{2}{N^2} \sum_{1 \le n < n' \le N} \operatorname{E} x_n x_{n'}$$

$$= -\frac{1}{N^2} N \left( \mu^2 + \sigma^2 \right) - \frac{2}{N^2} \frac{N(N-1)}{2} \mu^2.$$
(1.98)

The right hand side can be written as

$$-\frac{1}{N}(\mu^2 + \sigma^2) - \frac{N-1}{N}\mu^2 = -\mu^2 - \frac{1}{N}\sigma^2.$$
 (1.99)

Then,

$$E \sigma_{ML}^2 = \mu^2 + \sigma^2 - \mu^2 - \frac{1}{N} \sigma^2.$$
 (1.100)

$$E \sigma_{\rm ML}^2 = \frac{N-1}{N} \sigma^2. \tag{1.101}$$

Let  $x_1, \dots, x_N$  be variables such that

$$\begin{aligned}
\mathbf{E} \, x_n &= \mu, \\
\operatorname{var} x_n &= \sigma^2.
\end{aligned} \tag{1.102}$$

We have

$$E\left(\frac{1}{N}\sum_{n=1}^{N}(x_n-\mu)^2\right) = \frac{1}{N^2}\sum_{n=1}^{N}E(x_n-\mu)^2.$$
 (1.103)

The right hand side can be writen as

$$\frac{1}{N^2} \sum_{n=1}^{N} \operatorname{var} x_n = \frac{\sigma^2}{N}.$$
 (1.104)

Therefore,

$$E\left(\frac{1}{N}\sum_{n=1}^{N}(x_{n}-\mu)^{2}\right) = \frac{\sigma^{2}}{N}.$$
(1.105)

#### 1.14

Let

$$w_{dd'}^{S} = \frac{1}{2}(w_{dd'} + w_{d'd}),$$

$$w_{dd'}^{A} = \frac{1}{2}(w_{dd'} - w_{d'd}).$$
(1.106)

(a)

We have

$$w_{dd'} = w_{dd'}^{S} + w_{dd'}^{A},$$

$$w_{dd'}^{S} = w_{d'd}^{S},$$

$$w_{dd'}^{A} = -w_{d'd}^{A}.$$
(1.107)

(b)

We have

$$\sum_{d=1}^{D} \sum_{d'=1}^{D} w_{dd'}^{A} x_{d} x_{d'} = \frac{1}{2} \sum_{d=1}^{D} \sum_{d'=1}^{D} (w_{dd'} - w_{d'd}) x_{d} x_{d'}.$$
 (1.108)

The right hand side can be written as

$$\frac{1}{2} \left( \sum_{d=1}^{D} \sum_{d'=1}^{D} w_{dd'} x_d x_{d'} - \sum_{d=1}^{D} \sum_{d'=1}^{D} w_{d'd} x_d x_{d'} \right) = 0.$$
 (1.109)

Therefore,

$$\sum_{d=1}^{D} \sum_{d'=1}^{D} w_{dd'}^{A} x_{d} x_{d'} = 0.$$
 (1.110)

(c)

We have

$$\sum_{d=1}^{D} \sum_{d'=1}^{D} w_{dd'} x_d x_{d'} = \sum_{d=1}^{D} \sum_{d'=1}^{D} \left( w_{dd'}^{S} + w_{dd'}^{A} \right) x_d x_{d'}.$$
 (1.111)

By (b), the right hand side can be written as

$$\sum_{d=1}^{D} \sum_{d'=1}^{D} w_{dd'}^{S} x_{d} x_{d'} + \sum_{d=1}^{D} \sum_{d'=1}^{D} w_{dd'}^{A} x_{d} x_{d'} = \sum_{d=1}^{D} \sum_{d'=1}^{D} w_{dd'}^{S} x_{d} x_{d'}, \qquad (1.112)$$

Therefore,

$$\sum_{d=1}^{D} \sum_{d'=1}^{D} w_{dd'} x_d x_{d'} = \sum_{d=1}^{D} \sum_{d'=1}^{D} w_{dd'}^{S} x_d x_{d'}.$$
 (1.113)

(d)

Since  $\mathbf{W}^{\mathrm{S}}$  is a  $D \times D$  symmetric matrix, its number of independent parameters is  $\frac{D(D+1)}{2}$ .

#### 1.15

(a)

Let n(D, M) be the number of independent parameters of a polynomial in D dimensions and M orders. Then

$$n(1, M) = n(1, M - 1) = 1.$$
 (1.114)

Let us assume that

$$n(D,M) = \sum_{d=1}^{D} n(d, M-1).$$
(1.115)

The independent terms of a polynomial in D+1 dimensions and M orders can be split into 1. the ones of a polynomial in D dimensions and M orders and 2. the ones generated by multiplying the ones in D+1 dimensions and M orders by the D+1th variable. Then,

$$n(D+1,M) = n(D,M) + n(D+1,M-1), (1.116)$$

so that

$$n(D+1,M) = \sum_{d=1}^{D+1} n(d,M-1).$$
 (1.117)

Therefore, the assumption is proved by induction on D.

(b)

We have

$$\sum_{d=1}^{1} \frac{(d+M-2)!}{(d-1)!(M-1)!} = 1.$$
 (1.118)

Let us assume that

$$\sum_{d=1}^{D} \frac{(d+M-2)!}{(d-1)!(M-1)!} = \frac{(D+M-1)!}{(D-1)!M!}.$$
 (1.119)

Then,

$$\sum_{d=1}^{D+1} \frac{(d+M-2)!}{(d-1)!(M-1)!} = \frac{(D+M-1)!}{(D-1)!M!} + \frac{(D+M-1)!}{D!(M-1)!}.$$
 (1.120)

The right hand side can be written as

$$\frac{D(D+M-1)! + M(D+M-1)!}{D!M!} = \frac{(D+M)!}{D!M!}.$$
 (1.121)

Therefore, the assumption is proved by induction on D.

(c)

By 1.14(d),

$$n(D,2) = \frac{D(D+1)}{2}. (1.122)$$

Let us assume that

$$n(D,M) = \frac{(D+M-1)!}{(D-1)!M!}.$$
(1.123)

By (a),

$$n(D, M+1) = \sum_{d=1}^{D} n(d, M).$$
(1.124)

By the assumption and (b), the right hand side can be written as

$$\sum_{d=1}^{D} \frac{(d+M-1)!}{(d-1)!M!} = \frac{(D+M)!}{(D-1)!(M+1)!}.$$
 (1.125)

Therefore, the assumption is proved by induction on M.

#### 1.16

(a)

Let N(D, M) be the number of independent parameters in all of the terms up to and including the ones of D dimensions and M orders. By 1.15,

$$N(D, M) = \sum_{m=0}^{M} n(D, m), \qquad (1.126)$$

where

$$n(D,m) = \frac{(D+m-1)!}{(D-1)!m!}. (1.127)$$

(b)

By (a),

$$N(D,0) = 1. (1.128)$$

Let us assume that

$$\sum_{m=0}^{M} n(D,m) = \frac{(D+M)!}{D!M!}.$$
(1.129)

Then,

$$\sum_{m=0}^{M+1} n(D,m) = \frac{(D+M)!}{D!M!} + \frac{(D+M)!}{(D-1)!(M+1)!}.$$
 (1.130)

The right hand side can be written as

$$\frac{(M+1)(D+M)! + D(D+M)!}{D!(M+1)!} = \frac{(D+M+1)!}{D!(M+1)!}.$$
 (1.131)

Then, the assumption is proved by induction on M. Therefore,

$$N(D,M) = \frac{(D+M)!}{D!M!}.$$
(1.132)

(c)

By the approximation

$$n! \simeq n^n \exp(-n),\tag{1.133}$$

we have

$$\frac{(D+M)!}{D!M!} \simeq \frac{(D+M)^{D+M}}{D^D M^M}.$$
 (1.134)

The right hand side can be written as

$$D^{M}\left(1+\frac{M}{D}\right)^{D}\left(\frac{1}{M}+\frac{1}{D}\right)^{M}=M^{D}\left(1+\frac{D}{M}\right)^{M}\left(\frac{1}{D}+\frac{1}{M}\right)^{D}. \quad (1.135)$$

Therefore,

$$N(D,M) \simeq \begin{cases} D^M, & D \gg M, \\ M^D, & M \gg D. \end{cases}$$
 (1.136)

(d)

By (b),

$$N(10,3) = 286,$$
  
 $N(100,3) = 176851,$  (1.137)  
 $N(1000,3) = 167668501.$ 

#### 1.17

Let

$$\Gamma(x) = \int_0^\infty u^{x-1} \exp(-u) du. \tag{1.138}$$

(a)

We have

$$\Gamma(x+1) = \int_0^\infty u^x \exp(-u) du. \tag{1.139}$$

The right hand side can be written as

$$[-u^x \exp(-u)]_{u=0}^{u=\infty} + \int_0^\infty x u^{x-1} \exp(-u) du = x\Gamma(x).$$
 (1.140)

Therefore,

$$\Gamma(x+1) = x\Gamma(x). \tag{1.141}$$

(b)

We have

$$\Gamma(1) = \int_0^\infty \exp(-u)du,\tag{1.142}$$

so that

$$\Gamma(1) = 0!. \tag{1.143}$$

For a positive integer x, let us assume that

$$\Gamma(x) = (x - 1)!. \tag{1.144}$$

By (a),

$$\Gamma(x+1) = x\Gamma(x), \tag{1.145}$$

so that

$$\Gamma(x+1) = x!. \tag{1.146}$$

Therefore, the assumption is proved by induction on x.

#### 1.18

(a)

Let

$$\prod_{d=1}^{D} \int_{-\infty}^{\infty} \exp(-x_d^2) dx_i = S_D \int_{0}^{\infty} \exp(-r^2) r^{D-1} dr, \qquad (1.147)$$

where  $S_D$  is the surface area of a sphere of unit raidus in D dimensions. By 1.7, the left hand side can be written as  $\pi^{\frac{D}{2}}$ . By the transformation

$$s = r^2, (1.148)$$

the right hand side can be written as

$$\frac{S_D}{2} \int_0^\infty \exp(-s) s^{\frac{D-1}{2}} s^{-\frac{1}{2}} ds = \frac{S_D}{2} \Gamma\left(\frac{D}{2}\right). \tag{1.149}$$

Therefore,

$$S_D = \frac{2\pi^{\frac{D}{2}}}{\Gamma\left(\frac{D}{2}\right)}. (1.150)$$

(b)

The volume of the sphere can can be written as

$$V_D = S_D \int_0^1 r^{D-1} dr. (1.151)$$

Therefore,

$$V_D = \frac{S_D}{D}. ag{1.152}$$

(c)

By (a) and (b),

$$S_2 = 2\pi,$$
  $V_2 = \pi.$  (1.153)

Similarly,

$$S_3 = 4\pi,$$

$$V_3 = \frac{4}{3}\pi.$$
(1.154)

#### 1.19

(a)

The volume of a cube of side 2 in D dimensions is  $2^{D}$ . By 1.18, the ratio of the volume of the cocentric sphere of radius 1 divided by the volume of the

cube is given by

$$\frac{V_D}{2^D} = \frac{\pi^{\frac{D}{2}}}{D2^{D-1}\Gamma(\frac{D}{2})}. (1.155)$$

(b)

By (a) and the Stering's formula

$$\Gamma(x+1) \simeq (2\pi)^{\frac{1}{2}} \exp(-x)x^{\frac{x+1}{2}},$$
 (1.156)

we have

$$\frac{V_D}{2^D} \simeq \frac{\pi^{\frac{D}{2}}}{D2^{D-1}(2\pi)^{\frac{1}{2}} \exp\left(1 - \frac{D}{2}\right) \left(\frac{D}{2} - 1\right)^{\frac{D}{4}}}.$$
 (1.157)

The right hand side can be written as

$$\frac{1}{2e(2\pi)^{\frac{1}{2}}} \frac{1}{D} \left( \frac{e^2 \pi^2}{8D - 16} \right)^{\frac{D}{4}}.$$
 (1.158)

Therefore,

$$\lim_{D \to \infty} \frac{V_D}{2^D} = 0. \tag{1.159}$$

(c)

The ratio of the distance from the center of the cube to one of the corners divided by the perpendicular distance to one of the sides is given by

$$\frac{\sqrt{\sum_{i=1}^{D} 1^2}}{1} = \sqrt{D}.\tag{1.160}$$

Therefore, the ratio goes to  $\infty$  as  $D \to \infty$ .

#### 1.20

Let  $\mathbf{x}$  be a variable in D dimensions such that

$$p(\mathbf{x}) = (2\pi\sigma^2)^{-\frac{D}{2}} \exp\left(-\frac{\|\mathbf{x}\|^2}{2\sigma^2}\right). \tag{1.161}$$

(a)

We have

$$\int_{r \le \|\mathbf{x}\| \le r + \epsilon} p(\mathbf{x}) d\mathbf{x} = \int_{r}^{r + \epsilon} \int (2\pi\sigma^{2})^{-\frac{D}{2}} \exp\left(-\frac{r'^{2}}{2\sigma^{2}}\right) J dr' d\boldsymbol{\phi}, \qquad (1.162)$$

where  $\phi$  is the vector of the angular components of the polar coordinate and J is the Jacobian of the transformation from the Cartesian to polar coordinate. For a sufficiently small  $\epsilon$ , the right hand side can be approximated as

$$(2\pi\sigma^2)^{-\frac{D}{2}} \exp\left(-\frac{r^2}{2\sigma^2}\right) \int_r^{r+\epsilon} \int J dr' d\phi$$

$$= (2\pi\sigma^2)^{-\frac{D}{2}} \exp\left(-\frac{r^2}{2\sigma^2}\right) \int_{r<\|\mathbf{x}\| \le r+\epsilon} d\mathbf{x}.$$
(1.163)

Therefore,

$$\int_{r < \|\mathbf{x}\| \le r + \epsilon} p(\mathbf{x}) d\mathbf{x} \simeq p(r) \epsilon, \qquad (1.164)$$

where

$$p(r) = (2\pi\sigma^2)^{-\frac{D}{2}} S_D r^{D-1} \exp\left(-\frac{r^2}{2\sigma^2}\right),$$
 (1.165)

and  $S_D$  is the surface area of a unit sphere in D dimensions.

(b)

Setting the derivative of p(r) to zero gives

$$0 = (2\pi\sigma^2)^{-\frac{D}{2}} S_D \left( (D-1)r^{D-2} - \frac{r^D}{\sigma^2} \right) \exp\left( -\frac{r^2}{2\sigma^2} \right).$$
 (1.166)

Therefore, p(r) is maximised at a sigle stationary point

$$\hat{r} = \sqrt{D - 1}\sigma. \tag{1.167}$$

(c)

We have

$$\frac{p\left(\hat{r}+\epsilon\right)}{p(\hat{r})} = \left(\frac{\hat{r}+\epsilon}{\hat{r}}\right)^{D-1} \exp\left(-\frac{2\hat{r}\epsilon+\epsilon^2}{2\sigma^2}\right). \tag{1.168}$$

By (b), the right hand side can be written as

$$\exp\left((D-1)\ln\left(1+\frac{\epsilon}{\hat{r}}\right) - \frac{2\hat{r}\epsilon + \epsilon^2}{2\sigma^2}\right)$$

$$= \exp\left(\frac{\hat{r}^2}{\sigma^2}\ln\left(1+\frac{\epsilon}{\hat{r}}\right) - \frac{2\hat{r}\epsilon + \epsilon^2}{2\sigma^2}\right). \tag{1.169}$$

By the Taylor series

$$\ln(1+x) = x - \frac{1}{2}x^2 + o(x^3), \qquad (1.170)$$

the right hand side can be approximated as

$$\exp\left(\frac{\hat{r}^2}{\sigma^2}\left(\frac{\epsilon}{\hat{r}} - \frac{\epsilon^2}{2\hat{r}^2}\right) - \frac{2\hat{r}\epsilon + \epsilon^2}{2\sigma^2}\right) = \exp\left(-\frac{\epsilon^2}{\sigma^2}\right). \tag{1.171}$$

Therefore,

$$p(\hat{r} + \epsilon) \simeq p(\hat{r}) \exp\left(-\frac{\epsilon^2}{\sigma^2}\right).$$
 (1.172)

(d)

Let  $\hat{\mathbf{r}}$  be a vector of length  $\hat{r}$ . We have

$$\frac{p(\mathbf{0})}{p(\hat{\mathbf{r}})} = \exp\left(\frac{\hat{r}^2}{2\sigma^2}\right). \tag{1.173}$$

By (b), the right hand side can be written as

$$\exp\left(\frac{D-1}{2}\right). \tag{1.174}$$

Therefore,

$$\frac{p(\mathbf{0})}{p(\hat{\mathbf{r}})} = \exp\left(\frac{D-1}{2}\right). \tag{1.175}$$

#### 1.21

(a)

If  $0 \le a \le b$ , then

$$0 \le a(b-a). \tag{1.176}$$

$$a \le (ab)^{\frac{1}{2}}. (1.177)$$

(b)

For a two-class classification problem of  $\mathbf{x}$ , let the classes be  $\mathcal{C}_1$  and  $\mathcal{C}_2$  and let the decision regions be  $\mathcal{R}_1$  and  $\mathcal{R}_2$ . Let us choose the decision regions to minimise the probability of misclassification. Then,

$$p(\mathbf{x}, C_1) > p(\mathbf{x}, C_2) \Rightarrow \mathbf{x} \in C_1,$$
  

$$p(\mathbf{x}, C_2) > p(\mathbf{x}, C_1) \Rightarrow \mathbf{x} \in C_2.$$
(1.178)

By (a),

$$\int_{\mathcal{R}_{1}} p(\mathbf{x}, \mathcal{C}_{2}) d\mathbf{x} \leq \int_{\mathcal{R}_{1}} \left( p(\mathbf{x}, \mathcal{C}_{1}) p(\mathbf{x}, \mathcal{C}_{2}) \right)^{\frac{1}{2}} d\mathbf{x}, 
\int_{\mathcal{R}_{2}} p(\mathbf{x}, \mathcal{C}_{1}) d\mathbf{x} \leq \int_{\mathcal{R}_{2}} \left( p(\mathbf{x}, \mathcal{C}_{1}) p(\mathbf{x}, \mathcal{C}_{2}) \right)^{\frac{1}{2}} d\mathbf{x}.$$
(1.179)

Therefore,

$$\int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x} \le \int (p(\mathbf{x}, \mathcal{C}_1) p(\mathbf{x}, \mathcal{C}_2))^{\frac{1}{2}} d\mathbf{x}.$$
 (1.180)

#### 1.22

Let

$$EL = \sum_{k} \sum_{j} \int_{\mathcal{R}_{j}} L_{kj} p(\mathbf{x}, \mathcal{C}_{k}) d\mathbf{x}.$$
 (1.181)

If

$$L_{ki} = 1 - I_{ki}, (1.182)$$

then the right hand side can be written as

$$\sum_{k} \sum_{j} \int_{\mathcal{R}_{j}} (p(\mathbf{x}, \mathcal{C}_{k}) - p(\mathbf{x}, \mathcal{C}_{j})) d\mathbf{x} = \sum_{j} \int_{\mathcal{R}_{j}} \left( \sum_{k} p(\mathbf{x}, \mathcal{C}_{k}) - p(\mathbf{x}, \mathcal{C}_{j}) \right) d\mathbf{x}.$$
(1.183)

The right hand side can be written as

$$\sum_{i} \int_{\mathcal{R}_{j}} (p(\mathbf{x}) - p(\mathbf{x}, \mathcal{C}_{j})) d\mathbf{x} = 1 - \sum_{i} \int_{\mathcal{R}_{j}} p(\mathbf{x}, \mathcal{C}_{j}) d\mathbf{x}.$$
 (1.184)

Then,

$$EL = 1 - \sum_{j} \int_{\mathcal{R}_{j}} p(\mathcal{C}_{j}|\mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$
 (1.185)

Therefore, minimising E L reduces to choosing the criterion to maximise the posterior probability  $p(C_i|\mathbf{x})$ .

#### 1.23

Let

$$EL = \sum_{k} \sum_{j} \int_{\mathcal{R}_{j}} L_{kj} p(\mathbf{x}, \mathcal{C}_{k}) d\mathbf{x}.$$
 (1.186)

The right hand side can be written as

$$\sum_{j} \int_{\mathcal{R}_{j}} \sum_{k} L_{kj} p(\mathbf{x}, \mathcal{C}_{k}) d\mathbf{x} = \sum_{j} \int_{\mathcal{R}_{j}} \left( \sum_{k} L_{kj} p(\mathcal{C}_{k} | \mathbf{x}) \right) p(\mathbf{x}) d\mathbf{x}.$$
 (1.187)

Then,

$$EL = \sum_{j} \int_{\mathcal{R}_{j}} \left( \sum_{k} L_{kj} p(\mathcal{C}_{k} | \mathbf{x}) \right) p(\mathbf{x}) d\mathbf{x}.$$
 (1.188)

Therefore, minimising EL reduces to minimising  $\sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$ .

#### 1.24 (Incomplete)

Let

$$EL = \sum_{k} \sum_{j} \int_{\mathcal{R}_{j}} L_{kj} p(\mathbf{x}, \mathcal{C}_{k}) d\mathbf{x} + \lambda \int_{\forall kp(\mathcal{C}_{k}|\mathbf{x}) < \theta} p(\mathbf{x}) d\mathbf{x}.$$
 (1.189)

#### 1.25

Let

$$E L(\mathbf{t}, \mathbf{y}(\mathbf{x})) = \int \int \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}.$$
 (1.190)

Setting the derivative with respect to  $\mathbf{y}(\mathbf{x})$  to zero gives

$$\mathbf{0} = 2 \int (\mathbf{y}(\mathbf{x}) - \mathbf{t}) p(\mathbf{x}, \mathbf{t}) d\mathbf{t}. \tag{1.191}$$

The integral of the right hand side can be written as

$$\mathbf{y}(\mathbf{x}) \int p(\mathbf{x}, \mathbf{t}) d\mathbf{t} - \int \mathbf{t} p(\mathbf{x}, \mathbf{t}) d\mathbf{t} = \mathbf{y}(\mathbf{x}) p(\mathbf{x}) - p(\mathbf{x}) \int \mathbf{t} p(\mathbf{t}|\mathbf{x}) d\mathbf{t}. \quad (1.192)$$

The integral in the second term of the right hand side can be written as  $E_{\mathbf{t}}(\mathbf{t}|\mathbf{x})$ . Then, the right hand side can be written as

$$\mathbf{0} = p(\mathbf{x}) \left( \mathbf{y}(\mathbf{x}) - \mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}) \right). \tag{1.193}$$

Therefore,

$$\underset{\mathbf{y}(\mathbf{x})}{\operatorname{argmin}} E L(\mathbf{t}, \mathbf{y}(\mathbf{x})) = E_{\mathbf{t}}(\mathbf{t}|\mathbf{x}). \tag{1.194}$$

For a single target variable t, it reduces to

$$\underset{\mathbf{y}(\mathbf{x})}{\operatorname{argmin}} E L(t, \mathbf{y}(\mathbf{x})) = E_t(t|\mathbf{x}). \tag{1.195}$$

#### 1.26

Let

$$E L(\mathbf{t}, \mathbf{y}(\mathbf{x})) = \iint \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}.$$
 (1.196)

The right hand side can be written as

$$\iint \|\mathbf{y}(\mathbf{x}) - \mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}) + \mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}) - \mathbf{t}\|^{2} p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}$$

$$= \iint \|\mathbf{y}(\mathbf{x}) - \mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x})\|^{2} p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}$$

$$+ 2 \iint (\mathbf{y}(\mathbf{x}) - \mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}))^{\mathsf{T}} (\mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}) - \mathbf{t}) p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}$$

$$+ \iint \|\mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}) - \mathbf{t}\|^{2} p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}.$$
(1.197)

Let us look at each term of the right hand side. The first term can be written as

$$\int \|\mathbf{y}(\mathbf{x}) - \mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x})\|^{2} \left( \int p(\mathbf{x}, \mathbf{t}) d\mathbf{t} \right) d\mathbf{x} = \int \|\mathbf{y}(\mathbf{x}) - \mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x})\|^{2} p(\mathbf{x}) d\mathbf{x}.$$
(1.198)

The integral of the second term can be written as

$$\int (\mathbf{y}(\mathbf{x}) - \mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}))^{\mathsf{T}} \left( \int (\mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}) - \mathbf{t}) p(\mathbf{t}|\mathbf{x}) d\mathbf{t} \right) p(\mathbf{x}) d\mathbf{x}.$$
 (1.199)

Since

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

$$\int \mathbf{t}p(\mathbf{t}|\mathbf{x})d\mathbf{t} = E_{\mathbf{t}}(\mathbf{t}|\mathbf{x}),$$
(1.200)

the second term is zero. The third term can be written as

$$\int \left( \int \|\mathbf{E}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{t}|\mathbf{x}) d\mathbf{t} \right) p(\mathbf{x}) d\mathbf{x} = \int \operatorname{var}_{\mathbf{t}}(\mathbf{t}|\mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$
 (1.201)

Then,

$$EL(\mathbf{t}, \mathbf{y}(\mathbf{x})) = \int ||\mathbf{y}(\mathbf{x}) - E_{\mathbf{t}}(\mathbf{t}|\mathbf{x})||^{2} p(\mathbf{x}) d\mathbf{x} + \int var_{\mathbf{t}}(\mathbf{t}|\mathbf{x}) p(\mathbf{x}) d\mathbf{x}. \quad (1.202)$$

Therefore,

$$\underset{\mathbf{y}(\mathbf{x})}{\operatorname{argmin}} E L(\mathbf{t}, \mathbf{y}(\mathbf{x})) = E_{\mathbf{t}}(\mathbf{t}|\mathbf{x}). \tag{1.203}$$

### 1.27 (Incomplete)

(a)

Let

$$EL_q = \int \int |y(\mathbf{x}) - t|^q p(\mathbf{x}, t) d\mathbf{x} dt.$$
 (1.204)

Setting the derivative with respect to  $y(\mathbf{x})$  to zero gives

$$0 = qp(\mathbf{x}) \int |y(\mathbf{x}) - t|^{q-1} \operatorname{sign}(y(\mathbf{x}) - t)p(t|\mathbf{x})dt.$$
 (1.205)

$$\underset{y(\mathbf{x})}{\operatorname{argmin}} E L_q = \left\{ y(\mathbf{x}) \mid \int |y(\mathbf{x}) - t|^{q-1} \operatorname{sign}(y(\mathbf{x}) - t) p(t|\mathbf{x}) dt = 0 \right\}.$$
(1.206)

(b)

We have

$$EL_1 = \int \left( \int \operatorname{sign}(y(\mathbf{x}) - t) p(t|\mathbf{x}) dt \right) p(\mathbf{x}) d\mathbf{x}.$$
 (1.207)

The integral of the right hand side with respect to t can be written as

$$\int_{y(\mathbf{x})}^{\infty} p(t|\mathbf{x})dt - \int_{-\infty}^{y(\mathbf{x})} p(t|\mathbf{x})dt.$$
 (1.208)

Therefore,

$$\underset{y(\mathbf{x})}{\operatorname{argmin}} E L_1 = \operatorname{median}(t|\mathbf{x}). \tag{1.209}$$

(c)

We have

$$\lim_{q \to 0} \left( \underset{y(\mathbf{x})}{\operatorname{argmin}} \, \mathbf{E} \, L_q \right) = \operatorname{mode}(t|\mathbf{x})? \tag{1.210}$$

#### 1.28

(a)

Let us assume that

$$p(x,y) = p(x)p(y) \Rightarrow h(x,y) = h(x) + h(y).$$
 (1.211)

Then,

$$h\left(p^2\right) = 2h(p). \tag{1.212}$$

Let us assume that, for a positive integer n,

$$h\left(p^{n}\right) = nh(p). \tag{1.213}$$

Then, by the first assumption,

$$h(p^{n+1}) = h(p^n) + h(p),$$
 (1.214)

so that

$$h(p^{n+1}) = (n+1)h(p).$$
 (1.215)

Therefore, the second assumption is proved by induction on n.

(b)

For positive integers m and n,

$$h\left(p^{n}\right) = h\left(p^{\frac{n}{m}m}\right). \tag{1.216}$$

By the second assumption in (a), the left hand side can be written as nh(p). By the first assumption in (a), the right hand side can be written as  $mh\left(p^{\frac{n}{m}}\right)$ . Therefore,

$$h\left(p^{\frac{n}{m}}\right) = \frac{n}{m}h(p). \tag{1.217}$$

(c)

By the continuity, for a positive real number a,

$$h\left(p^{a}\right) = ah(p). \tag{1.218}$$

Taking the derivative with respect to a and substituting a = 1 gives

$$(p \ln p)h'(p) = h(p).$$
 (1.219)

Then,

$$\int \frac{h'(p)}{h(p)} dp = \int \frac{1}{p \ln p} dp + \text{const}.$$
 (1.220)

Ignorting the constants, the left hand side can be written as  $\ln h(p)$  and the right hand side can be written as  $\ln(\ln p)$ . Therefore,

$$h(p) \propto \ln p. \tag{1.221}$$

#### 1.29

Let x be an M-state discrete random variable. Then, the entropy is given by

$$H(x) = -\sum_{m=1}^{M} p(x_m) \ln p(x_m), \qquad (1.222)$$

where

$$\sum_{m=1}^{M} p(x_m) = 1. (1.223)$$

By the Jensen's inequality,

$$\sum_{m=1}^{M} p(x_i) \ln \frac{1}{p(x_m)} \le \ln \left( \sum_{m=1}^{M} 1 \right).$$
 (1.224)

Therefore,

$$H(x) \le \ln M. \tag{1.225}$$

#### 1.30

Let

$$p(x) = \mathcal{N}(x|\mu, \sigma^2),$$
  

$$q(x) = \mathcal{N}(x|m, s^2).$$
(1.226)

Then, the Kullback-Leibler divergence is given by

$$KL(p||q) = -\int p(x) \ln \frac{q(x)}{p(x)} dx. \qquad (1.227)$$

The right hand side can be written as

$$-\int_{-\infty}^{\infty} p(x) \ln \frac{(2\pi s^2)^{-\frac{1}{2}} \exp\left(-\frac{(x-m)^2}{2s^2}\right)}{(2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)} dx$$

$$= -\int_{-\infty}^{\infty} p(x) \left(-\frac{1}{2} \ln \frac{s^2}{\sigma^2} - \frac{(x-m)^2}{2s^2} + \frac{(x-\mu)^2}{2\sigma^2}\right) dx.$$
(1.228)

The right hand side can be written as

$$\ln \frac{s}{\sigma} \int_{-\infty}^{\infty} p(x)dx + \frac{1}{2s^2} \int_{-\infty}^{\infty} (x-m)^2 p(x)dx - \frac{1}{2\sigma^2} \int_{-\infty}^{\infty} (x-\mu)^2 p(x)dx. \quad (1.229)$$

The integral of the second term can be written as

$$\int_{-\infty}^{\infty} (x - \mu + \mu - m)^2 p(x) dx$$

$$= \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx + 2(\mu - m) \int_{-\infty}^{\infty} (x - \mu) p(x) dx$$

$$+ (\mu - m)^2 \int_{-\infty}^{\infty} p(x) dx.$$
(1.230)

$$KL(p||q) = \ln \frac{s}{\sigma} + \frac{\sigma^2 + (\mu - m)^2}{2s^2} - \frac{1}{2}.$$
 (1.231)

Let  $\mathbf{x}$  and  $\mathbf{y}$  be two variables. We have

$$H(\mathbf{x}) = -\int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x},$$

$$H(\mathbf{y}) = -\int p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{y},$$

$$H(\mathbf{x}, \mathbf{y}) = -\int \int p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}.$$
(1.232)

Since

$$H(\mathbf{x}) = -\int \left(\int p(\mathbf{x}, \mathbf{y}) d\mathbf{y}\right) \ln p(\mathbf{x}) d\mathbf{x},$$

$$H(\mathbf{y}) = -\int \left(\int p(\mathbf{x}, \mathbf{y}) d\mathbf{x}\right) \ln p(\mathbf{y}) d\mathbf{y},$$
(1.233)

we have

$$H(\mathbf{x}) + H(\mathbf{y}) - H(\mathbf{x}, \mathbf{y}) = -\iint p(\mathbf{x}, \mathbf{y}) \ln \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} d\mathbf{x} d\mathbf{y}.$$
 (1.234)

Since

$$\iint p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} = 1, \tag{1.235}$$

the Jensen's inequality can be used to have

$$-\iint p(\mathbf{x}, \mathbf{y}) \ln \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} d\mathbf{x} d\mathbf{y} \ge -\ln \left( \iint p(\mathbf{x})p(\mathbf{y}) d\mathbf{x} d\mathbf{y} \right). \tag{1.236}$$

The right hand side can be written as

$$-\ln\left(\int p(\mathbf{x})d\mathbf{x}\int p(\mathbf{y})d\mathbf{y}\right) = 0. \tag{1.237}$$

$$H(\mathbf{x}, \mathbf{y}) \le H(\mathbf{x}) + H(\mathbf{y}). \tag{1.238}$$

Let  $\mathbf{x}$  and  $\mathbf{y}$  be variables such that

$$\mathbf{y} = \mathbf{A}\mathbf{x},\tag{1.239}$$

where A is a nonsingular matrix. We have

$$\int p_x(\mathbf{x})d\mathbf{x} = \int p_x\left(\mathbf{A}^{-1}\mathbf{y}\right) \left| \det \mathbf{A}^{-1} \right| d\mathbf{y}.$$
 (1.240)

Then,

$$p_y(\mathbf{y}) = p_x \left( \mathbf{A}^{-1} \mathbf{y} \right) \left| \det \mathbf{A}^{-1} \right|. \tag{1.241}$$

We have

$$H(\mathbf{y}) = -\int p_y(\mathbf{y}) \ln p_y(\mathbf{y}) d\mathbf{y}.$$
 (1.242)

The right hand side can be written as

$$-\int p_{y}(\mathbf{y}) \ln \left( p_{x} \left( \mathbf{A}^{-1} \mathbf{y} \right) \left| \det \mathbf{A}^{-1} \right| \right) d\mathbf{y}$$

$$= -\int p_{y}(\mathbf{y}) \ln p_{x} \left( \mathbf{A}^{-1} \mathbf{y} \right) d\mathbf{y} + \ln \left| \det \mathbf{A} \right| \int p_{y}(\mathbf{y}) d\mathbf{y}.$$
(1.243)

The first term can be written as

$$-\left|\det \mathbf{A}^{-1}\right| \int p_x \left(\mathbf{A}^{-1} \mathbf{y}\right) \ln p_x \left(\mathbf{A}^{-1} \mathbf{y}\right) d\mathbf{y}. \tag{1.244}$$

By the transformation

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{y},\tag{1.245}$$

it can be written as

$$-\int p_x(\mathbf{x}) \ln p_x(\mathbf{x}) d\mathbf{x} = \mathbf{H}(\mathbf{x}). \tag{1.246}$$

$$H(\mathbf{y}) = H(\mathbf{x}) + \ln|\det \mathbf{A}|. \tag{1.247}$$

Let x and y be two discrete variables with K and L states. Then,

$$H(y|x) = -\sum_{k=1}^{K} \sum_{l=1}^{L} p(x_k, y_l) \ln p(y_l|x_k).$$
 (1.248)

If

$$H(y|x) = 0, (1.249)$$

then

$$0 = -\sum_{k=1}^{K} p(x_k) \sum_{l=1}^{L} p(y_l|x_k) \ln p(y_l|x_k).$$
 (1.250)

Since

$$p(x_k) \ge 0,$$
  
 $p(y_l|x_k) \ln p(y_l|x_k) \le 0,$  (1.251)

the equation reduces to

$$p(y_l|x_k) \ln p(y_l|x_k) = 0. (1.252)$$

Then,

$$p(y_l|x_k) = \begin{cases} 1, \\ 0. \end{cases}$$
 (1.253)

Since

$$\sum_{l=1}^{L} p(y_l|x_k) = 1, (1.254)$$

we have

$$p(y_l|x_k) = \begin{cases} 1, & \text{if } K = 1, \\ 0, & \text{otherwise.} \end{cases}$$
 (1.255)

#### 1.34

Let x be a variable. We have

$$H(x) = -\int_{-\infty}^{\infty} p(x) \ln p(x) dx. \tag{1.256}$$

In order to maximise H(x) with the constratints

$$\int_{-\infty}^{\infty} p(x)dx = 1,$$

$$\int_{-\infty}^{\infty} xp(x)dx = \mu,$$

$$\int_{-\infty}^{\infty} (x - \mu)^2 p(x)dx = \sigma^2,$$
(1.257)

let

$$L(p) = H(x) + \lambda_1 \left( \int_{-\infty}^{\infty} p(x)dx - 1 \right) + \lambda_2 \left( \int_{-\infty}^{\infty} x p(x)dx - \mu \right)$$

$$+ \lambda_3 \left( \int_{-\infty}^{\infty} (x - \mu)^2 p(x)dx - \sigma^2 \right).$$
(1.258)

Setting the variation with respect to p to zero gives

$$0 = -\ln p - 1 + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2. \tag{1.259}$$

Then,

$$p(x) = \exp(-1 + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2),$$
 (1.260)

so that

$$p(x) = c \exp\left(\lambda_3 \left(x - \left(\mu - \frac{\lambda_2}{2\lambda_3}\right)\right)^2\right),$$
 (1.261)

where

$$c = \exp\left(-1 + \lambda_1 - \frac{\lambda_2^2}{4\lambda_3}\right). \tag{1.262}$$

Substituting it to the constraints gives

$$c \int_{-\infty}^{\infty} \exp\left(\lambda_3 \left(x - \left(\mu - \frac{\lambda_2}{2\lambda_3}\right)\right)^2\right) dx = 1,$$

$$c \int_{-\infty}^{\infty} x \exp\left(\lambda_3 \left(x - \left(\mu - \frac{\lambda_2}{2\lambda_3}\right)\right)^2\right) dx = \mu,$$

$$c \int_{-\infty}^{\infty} (x - \mu)^2 \exp\left(\lambda_3 \left(x - \left(\mu - \frac{\lambda_2}{2\lambda_3}\right)\right)^2\right) dx = \sigma^2.$$

$$(1.263)$$

By the transformation

$$y = (-\lambda_3)^{\frac{1}{2}} \left( x - \left( \mu - \frac{\lambda_2}{2\lambda_3} \right) \right),$$
 (1.264)

they can be written as

$$c \int_{-\infty}^{\infty} \exp(-y^{2}) (-\lambda_{3})^{-\frac{1}{2}} dy = 1,$$

$$c \int_{-\infty}^{\infty} \left( (-\lambda_{3})^{-\frac{1}{2}} y + \mu - \frac{\lambda_{2}}{2\lambda_{3}} \right) \exp(-y^{2}) (-\lambda_{3})^{-\frac{1}{2}} dy = \mu,$$

$$c \int_{-\infty}^{\infty} \left( (-\lambda_{3})^{-\frac{1}{2}} y - \frac{\lambda_{2}}{2\lambda_{3}} \right)^{2} \exp(-y^{2}) (-\lambda_{3})^{-\frac{1}{2}} dy = \sigma^{2}.$$
(1.265)

Since

$$\int_{-\infty}^{\infty} \exp(-y^2) dy = \Gamma\left(\frac{1}{2}\right),$$

$$\int_{-\infty}^{\infty} y \exp(-y^2) dy = 0,$$

$$\int_{-\infty}^{\infty} y^2 \exp(-y^2) dy = \Gamma\left(\frac{3}{2}\right),$$
(1.266)

they can be written as

$$c(-\lambda_3)^{-\frac{1}{2}}\Gamma\left(\frac{1}{2}\right) = 1,$$

$$c\left(\mu - \frac{\lambda_2}{2\lambda_3}\right)(-\lambda_3)^{-\frac{1}{2}}\Gamma\left(\frac{1}{2}\right) = \mu,$$

$$c\left((-\lambda_3)^{-\frac{3}{2}}\Gamma\left(\frac{3}{2}\right) + (-\lambda_3)^{-\frac{1}{2}}\frac{\lambda_2^2}{4\lambda_3^2}\Gamma\left(\frac{1}{2}\right)\right) = \sigma^2.$$

$$(1.267)$$

Then,

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

$$\lambda_2 = 0,$$

$$\lambda_3 = -\frac{1}{2\sigma^2}.$$
(1.268)

$$p(x) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right).$$
 (1.269)

Let x be a variable such that

$$p(x) = \mathcal{N}\left(x|\mu, \sigma^2\right). \tag{1.270}$$

Then,

$$H(x) = -\int_{-\infty}^{\infty} \mathcal{N}\left(x|\mu,\sigma^2\right) \ln \mathcal{N}\left(x|\mu,\sigma^2\right) dx. \tag{1.271}$$

The right hand side can be written as

$$-\int_{-\infty}^{\infty} \mathcal{N}\left(x|\mu,\sigma^2\right) \left(-\frac{1}{2}\ln\left(2\pi\sigma^2\right) - \frac{1}{2\sigma^2}(x-\mu)^2\right) dx$$

$$= \frac{1}{2}\ln\left(2\pi\sigma^2\right) \int_{-\infty}^{\infty} \mathcal{N}\left(x|\mu,\sigma^2\right) dx + \frac{1}{2\sigma^2} \int_{-\infty}^{\infty} (x-\mu)^2 \mathcal{N}\left(x|\mu,\sigma^2\right) dx.$$
(1.272)

Therefore,

$$H(x) = \frac{1}{2} (1 + \ln(2\pi\sigma^2)).$$
 (1.273)

# 1.36 (Incomplete)

Let f be a strictly convex function. Then,

$$f(\lambda a + (1 - \lambda)b) \le \lambda f(a) + (1 - \lambda)f(b), \tag{1.274}$$

where  $a \leq b$  and  $0 \leq \lambda \leq 1$ . Let

$$x = \lambda a + (1 - \lambda)b. \tag{1.275}$$

Then, the inequality can be written as

$$f(x) \le \frac{b-x}{b-a}f(a) + \frac{x-a}{b-a}f(b).$$
 (1.276)

Let

$$g(x) = \frac{b-x}{b-a}f(a) + \frac{x-a}{b-a}f(b) - f(x). \tag{1.277}$$

Then,

$$g(x) \ge 0. \tag{1.278}$$

Additionally, for x > a,

$$g(x) = (x - a) \left( \frac{f(b) - f(a)}{b - a} - \frac{f(x) - f(a)}{x - a} \right).$$
 (1.279)

By the mean value theorem, there exists c and y such that  $a \leq c \leq b$ ,  $a \leq y \leq x$  and

$$f'(c) = \frac{f(b) - f(a)}{b - a},$$
  

$$f'(y) = \frac{f(x) - f(a)}{x - a}.$$
(1.280)

Then, for x > a, the inequality reduces to

$$f'(y) \le f'(c). \tag{1.281}$$

### 1.37

Let  $\mathbf{x}$  and  $\mathbf{y}$  be two variables. We have

$$H(\mathbf{x}, \mathbf{y}) = -\iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}.$$
 (1.282)

The right hand side can be written as

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

$$= - \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{x} d\mathbf{y} - \int \left( \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \right) \ln p(\mathbf{x}) d\mathbf{x}.$$
(1.283)

By the definition, the first and second terms of the right hand side can be written as H(y|x) and H(x). Therefore,

$$H(\mathbf{x}, \mathbf{y}) = H(\mathbf{y}|\mathbf{x}) + H(\mathbf{x}). \tag{1.284}$$

#### 1.38

Let f be a strictly convex function. Then,

$$f(\lambda x_1 + (1 - \lambda)x_2) \le \lambda f(x_1) + (1 - \lambda)f(x_2),$$
 (1.285)

where

$$0 \le \lambda \le 1. \tag{1.286}$$

Let us assume that

$$f\left(\sum_{m=1}^{M} \lambda_m x_m\right) \le \sum_{m=1}^{M} \lambda_m f(x_m), \tag{1.287}$$

where

$$\sum_{m=1}^{M} \lambda_m = 1,$$

$$\lambda_m \ge 0.$$
(1.288)

Since f is strictly convex,

$$f\left(\sum_{m=1}^{M+1} \lambda_m x_m\right) \le \lambda_{M+1} f(x_{M+1}) + (1 - \lambda_{M+1}) f\left(\sum_{m=1}^{M} \frac{\lambda_m}{1 - \lambda_{M+1}} x_m\right),\tag{1.289}$$

where

$$\sum_{m=1}^{M+1} \lambda_m = 1,$$

$$\lambda_m \ge 0.$$
(1.290)

By the assumption,

$$f\left(\sum_{m=1}^{M} \frac{\lambda_m}{1 - \lambda_{M+1}} x_m\right) \le \sum_{m=1}^{M} \frac{\lambda_m}{1 - \lambda_{M+1}} f(x_m).$$
 (1.291)

Then,

$$f\left(\sum_{m=1}^{M+1} \lambda_m x_m\right) \le \lambda_{M+1} f(x_{M+1}) + (1 - \lambda_{M+1}) \sum_{m=1}^{M} \frac{\lambda_m}{1 - \lambda_{M+1}} f(x_m), (1.292)$$

so that

$$f\left(\sum_{m=1}^{M+1} \lambda_m x_m\right) \le \sum_{m=1}^{M+1} \lambda_m f(x_m). \tag{1.293}$$

Therefore, the assumption is proved by induction on M.

Let x and y be two binary variables where

$$p(x = 0, y = 0) = \frac{1}{3},$$

$$p(x = 0, y = 1) = \frac{1}{3},$$

$$p(x = 1, y = 0) = 0,$$

$$p(x = 1, y = 1) = \frac{1}{3}.$$
(1.294)

(a)

We have

$$H(x) = -\sum_{x} p(x) \ln p(x).$$
 (1.295)

We have

$$p(x = 0) = p(x = 0, y = 0) + p(x = 0, y = 1),$$
  

$$p(x = 1) = p(x = 1, y = 0) + p(x = 1, y = 1).$$
(1.296)

Then,

$$p(x = 0) = \frac{2}{3},$$

$$p(x = 1) = \frac{1}{3}.$$
(1.297)

Therefore,

$$H(x) = \ln 3 - \frac{2}{3} \ln 2. \tag{1.298}$$

(b)

We have

$$H(y) = -\sum_{y} p(y) \ln p(y).$$
 (1.299)

$$p(y=0) = p(x=0, y=0) + p(x=1, y=0),$$
  

$$p(y=1) = p(x=0, y=1) + p(x=1, y=1).$$
(1.300)

Then,

$$p(y = 0) = \frac{1}{3},$$

$$p(y = 1) = \frac{2}{3}.$$
(1.301)

Therefore,

$$H(y) = \ln 3 - \frac{2}{3} \ln 2. \tag{1.302}$$

(c)

We have

$$H(y|x) = -\sum_{x,y} p(x,y) \ln p(y|x).$$
 (1.303)

By the Bayes' theorem,

$$p(y = 0|x = 0) = \frac{p(x = 0, y = 0)}{p(x = 0)},$$

$$p(y = 0|x = 1) = \frac{p(x = 1, y = 0)}{p(x = 1)},$$

$$p(y = 1|x = 0) = \frac{p(x = 0, y = 1)}{p(x = 0)},$$

$$p(y = 1|x = 1) = \frac{p(x = 1, y = 1)}{p(x = 1)}.$$

$$(1.304)$$

Then,

$$p(y = 0|x = 0) = \frac{1}{2},$$

$$p(y = 0|x = 1) = 0,$$

$$p(y = 1|x = 0) = \frac{1}{2},$$

$$p(y = 1|x = 1) = 1.$$
(1.305)

$$H(y|x) = \frac{2}{3}\ln 2. \tag{1.306}$$

(d)

We have

$$H(x|y) = -\sum_{x,y} p(x,y) \ln p(x|y).$$
 (1.307)

By the Bayes' theorem,

$$p(x = 0|y = 0) = \frac{p(x = 0, y = 0)}{p(y = 0)},$$

$$p(x = 0|y = 1) = \frac{p(x = 0, y = 1)}{p(y = 1)},$$

$$p(x = 1|y = 0) = \frac{p(x = 1, y = 0)}{p(y = 0)},$$

$$p(x = 1|y = 1) = \frac{p(x = 1, y = 1)}{p(y = 1)}.$$
(1.308)

Then,

$$p(x = 0|y = 0) = 1,$$

$$p(x = 0|y = 1) = \frac{1}{2},$$

$$p(x = 1|y = 0) = 0,$$

$$p(x = 1|y = 1) = \frac{1}{2}.$$
(1.309)

Therefore,

$$H(x|y) = \frac{2}{3} \ln 2. \tag{1.310}$$

(e)

We have

$$H(x,y) = -\sum_{x,y} p(x,y) \ln p(x,y).$$
 (1.311)

$$H(x, y) = \ln 3.$$
 (1.312)

(f)

We have

$$I(x,y) = -\sum_{x,y} p(x,y) \ln \frac{p(x)p(y)}{p(x,y)}.$$
 (1.313)

The right hand side can be written as

$$H(x) + H(y) - H(x, y).$$
 (1.314)

Therefore,

$$I(x,y) = \ln 3 - \frac{4}{3} \ln 2. \tag{1.315}$$

#### 1.40

Let  $\lambda_1, \dots, \lambda_M$  and  $x_1, \dots, x_M$  be numbers such that

$$\sum_{m=1}^{M} \lambda_m = 1,$$

$$\lambda_m \ge 0,$$

$$x_m > 0.$$
(1.316)

By the Jensen's inequality,

$$\sum_{m=1}^{M} \lambda_m \ln x_m \le \ln \left( \sum_{m=1}^{M} \lambda_m x_m \right), \tag{1.317}$$

so that

$$\prod_{m=1}^{M} x_m^{\lambda_m} \le \sum_{m=1}^{M} \lambda_m x_m. \tag{1.318}$$

Substituting

$$\lambda_m = \frac{1}{M} \tag{1.319}$$

to the inequality gives

$$\left(\prod_{m=1}^{M} x_m\right)^{\frac{1}{M}} \le \frac{1}{M} \sum_{m=1}^{M} x_m. \tag{1.320}$$

Let  $\mathbf{x}$  and  $\mathbf{y}$  be variables.

(a)

We have

$$I(\mathbf{x}, \mathbf{y}) = -\iint p(\mathbf{x}, \mathbf{y}) \ln \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} d\mathbf{x} d\mathbf{y}.$$
 (1.321)

By the Bayes' theorem, the right hand side can be written as

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

$$= - \iint p(\mathbf{x}, \mathbf{y}) \left( \ln p(\mathbf{x}) - \ln p(\mathbf{x}|\mathbf{y}) \right) d\mathbf{x} d\mathbf{y}.$$
(1.322)

The right hand side can be written as

$$-\int \left(\int p(\mathbf{x}, \mathbf{y}) d\mathbf{y}\right) \ln p(\mathbf{x}) d\mathbf{x} + \int \int p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}|\mathbf{y}) d\mathbf{x} d\mathbf{y}.$$
(1.323)

Therefore,

$$I(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) - H(\mathbf{x}|\mathbf{y}). \tag{1.324}$$

(b)

We have

$$I(\mathbf{x}, \mathbf{y}) = I(\mathbf{y}, \mathbf{x}). \tag{1.325}$$

By (a), the right hand side can be written as

$$H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x}). \tag{1.326}$$

$$I(\mathbf{x}, \mathbf{y}) = H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x}). \tag{1.327}$$

# 2 Probability Distributions

# 2.1

Let x be a variable such that

$$x \in \{0, 1\},\$$
  
$$p(x) = \mu^{x} (1 - \mu)^{1 - x}.$$
 (2.1)

(a)

We have

$$\sum_{x} p(x) = 1 - \mu + \mu. \tag{2.2}$$

Therefore,

$$\sum_{x} p(x) = 1. \tag{2.3}$$

(b)

We have

$$E x = \sum_{x} xp(x),$$

$$E x^{2} = \sum_{x} x^{2}p(x),$$
(2.4)

Then,

$$E x = \mu,$$
  

$$E x^2 = \mu.$$
(2.5)

We have

$$\operatorname{var} x = \operatorname{E} x^{2} - (\operatorname{E} x)^{2}.$$
 (2.6)

Therefore,

$$\operatorname{var} x = \mu(1 - \mu). \tag{2.7}$$

(c)

$$H(x) = -\sum_{x} p(x) \ln p(x). \tag{2.8}$$

Therefore,

$$H(x) = -\mu \ln \mu - (1 - \mu) \ln(1 - \mu). \tag{2.9}$$

# 2.2

Let x be a variable such that

$$x \in \{-1, 1\},\$$

$$p(x) = \left(\frac{1-\mu}{2}\right)^{\frac{1-x}{2}} \left(\frac{1+\mu}{2}\right)^{\frac{1+x}{2}}.$$
(2.10)

(a)

We have

$$\sum_{x} p(x) = \frac{1-\mu}{2} + \frac{1+\mu}{2}.$$
 (2.11)

Therefore,

$$\sum_{x} p(x) = 1. {(2.12)}$$

(b)

We have

$$E x = \sum_{x} xp(x),$$

$$E x^{2} = \sum_{x} x^{2}p(x).$$
(2.13)

The right hand sides can be written as

$$-\frac{1-\mu}{2} + \frac{1+\mu}{2} = \mu,$$

$$\frac{1-\mu}{2} + \frac{1+\mu}{2} = 1.$$
(2.14)

Then,

$$E x = \mu,$$
  

$$E x^2 = 1.$$
(2.15)

$$\operatorname{var} x = \operatorname{E} x^2 - (\operatorname{E} x)^2.$$
 (2.16)

Therefore,

$$var x = 1 - \mu^2. (2.17)$$

(c)

We have

$$H(x) = -\sum_{x} p(x|\mu) \ln p(x|\mu).$$
 (2.18)

Therefore,

$$H(x) = -\frac{1-\mu}{2} \ln \frac{1-\mu}{2} - \frac{1+\mu}{2} \ln \frac{1+\mu}{2}.$$
 (2.19)

### 2.3

(a)

We have

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

$$\binom{N}{n-1} = \frac{N!}{(n-1)!(N-n+1)!}$$
(2.20)

Then,

$$\binom{N}{n} + \binom{N}{n-1} = \frac{(N-n+1)N! + nN!}{n!(N-n+1)!}.$$
 (2.21)

The right hand side can be written as

$$\frac{(N+1)!}{n!(N+1-n)!} = \binom{N+1}{n}.$$
 (2.22)

Therefore,

$$\binom{N}{n} + \binom{N}{n-1} = \binom{N+1}{n}. \tag{2.23}$$

(b)

$$1 + x = \sum_{n=0}^{1} {1 \choose n} x^{n}.$$
 (2.24)

Let us assume that

$$(1+x)^N = \sum_{n=0}^N \binom{N}{n} x^n.$$
 (2.25)

Then,

$$(1+x)^{N+1} = \sum_{n=0}^{N} {N \choose n} x^n + \sum_{n=0}^{N} {N \choose n} x^{n+1}.$$
 (2.26)

By (a), the right hand side can be written as

$$\sum_{n=0}^{N} {N \choose n} x^n + \sum_{n=1}^{N+1} {N \choose n-1} x^n = 1 + x^{N+1} + \sum_{n=1}^{N} {N+1 \choose n} x^n.$$
 (2.27)

Then,

$$(1+x)^{N+1} = \sum_{n=0}^{N+1} {N+1 \choose n} x^n.$$
 (2.28)

Therefore, the assumption is proved by induction on N.

(c)

Let n be a variable such that

$$p(n) = \binom{N}{n} \mu^n (1 - \mu)^{N-n}.$$
 (2.29)

Then,

$$\sum_{n=0}^{N} p(n) = \sum_{n=0}^{N} {N \choose n} \mu^{n} (1-\mu)^{N-n}.$$
 (2.30)

By (b), the right hand side can be written as

$$(1-\mu)^N \sum_{n=0}^N \binom{N}{n} \left(\frac{\mu}{1-\mu}\right)^n = (1-\mu)^N \left(1 + \frac{\mu}{1-\mu}\right)^N.$$
 (2.31)

$$\sum_{n=0}^{N} p(n) = 1. (2.32)$$

Let n be a variable such that

$$p(n) = \binom{N}{n} \mu^n (1 - \mu)^{N-n}.$$
 (2.33)

(a)

We have

$$E n = \sum_{n=0}^{N} n \binom{N}{n} \mu^n (1 - \mu)^{N-n}.$$
 (2.34)

By 2.3(c),

$$\sum_{n=0}^{N} \binom{N}{n} \mu^n (1-\mu)^{N-n} = 1.$$
 (2.35)

Taking the derivative with respect to  $\mu$  gives

$$\sum_{n=0}^{N} n \binom{N}{n} \mu^{n-1} (1-\mu)^{N-n} - \sum_{n=0}^{N} (N-n) \binom{N}{n} \mu^{n} (1-\mu)^{N-n-1} = 0. \quad (2.36)$$

The first term of the left hand side can be written as

$$\frac{1}{\mu} \sum_{n=0}^{N} np(n) = \frac{1}{\mu} \operatorname{E} n.$$
 (2.37)

Since

$$(N-n)\binom{N}{n} = N\binom{N-1}{n},\tag{2.38}$$

the second term can be written as

$$-N\sum_{n=0}^{N-1} \binom{N-1}{n} \mu^n (1-\mu)^{N-n-1} = -N.$$
 (2.39)

$$E n = N\mu. (2.40)$$

(b)

By 2.3(c),

$$\sum_{n=0}^{N} \binom{N}{n} \mu^n (1-\mu)^{N-n} = 1.$$
 (2.41)

Taking the second derivative with respect to  $\mu$  gives

$$\sum_{n=0}^{N} n(n-1) \binom{N}{n} \mu^{n-2} (1-\mu)^{N-n}$$

$$-2 \sum_{n=0}^{N} n(N-n) \binom{N}{n} \mu^{n-1} (1-\mu)^{N-n-1}$$

$$+ \sum_{n=0}^{N} (N-n)(N-n-1) \binom{N}{n} \mu^{n} (1-\mu)^{N-n-2} = 0.$$
(2.42)

The first term of the left hand side can be written as

$$\frac{1}{\mu^2} \sum_{n=0}^{N} n(n-1)p(n) = \frac{1}{\mu^2} \operatorname{E} n(n-1).$$
 (2.43)

Since

$$n(N-n)\binom{N}{n} = N(N-1)\binom{N-2}{n-1},$$
  

$$(N-n)(N-n-1)\binom{N}{n} = N(N-1)\binom{N-2}{n},$$
(2.44)

the second and third terms can be written as

$$-2N(N-1)\sum_{n=1}^{N-1} {N-2 \choose n-1} \mu^{n-1} (1-\mu)^{N-n-1} = -2N(N-1),$$

$$N(N-1)\sum_{n=0}^{N} {N-2 \choose n} \mu^{n} (1-\mu)^{N-n-2} = N(N-1).$$
(2.45)

Then,

$$E n(n-1) = N(N-1)\mu^{2}.$$
 (2.46)

We have

$$var n = E n(n-1) + E n - (E n)^{2}.$$
 (2.47)

$$var n = N\mu(1-\mu). \tag{2.48}$$

We have

$$\Gamma(a)\Gamma(b) = \int_0^\infty x^{a-1} \exp(-x) dx \int_0^\infty y^{b-1} \exp(-y) dy. \tag{2.49}$$

By the transformation

$$t = x + y, (2.50)$$

the right hand side can be written as

$$\int_{0}^{\infty} x^{a-1} \left( \int_{x}^{\infty} (t-x)^{b-1} \exp(-t) dt \right) dx$$

$$= \int_{0}^{\infty} \left( \int_{0}^{t} x^{a-1} (t-x)^{b-1} dx \right) \exp(-t) dt.$$
(2.51)

By the transformation

$$x = t\mu, \tag{2.52}$$

the right hand side can be written as

$$\int_{0}^{\infty} \left( \int_{0}^{1} (t\mu)^{a-1} t^{b-1} (1-\mu)^{b-1} t d\mu \right) \exp(-t) dt$$

$$= \int_{0}^{1} \mu^{a-1} (1-\mu)^{b-1} d\mu \int_{0}^{\infty} t^{a+b-1} \exp(-t) dt.$$
(2.53)

Then,

$$\Gamma(a)\Gamma(b) = \Gamma(a+b) \int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu.$$
 (2.54)

Therefore,

$$\int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}.$$
 (2.55)

#### 2.6

Let  $\mu$  be a variable such that

$$p(\mu) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}.$$
 (2.56)

(a)

We have

By 2.5,

$$\int_{0}^{1} \mu^{a} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)},$$

$$\int_{0}^{1} \mu^{a+1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a+2)\Gamma(b)}{\Gamma(a+b+2)}.$$
(2.58)

Therefore,

$$E \mu = \frac{a}{a+b},$$

$$E \mu^2 = \frac{a(a+1)}{(a+b)(a+b+1)}.$$
(2.59)

(b)

We have

$$\operatorname{var} \mu = \operatorname{E} \mu^2 - (\operatorname{E} \mu)^2.$$
 (2.60)

By (a), the right hand side can be written as

$$\frac{a(a+1)}{(a+b)(a+b+1)} - \left(\frac{a}{a+b}\right)^2 = \frac{a(a+1)(a+b) - a^2(a+b+1)}{(a+b)^2(a+b+1)}.$$
 (2.61)

Therefore,

$$var \mu = \frac{ab}{(a+b)^2(a+b+1)}.$$
 (2.62)

(c)

Setting the derivative of p with respect to  $\mu$  to zero gives

$$0 = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \left( \frac{a-1}{\mu} - \frac{b-1}{1-\mu} \right). \tag{2.63}$$

$$\operatorname{mode} \mu = \frac{a - 1}{a + b - 2}.$$
 (2.64)

Let m and l be variables such that

$$p(m, l | \mu) = {m+l \choose m} \mu^m (1-\mu)^l,$$

$$p(\mu) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}.$$
(2.65)

By 2.6,

$$E\mu = \frac{a}{a+b}. (2.66)$$

Setting the derivative of  $p(m, l|\mu)$  with respect to  $\mu$  to zero gives

$$0 = {m+l \choose m} \mu^m (1-\mu)^l \left(\frac{m}{\mu} + \frac{l}{1-\mu}\right).$$
 (2.67)

Then, the maximum likelihood solution for  $\mu$  is given by

$$\mu_{\rm ML} = \frac{m}{m+l}.\tag{2.68}$$

By the Bayes' thereorem,

$$p(\mu|m, l)p(m, l) = p(m, l|\mu)p(\mu). \tag{2.69}$$

Then, by 2.5,

$$p(\mu|m,l) = \frac{\Gamma(m+l+a+b)}{\Gamma(m+a)\Gamma(l+b)} \mu^{m+a-1} (1-\mu)^{l+b-1}.$$
 (2.70)

The, by 2.6,

$$E(\mu|m,l) = \frac{m+a}{m+l+a+b}.$$
 (2.71)

Therefore,

$$E(\mu|m,l) = \lambda \mu_{\rm ML} + (1-\lambda) E \mu, \qquad (2.72)$$

where

$$\lambda = \frac{m+l}{m+l+a+b}. (2.73)$$

### 2.8

Let x and y be variables.

(a)

By the definition,

$$\mathbf{E} x = \int x p(x) dx. \tag{2.74}$$

The right hand side can be written as

$$\int x \left( \int p(x,y) dy \right) dx = \int \left( \int x p(x|y) dx \right) p(y) dy. \tag{2.75}$$

Therefore,

$$\mathbf{E} x = \mathbf{E}_y \left( \mathbf{E}_x(x|y) \right). \tag{2.76}$$

(b)

By the definition,

$$\operatorname{var} x = \operatorname{E} (x - \operatorname{E} x)^{2}. \tag{2.77}$$

By (a), the right hand side can be written as

$$E_y\left(E_x\left((x-E_x)^2|y\right)\right) = E_y\left(E_x\left((x-E_x(x|y) + E_x(x|y) - E_x)^2|y\right)\right). \tag{2.78}$$

The right hand side can be written as

$$E_{y} (E_{x} ((x - E_{x}(x|y))^{2}|y)) + 2 E_{y} (((E_{x}(x|y) - E_{x}) E_{x} (x - E_{x}(x|y))|y)) + E_{y} ((E_{x}(x|y) - E_{x})^{2}|y).$$
(2.79)

Let us look at each term of the right hand side. By the definition, the first term can be written as  $E_y(\operatorname{var}_x(x|y))$ . The second term can be written as

$$2 E_y ((E_x(x|y) - E_x) (E_x(x|y) - E_x(x|y))) = 0.$$
 (2.80)

By (a), the third term can be written as

$$E_y (E_x(x|y) - E_y (E_x(x|y)))^2 = var_y (E_x(x|y)).$$
 (2.81)

$$\operatorname{var} x = \operatorname{E}_{y} \left( \operatorname{var}_{x}(x|y) \right) + \operatorname{var}_{y} \left( \operatorname{E}_{x}(x|y) \right). \tag{2.82}$$

### 2.9 (Incomplete)

For a vector  $\boldsymbol{\mu}$  in 2 dimensions, by 2.5,

$$\int_{\substack{\mu_1 + \mu_2 = 1 \\ \mu_1 > 0, \mu_2 > 0}} \mu_1^{\alpha_1 - 1} \mu_2^{\alpha_2 - 1} d\boldsymbol{\mu} = \frac{\Gamma(\alpha_1) \Gamma(\alpha_2)}{\Gamma(\alpha_1 + \alpha_2)}.$$

For a vector  $\boldsymbol{\mu}$  in M dimensions, let us assume that

$$\int_{\substack{\sum_{m=1}^{M} \mu_m > 0}} \prod_{m=1}^{M} \mu_m^{\alpha_m - 1} d\boldsymbol{\mu} = \frac{\prod_{m=1}^{M} \Gamma(\alpha_m)}{\Gamma(\sum_{m=1}^{M} \alpha_m)}.$$

Under the constraint

$$\sum_{m=1}^{M+1} \mu_m = 1, \tag{2.83}$$

we have

$$\int_{0}^{c} \prod_{m=1}^{M+1} \mu_{m}^{\alpha_{m}-1} d\mu_{M+1} = \left(\prod_{m=1}^{M-1} \mu_{m}^{\alpha_{m}-1}\right) \int_{0}^{c} \mu_{M+1}^{\alpha_{M+1}-1} \left(c - \mu_{M+1}\right)^{\alpha_{M}-1} d\mu_{M+1},$$
(2.84)

where

$$c = 1 - \sum_{m=1}^{M-1} \mu_m. (2.85)$$

By the transformation

$$\mu'_{M+1} = \frac{\mu_{M+1}}{c},\tag{2.86}$$

the integral of the right hand side can be written as

$$\int_{0}^{1} (c\mu'_{M+1})^{\alpha_{M+1}-1} \left( c(1-\mu'_{M+1}) \right)^{\alpha_{M}-1} cd\mu'_{M+1} 
= c^{\alpha_{M}+\alpha_{M+1}-1} \int_{0}^{1} {\mu'_{M+1}}^{\alpha_{M+1}-1} (1-\mu'_{M+1})^{\alpha_{M}-1} d\mu'_{M+1}.$$
(2.87)

By 2.5, the integral of the right hand side can be written as

$$\frac{\Gamma(\alpha_M)\Gamma(\alpha_{M+1})}{\Gamma(\alpha_M + \alpha_{M+1})}. (2.88)$$

Then,

$$\int_{0}^{c} \prod_{m=1}^{M+1} \mu_{m}^{\alpha_{m}-1} d\mu_{M+1} = \left( \prod_{m=1}^{M-1} \mu_{m}^{\alpha_{m}-1} \right) c^{\alpha_{M} + \alpha_{M+1} - 1} \frac{\Gamma(\alpha_{M}) \Gamma(\alpha_{M+1})}{\Gamma(\alpha_{M} + \alpha_{M+1})}. \quad (2.89)$$

For a vector  $\boldsymbol{\mu}$  in M dimensions, by the assumption,

$$\int_{\sum_{m=1}^{M} \mu_m = 1} \left( \prod_{m=1}^{M-1} \mu_m^{\alpha_m - 1} \right) \mu_M^{\alpha_M + \alpha_{M+1} - 1} d\boldsymbol{\mu} = \frac{\left( \prod_{m=1}^{M-1} \Gamma(\alpha_m) \right) \Gamma(\alpha_M + \alpha_{M+1})}{\Gamma(\sum_{m=1}^{M+1} \alpha_m)}.$$

Then, for a vector  $\boldsymbol{\mu}$  in M+1 dimensions,

$$\int_{\sum_{m=1}^{M+1} \mu_m = 1} \prod_{m=1}^{M+1} \mu_m^{\alpha_m - 1} d\boldsymbol{\mu} = \frac{\Gamma(\alpha_M) \Gamma(\alpha_{M+1})}{\Gamma(\alpha_M + \alpha_{M+1})} \frac{\left(\prod_{m=1}^{M-1} \Gamma(\alpha_k)\right) \Gamma(\alpha_M + \alpha_{M+1})}{\Gamma(\sum_{m=1}^{M+1} \alpha_m)}?$$

The right hand side can be written as

$$\frac{\prod_{m=1}^{M+1} \Gamma(\alpha_m)}{\Gamma(\sum_{m=1}^{M+1} \alpha_m)}.$$
(2.90)

Therefore, the assumption is proved by induction on M.

#### 2.10

Let  $\mu$  be a vector such that

$$p(\boldsymbol{\mu}) = \frac{\Gamma(\sum_{k=1}^{K} \alpha_k)}{\prod_{k=1}^{K} \Gamma(\alpha_k)} \prod_{k=1}^{K} \mu_k^{\alpha_k - 1}.$$
 (2.91)

Then, by the definition,

$$E \mu_k = \int \mu_k p(\boldsymbol{\mu}) d\boldsymbol{\mu},$$

$$E \mu_k^2 = \int \mu_k^2 p(\boldsymbol{\mu}) d\boldsymbol{\mu},$$

$$E \mu_k \mu_{k'} = \int \mu_k \mu_{k'} p(\boldsymbol{\mu}) d\boldsymbol{\mu}.$$
(2.92)

Let  $k \neq k'$ . Then, by 2.9, the right hand sides can be written as

$$\frac{\Gamma\left(\sum_{k=1}^{K} \alpha_{k}\right)}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \frac{\frac{\Gamma(\alpha_{k}+1)}{\Gamma(\alpha_{k})} \prod_{k=1}^{K} \Gamma(\alpha_{k})}{\Gamma\left(\sum_{k=1}^{K} \alpha_{k}+1\right)} = \frac{\alpha_{k}}{\sum_{k=1}^{K} \alpha_{k}},$$

$$\frac{\Gamma\left(\sum_{k=1}^{K} \alpha_{k}\right)}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \frac{\frac{\Gamma(\alpha_{k}+2)}{\Gamma(\alpha_{k})} \prod_{k=1}^{K} \Gamma(\alpha_{k})}{\Gamma\left(\sum_{k=1}^{K} \alpha_{k}+2\right)} = \frac{\alpha_{k}(\alpha_{k}+1)}{\sum_{k=1}^{K} \alpha_{k}(\sum_{k=1}^{K} \alpha_{k}+1)}, \quad (2.93)$$

$$\frac{\Gamma\left(\sum_{k=1}^{K} \alpha_{k}\right)}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \frac{\frac{\Gamma(\alpha_{k}+1)\Gamma(\alpha_{k'}+1)}{\Gamma(\alpha_{k})\Gamma(\alpha_{k'})} \prod_{k=1}^{K} \Gamma(\alpha_{k})}{\Gamma\left(\sum_{k=1}^{K} \alpha_{k}+2\right)} = \frac{\alpha_{k}\alpha_{k'}}{\sum_{k=1}^{K} \alpha_{k}(\sum_{k=1}^{K} \alpha_{k}+1)}.$$

Then,

$$E \mu_k = \frac{\alpha_k}{\sum_{k=1}^K \alpha_k}.$$

$$E \mu_k^2 = \frac{\alpha_k(\alpha_k + 1)}{\sum_{k=1}^K \alpha_k \left(\sum_{k=1}^K \alpha_k + 1\right)},$$

$$E \mu_k \mu_{k'} = \frac{\alpha_k \alpha_{k'}}{\sum_{k=1}^K \alpha_k \left(\sum_{k=1}^K \alpha_k + 1\right)}.$$
(2.94)

Since

$$\operatorname{var} \mu_{k} = \operatorname{E} \mu_{k}^{2} - (\operatorname{E} \mu_{k})^{2}, \operatorname{cov} (\mu_{k}, \mu_{k'}) = \operatorname{E} \mu_{k} \mu_{k'} - \operatorname{E} \mu_{k} \operatorname{E} \mu_{k'},$$
(2.95)

we have

$$\operatorname{var} \mu_{k} = \frac{\alpha_{k} \left( \left( \sum_{k=1}^{K} \alpha_{k} \right) - \alpha_{k} \right)}{\left( \sum_{k=1}^{K} \alpha_{k} \right)^{2} \left( \sum_{k=1}^{K} \alpha_{k} + 1 \right)},$$

$$\operatorname{cov}(\mu_{k}, \mu_{k'}) = -\frac{\alpha_{k} \alpha_{k'}}{\left( \sum_{k=1}^{K} \alpha_{k} \right)^{2} \left( \sum_{k=1}^{K} \alpha_{k} + 1 \right)}.$$
(2.96)

#### 2.11

Let  $\mu$  be a variable such that

$$p(\boldsymbol{\mu}) = \frac{\Gamma\left(\sum_{k=1}^{K} \alpha_k\right)}{\prod_{k=1}^{K} \Gamma(\alpha_k)} \prod_{k=1}^{K} \mu_k^{\alpha_k - 1}.$$
 (2.97)

Then, by the definition,

$$E \ln \mu_k = \int (\ln \mu_k) \, p(\boldsymbol{\mu}) d\boldsymbol{\mu}. \tag{2.98}$$

Since

$$\frac{\partial}{\partial \alpha_k} p(\boldsymbol{\mu}) = \left( \frac{\Gamma'\left(\sum_{k=1}^K \alpha_k\right)}{\Gamma\left(\sum_{k=1}^K \alpha_k\right)} - \frac{\Gamma'(\alpha_k)}{\Gamma(\alpha_k)} + \ln \mu_k \right) p(\boldsymbol{\mu}), \tag{2.99}$$

we have

$$E \ln \mu_k = \frac{\partial}{\partial \alpha_k} \int p(\boldsymbol{\mu}) d\boldsymbol{\mu} + \left( \psi(\alpha_k) - \psi\left(\sum_{k=1}^K \alpha_k\right) \right) \int p(\boldsymbol{\mu}) d\boldsymbol{\mu}, \quad (2.100)$$

where

$$\psi(a) = \frac{d}{da} \ln \Gamma(a). \tag{2.101}$$

Therefore,

$$\operatorname{E} \ln \mu_k = \psi(\alpha_k) - \psi\left(\sum_{k=1}^K \alpha_k\right). \tag{2.102}$$

#### 2.12

Let x be a variable such that

$$p(x) = \frac{1}{b - a},\tag{2.103}$$

where a < b. Then

$$\int_{a}^{b} p(x)dx = 1. (2.104)$$

Then, by the definition,

$$E x = \frac{1}{b-a} \int_a^b x dx,$$

$$E x^2 = \frac{1}{b-a} \int_a^b x^2 dx.$$
(2.105)

Then,

$$E x = \frac{1}{2}(a+b),$$

$$E x^{2} = \frac{1}{3}(a^{2} + ab + b^{2}).$$
(2.106)

Since

$$var x = E x^2 - (E x)^2, (2.107)$$

we have

$$var x = \frac{1}{12}(b-a)^2. (2.108)$$

#### 2.13

Let  $\mathbf{x}$  be a variable in D dimensions and let

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

$$q(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{L}).$$
(2.109)

Then, by the definition, the Kulleback-Leibler divergence is given by

$$KL(p||q) = -\int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \ln \frac{\mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{L})}{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} d\mathbf{x}.$$
 (2.110)

Note that

$$\ln \frac{\mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{L})}{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})} = \ln \frac{(2\pi)^{-\frac{D}{2}} |\det \mathbf{L}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^{\mathsf{T}} \mathbf{L}^{-1}(\mathbf{x} - \mathbf{m})\right)}{(2\pi)^{-\frac{D}{2}} |\det \boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)}.$$
(2.111)

The right hand side can be written as

$$\frac{1}{2} \ln \left| \frac{\det \mathbf{\Sigma}}{\det \mathbf{L}} \right| + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2} (\mathbf{x} - \mathbf{m})^{\mathsf{T}} \mathbf{L}^{-1} (\mathbf{x} - \mathbf{m}). \quad (2.112)$$

Then, the integral can be written as

$$\frac{1}{2} \ln \left| \frac{\det \mathbf{\Sigma}}{\det \mathbf{L}} \right| \int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} 
+ \frac{1}{2} \int (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} 
- \frac{1}{2} \int (\mathbf{x} - \mathbf{m})^{\mathsf{T}} \mathbf{L}^{-1} (\mathbf{x} - \mathbf{m}) \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}.$$
(2.113)

Let us look at the integral of each term. The integral of the first term is 1. Since

$$\int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = \boldsymbol{\Sigma}, \tag{2.114}$$

we have

$$\int (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} (\mathbf{x} - \boldsymbol{\mu}) \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = \operatorname{tr} \boldsymbol{\Sigma}.$$
 (2.115)

Then, the integral of the second term can be written as

$$\operatorname{tr}\left(\mathbf{\Sigma}^{-1}\mathbf{\Sigma}\right) = D. \tag{2.116}$$

Since

$$(\mathbf{x} - \mathbf{m})^{\mathsf{T}} \mathbf{L}^{-1} (\mathbf{x} - \mathbf{m}) = (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - \mathbf{m})^{\mathsf{T}} \mathbf{L}^{-1} (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - \mathbf{m}), \quad (2.117)$$

the integral of the third term can be written as

$$\int (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{L}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}$$

$$+ 2(\boldsymbol{\mu} - \mathbf{m})^{\mathsf{T}} \mathbf{L}^{-1} \int (\mathbf{x} - \boldsymbol{\mu}) \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}$$

$$+ (\boldsymbol{\mu} - \mathbf{m})^{\mathsf{T}} \mathbf{L}^{-1} (\boldsymbol{\mu} - \mathbf{m}) \int \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}$$

$$= \operatorname{tr} (\mathbf{L}^{-1} \boldsymbol{\Sigma}) + (\boldsymbol{\mu} - \mathbf{m})^{\mathsf{T}} \mathbf{L}^{-1} (\boldsymbol{\mu} - \mathbf{m}).$$
(2.118)

Therefore,

$$KL(p||q) = \frac{1}{2} \left( \ln \left| \frac{\det \mathbf{L}}{\det \mathbf{\Sigma}} \right| - D + \operatorname{tr} \left( \mathbf{L}^{-1} \mathbf{\Sigma} \right) + (\boldsymbol{\mu} - \mathbf{m})^{\mathsf{T}} \mathbf{L}^{-1} (\boldsymbol{\mu} - \mathbf{m}) \right).$$
(2.119)

#### 2.14

Let  $\mathbf{x}$  be a variable in D dimensions. By the definition, the entropy is given by

$$H(\mathbf{x}) = -\int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x}.$$
 (2.120)

In order to maximise H(x) with the constratints

$$\int p(\mathbf{x})d\mathbf{x} = 1,$$

$$\int \mathbf{x}p(\mathbf{x})d\mathbf{x} = \boldsymbol{\mu},$$

$$\int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}}p(\mathbf{x})d\mathbf{x} = \boldsymbol{\Sigma},$$
(2.121)

let

$$L(p) = H(\mathbf{x}) + \lambda \left( \int p(\mathbf{x}) d\mathbf{x} - 1 \right) + \mathbf{l}^{\mathsf{T}} \left( \int \mathbf{x} p(\mathbf{x}) d\mathbf{x} - \boldsymbol{\mu} \right)$$
$$+ \mathbf{m}^{\mathsf{T}} \left( \int (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} p(\mathbf{x}) d\mathbf{x} - \boldsymbol{\Sigma} \right) \mathbf{m}.$$
(2.122)

Setting the variation with respect to p to zero gives

$$0 = -\ln p(\mathbf{x}) - 1 + \lambda + \mathbf{l}^{\mathsf{T}}\mathbf{x} + \mathbf{m}^{\mathsf{T}}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}}\mathbf{m}. \tag{2.123}$$

Then,

$$p(\mathbf{x}) = \exp\left(-1 + \lambda + \mathbf{l}^{\mathsf{T}}\mathbf{x} + \mathbf{m}^{\mathsf{T}}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}}\mathbf{m}\right), \tag{2.124}$$

so that

$$p(\mathbf{x}) = c \exp\left(-\left(\mathbf{x} - \boldsymbol{\mu} - \mathbf{M}\mathbf{I}\right)^{\mathsf{T}} \mathbf{M}^{-1} \left(\mathbf{x} - \boldsymbol{\mu} - \mathbf{M}\mathbf{I}\right)\right), \tag{2.125}$$

where

$$c = \exp(-1 + \lambda - \mathbf{l}^{\mathsf{T}}\mathbf{M}\mathbf{l}),$$
  

$$\mathbf{M} = -(\mathbf{m}\mathbf{m}^{\mathsf{T}})^{-1}.$$
(2.126)

Substituting it to the constraints and the transformation

$$\mathbf{y} = \mathbf{x} - \boldsymbol{\mu} - \mathbf{Ml} \tag{2.127}$$

gives

$$c \int \exp(-\mathbf{y}^{\mathsf{T}} \mathbf{M}^{-1} \mathbf{y}) d\mathbf{y} = 1,$$

$$c \int (\mathbf{y} + \boldsymbol{\mu} + \mathbf{M} \mathbf{l}) \exp(-\mathbf{y}^{\mathsf{T}} \mathbf{M}^{-1} \mathbf{y}) d\mathbf{y} = \boldsymbol{\mu},$$

$$c \int (\mathbf{y} + \mathbf{M} \mathbf{l}) (\mathbf{y} + \mathbf{M} \mathbf{l})^{\mathsf{T}} \exp(-\mathbf{y}^{\mathsf{T}} \mathbf{M}^{-1} \mathbf{y}) d\mathbf{y} = \boldsymbol{\Sigma}.$$
(2.128)

Since

$$\int \exp(-\mathbf{y}^{\mathsf{T}}\mathbf{y}) d\mathbf{y} = \left(\Gamma\left(\frac{1}{2}\right)\right)^{D},$$

$$\int \mathbf{y} \exp(-\mathbf{y}^{\mathsf{T}}\mathbf{y}) d\mathbf{y} = \mathbf{0},$$

$$\int \mathbf{y} \mathbf{y}^{\mathsf{T}} \exp(-\mathbf{y}^{\mathsf{T}}\mathbf{y}) d\mathbf{y} = \Gamma\left(\frac{3}{2}\right) \left(\Gamma\left(\frac{1}{2}\right)\right)^{D-1} \mathbf{I},$$
(2.129)

they can be written as

$$c\left(\Gamma\left(\frac{1}{2}\right)\right)^{D} |\det \mathbf{M}|^{\frac{1}{2}} = 1,$$

$$c(\boldsymbol{\mu} + \mathbf{Ml}) \left(\Gamma\left(\frac{1}{2}\right)\right)^{D} |\det \mathbf{M}|^{\frac{1}{2}} = \boldsymbol{\mu},$$

$$c\left(\Gamma\left(\frac{3}{2}\right) \left(\Gamma\left(\frac{1}{2}\right)\right)^{D-1} \mathbf{M} + \mathbf{Ml}(\mathbf{Ml})^{\mathsf{T}} \left(\Gamma\left(\frac{1}{2}\right)\right)^{D} |\det \mathbf{M}|^{\frac{1}{2}} = \boldsymbol{\Sigma}.$$

$$(2.130)$$

Then,

$$\lambda = 1 - \frac{D}{2} \ln \pi - \frac{1}{2} \ln |\det \mathbf{M}|,$$

$$\mathbf{l} = \mathbf{0},$$

$$\mathbf{M} = 2\Sigma.$$
(2.131)

Therefore,

$$p(\mathbf{x}) = (2\pi)^{-\frac{D}{2}} |\det \mathbf{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \tag{2.132}$$

### 2.15

Let  $\mathbf{x}$  be a variable in D dimensions such that

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}). \tag{2.133}$$

Then, by the definition, the entropy is given by

$$H(\mathbf{x}) = -\int \mathcal{N}(\mathbf{x}|\mu, \mathbf{\Sigma}) \ln \mathcal{N}(\mathbf{x}|\mu, \mathbf{\Sigma}) d\mathbf{x}.$$
 (2.134)

The right hand side can be written as

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

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

$$+ \frac{1}{2} \int (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}.$$
(2.135)

Let us look at each integral of the right hand side. The first integral is 1. Since

$$\int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = \boldsymbol{\Sigma}, \tag{2.136}$$

we have

$$\int (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} (\mathbf{x} - \boldsymbol{\mu}) \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = \operatorname{tr} \boldsymbol{\Sigma}.$$
 (2.137)

Then, the second integral can be written as

$$\operatorname{tr}\left(\mathbf{\Sigma}^{-1}\mathbf{\Sigma}\right) = D. \tag{2.138}$$

Therefore,

$$H(\mathbf{x}) = \frac{D}{2} (1 + \ln(2\pi)) + \frac{1}{2} \ln|\det \Sigma|.$$
 (2.139)

### 2.16

Let x be a variable such that

$$x = x_1 + x_2, (2.140)$$

where

$$p(x_1) = \mathcal{N}\left(x_1|\mu_1, \tau_1^{-1}\right), p(x_2) = \mathcal{N}\left(x_2|\mu_2, \tau_2^{-1}\right).$$
 (2.141)

By marginalisation,

$$p(x) = \int_{-\infty}^{\infty} p(x|x_2)p(x_2)dx_2. \tag{2.142}$$

The right hand side can be written as

$$\int_{-\infty}^{\infty} \mathcal{N}\left(x|\mu_1 + x_2, \tau_1^{-1}\right) \mathcal{N}\left(x_2|\mu_2, \tau_2^{-1}\right) dx_2$$

$$= \int_{-\infty}^{\infty} \left(\frac{\tau_1}{2\pi}\right)^{\frac{1}{2}} \exp\left(-\frac{\tau_1}{2}(x - \mu_1 - x_2)^2\right) \left(\frac{\tau_2}{2\pi}\right)^{\frac{1}{2}} \exp\left(-\frac{\tau_2}{2}(x_2 - \mu_2)^2\right) dx_2.$$
(2.143)

The logarithm of the integrand except the terms independent of x and z is given by

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

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

$$= -\frac{\tau_1 + \tau_2}{2} \left( x_2 - \frac{\tau_1(x - \mu_1) + \tau_2 \mu_2}{\tau_1 + \tau_2} \right)^2 - \frac{\tau_1 \tau_2}{2(\tau_1 + \tau_2)} (x - \mu_1 - \mu_2)^2.$$
(2.144)

Then,

$$p(x) = \mathcal{N}\left(x|\mu_1 + \mu_2, \tau_1^{-1} + \tau_2^{-1}\right).$$
 (2.145)

Therefore, by 1.35,

$$H(x) = \frac{1}{2} \left( 1 + \ln(2\pi) + \ln\left(\tau_1^{-1} + \tau_2^{-1}\right) \right). \tag{2.146}$$

#### 2.17

Let  $\Sigma$  be a matrix and

$$\mathbf{S} = \frac{1}{2} \left( \mathbf{\Sigma}^{-1} + \left( \mathbf{\Sigma}^{-1} \right)^{\mathsf{T}} \right),$$

$$\mathbf{A} = \frac{1}{2} \left( \mathbf{\Sigma}^{-1} - \left( \mathbf{\Sigma}^{-1} \right)^{\mathsf{T}} \right).$$
(2.147)

Then,

$$\Sigma^{-1} = \mathbf{S} + \mathbf{A},\tag{2.148}$$

so that

$$(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{S} (\mathbf{x} - \boldsymbol{\mu}) + (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{A} (\mathbf{x} - \boldsymbol{\mu}). \quad (2.149)$$

The second term of the right hand side can be written as

$$\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \left(\boldsymbol{\Sigma}^{-1}\right)^{\mathsf{T}} (\mathbf{x} - \boldsymbol{\mu}). \tag{2.150}$$

The second term of the right hand side can be written as

$$-\frac{1}{2} \left( \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)^{\mathsf{T}} (\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}). \tag{2.151}$$

Then,

$$(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{A} (\mathbf{x} - \boldsymbol{\mu}) = 0. \tag{2.152}$$

Therefore,

$$(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{S} (\mathbf{x} - \boldsymbol{\mu}). \tag{2.153}$$

### 2.18

(a)

Let  $\Sigma$  be a  $D \times D$  real symmetric matrix such that

$$\Sigma \mathbf{u}_d = \lambda_d \mathbf{u}_d, \tag{2.154}$$

where  $\mathbf{u}_1, \dots, \mathbf{u}_D$  are unit vectors. Then,

$$\overline{\mathbf{u}_d}^{\mathsf{T}} \mathbf{\Sigma} \mathbf{u}_d = \lambda_d, \tag{2.155}$$

where  $\overline{\mathbf{u}_d}$  is the conjugate of  $\mathbf{u}_d$ . Since  $\Sigma$  is real and symmetric, the left hand side can be written as

$$\overline{\mathbf{u}_d}^{\mathsf{T}} \overline{\mathbf{\Sigma}}^{\mathsf{T}} \mathbf{u}_d = \left( \overline{\mathbf{\Sigma}} \overline{\mathbf{u}_d} \right)^{\mathsf{T}} \mathbf{u}_d. \tag{2.156}$$

The right hand side can be written as

$$\overline{\lambda_d} \overline{\mathbf{u}_d}^{\mathsf{T}} \mathbf{u}_d = \overline{\lambda_d}. \tag{2.157}$$

$$\lambda_d = \overline{\lambda_d}.\tag{2.158}$$

(b)

For  $d \neq d'$ , taking the inner product with  $\mathbf{u}'_d$  on both sides of

$$\Sigma \mathbf{u}_d = \lambda_d \mathbf{u}_d \tag{2.159}$$

gives

$$\mathbf{u}_{d'}^{\mathsf{T}} \mathbf{\Sigma} \mathbf{u}_d = \lambda_d \mathbf{u}_{d'}^{\mathsf{T}} \mathbf{u}_d. \tag{2.160}$$

Since  $\Sigma$  is symmetric, the left hand side can be written as

$$\mathbf{u}_{d'}^{\mathsf{T}} \mathbf{\Sigma}^{\mathsf{T}} \mathbf{u}_d = (\mathbf{\Sigma} \mathbf{u}_{d'})^{\mathsf{T}} \mathbf{u}_d. \tag{2.161}$$

The right hand side can be written as  $\lambda_{d'} \mathbf{u}_{d'}^{\mathsf{T}} \mathbf{u}_{d}$ . Then,

$$\lambda_d \mathbf{u}_{d'}^{\mathsf{T}} \mathbf{u}_d = \lambda_{d'} \mathbf{u}_{d'}^{\mathsf{T}} \mathbf{u}_d. \tag{2.162}$$

Therefore, if  $\lambda_d \neq \lambda_{d'}$ , then

$$\mathbf{u}_{d'}^{\mathsf{T}}\mathbf{u}_d = 0. \tag{2.163}$$

# 2.19

Let  $\Sigma$  be a  $D \times D$  real symmetric matrix such that

$$\Sigma \mathbf{u}_d = \lambda_d \mathbf{u}_d, \tag{2.164}$$

where  $\mathbf{u}_1, \dots, \mathbf{u}_D$  are unit vectors. Let

$$\Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_D), 
\mathbf{U} = [\mathbf{u}_1 \dots \mathbf{u}_D].$$
(2.165)

Then

$$\Sigma \mathbf{U} = \mathbf{U} \mathbf{\Lambda}. \tag{2.166}$$

By 2.18,

$$\mathbf{U}^{\mathsf{T}}\mathbf{U} = \mathbf{I}.\tag{2.167}$$

Then,

$$\Sigma = \mathbf{U}\Lambda\mathbf{U}^{\mathsf{T}}, \Sigma^{-1} = \mathbf{U}\Lambda^{-1}\mathbf{U}^{\mathsf{T}},$$
(2.168)

$$\Sigma = \sum_{d=1}^{D} \lambda_d \mathbf{u}_d \mathbf{u}_d^{\mathsf{T}},$$

$$\Sigma^{-1} = \sum_{d=1}^{D} \frac{1}{\lambda_d} \mathbf{u}_d \mathbf{u}_d^{\mathsf{T}}.$$
(2.169)

Let  $\Sigma$  be a  $D \times D$  real symmetric matrix such that

$$\Sigma \mathbf{u}_d = \lambda_d \mathbf{u}_d, \tag{2.170}$$

where  $u_1, \dots, u_D$  are unit vectors. Let

$$\Lambda = \operatorname{diag}(\lambda_1, \cdots, \lambda_D), 
\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_D].$$
(2.171)

By 2.19,

$$\mathbf{a}^{\mathsf{T}} \mathbf{\Sigma} \mathbf{a} = \mathbf{b}^{\mathsf{T}} \mathbf{\Lambda} \mathbf{b},\tag{2.172}$$

where

$$\mathbf{b} = \mathbf{U}^{\mathsf{T}} \mathbf{a}.\tag{2.173}$$

The right hand side can be written as  $\sum_{d=1}^{D} \lambda_d b_d^2$ . Therefore, the necessary and sufficient condition for

$$\mathbf{a}^{\mathsf{T}} \mathbf{\Sigma} \mathbf{a} > 0 \tag{2.174}$$

for any real vector  $\mathbf{a}$  is

$$\lambda_d > 0. \tag{2.175}$$

### 2.21

Let  $\Sigma$  be a  $D \times D$  real symmetric matrix. Then the number of independent parameters is  $\frac{D(D+1)}{2}$ .

### 2.22

Let  $\Sigma$  be a  $D \times D$  symmetric matrix and

$$\Sigma \Lambda = I. \tag{2.176}$$

Taking the transpose of the both sides gives

$$\mathbf{\Lambda}^{\mathsf{T}} \mathbf{\Sigma} = \mathbf{I}.\tag{2.177}$$

$$\mathbf{\Lambda}^{\mathsf{T}} = \mathbf{\Lambda}.\tag{2.178}$$

Let  $\Sigma$  be a  $D \times D$  real symmetric matrix such that

$$\Sigma \mathbf{u}_d = \lambda_d \mathbf{u}_d, \tag{2.179}$$

where  $u_1, \dots, u_D$  are unit vectors. Let

$$\mathbf{\Lambda}' = \operatorname{diag}\left(\lambda_1^{-\frac{1}{2}}, \cdots, \lambda_D^{-\frac{1}{2}}\right),$$

$$\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_D].$$
(2.180)

By 2.19,

$$\int_{(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu}) = \Delta} d\mathbf{x} = \int_{(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}} \mathbf{U} \boldsymbol{\Lambda}' \boldsymbol{\Lambda}'^{\mathsf{T}} \mathbf{U}^{\mathsf{T}} (\mathbf{x}-\boldsymbol{\mu}) = \Delta} d\mathbf{x}.$$
 (2.181)

By the transformation

$$\mathbf{y} = \mathbf{\Lambda}'^{\mathsf{T}} \mathbf{U}^{\mathsf{T}} (\mathbf{x} - \boldsymbol{\mu}) \tag{2.182}$$

and the property

$$\mathbf{U}^{\mathsf{T}}\mathbf{U} = \mathbf{I},\tag{2.183}$$

the right hand side can be written as

$$\int_{\|\mathbf{v}\|^2 = \Delta} \left| \det \left( \mathbf{U} \mathbf{\Lambda}'^{-1} \right) \right| d\mathbf{y} = \left| \det \mathbf{\Sigma} \right|^{\frac{1}{2}} \int_{\|\mathbf{v}\|^2 = \Delta} d\mathbf{y}. \tag{2.184}$$

Therefore,

$$\int_{(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu}) = \Delta} d\mathbf{x} = |\det \mathbf{\Sigma}|^{\frac{1}{2}} \Delta^D V_D, \qquad (2.185)$$

where

$$V_D = \int_{\|\mathbf{x}\| = 1} d\mathbf{x}.$$
 (2.186)

# 2.24

Let **A** be a square matrix and **D** be an invertible matrix. We have

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{B}\mathbf{D}^{-1} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{O} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}. \tag{2.187}$$

The right hand side can be written as

$$\begin{bmatrix} \mathbf{I} & \mathbf{B}\mathbf{D}^{-1} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{O} \\ \mathbf{O} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{D}^{-1}\mathbf{C} & \mathbf{I} \end{bmatrix}. \tag{2.188}$$

Then,

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ -\mathbf{D}^{-1}\mathbf{C} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{M} & \mathbf{O} \\ \mathbf{O} & \mathbf{D}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\mathbf{B}\mathbf{D}^{-1} \\ \mathbf{O} & \mathbf{I} \end{bmatrix}, \tag{2.189}$$

where

$$\mathbf{M} = \left(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}\right)^{-1}.\tag{2.190}$$

Therefore,

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{M} & -\mathbf{B}\mathbf{D}^{-1}\mathbf{M} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{bmatrix}.$$
 (2.191)

#### 2.25

Let  $\mathbf{x}$  be a variable such that

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}), \tag{2.192}$$

where

$$\mathbf{x} = egin{bmatrix} \mathbf{x}_a \ \mathbf{x}_b \ \mathbf{x}_c \end{bmatrix}, oldsymbol{\mu} = egin{bmatrix} oldsymbol{\mu}_a \ oldsymbol{\mu}_c \end{bmatrix}, oldsymbol{\Sigma} = egin{bmatrix} oldsymbol{\Sigma}_{aa} & oldsymbol{\Sigma}_{ab} & oldsymbol{\Sigma}_{ac} \ oldsymbol{\Sigma}_{ca} & oldsymbol{\Sigma}_{cb} & oldsymbol{\Sigma}_{cc} \end{bmatrix}.$$

Let

$$\Lambda = \Sigma^{-1},\tag{2.193}$$

where

$$oldsymbol{\Lambda} = egin{bmatrix} oldsymbol{\Lambda}_{aa} & oldsymbol{\Lambda}_{ab} & oldsymbol{\Lambda}_{ac} \ oldsymbol{\Lambda}_{ba} & oldsymbol{\Lambda}_{bb} & oldsymbol{\Lambda}_{bc} \ oldsymbol{\Lambda}_{ca} & oldsymbol{\Lambda}_{cb} & oldsymbol{\Lambda}_{cc} \end{bmatrix}.$$

Then, the logarithm of  $p(\mathbf{x})$  except the terms independent of  $\mathbf{x}_a$  can be written as

$$-\frac{1}{2}(\mathbf{x}_{a}-\boldsymbol{\mu}_{a})^{\mathsf{T}}\boldsymbol{\Lambda}_{aa}(\mathbf{x}_{a}-\boldsymbol{\mu}_{a})-\frac{1}{2}(\mathbf{x}_{a}-\boldsymbol{\mu}_{a})^{\mathsf{T}}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_{b}-\boldsymbol{\mu}_{b})$$

$$-\frac{1}{2}(\mathbf{x}_{a}-\boldsymbol{\mu}_{a})^{\mathsf{T}}\boldsymbol{\Lambda}_{ac}(\mathbf{x}_{c}-\boldsymbol{\mu}_{c})-\frac{1}{2}(\mathbf{x}_{b}-\boldsymbol{\mu}_{b})^{\mathsf{T}}\boldsymbol{\Lambda}_{ba}(\mathbf{x}_{a}-\boldsymbol{\mu}_{a})$$

$$-\frac{1}{2}(\mathbf{x}_{c}-\boldsymbol{\mu}_{c})^{\mathsf{T}}\boldsymbol{\Lambda}_{ca}(\mathbf{x}_{a}-\boldsymbol{\mu}_{a}).$$
(2.194)

Except the terms independent of  $\mathbf{x}_a$ , it can be written as

$$-\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_{a|b,c})^{\mathsf{T}} \boldsymbol{\Sigma}_{a|b,c}^{-1}(\mathbf{x}_a - \boldsymbol{\mu}_{a|b,c}), \qquad (2.195)$$

where

$$\boldsymbol{\mu}_{a|b,c} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ac} (\mathbf{x}_c - \boldsymbol{\mu}_c),$$

$$\boldsymbol{\Sigma}_{a|b,c} = \boldsymbol{\Lambda}_{aa}^{-1}.$$
(2.196)

Then,

$$p(\mathbf{x}_a|\mathbf{x}_b,\mathbf{x}_c) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b,c},\boldsymbol{\Sigma}_{a|b,c}). \tag{2.197}$$

By marginalisation,

$$p(\mathbf{x}_a|\mathbf{x}_b) = \int p(\mathbf{x}_a|\mathbf{x}_b, \mathbf{x}_c) p(\mathbf{x}_c) d\mathbf{x}_c.$$
 (2.198)

The integrand of the right hand side except the terms independent of  $\mathbf{x}_c$  can be written as

$$-\frac{1}{2} \left( \mathbf{x}_{a} - \boldsymbol{\mu}_{a|b,c} \right)^{\mathsf{T}} \boldsymbol{\Sigma}_{a|b,c}^{-1} \left( \mathbf{x}_{a} - \boldsymbol{\mu}_{a|b,c} \right) - \frac{1}{2} (\mathbf{x}_{c} - \boldsymbol{\mu}_{c})^{\mathsf{T}} \boldsymbol{\Lambda}_{cc} (\mathbf{x}_{c} - \boldsymbol{\mu}_{c})$$

$$= -\frac{1}{2} \mathbf{v}^{\mathsf{T}} \mathbf{M} \mathbf{v},$$
(2.199)

where

$$\mathbf{v} = \begin{bmatrix} \mathbf{x}_{c} - \boldsymbol{\mu}_{c} \\ \mathbf{x}_{a} - \boldsymbol{\mu}_{a} + \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_{b} - \boldsymbol{\mu}_{b}) \end{bmatrix},$$

$$\mathbf{M} = \begin{bmatrix} \boldsymbol{\Lambda}_{cc} + \boldsymbol{\Lambda}_{ac}^{\mathsf{T}} \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ac} & \boldsymbol{\Lambda}_{ac}^{\mathsf{T}} \\ \boldsymbol{\Lambda}_{ac} & \boldsymbol{\Lambda}_{aa} \end{bmatrix}.$$
(2.200)

By 2.24,

$$\mathbf{M}^{-1} = \begin{bmatrix} \mathbf{\Lambda}_{cc}^{-1} & -\mathbf{\Lambda}_{cc}^{-1} \mathbf{\Lambda}_{aa}^{\mathsf{T}} \mathbf{\Lambda}_{aa}^{-1} \\ -\mathbf{\Lambda}_{aa}^{-1} \mathbf{\Lambda}_{ac} \mathbf{\Lambda}_{cc}^{-1} & \mathbf{\Lambda}_{aa}^{-1} + \mathbf{\Lambda}_{aa}^{-1} \mathbf{\Lambda}_{ac} \mathbf{\Lambda}_{cc}^{-1} \mathbf{\Lambda}_{ac}^{\mathsf{T}} \mathbf{\Lambda}_{aa}^{-1} \end{bmatrix}.$$
(2.201)

Therefore,

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b}), \qquad (2.202)$$

where

$$\mu_{a|b} = \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab} \left( \mathbf{x}_b - \mu_b \right),$$

$$\Sigma_{a|b} = \Lambda_{aa}^{-1} + \Lambda_{aa}^{-1} \Lambda_{ac} \Lambda_{cc}^{-1} \Lambda_{ac}^{\mathsf{T}} \Lambda_{aa}^{-1}.$$
(2.203)

Let  $\mathbf{A}$  be a square matrix and  $\mathbf{C}$  be an invertible matrix. By 2.24,

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{D} & -\mathbf{C}^{-1} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{M} & \mathbf{MBC} \\ \mathbf{CDM} & -\mathbf{C} + \mathbf{CDMBC} \end{bmatrix}, \tag{2.204}$$

where

$$\mathbf{M} = (\mathbf{A} + \mathbf{BCD})^{-1}. \tag{2.205}$$

By 2.24,

$$\begin{bmatrix} -\mathbf{C}^{-1} & \mathbf{D} \\ \mathbf{B} & \mathbf{A} \end{bmatrix}^{-1} = \begin{bmatrix} -\mathbf{N} & \mathbf{N}\mathbf{D}\mathbf{A}^{-1} \\ \mathbf{A}^{-1}\mathbf{B}\mathbf{N} & \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}\mathbf{N}\mathbf{D}\mathbf{A}^{-1} \end{bmatrix}, (2.206)$$

where

$$\mathbf{N} = \left(\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B}\right)^{-1}.$$
 (2.207)

Therefore,

$$\mathbf{M} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} \mathbf{N} \mathbf{D} \mathbf{A}^{-1}. \tag{2.208}$$

#### 2.27

(a)

Let  $\mathbf{x}$  and  $\mathbf{z}$  be two variables. By the definition,

$$E(\mathbf{x} + \mathbf{z}) = \int \int (\mathbf{x} + \mathbf{z}) p(\mathbf{x}, \mathbf{z}) d\mathbf{x} d\mathbf{z}.$$
 (2.209)

The right hand side can be written as

$$\int \mathbf{x} \left( \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \right) d\mathbf{x} + \int \mathbf{z} \left( \int p(\mathbf{x}, \mathbf{z}) d\mathbf{x} \right) d\mathbf{z} = \int \mathbf{x} p(\mathbf{x}) d\mathbf{x} + \int \mathbf{z} p(\mathbf{z}) d\mathbf{z}.$$
(2.210)

$$E(\mathbf{x} + \mathbf{z}) = E \mathbf{x} + E \mathbf{z}. \tag{2.211}$$

(b)

Let  $\mathbf{x}$  and  $\mathbf{z}$  be two independent variables. By the definition,

$$cov(\mathbf{x} + \mathbf{z}) = \int \int (\mathbf{x} + \mathbf{z} - E(\mathbf{x} + \mathbf{z})) (\mathbf{x} + \mathbf{z} - E(\mathbf{x} + \mathbf{z}))^{\mathsf{T}} p(\mathbf{x}, \mathbf{z}) d\mathbf{x} d\mathbf{z}.$$
(2.212)

The right hand side can be written as

$$\int \int (\mathbf{x} - \mathbf{E} \mathbf{x}) (\mathbf{x} - \mathbf{E} \mathbf{x})^{\mathsf{T}} p(\mathbf{x}, \mathbf{z}) d\mathbf{x} d\mathbf{z} + \int \int (\mathbf{x} - \mathbf{E} \mathbf{x}) (\mathbf{z} - \mathbf{E} \mathbf{z})^{\mathsf{T}} p(\mathbf{x}, \mathbf{z}) d\mathbf{x} d\mathbf{z} 
+ \int \int (\mathbf{z} - \mathbf{E} \mathbf{z}) (\mathbf{x} - \mathbf{E} \mathbf{x})^{\mathsf{T}} p(\mathbf{x}, \mathbf{z}) d\mathbf{x} d\mathbf{z} + \int \int (\mathbf{z} - \mathbf{E} \mathbf{z}) (\mathbf{z} - \mathbf{E} \mathbf{z})^{\mathsf{T}} p(\mathbf{x}, \mathbf{z}) d\mathbf{x} d\mathbf{z}.$$
(2.213)

Each term can be written as

$$\int (\mathbf{x} - \mathbf{E} \, \mathbf{x}) (\mathbf{x} - \mathbf{E} \, \mathbf{x})^{\mathsf{T}} \left( \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \right) d\mathbf{x} = \int (\mathbf{x} - \mathbf{E} \, \mathbf{x}) (\mathbf{x} - \mathbf{E} \, \mathbf{x})^{\mathsf{T}} p(\mathbf{x}) d\mathbf{x},$$

$$\int (\mathbf{x} - \mathbf{E} \, \mathbf{x}) p(\mathbf{x}) d\mathbf{x} \int (\mathbf{z} - \mathbf{E} \, \mathbf{z})^{\mathsf{T}} p(\mathbf{z}) d\mathbf{z} = (\mathbf{E} \, \mathbf{x} - \mathbf{E} \, \mathbf{x}) (\mathbf{E} \, \mathbf{z} - \mathbf{E} \, \mathbf{z})^{\mathsf{T}},$$

$$\int (\mathbf{z} - \mathbf{E} \, \mathbf{z}) p(\mathbf{z}) d\mathbf{z} \int (\mathbf{x} - \mathbf{E} \, \mathbf{x})^{\mathsf{T}} p(\mathbf{x}) d\mathbf{x} = (\mathbf{E} \, \mathbf{z} - \mathbf{E} \, \mathbf{z}) (\mathbf{E} \, \mathbf{x} - \mathbf{E} \, \mathbf{x})^{\mathsf{T}},$$

$$\int (\mathbf{z} - \mathbf{E} \, \mathbf{z}) (\mathbf{z} - \mathbf{E} \, \mathbf{z})^{\mathsf{T}} \left( \int p(\mathbf{x}, \mathbf{z}) d\mathbf{x} \right) d\mathbf{z} = \int (\mathbf{z} - \mathbf{E} \, \mathbf{z}) (\mathbf{z} - \mathbf{E} \, \mathbf{z})^{\mathsf{T}} p(\mathbf{z}) d\mathbf{z}.$$

$$(2.214)$$

Therefore,

$$cov(\mathbf{x} + \mathbf{z}) = cov \,\mathbf{x} + cov \,\mathbf{z}.\tag{2.215}$$

#### 2.28

Let  $\mathbf{z}$  be a variable such that

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix},$$

$$\mathbf{E} \mathbf{z} = \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \end{bmatrix},$$

$$\mathbf{cov} \mathbf{z} = \begin{bmatrix} \boldsymbol{\Lambda}^{-1} & \boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathsf{T}} \\ \mathbf{A}\boldsymbol{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathsf{T}} \end{bmatrix},$$
(2.216)

where  $\mathbf{x}$  and  $\mathbf{y}$  are Gaussian variables. By 2.29,

$$(\cos \mathbf{z})^{-1} = \begin{bmatrix} \mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A} & -\mathbf{A}^{\mathsf{T}} \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{bmatrix}.$$

Then,  $\ln p(\mathbf{z})$  except the terms independent of  $\mathbf{x}$  and  $\mathbf{y}$  is given by

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

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

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

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

$$+ \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A} (\mathbf{x} - \boldsymbol{\mu}).$$
(2.217)

The right hand side can be written as

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2} (\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})^{\mathsf{T}} \mathbf{L} (\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}). \tag{2.218}$$

Therefore,

$$p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}\right),$$
  

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}\left(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}\right).$$
(2.219)

#### 2.29

Let

$$\mathbf{R} = \begin{bmatrix} \mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A} & -\mathbf{A}^{\mathsf{T}} \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{bmatrix}. \tag{2.220}$$

By 2.24,

$$\mathbf{R}^{-1} = \begin{bmatrix} \mathbf{\Lambda}^{-1} & \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \\ \mathbf{A}\mathbf{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \end{bmatrix}.$$
 (2.221)

#### 2.30

Let

$$\mathbf{R}^{-1} = \begin{bmatrix} \mathbf{\Lambda}^{-1} & \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \\ \mathbf{A}\mathbf{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \end{bmatrix}. \tag{2.222}$$

Then,

$$\mathbf{R}^{-1} \begin{bmatrix} \mathbf{\Lambda} \boldsymbol{\mu} - \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{A} \boldsymbol{\mu} + \mathbf{b} \end{bmatrix}. \tag{2.223}$$

# 2.31

Let y be a variable such that

$$\mathbf{y} = \mathbf{x} + \mathbf{z},\tag{2.224}$$

where

$$p(\mathbf{x}) = \mathcal{N} \left( \mathbf{x} | \boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{x}} \right),$$
  

$$p(\mathbf{z}) = \mathcal{N} \left( \mathbf{z} | \boldsymbol{\mu}_{\mathbf{z}}, \boldsymbol{\Sigma}_{\mathbf{z}} \right).$$
(2.225)

By marginalisation,

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x}.$$
 (2.226)

The right hand side can be written as

$$\int \mathcal{N}(\mathbf{y}|\mathbf{x} + \boldsymbol{\mu}_{\mathbf{z}}, \boldsymbol{\Sigma}_{\mathbf{z}}) \,\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{x}}) \, d\mathbf{x}. \tag{2.227}$$

The logarithm of the integrand except the terms independent of  ${\bf x}$  and  ${\bf y}$  is given by

$$-\frac{1}{2}(\mathbf{y} - \mathbf{x} - \boldsymbol{\mu}_{\mathbf{z}})^{\mathsf{T}} \boldsymbol{\Sigma}_{\mathbf{z}}^{-1} (\mathbf{y} - \mathbf{x} - \boldsymbol{\mu}_{\mathbf{z}}) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}})^{\mathsf{T}} \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}}). \quad (2.228)$$

The terms except the ones independent of  $\mathbf{x}$  and  $\mathbf{y}$  can be written as

$$-\frac{1}{2}\mathbf{u}^{\mathsf{T}}\mathbf{R}\mathbf{u} + \mathbf{u}^{\mathsf{T}}\mathbf{v} = -\frac{1}{2}\left(\mathbf{u} - \mathbf{R}^{-1}\mathbf{v}\right)^{\mathsf{T}}\mathbf{R}\left(\mathbf{u} - \mathbf{R}^{-1}\mathbf{v}\right) + \frac{1}{2}\mathbf{v}^{\mathsf{T}}\mathbf{R}^{-1}\mathbf{v}, \quad (2.229)$$

where

$$\mathbf{u} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix},$$

$$\mathbf{R} = \begin{bmatrix} \mathbf{\Sigma}_{\mathbf{x}}^{-1} + \mathbf{\Sigma}_{\mathbf{z}}^{-1} & -\mathbf{\Sigma}_{\mathbf{z}}^{-1} \\ -\mathbf{\Sigma}_{\mathbf{z}}^{-1} & \mathbf{\Sigma}_{\mathbf{z}}^{-1} \end{bmatrix},$$

$$\mathbf{v} = \begin{bmatrix} \mathbf{\Sigma}_{\mathbf{x}}^{-1} \boldsymbol{\mu}_{\mathbf{x}} - \mathbf{\Sigma}_{\mathbf{z}}^{-1} \boldsymbol{\mu}_{\mathbf{z}} \\ \mathbf{\Sigma}_{\mathbf{z}}^{-1} \boldsymbol{\mu}_{\mathbf{z}} \end{bmatrix},$$
(2.230)

By 2.29 and 2.30,

$$\mathbf{R}^{-1} = \begin{bmatrix} \mathbf{\Sigma}_{\mathbf{x}} & \mathbf{\Sigma}_{\mathbf{x}} \\ \mathbf{\Sigma}_{\mathbf{x}} & \mathbf{\Sigma}_{\mathbf{x}} + \mathbf{\Sigma}_{\mathbf{z}} \end{bmatrix},$$

$$\mathbf{R}^{-1}\mathbf{v} = \begin{bmatrix} \boldsymbol{\mu}_{\mathbf{x}} \\ \boldsymbol{\mu}_{\mathbf{x}} + \boldsymbol{\mu}_{\mathbf{z}} \end{bmatrix}.$$
(2.231)

Therefore,

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}_{\mathbf{x}} + \boldsymbol{\mu}_{\mathbf{z}}, \boldsymbol{\Sigma}_{\mathbf{x}} + \boldsymbol{\Sigma}_{\mathbf{z}}).$$
 (2.232)

#### 2.32

Let  $\mathbf{x}$  and  $\mathbf{y}$  be variables such that

$$p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}\right),$$
  

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}\left(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}\right).$$
(2.233)

By the Bayes' theorem,

$$p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y}). \tag{2.234}$$

The logarithm of the left hand side except the terms independent of  $\mathbf{x}$  and  $\mathbf{y}$  can be written as

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

$$= -\frac{1}{2} (\mathbf{y} - \mathbf{A}(\mathbf{x} - \boldsymbol{\mu}) - \mathbf{A}\boldsymbol{\mu} - \mathbf{b})^{\mathsf{T}} \mathbf{L} (\mathbf{y} - \mathbf{A}(\mathbf{x} - \boldsymbol{\mu}) - \mathbf{A}\boldsymbol{\mu} - \mathbf{b}) \qquad (2.235)$$

$$-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}).$$

The right hand side can be written as

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

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

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

$$-\frac{1}{2} (\mathbf{y} - \mathbf{A} \boldsymbol{\mu} - \mathbf{b})^{\mathsf{T}} \mathbf{L} (\mathbf{y} - \mathbf{A} \boldsymbol{\mu} - \mathbf{b}),$$
(2.236)

where

$$\mathbf{z} = (\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A})^{-1} \mathbf{A}^{\mathsf{T}} \mathbf{L} (\mathbf{y} - \mathbf{A} \boldsymbol{\mu} - \mathbf{b}). \tag{2.237}$$

The right hand side can be written as

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu} - \mathbf{z})^{\mathsf{T}} (\boldsymbol{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A}) (\mathbf{x} - \boldsymbol{\mu} - \mathbf{z})$$

$$-\frac{1}{2}(\mathbf{y} - \mathbf{A} \boldsymbol{\mu} - \mathbf{b})^{\mathsf{T}} \mathbf{M} (\mathbf{y} - \mathbf{A} \boldsymbol{\mu} - \mathbf{b}),$$
(2.238)

where

$$\mathbf{M} = \mathbf{L} - \mathbf{L}\mathbf{A} \left(\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}}\mathbf{L}\mathbf{A}\right)^{-1} \mathbf{A}^{\mathsf{T}}\mathbf{L}. \tag{2.239}$$

We have

$$\mu + z = (\Lambda + A^{\mathsf{T}}LA)^{-1} ((\Lambda + A^{\mathsf{T}}LA) \mu + A^{\mathsf{T}}L(y - A\mu - b)).$$
 (2.240)

Then,

$$\mu + \mathbf{z} = (\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A})^{-1} (\mathbf{A}^{\mathsf{T}} \mathbf{L} (\mathbf{y} - \mathbf{b}) + \mathbf{\Lambda} \mu).$$
 (2.241)

By 2.26,

$$(\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A})^{-1} = \mathbf{\Lambda}^{-1} - \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \left( \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \right)^{-1} \mathbf{A} \mathbf{\Lambda}^{-1}. \tag{2.242}$$

Then,

$$\mathbf{M} = \mathbf{L} - \mathbf{L} \mathbf{A} \left( \mathbf{\Lambda}^{-1} - \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \left( \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \right)^{-1} \mathbf{A} \mathbf{\Lambda}^{-1} \right) \mathbf{A}^{\mathsf{T}} \mathbf{L}. \quad (2.243)$$

The right hand side can be written as

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

$$= \mathbf{L} - \mathbf{L} \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \left( \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \right)^{-1} \left( \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} - \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \right) \mathbf{L}.$$
(2.244)

The right hand side can be written as

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

$$= \mathbf{L} \left( \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} - \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \right) \left( \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \right)^{-1}.$$
(2.245)

Then,

$$\mathbf{M} = \left(\mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathsf{T}}\right)^{-1}.\tag{2.246}$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}\left(\mathbf{x}|\left(\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}}\mathbf{L}\mathbf{A}\right)^{-1}\left(\mathbf{A}^{\mathsf{T}}\mathbf{L}(\mathbf{y} - \mathbf{b}) + \mathbf{\Lambda}\boldsymbol{\mu}\right), \left(\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}}\mathbf{L}\mathbf{A}\right)^{-1}\right),$$
  

$$p(\mathbf{y}) = \mathcal{N}\left(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1}\mathbf{A}^{\mathsf{T}}\right).$$
(2.247)

#### 2.33

Refer to 2.32, while a different approach is presented below.

Let  $\mathbf{x}$  and  $\mathbf{y}$  be variables such that

$$p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}\right),$$
  

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}\left(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}\right).$$
(2.248)

By the Bayes' theorem,

$$p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y}). \tag{2.249}$$

The logarithm of the left hand side except the terms independent of  $\mathbf{x}$  and  $\mathbf{y}$  can be written as

$$-\frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})^{\mathsf{T}} \mathbf{L} (\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}). \tag{2.250}$$

The terms except the ones independent of  $\mathbf{x}$  and  $\mathbf{y}$  can be written as

$$-\frac{1}{2} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A} & -\mathbf{A}^{\mathsf{T}} \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} + \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} -\mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{b} + \mathbf{\Lambda} \boldsymbol{\mu} \\ \mathbf{L} \mathbf{b} \end{bmatrix}. \quad (2.251)$$

By 2.24,

$$\begin{bmatrix} \mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A} & -\mathbf{A}^{\mathsf{T}} \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{\Lambda}^{-1} & \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \\ \mathbf{A} \mathbf{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^{\mathsf{T}} \end{bmatrix}, \quad (2.252)$$

so that

$$\begin{bmatrix} \mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A} & -\mathbf{A}^{\mathsf{T}} \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{bmatrix}^{-1} \begin{bmatrix} -\mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{b} + \mathbf{\Lambda} \boldsymbol{\mu} \\ \mathbf{L} \mathbf{b} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{A} \boldsymbol{\mu} + \mathbf{b} \end{bmatrix}. \tag{2.253}$$

Then,

$$p(\mathbf{y}) = \mathcal{N}\left(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathsf{T}}\right). \tag{2.254}$$

By 2.25,

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu} + (\boldsymbol{\Lambda} + \mathbf{A}^{\mathsf{T}}\mathbf{L}\mathbf{A})^{-1}\mathbf{A}^{\mathsf{T}}\mathbf{L}\left(\mathbf{y} - \mathbf{A}\boldsymbol{\mu} - \mathbf{b}\right), (\boldsymbol{\Lambda} + \mathbf{A}^{\mathsf{T}}\mathbf{L}\mathbf{A})^{-1}\right).$$
(2.255)

We have

$$\mu + (\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A})^{-1} \mathbf{A}^{\mathsf{T}} \mathbf{L} (\mathbf{y} - \mathbf{A} \boldsymbol{\mu} - \mathbf{b})$$

$$= (\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}} \mathbf{L} \mathbf{A})^{-1} (\mathbf{\Lambda} \boldsymbol{\mu} + \mathbf{A}^{\mathsf{T}} \mathbf{L} (\mathbf{y} - \mathbf{b})).$$
(2.256)

Therefore,

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}\left(\mathbf{x}|\left(\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}}\mathbf{L}\mathbf{A}\right)^{-1}\left(\mathbf{\Lambda}\boldsymbol{\mu} + \mathbf{A}^{\mathsf{T}}\mathbf{L}(\mathbf{y} - \mathbf{b})\right), \left(\mathbf{\Lambda} + \mathbf{A}^{\mathsf{T}}\mathbf{L}\mathbf{A}\right)^{-1}\right).$$
(2.257)

# 2.34

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be variables such that

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}). \tag{2.258}$$

Then,

$$\ln\left(\prod_{n=1}^{N} p(\mathbf{x}_n)\right) = -\frac{ND}{2}\ln(2\pi) - \frac{N}{2}\ln|\det \mathbf{\Sigma}| - \frac{1}{2}\sum_{n=1}^{N}(\mathbf{x}_n - \boldsymbol{\mu})^{\mathsf{T}}\mathbf{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}).$$
(2.259)

By 3.21(a), setting the derivatives with respect to  $\mu$  and  $\Sigma$  to zero gives

$$\mathbf{0} = \sum_{n=1}^{N} \left( \mathbf{\Sigma}^{-1} + \left( \mathbf{\Sigma}^{-1} \right)^{\mathsf{T}} \right) (\mathbf{x}_{n} - \boldsymbol{\mu}),$$

$$\mathbf{O} = -\frac{N}{2} \left( \mathbf{\Sigma}^{-1} \right)^{\mathsf{T}} + \frac{1}{2} \left( \mathbf{\Sigma}^{-1} \right)^{2} \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathsf{T}}.$$

$$(2.260)$$

Therefore, the maximum likelihood solutions for  $\mu$  and  $\Sigma$  are given by

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

$$\boldsymbol{\Sigma}_{\mathrm{ML}} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}_{\mathrm{ML}}) (\mathbf{x}_{n} - \boldsymbol{\mu}_{\mathrm{ML}})^{\mathsf{T}}.$$
(2.261)

# 2.35

(a)

Let  $\mathbf{x}$  be a variable such that

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}).$$
 (2.262)

Then,

$$\mathbf{E} \mathbf{x} \mathbf{x}^{\mathsf{T}} = \mathbf{E} (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu})^{\mathsf{T}}. \tag{2.263}$$

The right hand side can be written as

$$E(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} + \boldsymbol{\mu} E(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} + E(\mathbf{x} - \boldsymbol{\mu})\boldsymbol{\mu}^{\mathsf{T}} + \boldsymbol{\mu}\boldsymbol{\mu}^{\mathsf{T}}. \tag{2.264}$$

Since

$$\begin{aligned}
\mathbf{E} \mathbf{x} &= \boldsymbol{\mu}, \\
\cos \mathbf{x} &= \boldsymbol{\Sigma},
\end{aligned} (2.265)$$

The right hand side can be written as  $\Sigma + \mu \mu^{\dagger}$ . Therefore,

$$\mathbf{E} \mathbf{x} \mathbf{x}^{\mathsf{T}} = \mathbf{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^{\mathsf{T}}. \tag{2.266}$$

(b)

Let  $\mathbf{x}_n$  and  $\mathbf{x}_m$  be independent variables such that

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

$$p(\mathbf{x}_m) = \mathcal{N}(\mathbf{x}_m | \boldsymbol{\mu}, \boldsymbol{\Sigma}).$$
(2.267)

If  $n \neq m$ , then

$$\mathbf{E}\,\mathbf{x}_n\mathbf{x}_m^{\mathsf{T}} = \mathbf{E}\,\mathbf{x}_n\,\mathbf{E}\,\mathbf{x}_m^{\mathsf{T}}.\tag{2.268}$$

The right hand side can be written as  $\mu\mu^{\dagger}$ . Therefore, by (a),

$$\mathbf{E} \, \mathbf{x}_n \mathbf{x}_m^{\mathsf{T}} = I_{nm} \mathbf{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^{\mathsf{T}}. \tag{2.269}$$

(c)

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be variables such that

$$p(\mathbf{x}_n) = \mathcal{N}\left(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}\right). \tag{2.270}$$

By 2.34, the maximum likelihood solutions for  $\mu$  and  $\Sigma$  are given by

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

$$\boldsymbol{\Sigma}_{\mathrm{ML}} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}_{\mathrm{ML}}) (\mathbf{x}_{n} - \boldsymbol{\mu}_{\mathrm{ML}})^{\mathsf{T}}.$$
(2.271)

Then,

$$E \Sigma_{\mathrm{ML}} = \frac{1}{N} \sum_{n=1}^{N} E(\mathbf{x}_{n} - \boldsymbol{\mu}_{\mathrm{ML}}) (\mathbf{x}_{n} - \boldsymbol{\mu}_{\mathrm{ML}})^{\mathsf{T}}.$$
 (2.272)

The right hand side can be written as

$$\frac{1}{N} \sum_{n=1}^{N} \mathbf{E} \mathbf{x}_{n} \mathbf{x}_{n}^{\mathsf{T}} - \frac{1}{N^{2}} \sum_{n=1}^{N} \mathbf{E} \left( \sum_{n=1}^{N} \mathbf{x}_{n} \right) \mathbf{x}_{n}^{\mathsf{T}} - \frac{1}{N^{2}} \sum_{n=1}^{N} \mathbf{E} \mathbf{x}_{n} \left( \sum_{n=1}^{N} \mathbf{x}_{n} \right)^{\mathsf{T}} + \frac{1}{N^{3}} \sum_{n=1}^{N} \mathbf{E} \left( \sum_{n=1}^{N} \mathbf{x}_{n} \right) \left( \sum_{n=1}^{N} \mathbf{x}_{n} \right)^{\mathsf{T}}.$$
(2.273)

By (b), the first term can be written as  $\Sigma + \mu \mu^{\dagger}$ . By (b), the second and third terms can be written as

$$-\frac{1}{N}\left((\mathbf{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^{\mathsf{T}}) + (N-1)\boldsymbol{\mu}\boldsymbol{\mu}^{\mathsf{T}}\right) = -\frac{1}{N}\mathbf{\Sigma} - \boldsymbol{\mu}\boldsymbol{\mu}^{\mathsf{T}}.$$
 (2.274)

By (b), the fourth term can be written as

$$\frac{1}{N^2} \left( N \left( \mathbf{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^{\mathsf{T}} \right) + N (N - 1) \boldsymbol{\mu} \boldsymbol{\mu}^{\mathsf{T}} \right) = \frac{1}{N} \mathbf{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^{\mathsf{T}}. \tag{2.275}$$

Then,

$$E \Sigma_{ML} = (\Sigma + \mu \mu^{\mathsf{T}}) + 2\left(-\frac{1}{N}\Sigma - \mu \mu^{\mathsf{T}}\right) + \frac{1}{N}\Sigma + \mu \mu^{\mathsf{T}}.$$
 (2.276)

Therefore,

$$E \Sigma_{\rm ML} = \frac{N-1}{N} \Sigma. \tag{2.277}$$

# 2.36

Let  $x_1, \dots, x_N$  be variables such that

$$p(x_n) = \mathcal{N}\left(x_n|\mu, \sigma^2\right). \tag{2.278}$$

Let us assume that  $\mu$  is known. By 2.34, the maximum likelihood solution for  $\sigma^2$  is given by

$$\sigma_{\rm ML}^{2(N)} = \frac{1}{N} \sum_{n=1}^{N} (x_n - \mu)^2.$$
 (2.279)

The right hand side can be written as

$$\frac{1}{N}(x_N - \mu)^2 + \frac{1}{N} \sum_{n=1}^{N-1} (x_n - \mu)^2 = \frac{1}{N}(x_N - \mu)^2 + \frac{N-1}{N} \sigma_{ML}^{2(N-1)}. \quad (2.280)$$

Then,

$$\sigma_{\rm ML}^{2(N)} = \sigma_{\rm ML}^{2(N-1)} + \frac{1}{N} \left( (x_N - \mu)^2 - \sigma_{\rm ML}^{2(N-1)} \right). \tag{2.281}$$

We have

$$\frac{\partial}{\partial \sigma^2} \left( -\ln p(x_n) \right) = \frac{1}{2\sigma^2} - \frac{1}{2(\sigma^2)^2} (x_n - \mu)^2.$$
 (2.282)

Therefore,

$$\sigma_{\rm ML}^{2(N)} = \sigma_{\rm ML}^{2(N-1)} - \frac{2\left(\sigma_{\rm ML}^{2(N-1)}\right)^{2}}{N} \left(\frac{\partial}{\partial \sigma^{2}} \left(-\ln p(x_{N})\right)\right) \bigg|_{\sigma^{2} = \sigma_{\rm ML}^{2(N-1)}}.$$
(2.283)

# 2.37

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be variables such that

$$p(\mathbf{x}_n) = \mathcal{N}\left(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}\right). \tag{2.284}$$

Let us assume that  $\mu$  is known. By 2.34, the maximum likelihood solution for  $\Sigma$  is given by

$$\Sigma_{\mathrm{ML}}^{(N)} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^{\mathsf{T}}.$$
 (2.285)

The right hand side can be written as

$$\frac{1}{N}(\mathbf{x}_N - \boldsymbol{\mu})(\mathbf{x}_N - \boldsymbol{\mu})^{\mathsf{T}} + \frac{1}{N} \sum_{n=1}^{N-1} (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^{\mathsf{T}}$$

$$= \frac{1}{N} (\mathbf{x}_N - \boldsymbol{\mu})(\mathbf{x}_N - \boldsymbol{\mu})^{\mathsf{T}} + \frac{N-1}{N} \boldsymbol{\Sigma}_{\mathrm{ML}}^{(N-1)}.$$
(2.286)

Then,

$$\Sigma_{\mathrm{ML}}^{(N)} = \Sigma_{\mathrm{ML}}^{(N-1)} + \frac{1}{N} \left( (\mathbf{x}_N - \boldsymbol{\mu}) (\mathbf{x}_N - \boldsymbol{\mu})^{\mathsf{T}} - \Sigma_{\mathrm{ML}}^{(N-1)} \right). \tag{2.287}$$

By 3.21(a), we have

$$\frac{\partial}{\partial \mathbf{\Sigma}} \left( -\ln p(x_n) \right) = -\frac{1}{2} \left( \mathbf{\Sigma}^{-1} \right)^{\mathsf{T}} + \frac{1}{2} \left( \mathbf{\Sigma}^{-1} \right)^2 (\mathbf{x}_N - \boldsymbol{\mu}) (\mathbf{x}_N - \boldsymbol{\mu})^{\mathsf{T}}. \tag{2.288}$$

Therefore,

$$\Sigma_{\mathrm{ML}}^{(N)} = \Sigma_{\mathrm{ML}}^{(N-1)} - \frac{\left(\Sigma_{\mathrm{ML}}^{(N-1)}\right)^{2}}{N} \left(\frac{\partial}{\partial \Sigma} \left(-\ln p\left(\mathbf{x}_{N}\right)\right)\right) \Big|_{\Sigma = \Sigma_{\mathrm{MI}}^{(N-1)}}.$$
 (2.289)

# 2.38

Let  $x_1, \dots, x_N$  be variables such that

$$p(x_n|\mu) = \mathcal{N}\left(x_n|\mu, \sigma^2\right),$$
  

$$p(\mu) = \mathcal{N}\left(\mu|\mu_0, \sigma_0^2\right).$$
(2.290)

By the Bayes' theorem,

$$p(\mu|\mathbf{x})p(\mathbf{x}) = p(\mathbf{x}|\mu)p(\mu). \tag{2.291}$$

The logarithm of the right hand side excpt the terms independent of  $\mathbf{x}$  and  $\mu$  can be written as

$$-\frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu)^2 - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2.$$
 (2.292)

By 2.34, the maximum likelihood solution for  $\mu$  is given by

$$\mu_{\rm ML} = \frac{1}{N} \sum_{n=1}^{N} x_n. \tag{2.293}$$

Then, the first term can be written as

$$-\frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu_{\text{ML}} + \mu_{\text{ML}} - \mu)^2$$

$$= -\frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu_{\text{ML}})^2 - \frac{\mu_{\text{ML}} - \mu}{\sigma^2} \sum_{n=1}^{N} (x_n - \mu_{\text{ML}}) - \frac{N}{2\sigma^2} (\mu_{\text{ML}} - \mu)^2.$$
(2.294)

Since the second term of the right hand side is zero, the logarithm except the terms independent of  $\mathbf{x}$  and  $\mu$  can be written as

$$-\frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu_{\rm ML})^2 - \frac{N}{2\sigma^2} (\mu_{\rm ML} - \mu)^2 - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2$$

$$= -\frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu_{\rm ML})^2 - \frac{1}{2\sigma_N^2} (\mu - \mu_N)^2 + \frac{\mu_N^2}{2\sigma_N^2} - \frac{N\mu_{\rm ML}^2}{2\sigma^2} - \frac{\mu_0^2}{2\sigma_0^2},$$
(2.295)

where

$$\mu_{N} = \frac{N\sigma_{0}^{2}}{N\sigma_{0}^{2} + \sigma^{2}}\mu_{ML} + \frac{\sigma^{2}}{N\sigma_{0}^{2} + \sigma^{2}}\mu_{0},$$

$$\sigma_{N}^{2} = \frac{\sigma^{2}\sigma_{0}^{2}}{N\sigma_{0}^{2} + \sigma^{2}}.$$
(2.296)

Therefore,

$$p(\mu|\mathbf{x}) = \mathcal{N}\left(\mu|\mu_N, \sigma_N^2\right). \tag{2.297}$$

#### 2.39

Let  $x_1, \dots, x_N$  be variables such that

$$p(x_n|\mu) = \mathcal{N}\left(x_n|\mu, \sigma^2\right),$$
  

$$p(\mu) = \mathcal{N}\left(\mu|\mu_0, \sigma_0^2\right).$$
(2.298)

(a)

By 2.38,

$$p(\mu|\mathbf{x}) = \mathcal{N}\left(\mu|\mu_N, \sigma_N^2\right),$$
 (2.299)

where

$$\mu_N = \frac{\sigma_0^2}{N\sigma_0^2 + \sigma^2} \sum_{n=1}^N x_n + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0,$$

$$\sigma_N^2 = \frac{\sigma^2 \sigma_0^2}{N\sigma_0^2 + \sigma^2}.$$
(2.300)

Then,

$$\frac{1}{\sigma_N^2} = \frac{(N-1)\sigma_0^2 + \sigma^2}{\sigma^2 \sigma_0^2} + \frac{1}{\sigma^2}.$$
 (2.301)

Therefore,

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_{N-1}^2} + \frac{1}{\sigma^2}. (2.302)$$

We have

$$\frac{\mu_N}{\sigma_N^2} = \frac{1}{\sigma^2} \sum_{n=1}^N x_n + \frac{\mu_0}{\sigma_0^2},\tag{2.303}$$

so that

$$\frac{\mu_{N-1}}{\sigma_{N-1}^2} = \frac{1}{\sigma^2} \sum_{n=1}^{N-1} x_n + \frac{\mu_0}{\sigma_0^2}.$$
 (2.304)

Therefore,

$$\frac{\mu_N}{\sigma_N^2} = \frac{\mu_{N-1}}{\sigma_{N-1}^2} + \frac{x_N}{\sigma^2}.$$
 (2.305)

(b)

By the Bayes' theorem,

$$p(\mu|\mathbf{x}_N)p(\mathbf{x}_N) = p(\mathbf{x}_N|\mu)p(\mu). \tag{2.306}$$

Since  $x_N$  and  $\mathbf{x}_{N-1}$  are independent, it can be written as

$$p(\mu|\mathbf{x}_N)p(x_N)p(\mathbf{x}_{N-1}) = p(x_N|\mu)p(\mathbf{x}_{N-1}|\mu)p(\mu). \tag{2.307}$$

By the Bayes' theorem, the right hand side can be written as

$$p(x_N|\mu)p(\mu|\mathbf{x}_{N-1})p(\mathbf{x}_{N-1}).$$
 (2.308)

Then,

$$p(\mu|\mathbf{x}_N)p(x_N) = p(\mu|\mathbf{x}_{N-1})p(x_N|\mu). \tag{2.309}$$

The logarithm of the integrand of the right hand side except the terms independent of  $\mu$  or  $x_N$  can be written as

$$-\frac{1}{2\sigma_{N-1}^{2}}(\mu-\mu_{N-1})^{2} - \frac{1}{2\sigma^{2}}(x_{N}-\mu)^{2}$$

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

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

$$(2.310)$$

Then,

$$\mu_N = \frac{\frac{1}{\sigma_{N-1}^2}}{\frac{1}{\sigma_{N-1}^2} + \frac{1}{\sigma^2}} \mu_{N-1} + \frac{\frac{1}{\sigma^2}}{\frac{1}{\sigma_{N-1}^2} + \frac{1}{\sigma^2}} x_N,$$

$$\sigma_N^2 = \frac{1}{\frac{1}{\sigma_{N-1}^2} + \frac{1}{\sigma^2}}.$$
(2.311)

Therefore,

$$\frac{\mu_N}{\sigma_N^2} = \frac{\mu_{N-1}}{\sigma_{N-1}^2} + \frac{x_N}{\sigma^2}, 
\frac{1}{\sigma_N^2} = \frac{1}{\sigma_{N-1}^2} + \frac{1}{\sigma^2}.$$
(2.312)

# 2.40

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be variables such that

$$p(\mathbf{x}_n|\boldsymbol{\mu}) = \mathcal{N}\left(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma}\right),$$
  

$$p(\boldsymbol{\mu}) = \mathcal{N}\left(\boldsymbol{\mu}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0\right).$$
(2.313)

By 2.34, the maximum likelihood solution for  $\mu$  is given by

$$\boldsymbol{\mu}_{\mathrm{ML}} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_{n}, \qquad (2.314)$$

By the Bayes' theorem,

$$p(\boldsymbol{\mu}|\mathbf{X})p(\mathbf{X}) = p(\mathbf{X}|\boldsymbol{\mu})p(\boldsymbol{\mu}). \tag{2.315}$$

The logarithm of the right hand side excpt the terms independent of  ${\bf X}$  and  ${\boldsymbol \mu}$  can be written as

$$-\frac{1}{2}\sum_{n=1}^{N}(\mathbf{x}_{n}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}_{n}-\boldsymbol{\mu})-\frac{1}{2}(\boldsymbol{\mu}-\boldsymbol{\mu}_{0})\boldsymbol{\Sigma}_{0}^{-1}(\boldsymbol{\mu}-\boldsymbol{\mu}_{0})^{\mathsf{T}}.$$
 (2.316)

The first term can be written as

$$-\frac{1}{2}\sum_{n=1}^{N}(\mathbf{x}_{n}-\boldsymbol{\mu}_{\mathrm{ML}}+\boldsymbol{\mu}_{\mathrm{ML}}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}_{n}-\boldsymbol{\mu}_{\mathrm{ML}}+\boldsymbol{\mu}_{\mathrm{ML}}-\boldsymbol{\mu})$$

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

$$-\frac{N}{2}(\boldsymbol{\mu}_{\mathrm{ML}}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_{\mathrm{ML}}-\boldsymbol{\mu}).$$
(2.317)

The second term of the right hand side is zero. Then, the logarithm except the terms independent of X and  $\mu$  can be written as

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

$$-\frac{1}{2}(\boldsymbol{\mu}-\boldsymbol{\mu}_{0})\boldsymbol{\Sigma}_{0}^{-1}(\boldsymbol{\mu}-\boldsymbol{\mu}_{0})^{\mathsf{T}}$$

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

$$+\frac{1}{2}\boldsymbol{\mu}_{N}^{\mathsf{T}}\boldsymbol{\Sigma}_{N}^{-1}\boldsymbol{\mu}_{N} - \frac{N}{2}\boldsymbol{\mu}_{\mathrm{ML}}^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_{\mathrm{ML}},$$

$$(2.318)$$

where

$$\mu_{N} = (N\Sigma^{-1} + \Sigma_{0}^{-1})^{-1} (N\Sigma^{-1}\mu_{ML} + \Sigma_{0}^{-1}\mu_{0}),$$
  

$$\Sigma_{N} = (N\Sigma^{-1} + \Sigma_{0}^{-1})^{-1}.$$
(2.319)

Therefore,

$$p(\boldsymbol{\mu}|\mathbf{X}) = \mathcal{N}(\boldsymbol{\mu}|\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N). \tag{2.320}$$

#### 2.41

By the definition,

$$Gam(\lambda|a,b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda). \tag{2.321}$$

Then,

$$\int_0^\infty \operatorname{Gam}(\lambda|a,b)d\lambda = \frac{b^a}{\Gamma(a)} \int_0^\infty \lambda^{a-1} \exp(-b\lambda)d\lambda. \tag{2.322}$$

By the transformation

$$\lambda' = b\lambda, \tag{2.323}$$

the right hand side can be written as

$$\frac{b^a}{\Gamma(a)} \int_0^\infty \left(\frac{\lambda'}{b}\right)^{a-1} \exp(-\lambda') \frac{1}{b} d\lambda' = \frac{1}{\Gamma(a)} \int_0^\infty {\lambda'}^{a-1} \exp(-\lambda') d\lambda'. \quad (2.324)$$

The right hand side can be written as

$$\frac{1}{\Gamma(a)}\Gamma(a) = 1. \tag{2.325}$$

Therefore,

$$\int_0^\infty \operatorname{Gam}(\lambda|a,b)d\lambda = 1. \tag{2.326}$$

# 2.42

Let  $\lambda$  be a variable such that

$$p(\lambda) = \operatorname{Gam}(\lambda|a, b). \tag{2.327}$$

By the definition,

$$Gam(\lambda|a,b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda).$$
 (2.328)

(a)

We have

$$E \lambda = \frac{b^a}{\Gamma(a)} \int_0^\infty \lambda^a \exp\left(-\frac{\lambda}{b}\right) d\lambda. \tag{2.329}$$

By the transformation

$$\lambda' = b\lambda, \tag{2.330}$$

the right hand side can be written as

$$\frac{b^a}{\Gamma(a)} \int_0^\infty \left(\frac{\lambda'}{b}\right)^a \exp(-\lambda') \frac{1}{b} d\lambda' = \frac{1}{b\Gamma(a)} \int_0^\infty {\lambda'}^a \exp(-\lambda') d\lambda'. \tag{2.331}$$

The right hand side can be written as

$$\frac{1}{b\Gamma(a)}\Gamma(a+1) = \frac{a}{b}.$$
(2.332)

$$E\lambda = \frac{a}{b}. (2.333)$$

(b)

We have

$$E \lambda^{2} = \frac{b^{a}}{\Gamma(a)} \int_{0}^{\infty} \lambda^{a+1} \exp\left(-\frac{\lambda}{b}\right) d\lambda. \tag{2.334}$$

By the transformation

$$\lambda' = b\lambda, \tag{2.335}$$

the right hand side can be written as

$$\frac{b^a}{\Gamma(a)} \int_0^\infty \left(\frac{\lambda'}{b}\right)^{a+1} \exp(-\lambda') \frac{1}{b} d\lambda' = \frac{1}{b^2 \Gamma(a)} \int_0^\infty {\lambda'}^{a+1} \exp(-\lambda') d\lambda'. \quad (2.336)$$

The right hand side can be written as

$$\frac{1}{b^2\Gamma(a)}\Gamma(a+2) = \frac{a(a+1)}{b^2}.$$
 (2.337)

Then,

$$E \lambda^2 = \frac{a(a+1)}{b^2}.$$
 (2.338)

We have

$$\operatorname{var} \lambda = \operatorname{E} \lambda^2 - (\operatorname{E} \lambda)^2. \tag{2.339}$$

Therefore,

$$\operatorname{var} \lambda = \frac{a}{b^2}.\tag{2.340}$$

(c)

Setting the derivative of  $Gam(\lambda|a,b)$  with respect to  $\lambda$  to zero gives

$$0 = \frac{b^a}{\Gamma(a)} \left( \frac{a-1}{\lambda} - b \right) \lambda^{a-1} \exp\left( -\frac{\lambda}{b} \right). \tag{2.341}$$

Therefore,

$$\operatorname{mode} \lambda = \frac{a-1}{b}.\tag{2.342}$$

# 2.43

Let

$$p\left(x|\sigma^2,q\right) = \frac{q}{2\Gamma(\frac{1}{q})} \left(2\sigma^2\right)^{-\frac{1}{q}} \exp\left(-\frac{|x|^q}{2\sigma^2}\right). \tag{2.343}$$

(a)

We have

$$\int_{-\infty}^{\infty} p\left(x|\sigma^2, q\right) dx = \frac{q}{\Gamma(\frac{1}{q})} \left(2\sigma^2\right)^{-\frac{1}{q}} \int_{0}^{\infty} \exp\left(-\frac{x^q}{2\sigma^2}\right) dx. \tag{2.344}$$

By the transformation

$$x' = \frac{x^q}{2\sigma^2},\tag{2.345}$$

the right hand side can be written as

$$\frac{q}{\Gamma(\frac{1}{q})} \left(2\sigma^2\right)^{-\frac{1}{q}} \int_0^\infty \exp(-x') \left(2\sigma^2\right)^{\frac{1}{q}} \frac{1}{q} x^{\frac{1}{q}-1} dx'$$

$$= \frac{1}{\Gamma(\frac{1}{q})} \int_0^\infty x^{\frac{1}{q}-1} \exp(-x') dx'.$$
(2.346)

The right hand side can be written as

$$\frac{1}{\Gamma(\frac{1}{q})}\Gamma\left(\frac{1}{q}\right) = 1. \tag{2.347}$$

Therefore,

$$\int_{-\infty}^{\infty} p\left(x|\sigma^2, q\right) dx = 1. \tag{2.348}$$

(b)

We have

$$p\left(x|\sigma^2,2\right) = \frac{1}{\Gamma(\frac{1}{2})} \left(2\sigma^2\right)^{-\frac{1}{2}} \exp\left(-\frac{x^2}{2\sigma^2}\right). \tag{2.349}$$

Therefore,

$$p(x|\sigma^2, 2) = \mathcal{N}(x|0, \sigma^2). \tag{2.350}$$

(c)

Let  $\mathbf{t}=(t_1,\cdots,t_N)^\intercal$  and  $\mathbf{X}=\{\mathbf{x}_1,\cdots,\mathbf{x}_N\}$  such that

$$t_n = y(\mathbf{x}_n, \mathbf{w}) + \epsilon_n, \tag{2.351}$$

where

$$p(\epsilon_n) = p\left(\epsilon_n | \sigma^2, q\right). \tag{2.352}$$

Then, the logarithm of  $p(\epsilon_n)$  except the terms independent of **w** and  $\sigma^2$  can be written as

$$-\frac{|\epsilon_n|^q}{2\sigma^2} - \frac{1}{q}\ln\left(2\sigma^2\right). \tag{2.353}$$

Therefore, the logarithm of  $p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)$  except the terms independent of  $\mathbf{w}$  and  $\sigma^2$  can be written as

$$-\frac{1}{2\sigma^2} \sum_{n=1}^{N} \left| y(\mathbf{x}_n, \mathbf{w}) - t_n \right|^q - \frac{N}{q} \ln\left(2\sigma^2\right). \tag{2.354}$$

# 2.44

Let  $x_1, \dots, x_N$  be variables such that

$$p(x_n|\mu,\tau) = \mathcal{N}\left(x_n|\mu,\tau^{-1}\right),$$
  

$$p(\mu,\tau) = \mathcal{N}\left(\mu|\mu_0,(\beta\tau)^{-1}\right) \operatorname{Gam}(\tau|a,b).$$
(2.355)

By the Bayes' theorem,

$$p(\mu, \tau | \mathbf{x}) p(\mathbf{x}) = p(\mathbf{x} | \mu, \tau) p(\mu, \tau). \tag{2.356}$$

The logarithm of the right hand side except the terms independent of  $\mathbf{x}$ ,  $\mu$  and  $\tau$  can be written as

$$\frac{N}{2}\ln\tau - \frac{\tau}{2}\sum_{n=1}^{N}(x_n - \mu)^2 + \frac{1}{2}\ln\tau - \frac{\beta\tau}{2}(\mu - \mu_0)^2 + (a-1)\ln\tau - b\tau$$

$$= \left(a + \frac{N-1}{2}\right)\ln\tau - \frac{\tau}{2}\sum_{n=1}^{N}(x_n - \mu)^2 - \frac{\beta\tau}{2}(\mu - \mu_0)^2 - b\tau.$$
(2.357)

Let

$$\bar{x} = \frac{1}{N} \sum_{n=1}^{N} x_n. \tag{2.358}$$

Then,

$$\sum_{n=1}^{N} (x_n - \mu)^2 = \sum_{n=1}^{N} (x_n - \bar{x} + \bar{x} - \mu)^2.$$
 (2.359)

The right hand side can be written as

$$\sum_{n=1}^{N} (x_n - \bar{x})^2 + 2(\bar{x} - \mu) \sum_{n=1}^{N} (x_n - \bar{x}) + N(\bar{x} - \mu)^2$$

$$= \sum_{n=1}^{N} (x_n - \bar{x})^2 + N(\bar{x} - \mu)^2.$$
(2.360)

Then, the logarithm except the terms independent of  $\mathbf{x}$ ,  $\mu$  and  $\tau$  can be written as

$$\left(a + \frac{N-1}{2}\right) \ln \tau - \frac{N\tau}{2} (\bar{x} - \mu)^2 - \frac{\beta\tau}{2} (\mu - \mu_0)^2 - b\tau - \frac{\tau}{2} \sum_{n=1}^{N} (x_n - \bar{x})^2.$$
(2.361)

The second and third terms can be written as

$$-\frac{N\tau}{2}(\bar{x}-\mu)^2 - \frac{\beta\tau}{2}(\mu-\mu_0)^2$$

$$= -\frac{(N+\beta)\tau}{2}\left(\mu - \frac{N\bar{x}+\beta\mu_0}{N+\beta}\right)^2 + \frac{(N\bar{x}+\beta\mu_0)^2\tau}{2(N+\beta)} - \frac{N\tau}{2}\bar{x}^2 - \frac{\beta\tau}{2}\mu_0^2.$$
(2.362)

The second, third and forth terms on the right hand side can be written as

$$\frac{(N\bar{x} + \beta\mu_0)^2\tau - (N+\beta)N\tau\bar{x}^2 - (N+\beta)\beta\tau\mu_0^2}{2(N+\beta)} = -\frac{N\beta\tau(\bar{x} - \mu_0)^2}{2(N+\beta)}.$$
(2.363)

Then, the logrithm except the terms independent of  $\mathbf{x}$ ,  $\mu$  and  $\tau$  can be written as

$$-\frac{(N+\beta)\tau}{2} \left(\mu - \frac{N\bar{x} + \beta\mu_0}{N+\beta}\right)^2 + \left(a + \frac{N-1}{2}\right) \ln \tau - \left(b + \frac{N\beta(\bar{x} - \mu_0)^2}{2(N+\beta)} + \frac{1}{2} \sum_{n=1}^{N} (x_n - \bar{x})^2\right) \tau.$$
 (2.364)

$$p(\mu, \tau | \mathbf{x}) = \mathcal{N}\left(\mu | \mu_N, \tau_N^{-1}\right) \operatorname{Gam}\left(\tau | a_N, b_N\right), \qquad (2.365)$$

where

$$\mu_{N} = \frac{N\bar{x} + \beta\mu_{0}}{N + \beta},$$

$$\tau_{N} = (N + \beta)\tau,$$

$$a_{N} = a + \frac{N+1}{2},$$

$$b_{N} = b + \frac{N\beta(\bar{x} - \mu_{0})^{2}}{2(N+\beta)} + \frac{1}{2}\sum_{n=1}^{N}(x_{n} - \bar{x})^{2}.$$
(2.366)

# 2.45

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be variables in D dimensions such that

$$p(\mathbf{x}_n|\mathbf{\Lambda}) = \mathcal{N}\left(\mathbf{x}_n|\boldsymbol{\mu}, \mathbf{\Lambda}^{-1}\right),$$
  

$$p(\mathbf{\Lambda}) = \mathcal{W}(\mathbf{\Lambda}|\mathbf{W}, \nu),$$
(2.367)

where

$$W(\mathbf{\Lambda}|\mathbf{W},\nu) = B(\mathbf{W},\nu)|\det\mathbf{\Lambda}|^{\frac{\nu-D-1}{2}}\exp\left(-\frac{1}{2}\operatorname{tr}\left(\mathbf{W}^{-1}\mathbf{\Lambda}\right)\right). \tag{2.368}$$

By the Bayes' theorem,

$$p(\mathbf{\Lambda}|\mathbf{X})p(\mathbf{X}) = p(\mathbf{X}|\mathbf{\Lambda})p(\mathbf{\Lambda}). \tag{2.369}$$

The logarithm of the right hand side except the terms independent of  $\Lambda$  can be written as

$$-\frac{N}{2}\ln\left|\det\left(\mathbf{\Lambda}^{-1}\right)\right| - \frac{1}{2}\sum_{n=1}^{N}(\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Lambda}(\mathbf{x}_{n} - \boldsymbol{\mu})$$

$$+\frac{\nu - D - 1}{2}\ln\left|\det\mathbf{\Lambda}\right| - \frac{1}{2}\operatorname{tr}\left(\mathbf{W}^{-1}\boldsymbol{\Lambda}\right)$$

$$=\frac{\nu + N - D - 1}{2}\ln\left|\det\mathbf{\Lambda}\right| - \frac{1}{2}\operatorname{tr}\left(\left(\mathbf{W}^{-1} + \mathbf{S}\right)\boldsymbol{\Lambda}\right),$$
(2.370)

where

$$\mathbf{S} = \sum_{n=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^{\mathsf{T}}.$$
 (2.371)

Then,

$$p(\mathbf{\Lambda}|\mathbf{X}) = \mathcal{W}\left(\mathbf{\Lambda}|\left(\mathbf{W}^{-1} + \mathbf{S}\right)^{-1}, \nu + N\right).$$
 (2.372)

Therefore, W is a conjugate prior distribution of  $\Lambda$ .

# 2.46

Let x be a variable such that

$$p(x|\tau) = \mathcal{N}\left(x|\mu, \tau^{-1}\right),$$
  

$$p(\tau) = \operatorname{Gam}(\tau|a, b).$$
(2.373)

By marginalisation,

$$p(x) = \int_0^\infty p(x|\tau)p(\tau)d\tau. \tag{2.374}$$

The right hand side can be written as

$$\int_{0}^{\infty} (2\pi\tau^{-1})^{-\frac{1}{2}} \exp\left(-\frac{\tau}{2}(x-\mu)^{2}\right) \frac{b^{a}}{\Gamma(a)} \tau^{a-1} \exp(-b\tau) d\tau$$

$$= (2\pi)^{-\frac{1}{2}} \frac{b^{a}}{\Gamma(a)} \int_{0}^{\infty} \tau^{a-\frac{1}{2}} \exp\left(-\left(b + \frac{(x-\mu)^{2}}{2}\right)\tau\right) d\tau.$$
(2.375)

By the transformation

$$\tau' = \left(b + \frac{(x-\mu)^2}{2}\right)\tau,\tag{2.376}$$

the integral of the right hand side can be written as

$$\int_{0}^{\infty} \left( \frac{\tau'}{b + \frac{(x-\mu)^{2}}{2}} \right)^{a-\frac{1}{2}} \exp(-\tau') \frac{1}{b + \frac{(x-\mu)^{2}}{2}} d\tau'$$

$$= \Gamma \left( a + \frac{1}{2} \right) \left( b + \frac{(x-\mu)^{2}}{2} \right)^{-a-\frac{1}{2}}.$$
(2.377)

Then,

$$p(x) = (2\pi)^{-\frac{1}{2}} \frac{\Gamma(a + \frac{1}{2})}{\Gamma(a)} b^a \left( b + \frac{(x - \mu)^2}{2} \right)^{-a - \frac{1}{2}}.$$
 (2.378)

By the transformation

$$\nu = 2a,$$

$$\lambda = \frac{a}{b},$$
(2.379)

the right hand side can be written as

$$(2\pi)^{-\frac{1}{2}} \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})} \left(\frac{\nu}{2\lambda}\right)^{-\frac{1}{2}} \left(1 + \frac{(x-\mu)^2}{\frac{\nu}{\lambda}}\right)^{-\frac{\nu+1}{2}}$$

$$= \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})} \left(\frac{\lambda}{\pi\nu}\right)^{\frac{1}{2}} \left(1 + \frac{\lambda}{\nu}(x-\mu)^2\right)^{-\frac{\nu+1}{2}}.$$
(2.380)

Therefore,

$$p(x) = \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})} \left(\frac{\lambda}{\pi\nu}\right)^{\frac{1}{2}} \left(1 + \frac{\lambda}{\nu}(x-\mu)^2\right)^{-\frac{\nu+1}{2}}.$$
 (2.381)

# 2.47

Let

$$\operatorname{St}(x|\mu,\lambda,\nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})} \left(\frac{\lambda}{\pi\nu}\right)^{\frac{1}{2}} \left(1 + \frac{\lambda}{\nu}(x-\mu)^2\right)^{-\frac{\nu+1}{2}}.$$
 (2.382)

By the transformation

$$\frac{1}{y} = \frac{\lambda}{\nu} (x - \mu)^2,$$
 (2.383)

the right hand side except the terms independent of x can be written as

$$\left(1 + \frac{1}{y}\right)^{-\frac{\lambda(x-\mu)^2}{2}y - \frac{1}{2}}.$$
(2.384)

By the property

$$\lim_{x \to \infty} \left( 1 + \frac{1}{x} \right)^x = e,\tag{2.385}$$

we have

$$\lim_{y \to \infty} \left( 1 + \frac{1}{y} \right)^{-\frac{\lambda(x-\mu)^2}{2}y - \frac{1}{2}} = \exp\left( -\frac{\lambda}{2}(x-\mu)^2 \right). \tag{2.386}$$

$$\lim_{\nu \to \infty} \operatorname{St}(x|\mu, \lambda, \nu) = \mathcal{N}(x|\mu, \lambda^{-1}). \tag{2.387}$$

#### 2.48

Let  $\mathbf{x}$  be a variable in D dimensions such that

$$p(\mathbf{x}|\eta) = \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu}, (\eta\boldsymbol{\Lambda})^{-1}\right),$$
  

$$p(\eta) = \operatorname{Gam}\left(\eta \mid \frac{\nu}{2}, \frac{\nu}{2}\right).$$
(2.388)

By marginalisation,

$$p(\mathbf{x}) = \int_0^\infty p(\mathbf{x}|\eta)p(\eta)d\eta. \tag{2.389}$$

The right hand side can be written as

$$\int_{0}^{\infty} (2\pi)^{-\frac{D}{2}} \left| \det(\eta \mathbf{\Lambda})^{-1} \right|^{-\frac{1}{2}} \exp\left(-\frac{\eta}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Lambda} (\mathbf{x} - \boldsymbol{\mu})\right) \frac{\left(\frac{\nu}{2}\right)^{\frac{\nu}{2}}}{\Gamma(\frac{\nu}{2})} \eta^{\frac{\nu}{2} - 1} \exp\left(-\frac{\nu}{2} \eta\right) d\eta$$

$$= (2\pi)^{-\frac{D}{2}} \frac{\left(\frac{\nu}{2}\right)^{\frac{\nu}{2}}}{\Gamma(\frac{\nu}{2})} \left| \det \mathbf{\Lambda} \right|^{\frac{1}{2}} \int_{0}^{\infty} \eta^{\frac{D+\nu}{2} - 1} \exp\left(-\frac{\nu + (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Lambda} (\mathbf{x} - \boldsymbol{\mu})}{2} \eta\right) d\eta.$$
(2.390)

By the transformation

$$\eta' = \frac{\nu + (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})}{2} \eta, \tag{2.391}$$

the integral of the right hand side can be written as

$$\int_{0}^{\infty} \left( \frac{2\eta'}{\nu + (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})} \right)^{\frac{D+\nu}{2} - 1} \exp(-\eta') \frac{2}{\nu + (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})} d\eta'$$

$$= \left( \frac{2}{\nu + (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})} \right)^{\frac{D+\nu}{2}} \int_{0}^{\infty} \eta'^{\frac{D+\nu}{2} - 1} \exp(-\eta') d\eta'.$$
(2.392)

Then,

$$p(\mathbf{x}) = (2\pi)^{-\frac{D}{2}} \frac{\left(\frac{\nu}{2}\right)^{\frac{\nu}{2}}}{\Gamma(\frac{\nu}{2})} |\det \mathbf{\Lambda}|^{\frac{1}{2}} \left( \frac{2}{\nu + (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Lambda} (\mathbf{x} - \boldsymbol{\mu})} \right)^{\frac{D+\nu}{2}} \Gamma\left(\frac{D+\nu}{2}\right). \tag{2.393}$$

The right hand side can be written as

$$(2\pi)^{-\frac{D}{2}} \frac{\Gamma(\frac{D+\nu}{2})}{\Gamma(\frac{\nu}{2})} |\det \mathbf{\Lambda}|^{\frac{1}{2}} \left(\frac{\nu}{2}\right)^{-\frac{D}{2}} \left(\frac{\nu}{2}\right)^{\frac{D+\nu}{2}} \left(\frac{2}{\nu + (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Lambda} (\mathbf{x} - \boldsymbol{\mu})}\right)^{\frac{D+\nu}{2}}$$

$$= (2\pi)^{-\frac{D}{2}} \frac{\Gamma(\frac{D+\nu}{2})}{\Gamma(\frac{\nu}{2})} |\det \mathbf{\Lambda}|^{\frac{1}{2}} \left(\frac{\nu}{2}\right)^{-\frac{D}{2}} \left(1 + \frac{(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Lambda} (\mathbf{x} - \boldsymbol{\mu})}{\nu}\right)^{-\frac{D+\nu}{2}}.$$

$$(2.394)$$

Therefore,

$$p(\mathbf{x}) = \frac{\Gamma(\frac{D+\nu}{2})}{\Gamma(\frac{\nu}{2})} \frac{|\det \mathbf{\Lambda}|^{\frac{1}{2}}}{(\pi\nu)^{\frac{D}{2}}} \left(1 + \frac{(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Lambda} (\mathbf{x} - \boldsymbol{\mu})}{\nu}\right)^{-\frac{D+\nu}{2}}.$$
 (2.395)

# 2.49

Let  $\mathbf{x}$  be a variable such that

$$p(\mathbf{x}) = \text{St}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\nu}), \tag{2.396}$$

where

$$\operatorname{St}(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Lambda},\nu) = \int \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu},(\eta\boldsymbol{\Lambda})^{-1}\right) \operatorname{Gam}\left(\eta \mid \frac{\nu}{2}, \frac{\nu}{2}\right) d\eta. \tag{2.397}$$

(a)

We have

$$\mathbf{E}\,\mathbf{x} = \int \mathbf{x} \mathrm{St}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) d\mathbf{x}. \tag{2.398}$$

The right hand side can be written as

$$\int \mathbf{x} \left( \int \mathcal{N} \left( \mathbf{x} | \boldsymbol{\mu}, (\eta \boldsymbol{\Lambda})^{-1} \right) \operatorname{Gam} \left( \eta \mid \frac{\nu}{2}, \frac{\nu}{2} \right) d\eta \right) d\mathbf{x}$$

$$= \int \left( \int \mathbf{x} \mathcal{N} \left( \mathbf{x} | \boldsymbol{\mu}, (\eta \boldsymbol{\Lambda})^{-1} \right) d\mathbf{x} \right) \operatorname{Gam} \left( \eta \mid \frac{\nu}{2}, \frac{\nu}{2} \right) d\eta.$$
(2.399)

The right hand side can be written as

$$\boldsymbol{\mu} \int \operatorname{Gam}\left(\eta \mid \frac{\nu}{2}, \frac{\nu}{2}\right) d\eta = \boldsymbol{\mu}. \tag{2.400}$$

Therefore,

$$\mathbf{E}\,\mathbf{x} = \boldsymbol{\mu}.\tag{2.401}$$

(b)

By (a), we have

$$\operatorname{cov} \mathbf{x} = \int (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \operatorname{St}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) d\mathbf{x}. \tag{2.402}$$

The right hand side can be written as

$$\int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \left( \int \mathcal{N} \left( \mathbf{x} | \boldsymbol{\mu}, (\eta \boldsymbol{\Lambda})^{-1} \right) \operatorname{Gam} \left( \eta \mid \frac{\nu}{2}, \frac{\nu}{2} \right) d\eta \right) d\mathbf{x}$$

$$= \int \left( \int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathcal{N} \left( \mathbf{x} | \boldsymbol{\mu}, (\eta \boldsymbol{\Lambda})^{-1} \right) d\mathbf{x} \right) \operatorname{Gam} \left( \eta \mid \frac{\nu}{2}, \frac{\nu}{2} \right) d\eta.$$
(2.403)

The right hand side can be written as

$$\int (\eta \mathbf{\Lambda})^{-1} \operatorname{Gam}\left(\eta \mid \frac{\nu}{2}, \frac{\nu}{2}\right) d\eta = \frac{\left(\frac{\nu}{2}\right)^{\frac{\nu}{2}}}{\Gamma\left(\frac{\nu}{2}\right)} \left(\int \eta^{\frac{\nu}{2}-2} \exp\left(-\frac{\nu}{2}\eta\right) d\eta\right) \mathbf{\Lambda}^{-1}.$$
(2.404)

By the transformation

$$\eta' = \frac{\nu}{2}\eta,\tag{2.405}$$

the integral of the right hand side can be written as

$$\int \left(\frac{2}{\nu}\eta'\right)^{\frac{\nu}{2}-2} \exp(-\eta') \frac{2}{\nu} d\eta' = \left(\frac{2}{\nu}\right)^{\frac{\nu}{2}-1} \Gamma\left(\frac{\nu}{2}-1\right). \tag{2.406}$$

Then,

$$\operatorname{cov} \mathbf{x} = \frac{\left(\frac{\nu}{2}\right)^{\frac{\nu}{2}}}{\Gamma\left(\frac{\nu}{2}\right)} \left(\frac{2}{\nu}\right)^{\frac{\nu}{2}-1} \Gamma\left(\frac{\nu}{2}-1\right) \mathbf{\Lambda}^{-1}. \tag{2.407}$$

Therefore,

$$\operatorname{cov} \mathbf{x} = \frac{\nu}{\nu - 2} \mathbf{\Lambda}^{-1}. \tag{2.408}$$

(c)

Setting the derivative of  $p(\mathbf{x})$  to zero gives

$$\mathbf{0} = -\frac{1}{2} \left( \mathbf{\Lambda} + \mathbf{\Lambda}^{\mathsf{T}} \right) \left( \mathbf{x} - \boldsymbol{\mu} \right) \int \eta \mathcal{N} \left( \mathbf{x} | \boldsymbol{\mu}, (\eta \mathbf{\Lambda})^{-1} \right) \operatorname{Gam} \left( \eta \mid \frac{\nu}{2}, \frac{\nu}{2} \right) d\eta. \quad (2.409)$$

$$mode \mathbf{x} = \boldsymbol{\mu}. \tag{2.410}$$

# 2.50

Let

$$\operatorname{St}(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Lambda},\nu) = \frac{\Gamma(\frac{D+\nu}{2})}{\Gamma(\frac{\nu}{2})} \frac{(\det \boldsymbol{\Lambda})^{\frac{1}{2}}}{(\pi\nu)^{\frac{D}{2}}} \left(1 + \frac{(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Lambda}(\mathbf{x}-\boldsymbol{\mu})}{\nu}\right)^{-\frac{D+\nu}{2}}. \quad (2.411)$$

By the transformation

$$y = \frac{\nu}{(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})},$$
 (2.412)

the right hand side except the terms independent of x can be written as

$$\left(1 + \frac{1}{y}\right)^{-\frac{(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})}{2} y - \frac{D}{2}}$$
(2.413)

By the property

$$\lim_{x \to \infty} \left( 1 + \frac{1}{x} \right)^x = e,\tag{2.414}$$

we have

$$\lim_{y \to \infty} \left( 1 + \frac{1}{y} \right)^{-\frac{(\mathbf{x} - \boldsymbol{\mu})^\mathsf{T} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})}{2} y - \frac{D}{2}} = \exp\left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\mathsf{T} \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}) \right). \quad (2.415)$$

Therefore,

$$\lim_{\nu \to \infty} \operatorname{St}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}). \tag{2.416}$$

# 2.51

(a)

We have

$$\exp(iA)\exp(-iA) = 1. \tag{2.417}$$

The left hand side can be written as

$$(\cos A + i\sin A)(\cos A - i\sin A) = \cos^2 A + \sin^2 A. \tag{2.418}$$

$$\cos^2 A + \sin^2 A = 1. \tag{2.419}$$

(b)

We have

$$\cos(A - B) = \operatorname{Re}\left(\exp\left(i(A - B)\right)\right). \tag{2.420}$$

The right hand side can be written as

$$\operatorname{Re}\left(\exp(iA)\exp(-iB)\right) = \operatorname{Re}\left((\cos A + i\sin A)(\cos B - i\sin B)\right). \quad (2.421)$$

Therefore,

$$\cos(A - B) = \cos A \cos B + \sin A \sin B. \tag{2.422}$$

(c)

We have

$$\sin(A - B) = \text{Im} (\exp(i(A - B))).$$
 (2.423)

The right hand side can be written as

$$\operatorname{Im}\left(\exp(iA)\exp(-iB)\right) = \left((\cos A + i\sin A)(\cos B - i\sin B)\right). \tag{2.424}$$

Therefore,

$$\sin(A - B) = \sin A \cos B - \cos A \sin B. \tag{2.425}$$

# 2.52

Let

$$f(\theta|\theta_0, m) = \frac{1}{2\pi I_0(m)} \exp(m\cos(\theta - \theta_0)),$$
 (2.426)

where

$$I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp(m\cos\theta) d\theta.$$
 (2.427)

By the Taylor series

$$\cos \alpha = 1 - \frac{1}{2}\alpha^2 + O\left(\alpha^4\right), \qquad (2.428)$$

the right hand side can be written as

$$\frac{\exp\left(m\left(1 - \frac{1}{2}(\theta - \theta_0)^2 + O\left((\theta - \theta_0)^4\right)\right)\right)}{\int_0^{2\pi} \exp\left(m\left(1 - \frac{1}{2}\theta^2 + O(\theta^4)\right)\right) d\theta} 
= \exp\left(-\frac{m}{2}(\theta - \theta_0)^2\right) \frac{\exp\left(mO\left((\theta - \theta_0)^4\right)\right)}{\int_0^{2\pi} \exp\left(m\left(-\frac{1}{2}\theta^2 + O(\theta^4)\right)\right) d\theta}.$$
(2.429)

$$\lim_{m \to \infty} f(\theta | \theta_0, m) = \mathcal{N}\left(\theta | \theta_0, m^{-1}\right). \tag{2.430}$$

2.53

Let

$$\sum_{n=1}^{N} \sin(\theta_n - \theta_0) = 0. \tag{2.431}$$

The left hand side can be written as

$$\sum_{n=1}^{N} (\sin \theta_n \cos \theta_0 - \cos \theta_n \sin \theta_0) = \cos \theta_0 \sum_{n=1}^{N} \sin \theta_n - \sin \theta_0 \sum_{n=1}^{N} \cos \theta_n. \quad (2.432)$$

Therefore,

$$\theta_0 = \arctan\left(\frac{\sum_{n=1}^N \sin \theta_n}{\sum_{n=1}^N \cos \theta_n}\right). \tag{2.433}$$

# 2.54

Let

$$f(\theta|\theta_0, m) = \frac{1}{2\pi I_0(m)} \exp(m\cos(\theta - \theta_0)),$$
 (2.434)

where

$$I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp(m\cos\theta) d\theta. \tag{2.435}$$

Setting the first and second derivatives with respect to  $\theta$  to zero gives

$$0 = -m\sin(\theta - \theta_0)f(\theta|\theta_0, m),$$
  

$$0 = (m^2\sin^2(\theta - \theta_0) - m\cos(\theta - \theta_0))f(\theta|\theta_0, m).$$
(2.436)

Therefore,

$$\underset{\theta}{\operatorname{argmax}} f(\theta|\theta_0, m) = \theta_0,$$

$$\underset{\theta}{\operatorname{argmin}} f(\theta|\theta_0, m) = \theta_0 - \pi \operatorname{sgn}(\theta_0 - \pi).$$
(2.437)

#### 2.55

(a)

Let  $\theta_1, \dots, \theta_N$  be variables such that

$$p(\theta_n) = f(\theta_n | \theta_0, m), \tag{2.438}$$

where

$$f(\theta|\theta_0, m) = \frac{1}{2\pi I_0(m)} \exp\left(m\cos(\theta - \theta_0)\right),$$

$$I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp(m\cos\theta) d\theta.$$
(2.439)

Then,

$$\ln\left(\prod_{n=1}^{N} p(\theta_n | \theta_0, m)\right) = -\frac{N}{2} \ln\left(2\pi I_0(m)\right) + m \sum_{n=1}^{N} \cos(\theta_n - \theta_0).$$
 (2.440)

Setting the derivative with respect to  $\theta_0$  to zero gives

$$0 = m \sum_{n=1}^{N} \sin(\theta_n - \theta_0). \tag{2.441}$$

Therefore, by 2.53, the maximum likelihood solution for  $\theta_0$  is given by

$$\theta_0^{\text{ML}} = \arctan\left(\frac{\sum_{n=1}^N \sin \theta_n}{\sum_{n=1}^N \cos \theta_n}\right). \tag{2.442}$$

(b)

Let

$$\bar{r}\cos\bar{\theta} = \frac{1}{N} \sum_{n=1}^{N} \cos\theta_n,$$

$$\bar{r}\sin\bar{\theta} = \frac{1}{N} \sum_{n=1}^{N} \sin\theta_n.$$
(2.443)

By (a), 
$$\bar{\theta} = \theta_0^{\text{ML}}. \tag{2.444}$$

We have

$$\frac{1}{N} \sum_{n=1}^{N} \cos\left(\theta_n - \theta_0^{\text{ML}}\right) = \left(\frac{1}{N} \sum_{n=1}^{N} \cos\theta_n\right) \cos\theta_0^{\text{ML}} + \left(\frac{1}{N} \sum_{n=1}^{N} \sin\theta_n\right) \sin\theta_0^{\text{ML}}.$$
(2.445)

The right hand side can be written as

$$\bar{r}\cos^2\bar{\theta} + \bar{r}\sin^2\bar{\theta} = \bar{r}.$$
 (2.446)

Therefore,

$$\frac{1}{N} \sum_{n=1}^{N} \cos\left(\theta_n - \theta_0^{\mathrm{ML}}\right) = \bar{r}.$$
 (2.447)

2.56

(a)

Let

Beta
$$(\mu|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}.$$
 (2.448)

The right hand side can be written as

$$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \exp((a-1)\ln\mu + (b-1)\ln(1-\mu))$$
 (2.449)

Therefore, the natural parameters are given by

$$\boldsymbol{\eta} = \begin{bmatrix} a-1 \\ b-1 \end{bmatrix}.$$

(b)

Let

$$Gam(\lambda|a,b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda). \tag{2.450}$$

The right hand side can be written as

$$\frac{b^a}{\Gamma(a)} \exp\left((a-1)\ln\lambda - b\lambda\right). \tag{2.451}$$

Therefore, the natural parameters are given by

$$\boldsymbol{\eta} = \begin{bmatrix} a-1 \\ -b \end{bmatrix}.$$

(c)

Let

$$f(\theta|\theta_0, m) = \frac{1}{2\pi I_0(m)} \exp\left(m\cos(\theta - \theta_0)\right), \qquad (2.452)$$

where

$$I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp(m\cos\theta) d\theta, \qquad (2.453)$$

the right hand side can be written as

$$\frac{1}{2\pi I_0(m)} \exp(m\cos\theta_0\cos\theta + m\sin\theta_0\sin\theta). \tag{2.454}$$

Therefore, the natural parameters are given by

$$\boldsymbol{\eta} = \begin{bmatrix} m\cos\theta_0 \\ m\sin\theta_0 \end{bmatrix}.$$

# 2.57

By the definition,

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} (\det \boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right). \quad (2.455)$$

Therefore,

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp\left(\boldsymbol{\eta}^{\mathsf{T}}\mathbf{u}(\mathbf{x})\right), \tag{2.456}$$

where

$$h(\mathbf{x}) = (2\pi)^{-\frac{D}{2}},$$

$$g(\boldsymbol{\eta}) = (\det(-2\boldsymbol{\eta}_2))^{-\frac{1}{2}} \exp\left(\frac{1}{4}\boldsymbol{\eta}_1^{\mathsf{T}}\boldsymbol{\eta}_2^{-1}\boldsymbol{\eta}_1\right),$$

$$\boldsymbol{\eta} = \begin{bmatrix} \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ -\frac{1}{2}\boldsymbol{\Sigma}^{-1} \end{bmatrix},$$

$$\mathbf{u}(\mathbf{x}) = \begin{bmatrix} \mathbf{x} \\ \mathbf{x}\mathbf{x}^{\mathsf{T}} \end{bmatrix}.$$

# 2.58

Let  $\mathbf{x}$  be a variable such that

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp\left(\boldsymbol{\eta}^{\mathsf{T}}\mathbf{u}(\mathbf{x})\right). \tag{2.457}$$

Then, taking the first derivative of

$$\int p(\mathbf{x}|\boldsymbol{\eta})d\mathbf{x} = 1 \tag{2.458}$$

with respect to  $\eta$  gives

$$\nabla g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp(\boldsymbol{\eta}^{\mathsf{T}} \mathbf{u}(\mathbf{x})) d\mathbf{x} + g(\boldsymbol{\eta}) \int \mathbf{u}(\mathbf{x}) h(\mathbf{x}) \exp(\boldsymbol{\eta}^{\mathsf{T}} \mathbf{u}(\mathbf{x})) d\mathbf{x} = \mathbf{0}.$$
(2.459)

The left hand side can be written as

$$\frac{\nabla g(\boldsymbol{\eta})}{g(\boldsymbol{\eta})} \int p(\mathbf{x}|\boldsymbol{\eta}) d\mathbf{x} + \int \mathbf{u}(\mathbf{x}) p(\mathbf{x}|\boldsymbol{\eta}) d\mathbf{x} = \frac{\nabla g(\boldsymbol{\eta})}{g(\boldsymbol{\eta})} + \mathbf{E} \mathbf{u}(\mathbf{x}).$$
(2.460)

Therefore,

$$\mathbf{E}\,\mathbf{u}(\mathbf{x}) = -\frac{\nabla g(\boldsymbol{\eta})}{g(\boldsymbol{\eta})}.\tag{2.461}$$

Thus,

$$\mathbf{E}\,\mathbf{u}(\mathbf{x}) = -\nabla \ln g(\boldsymbol{\eta}). \tag{2.462}$$

Taking the second derivative with respect to  $\eta$  gives

$$\nabla \nabla g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp(\boldsymbol{\eta}^{\mathsf{T}} \mathbf{u}(\mathbf{x})) d\mathbf{x} + 2\nabla g(\boldsymbol{\eta}) \int \mathbf{u}(\mathbf{x})^{\mathsf{T}} h(\mathbf{x}) \exp(\boldsymbol{\eta}^{\mathsf{T}} \mathbf{u}(\mathbf{x})) d\mathbf{x}$$
$$+ g(\boldsymbol{\eta}) \int \mathbf{u}(\mathbf{x}) \mathbf{u}(\mathbf{x})^{\mathsf{T}} h(\mathbf{x}) \exp(\boldsymbol{\eta}^{\mathsf{T}} \mathbf{u}(\mathbf{x})) d\mathbf{x} = \mathbf{O}.$$
(2.463)

The left hand side can be written as

$$\frac{\nabla \nabla g(\boldsymbol{\eta})}{g(\boldsymbol{\eta})} \int p(\mathbf{x}|\boldsymbol{\eta}) d\mathbf{x} + \frac{2\nabla g(\boldsymbol{\eta})}{g(\boldsymbol{\eta})} \int \mathbf{u}(\mathbf{x})^{\mathsf{T}} p(\mathbf{x}|\boldsymbol{\eta}) d\mathbf{x} + \int \mathbf{u}(\mathbf{x}) \mathbf{u}(\mathbf{x})^{\mathsf{T}} p(\mathbf{x}|\boldsymbol{\eta}) d\mathbf{x} 
= \frac{\nabla \nabla g(\boldsymbol{\eta})}{g(\boldsymbol{\eta})} - 2 \operatorname{E} \mathbf{u}(\mathbf{x}) \operatorname{E} \mathbf{u}(\mathbf{x})^{\mathsf{T}} + \operatorname{E} (\mathbf{u}(\mathbf{x}) \mathbf{u}(\mathbf{x})^{\mathsf{T}}).$$
(2.464)

Therefore,

$$E(\mathbf{u}(\mathbf{x})\mathbf{u}(\mathbf{x})^{\mathsf{T}}) = -\frac{\nabla \nabla g(\boldsymbol{\eta})}{g(\boldsymbol{\eta})} + \frac{2\nabla g(\boldsymbol{\eta})(\nabla g(\boldsymbol{\eta}))^{\mathsf{T}}}{g^2(\boldsymbol{\eta})}.$$
 (2.465)

By the definition,

$$\operatorname{cov} \mathbf{u}(\mathbf{x}) = \operatorname{E} (\mathbf{u}(\mathbf{x})\mathbf{u}(\mathbf{x})^{\mathsf{T}}) - \operatorname{E} \mathbf{u}(\mathbf{x})\operatorname{E} \mathbf{u}(\mathbf{x})^{\mathsf{T}}. \tag{2.466}$$

Thus,

$$\operatorname{cov} \mathbf{u}(\mathbf{x}) = -\frac{\nabla \nabla g(\boldsymbol{\eta})}{g(\boldsymbol{\eta})} + \frac{\nabla g(\boldsymbol{\eta})(\nabla g(\boldsymbol{\eta}))^{\mathsf{T}}}{g^2(\boldsymbol{\eta})}.$$
 (2.467)

Hence,

$$\operatorname{cov} \mathbf{u}(\mathbf{x}) = -\nabla \nabla \ln g(\boldsymbol{\eta}). \tag{2.468}$$

# 2.59

Let

$$p(x|\sigma) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right). \tag{2.469}$$

Then

$$\int p(x|\sigma)dx = \frac{1}{\sigma} \int f\left(\frac{x}{\sigma}\right) dx. \tag{2.470}$$

By the transformation

$$x' = \frac{x}{\sigma},\tag{2.471}$$

the right hand side can be written as

$$\frac{1}{\sigma} \int f(x')\sigma dx' = \int f(x')dx'. \tag{2.472}$$

Therefore,  $p(x|\sigma)$  will be normalised if f(x) is normalised.

# 2.60

Let  $\mathbf{x}$  be a variable such that

$$\mathbf{x} \in \mathcal{R}_i \Rightarrow p(\mathbf{x}) = h_i,$$
 (2.473)

where

$$\int_{\mathcal{R}_i} d\mathbf{x} = \Delta_i. \tag{2.474}$$

Since

$$\int p(\mathbf{x})d\mathbf{x} = 1, \tag{2.475}$$

we have

$$\sum_{i} h_i \Delta_i = 1. \tag{2.476}$$

Let N be the total number of observations and  $n_i$  be the number of observations which fall in  $\mathcal{R}_i$ . Then, the logarithm of the likelihood is given by

$$\ln\left(\prod_{i} h_i^{n_i}\right) = \sum_{i} n_i \ln h_i, \tag{2.477}$$

where

$$\sum_{i} n_i = N. \tag{2.478}$$

Setting the derivatives of

$$\sum_{i} n_{i} \ln h_{i} + \lambda \left( \sum_{i} h_{i} \Delta_{i} - 1 \right) \tag{2.479}$$

with respect to  $h_i$  and  $\lambda$  to zero gives

$$\frac{n_i}{h_i} + \lambda \Delta_i = 0,$$

$$\sum_i h_i \Delta_i - 1 = 0.$$
(2.480)

Then,

$$\lambda = -N,$$

$$h_i = \frac{n_i}{N\Delta_i}.$$
(2.481)

Therefore, the maximum likelihood estimator for the  $\{h_i\}$  is  $\frac{n_i}{N\Delta_i}$ .

# 2.61 (Incomplete)

Let  $\mathbf{x}$  be a variable and  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be observations. Let

$$p(\mathbf{x}) = \frac{K}{NV(\mathbf{x})},\tag{2.482}$$

where

$$V(\mathbf{x}) = \int_{\|\mathbf{x}' - \mathbf{x}\| \le \|\mathbf{x}_{(K)} - \mathbf{x}\|} d\mathbf{x}', \qquad (2.483)$$

K is a constant and  $\mathbf{x}_{(K)}$  is the Kth nearest observation from the point  $\mathbf{x}$ .

# 3 Linear Models for Regression

# 3.1

By the definition,

$$tanh a = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}.$$
(3.1)

The right hand side can be written as

$$\frac{1 - \exp(-2a)}{1 + \exp(-2a)} = \frac{2}{1 + \exp(-2a)} - 1. \tag{3.2}$$

Therefore,

$$tanh a = 2\sigma(2a) - 1,$$
(3.3)

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}. (3.4)$$

Let

$$y(x_n, \mathbf{w}) = w_0 + \sum_{m=1}^{M} w_j \sigma\left(\frac{x - \mu_j}{s}\right). \tag{3.5}$$

By the result above, the right hand side can be written as

$$w_0 + \sum_{m=1}^{M} w_m \frac{1 + \tanh\left(\frac{x - \mu_m}{2s}\right)}{2} = w_0 + \frac{1}{2} \sum_{m=1}^{M} w_m + \frac{1}{2} \sum_{m=1}^{M} w_m \tanh\left(\frac{x - \mu_m}{2s}\right).$$
(3.6)

Therefore,  $y(x_n, \mathbf{w})$  is equivalent to

$$y(x_n, \mathbf{u}) = u_0 + \sum_{m=1}^{M} u_m \tanh\left(\frac{x - \mu_m}{2s}\right), \tag{3.7}$$

where

$$u_0 = w_0 + \frac{1}{2} \sum_{m=1}^{M} w_m,$$

$$u_m = \frac{1}{2} w_m.$$
(3.8)

## 3.2 (Incomplete)

Let  $\Phi$  be an  $N \times M$  matarix. Then, for any vector  $\mathbf{v}$  in N dimensions,

$$\mathbf{\Phi} \left(\mathbf{\Phi}^{\mathsf{T}}\mathbf{\Phi}\right)^{-1} \mathbf{\Phi}^{\mathsf{T}}\mathbf{v} \tag{3.9}$$

is a projection of  $\mathbf{v}$  onto the space spanned by the columns of  $\mathbf{\Phi}$ ? Additionally, for a vector  $\mathbf{t}$  in N dimensions,

$$(\mathbf{\Phi}^{\mathsf{T}}\mathbf{\Phi})^{-1}\mathbf{\Phi}^{\mathsf{T}}\mathbf{t} \tag{3.10}$$

is an orthogonal projection of  ${\bf t}$  onto the space spanned by the columns of  ${\bf \Phi}$ ?

## 3.3

Let

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} r_n \left( t_n - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n) \right)^2.$$
 (3.11)

The right hand side can be written as

$$\frac{1}{2} \|\mathbf{t}' - \mathbf{\Phi}' \mathbf{w}\|^2, \tag{3.12}$$

where

$$\mathbf{t}' = egin{bmatrix} \sqrt{r_1}t_1 \ dots \ \sqrt{r_N}t_N \end{bmatrix}, \mathbf{\Phi}' = egin{bmatrix} \sqrt{r_1}oldsymbol{\phi}(\mathbf{x}_1)^\intercal \ dots \ \sqrt{r_N}oldsymbol{\phi}(\mathbf{x}_N)^\intercal \end{bmatrix}.$$

Setting the derivative with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = -\mathbf{\Phi}^{\prime\mathsf{T}}(\mathbf{t}^{\prime} - \mathbf{\Phi}^{\prime}\mathbf{w}). \tag{3.13}$$

Therefore,

$$\underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w}) = \left(\mathbf{\Phi}^{\prime \mathsf{T}} \mathbf{\Phi}^{\prime}\right)^{-1} \mathbf{\Phi}^{\prime \mathsf{T}} \mathbf{t}^{\prime}. \tag{3.14}$$

# 3.4 (Incomplete)

Let

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \left( y(\mathbf{x}_n, \mathbf{w}) - t_n \right)^2, \tag{3.15}$$

where

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{m=1}^{M} w_m(x_m + \epsilon_m),$$
  

$$p(\epsilon_m) = \mathcal{N}\left(\epsilon_m | 0, \sigma^2\right).$$
(3.16)

Setting the derivative with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = \sum_{n=1}^{N} \begin{bmatrix} 1 \\ \mathbf{x}_n + \boldsymbol{\epsilon}_n \end{bmatrix} (y(\mathbf{x}_n, \mathbf{w}) - t_n).$$

The right hand side can be written as

$$\sum_{n=1}^{N} \begin{bmatrix} 1 \\ \mathbf{x}_n \end{bmatrix} (y(\mathbf{x}_n, \mathbf{w}) - t_n) + \sum_{n=1}^{N} \begin{bmatrix} 0 \\ \boldsymbol{\epsilon}_n \end{bmatrix} (y(\mathbf{x}_n, \mathbf{w}) - t_n).$$

#### 3.5

Let

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} (t_n - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n))^2.$$
 (3.17)

Then, the minimisation of  $E(\mathbf{w})$  under the constraint

$$\sum_{m=1}^{M} \left| w_m \right|^q \le \eta \tag{3.18}$$

reduces to the minimisation of

$$E(\mathbf{w}) + \lambda \left( \sum_{m=1}^{M} |w_m|^q - \eta \right) \tag{3.19}$$

with respect to  $\mathbf{w}$  and  $\lambda$ . Then,

$$\eta = \sum_{m=1}^{M} |w_m^*(\lambda)|^q, \tag{3.20}$$

where

$$\mathbf{w}^*(\lambda) = \underset{\mathbf{w}}{\operatorname{argmin}} \left( E(\mathbf{w}) + \lambda \left( \sum_{m=1}^M |w_m|^q - \eta \right) \right). \tag{3.21}$$

#### 3.6

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables in D dimensions such that

$$p(\mathbf{t}_n|\mathbf{y}_n) = \mathcal{N}\left(\mathbf{t}_n|\mathbf{y}_n, \mathbf{\Sigma}\right), \tag{3.22}$$

where

$$\mathbf{y}_n = \mathbf{W}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n). \tag{3.23}$$

Then,

$$\ln p(\mathbf{T}|\mathbf{Y}) = -\frac{ND}{2}\ln(2\pi) - \frac{N}{2}\ln(\det \Sigma)$$

$$-\frac{1}{2}\sum_{n=1}^{N} (\mathbf{t}_n - \mathbf{W}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n))^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{W}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n)).$$
(3.24)

By 3.21(a), setting the derivatives with respect to **W** and  $\Sigma$  to zero gives

$$\mathbf{O} = -\frac{1}{2} \left( \mathbf{\Sigma}^{-1} + \left( \mathbf{\Sigma}^{-1} \right)^{\mathsf{T}} \right) \sum_{n=1}^{N} \left( \mathbf{t}_{n} - \mathbf{W}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_{n}) \right) \left( \boldsymbol{\phi}(\mathbf{x}_{n}) \right)^{\mathsf{T}},$$

$$\mathbf{O} = -\frac{N}{2} \left( \mathbf{\Sigma}^{-1} \right)^{\mathsf{T}} + \frac{1}{2} \left( \mathbf{\Sigma}^{-1} \right)^{2} \sum_{n=1}^{N} \left( \mathbf{t}_{n} - \mathbf{W}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_{n}) \right) \left( \mathbf{t}_{n} - \mathbf{W}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_{n}) \right)^{\mathsf{T}}.$$

$$(3.25)$$

Therefore, the maximum likelihood solutions for W and  $\Sigma$  are given by

$$\mathbf{W}_{\mathrm{ML}} = (\mathbf{\Phi}^{\mathsf{T}}\mathbf{\Phi})^{-1} \mathbf{\Phi}^{\mathsf{T}}\mathbf{t},$$

$$\mathbf{\Sigma}_{\mathrm{ML}} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{t}_{n} - \mathbf{W}_{\mathrm{ML}}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_{n})) (\mathbf{t}_{n} - \mathbf{W}_{\mathrm{ML}}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_{n}))^{\mathsf{T}}.$$
(3.26)

### 3.7

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}_n, \beta^{-1}\right),$$
  

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0).$$
(3.27)

By the Bayes' theorem,

$$p(\mathbf{w}|\mathbf{t})p(\mathbf{t}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}). \tag{3.28}$$

The logarithm of the right hand side except the terms independent of  ${\bf t}$  and  ${\bf w}$  can be written as

$$-\frac{\beta}{2} \sum_{n=1}^{N} (t_n - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n)^2 - \frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^{\mathsf{T}} \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0)$$

$$= -\frac{1}{2} \begin{bmatrix} \mathbf{w} \\ \mathbf{t} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} & -\beta \mathbf{\Phi}^{\mathsf{T}} \\ -\beta \mathbf{\Phi} & \beta \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{t} \end{bmatrix} + \begin{bmatrix} \mathbf{w} \\ \mathbf{t} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{S}_0^{-1} \mathbf{m}_0 \\ \mathbf{0} \end{bmatrix}.$$
(3.29)

Therefore,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S}), \tag{3.30}$$

where

$$\mathbf{m} = \mathbf{S} \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{t} \right), \mathbf{S} = \left( \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} \right)^{-1}.$$
 (3.31)

#### 3.8

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}_n, \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0).$$
 (3.32)

By 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S}), \tag{3.33}$$

where

$$\mathbf{m} = \mathbf{S} \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{t} \right), \mathbf{S} = \left( \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} \right)^{-1}.$$
 (3.34)

By the Bayes' theorem,

$$p(\mathbf{w}|\mathbf{t}')p(\mathbf{t}') = p(\mathbf{t}'|\mathbf{w})p(\mathbf{w}), \tag{3.35}$$

where

$$\mathbf{t}' = \begin{bmatrix} \mathbf{t} \\ t_{N+1} \end{bmatrix}. \tag{3.36}$$

Since  $\mathbf{t}$  and  $t_{N+1}$  are independent, it can be written as

$$p(\mathbf{w}|\mathbf{t}')p(t_{N+1})p(\mathbf{t}) = p(t_{N+1}|\mathbf{w})p(\mathbf{t}|\mathbf{w})p(\mathbf{w}). \tag{3.37}$$

By the Bayes' theorem, the right hand side can be written as

$$p(t_{N+1}|\mathbf{w})p(\mathbf{w}|\mathbf{t})p(\mathbf{t}). \tag{3.38}$$

Then,

$$p(\mathbf{w}|\mathbf{t}')p(t_{N+1}) = p(\mathbf{w}|\mathbf{t})p(t_{N+1}|\mathbf{w}). \tag{3.39}$$

The logarithm of the right hand side except the terms independent of  $\mathbf{w}$  can be written as

$$-\frac{1}{2} (\mathbf{w} - \mathbf{m})^{\mathsf{T}} \mathbf{S}^{-1} (\mathbf{w} - \mathbf{m}) - \frac{\beta}{2} (t_{N+1} - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_{N+1})^{2}$$

$$= -\frac{1}{2} \begin{bmatrix} \mathbf{w} \\ t_{N+1} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{S}^{-1} + \beta \boldsymbol{\phi}_{N+1} \boldsymbol{\phi}_{N+1}^{\mathsf{T}} & -\beta \boldsymbol{\phi}_{N+1} \\ -\beta \boldsymbol{\phi}_{N+1}^{\mathsf{T}} & \beta \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ t_{N+1} \end{bmatrix}$$

$$+ \begin{bmatrix} \mathbf{w} \\ t_{N+1} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{S}^{-1} \mathbf{m} \\ 0 \end{bmatrix} - \frac{1}{2} \mathbf{m}^{\mathsf{T}} \mathbf{S}^{-1} \mathbf{m}.$$
(3.40)

Then,

$$p(\mathbf{w}|\mathbf{t}') = \mathcal{N}(\mathbf{w}|\mathbf{m}', \mathbf{S}'), \qquad (3.41)$$

where

$$\mathbf{m}' = \mathbf{S}' \left( \mathbf{S}^{-1} \mathbf{m} + \beta t_{N+1} \boldsymbol{\phi}_{N+1} \right),$$
  
$$\mathbf{S}' = \left( \mathbf{S}^{-1} + \beta \boldsymbol{\phi}_{N+1} \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \right)^{-1}.$$
 (3.42)

We have

$$\mathbf{S}^{-1}\mathbf{m} + \beta t_{N+1} \boldsymbol{\phi}_{N+1} = \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \boldsymbol{\Phi}^{\prime \mathsf{T}} \mathbf{t}^{\prime}, \tag{3.43}$$

and

$$\mathbf{S}^{-1} + \beta \phi_{N+1} \phi_{N+1}^{\mathsf{T}} = \mathbf{S}_0^{-1} + \beta \Phi'^{\mathsf{T}} \Phi', \tag{3.44}$$

where

$$\mathbf{\Phi}' = \begin{bmatrix} \mathbf{\Phi} \\ \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \end{bmatrix}. \tag{3.45}$$

Therefore,

$$\mathbf{m}' = \mathbf{S}' \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{\Phi}'^{\mathsf{T}} \mathbf{t}' \right), \mathbf{S}' = \left( \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}'^{\mathsf{T}} \mathbf{\Phi}' \right)^{-1}.$$
 (3.46)

#### 3.9

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}_n, \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0).$$
 (3.47)

Then, by 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S}), \tag{3.48}$$

where

$$\mathbf{m} = \mathbf{S} \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{t} \right),$$

$$\mathbf{S} = \left( \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} \right)^{-1}.$$
(3.49)

Let

$$\mathbf{t}' = \begin{bmatrix} \mathbf{t} \\ t_{N+1} \end{bmatrix}, \mathbf{\Phi} = \begin{bmatrix} \mathbf{\Phi} \\ \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \end{bmatrix}. \tag{3.50}$$

Therefore,

$$p(\mathbf{w}|\mathbf{t}') = \mathcal{N}(\mathbf{w}|\mathbf{m}', \mathbf{t}'),$$
 (3.51)

where

$$\mathbf{m}' = \mathbf{S}' \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{\Phi}'^{\mathsf{T}} \mathbf{t}' \right),$$
  
$$\mathbf{S}' = \left( \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}'^{\mathsf{T}} \mathbf{\Phi}' \right)^{-1}.$$
 (3.52)

#### 3.10

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}_n, \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0\right).$$
 (3.53)

By 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S}),$$
 (3.54)

where

$$\mathbf{m} = \mathbf{S} \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{t} \right),$$
  
$$\mathbf{S} = \left( \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} \right)^{-1}.$$
 (3.55)

By marginalisation,

$$p(t_{N+1}|\mathbf{t}) = \int p(t_{N+1}|\mathbf{w})p(\mathbf{w}|\mathbf{t})d\mathbf{w}.$$
 (3.56)

The logarithm of the integrand of the right hand side except the terms independent of  $t_{N+1}$  and **w** can be written as

$$-\frac{\beta}{2} (t_{N+1} - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_{N+1})^{2} - \frac{1}{2} (\mathbf{w} - \mathbf{m})^{\mathsf{T}} \mathbf{S}^{-1} (\mathbf{w} - \mathbf{m})$$

$$= -\frac{1}{2} \begin{bmatrix} \mathbf{w} \\ t_{N+1} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{S}^{-1} + \beta \boldsymbol{\phi}_{N+1} \boldsymbol{\phi}_{N+1}^{\mathsf{T}} & -\beta \boldsymbol{\phi}_{N+1} \\ -\beta \boldsymbol{\phi}_{N+1}^{\mathsf{T}} & \beta \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ t_{N+1} \end{bmatrix}$$

$$+ \begin{bmatrix} \mathbf{w} \\ t_{N+1} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{S}^{-1} \mathbf{m} \\ 0 \end{bmatrix} - \frac{1}{2} \mathbf{m}^{\mathsf{T}} \mathbf{S}^{-1} \mathbf{m}.$$
(3.57)

By 2.24,

$$\begin{bmatrix} \mathbf{S}^{-1} + \beta \boldsymbol{\phi}_{N+1} \boldsymbol{\phi}_{N+1}^{\mathsf{T}} & -\beta \boldsymbol{\phi}_{N+1} \\ -\beta \boldsymbol{\phi}_{N+1}^{\mathsf{T}} & \beta \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{S} & \mathbf{S} \boldsymbol{\phi}_{N+1} \\ \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \mathbf{S} & \beta^{-1} + \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \mathbf{S} \boldsymbol{\phi}_{N+1} \end{bmatrix}.$$
(3.58)

Then,

$$\begin{bmatrix} \mathbf{S}^{-1} + \beta \boldsymbol{\phi}_{N+1} \boldsymbol{\phi}_{N+1}^{\mathsf{T}} & \beta \boldsymbol{\phi}_{N+1} \\ \beta \boldsymbol{\phi}_{N+1}^{\mathsf{T}} & \beta \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{S}^{-1} \mathbf{m} \\ 0 \end{bmatrix} = \begin{bmatrix} \mathbf{m} \\ \mathbf{m}^{\mathsf{T}} \boldsymbol{\phi}_{N+1} \end{bmatrix}.$$
(3.59)

Therefore,

$$p(t_{N+1}|\mathbf{t}) = \mathcal{N}\left(t_{N+1}|\mathbf{m}^{\mathsf{T}}\boldsymbol{\phi}_{N+1}, \sigma^{2}\right), \tag{3.60}$$

where

$$\sigma^2 = \beta^{-1} + \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \mathbf{S} \boldsymbol{\phi}_{N+1}. \tag{3.61}$$

## 3.11

Let  $t_1, \dots, t_N$  be a variable such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}_n, \beta^{-1}\right),$$
  

$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0\right).$$
(3.62)

By 3.10,

$$p(t_{N+1}|\mathbf{t}) = \mathcal{N}\left(t_{N+1}|\mathbf{m}^{\mathsf{T}}\boldsymbol{\phi}_{N+1}, \sigma^{2}\right), \tag{3.63}$$

where

$$\mathbf{m} = \mathbf{S} \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{t} \right),$$

$$\mathbf{S} = \left( \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} \right)^{-1},$$

$$\sigma^2 = \beta^{-1} + \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \mathbf{S} \boldsymbol{\phi}_{N+1}.$$
(3.64)

Then,

$$\sigma^2 - {\sigma'}^2 = \phi^{\mathsf{T}} (\mathbf{S} - \mathbf{S}') \phi. \tag{3.65}$$

We have

$$\mathbf{S}' = \left(\mathbf{S}^{-1} + \beta \phi_{N+1} \phi_{N+1}^{\mathsf{T}}\right)^{-1}.$$
 (3.66)

By 2.24,

$$\begin{bmatrix} \mathbf{S}^{-1} & \beta \boldsymbol{\phi}_{N+1} \\ \beta \boldsymbol{\phi}_{N+1}^{\mathsf{T}} & -\beta \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{S}' & \mathbf{S}' \boldsymbol{\phi}_{N+1} \\ \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \mathbf{S}' & -\beta^{-1} + \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \mathbf{S}' \boldsymbol{\phi}_{N+1} \end{bmatrix}, \quad (3.67)$$

and

$$\begin{bmatrix} -\beta & \beta \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \\ \beta \boldsymbol{\phi}_{N+1} & \mathbf{S}^{-1} \end{bmatrix}^{-1} = \begin{bmatrix} -c & \beta c \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \mathbf{S} \\ \beta c \mathbf{S} \boldsymbol{\phi}_{N+1} & \mathbf{S} - \beta^2 c \mathbf{S} \boldsymbol{\phi}_{N+1} \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \mathbf{S} \end{bmatrix}, \quad (3.68)$$

where

$$c = \left(\beta + \beta^2 \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \mathbf{S} \boldsymbol{\phi}_{N+1}\right)^{-1}. \tag{3.69}$$

Then,

$$\mathbf{S}' = \mathbf{S} - \beta^2 c \mathbf{S} \phi_{N+1} \phi_{N+1}^{\mathsf{T}} \mathbf{S}. \tag{3.70}$$

Then,

$$\phi(\mathbf{x})^{\mathsf{T}}(\mathbf{S}_{N} - \mathbf{S}_{N+1})\phi(\mathbf{x}) = \frac{\beta \left(\phi(\mathbf{x})^{\mathsf{T}}\mathbf{S}_{N}\phi(\mathbf{x}_{N+1})\right)^{2}}{1 + \beta\phi(\mathbf{x}_{N+1})\mathbf{S}_{N}\phi(\mathbf{x}_{N+1})^{\mathsf{T}}}.$$
(3.71)

Therefore,

$$\sigma_{N+1}^2(\mathbf{x}) \le \sigma_N^2(\mathbf{x}). \tag{3.72}$$

#### 3.12

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w},\beta) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n),\beta^{-1}\right),$$
  

$$p(\mathbf{w},\beta) = \mathcal{N}\left(\mathbf{w}|\mathbf{m}_0,\beta^{-1}\mathbf{S}_0\right)\operatorname{Gam}(\beta|a_0,b_0),$$
(3.73)

where **w** and  $\phi$  are vectors in M dimensions. By the Bayes' theorem,

$$p(\mathbf{w}, \beta | \mathbf{t}) p(\mathbf{t}) = p(\mathbf{t} | \mathbf{w}, \beta) p(\mathbf{w}, \beta). \tag{3.74}$$

The logarithm of the right hand side except the terms independent of  $\mathbf{t}$ ,  $\mathbf{w}$  and  $\beta$  can be written as

$$-\frac{N}{2}\ln\beta^{-1} - \frac{\beta}{2}\sum_{n=1}^{N}(t_n - \mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n))^2 - \frac{M}{2}\ln\beta^{-1} - \frac{\beta}{2}(\mathbf{w} - \mathbf{m}_0)^{\mathsf{T}}\mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0)$$

$$+ (a_0 - 1)\ln\beta - b_0\beta$$

$$= -\frac{M}{2}\ln\beta - \frac{\beta}{2}\mathbf{w}^{\mathsf{T}}\left(\mathbf{S}_0^{-1} + \mathbf{\Phi}^{\mathsf{T}}\mathbf{\Phi}\right)\mathbf{w} + \beta\mathbf{w}^{\mathsf{T}}\left(\mathbf{S}_0^{-1}\mathbf{m}_0 + \mathbf{\Phi}^{\mathsf{T}}\mathbf{t}\right) - \frac{\beta}{2}\|\mathbf{t}\|^2 - \frac{\beta}{2}\mathbf{m}_0^{\mathsf{T}}\mathbf{S}_0^{-1}\mathbf{m}_0$$

$$+ \left(a_0 + \frac{N}{2} - 1\right)\ln\beta - b_0\beta.$$
(3.75)

The right hand side can be written as

$$-\frac{M}{2}\ln\beta - \frac{\beta}{2}(\mathbf{w} - \mathbf{m}_N)^{\mathsf{T}}\mathbf{S}_N^{-1}(\mathbf{w} - \mathbf{m}_N) + (a_N - 1)\ln\beta - b_N\beta, \quad (3.76)$$

where

$$\mathbf{m}_{N} = \mathbf{S}_{N} \left( \mathbf{S}_{0}^{-1} \mathbf{m}_{0} + \mathbf{\Phi}^{\mathsf{T}} \mathbf{t} \right),$$

$$\mathbf{S}_{N}^{-1} = \mathbf{S}_{0}^{-1} + \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi},$$

$$a_{N} = a_{0} + \frac{N}{2},$$

$$b_{N} = b_{0} + \frac{1}{2} \|\mathbf{t}\|^{2} + \frac{1}{2} \mathbf{m}_{0}^{\mathsf{T}} \mathbf{S}_{0} \mathbf{m}_{0} - \frac{1}{2} \mathbf{m}_{N}^{\mathsf{T}} \mathbf{S}_{N}^{-1} \mathbf{m}_{N}.$$

$$(3.77)$$

Therefore,

$$p(\mathbf{w}, \beta | \mathbf{t}) = \mathcal{N}\left(\mathbf{w} | \mathbf{m}_N, \beta^{-1} \mathbf{S}_N\right) \operatorname{Gam}(\beta | a_N, b_N).$$
 (3.78)

Substituting it to the result of the Bayes' theorem above, we have

$$p(\mathbf{t}) = \frac{\mathcal{N}\left(\mathbf{t}|\mathbf{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I}\right) \mathcal{N}\left(\mathbf{w}|\mathbf{m}_{0}, \beta^{-1}\mathbf{S}_{0}\right) \operatorname{Gam}(\beta|a_{0}, b_{0})}{\mathcal{N}\left(\mathbf{w}|\mathbf{m}_{N}, \beta^{-1}\mathbf{S}_{N}\right) \operatorname{Gam}(\beta|a_{N}, b_{N})}.$$
(3.79)

The logarithm of the right hand side can be written as

$$-\frac{N}{2}\ln(2\pi) - \frac{N}{2}\ln\beta^{-1} - \frac{\beta}{2}(\mathbf{t} - \mathbf{\Phi}\mathbf{w})^{\mathsf{T}}(\mathbf{t} - \mathbf{\Phi}\mathbf{w})$$

$$-\frac{M}{2}\ln(2\pi) - \frac{M}{2}\ln\beta^{-1} - \frac{1}{2}\det\mathbf{S}_{0} - \frac{\beta}{2}(\mathbf{w} - \mathbf{m}_{0})^{\mathsf{T}}\mathbf{S}_{0}^{-1}(\mathbf{w} - \mathbf{m}_{0})$$

$$+ a_{0}\ln b_{0} - \ln\Gamma(a_{0}) + (a_{0} - 1)\ln\beta - b_{0}\beta$$

$$+ \frac{M}{2}\ln(2\pi) + \frac{M}{2}\ln\beta^{-1} + \frac{1}{2}\det\mathbf{S}_{N} + \frac{\beta}{2}(\mathbf{w} - \mathbf{m}_{N})^{\mathsf{T}}\mathbf{S}_{N}^{-1}(\mathbf{w} - \mathbf{m}_{0})$$

$$- a_{N}\ln b_{N} + \ln\Gamma(a_{N}) - (a_{N} - 1)\ln\beta + b_{N}\beta$$

$$= -\frac{N}{2}\ln(2\pi) - \frac{1}{2}\det\mathbf{S}_{0} + a_{0}\ln b_{0} - \ln\Gamma(a_{0}) + \frac{1}{2}\det\mathbf{S}_{N} - a_{N}\ln b_{N} + \ln\Gamma(a_{N}).$$
(3.80)

Therefore,

$$p(\mathbf{t}) = (2\pi)^{-\frac{N}{2}} \left( \frac{\det \mathbf{S}_N}{\det \mathbf{S}_0} \right)^{\frac{1}{2}} \frac{\Gamma(a_N)}{\Gamma(a_0)} \frac{b_0^{a_0}}{b_N^{a_N}}.$$
 (3.81)

#### 3.13

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w},\beta) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n),\beta^{-1}\right),$$
  

$$p(\mathbf{w},\beta) = \mathcal{N}\left(\mathbf{w}|\mathbf{m}_0,\beta^{-1}\mathbf{S}_0\right)\operatorname{Gam}(\beta|a_0,b_0),$$
(3.82)

where **w** and  $\phi$  are vectors in M dimensions. Then, by 3.12,

$$p(\mathbf{w}, \beta | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \beta^{-1} \mathbf{S}_N) \operatorname{Gam}(\beta | a_N, b_N),$$
 (3.83)

where

$$\mathbf{m}_{N} = \mathbf{S}_{N} \left( \mathbf{S}_{0}^{-1} \mathbf{m}_{0} + \mathbf{\Phi}^{\mathsf{T}} \mathbf{t} \right),$$

$$\mathbf{S}_{N}^{-1} = \mathbf{S}_{0}^{-1} + \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi},$$

$$a_{N} = a_{0} + \frac{N}{2},$$

$$b_{N} = b_{0} + \frac{1}{2} \|\mathbf{t}\|^{2} + \frac{1}{2} \mathbf{m}_{0}^{\mathsf{T}} \mathbf{S}_{0} \mathbf{m}_{0} - \frac{1}{2} \mathbf{m}_{N}^{\mathsf{T}} \mathbf{S}_{N}^{-1} \mathbf{m}_{N}.$$

$$(3.84)$$

By marginalisation,

$$p(t|\mathbf{t}) = \int \int p(t|\mathbf{w}, \beta)p(\mathbf{w}, \beta|\mathbf{t})d\mathbf{w}d\beta.$$
 (3.85)

The right hand side can be written as

$$\int \left( \int \mathcal{N} \left( t | \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}), \beta^{-1} \right) \mathcal{N} \left( \mathbf{w} | \mathbf{m}_{N}, \beta^{-1} \mathbf{S}_{N} \right) d\mathbf{w} \right) \operatorname{Gam}(\beta | a_{N}, b_{N}) d\beta.$$
(3.86)

The logarithm of the integrand with respect to  $\mathbf{w}$  except the terms indepndent of  $\mathbf{w}$  can be written as

$$-\frac{\beta}{2} \left(t - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x})\right)^{2} - \frac{\beta}{2} (\mathbf{w} - \mathbf{m}_{N})^{\mathsf{T}} \mathbf{S}_{N}^{-1} (\mathbf{w} - \mathbf{m}_{N}). \tag{3.87}$$

It can be written as

$$-\frac{\beta}{2} \begin{bmatrix} \mathbf{w} \\ t \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{S}_N^{-1} + \boldsymbol{\phi}(\mathbf{x}) \boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}} & -\boldsymbol{\phi}(\mathbf{x}) \\ -\boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}} & 1 \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ t \end{bmatrix} + \beta \begin{bmatrix} \mathbf{w} \\ t \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{S}_N^{-1} \mathbf{m}_N \\ 0 \end{bmatrix} - \frac{\beta}{2} \mathbf{m}_N^{\mathsf{T}} \mathbf{S}_N^{-1} \mathbf{m}_N.$$

By 2.24,

$$\begin{bmatrix} \mathbf{S}_N^{-1} + \boldsymbol{\phi}(\mathbf{x})\boldsymbol{\phi}(\mathbf{x})^\intercal & -\boldsymbol{\phi}(\mathbf{x}) \\ -\boldsymbol{\phi}(\mathbf{x})^\intercal & 1 \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{S}_N & \mathbf{S}_N\boldsymbol{\phi}(\mathbf{x}) \\ \boldsymbol{\phi}(\mathbf{x})^\intercal\mathbf{S}_N & 1 + \boldsymbol{\phi}(\mathbf{x})^\intercal\mathbf{S}_N\boldsymbol{\phi}(\mathbf{x}) \end{bmatrix}.$$

Then,

$$\begin{bmatrix} \mathbf{S}_N^{-1} + \boldsymbol{\phi}(\mathbf{x})\boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}} & -\boldsymbol{\phi}(\mathbf{x}) \\ -\boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}} & 1 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{S}_N^{-1}\mathbf{m}_N \\ 0 \end{bmatrix} = \begin{bmatrix} \mathbf{m}_N \\ \mathbf{m}_N^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}) \end{bmatrix}.$$

Therefore, the integral with respect to  $\mathbf{w}$  can be written as

$$\mathcal{N}\left(t|\mathbf{m}_{N}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}),\beta^{-1}\left(1+\boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}}\mathbf{S}_{N}\boldsymbol{\phi}(\mathbf{x})\right)\right). \tag{3.88}$$

Then, the logarithm of the integrand with respect to  $\beta$  except the terms independent of  $\beta$  can be written as

$$-\frac{1}{2}\ln\beta^{-1} - \frac{\beta}{2\left(1 + \boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}}\mathbf{S}_{N}\boldsymbol{\phi}(\mathbf{x})\right)} \left(t - \mathbf{m}_{N}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x})\right)^{2} + (a_{N} - 1)\ln\beta - b_{N}\beta$$

$$= \left(a_{N} + \frac{1}{2} - 1\right)\ln\beta - \left(b_{N} + \frac{\left(t - \mathbf{m}_{N}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x})\right)^{2}}{2\left(1 + \boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}}\mathbf{S}_{N}\boldsymbol{\phi}(\mathbf{x})\right)}\right)\beta.$$
(3.89)

Therefore, the integral with respect to  $\beta$  except the terms independent of t can be written as

$$\left(b_N + \frac{\left(t - \mathbf{m}_N^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x})\right)^2}{2\left(1 + \boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}} \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x})\right)}\right)^{-a_N - \frac{1}{2}}.$$
(3.90)

Thus,

$$p(t|\mathbf{x}, \mathbf{t}) = \operatorname{St}(t|\mu, \lambda, \nu), \tag{3.91}$$

where

$$\mu = \mathbf{m}_{N}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}),$$

$$\lambda = \frac{a_{N}}{b_{N}} (1 + \boldsymbol{\phi}(\mathbf{x})^{\mathsf{T}} \mathbf{S}_{N} \boldsymbol{\phi}(\mathbf{x}))^{-1},$$

$$\nu = 2a_{N}.$$
(3.92)

# 3.14 (Incomplete)

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}),$$
(3.93)

where **w** and  $\phi$  are vectors in M dimensions. Then, by 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N), \tag{3.94}$$

where

$$\mathbf{m}_{N} = \beta \mathbf{S}_{N} \mathbf{\Phi}^{\mathsf{T}} \mathbf{t}, \mathbf{S}_{N}^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi}.$$
 (3.95)

Let

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^{\mathsf{T}} \phi(\mathbf{x}). \tag{3.96}$$

Then,

$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^{N} k(\mathbf{x}, \mathbf{x}_n) t_n,$$
 (3.97)

where

$$k(\mathbf{x}, \mathbf{x}') = \beta \phi(\mathbf{x})^{\mathsf{T}} \mathbf{S}_N \phi(\mathbf{x}'). \tag{3.98}$$

Let us suppose that  $\phi_j(\mathbf{x})$  are linearly independent, N > M and

$$\phi_0(\mathbf{x}) = 1. \tag{3.99}$$

Then, we can construct a new basis set  $\psi_i(\mathbf{x})$  such that

$$\mathbf{\Psi}^{\mathsf{T}}\mathbf{\Psi} = \mathbf{I}?\tag{3.100}$$

$$\sum_{n=1}^{N} \psi_j(\mathbf{x}_n) \psi_k(\mathbf{x}_n) = I_{jk}? \tag{3.101}$$

where

$$oldsymbol{\Psi} = egin{bmatrix} oldsymbol{\psi}(\mathbf{x}_1)^\intercal \ dots \ oldsymbol{\psi}(\mathbf{x}_N)^\intercal \end{bmatrix}$$

and

$$\psi_0(\mathbf{x}) = 1. \tag{3.102}$$

Under the basis set, if  $\alpha = 0$ , then

$$\mathbf{S}_N^{-1} = \beta \mathbf{I},\tag{3.103}$$

so that

$$k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\psi}(\mathbf{x})^{\mathsf{T}} \boldsymbol{\psi}(\mathbf{x}'). \tag{3.104}$$

Then,

$$\sum_{n=1}^{N} k(\mathbf{x}, \mathbf{x}_n) = \sum_{n=1}^{N} \sum_{j=0}^{M-1} \psi_j(\mathbf{x}) \psi_j(\mathbf{x}_n) = 1?$$
 (3.105)

#### 3.15

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}),$$
(3.106)

where **w** and  $\phi$  are vectors in M dimensions. By 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N), \tag{3.107}$$

where

$$\mathbf{m}_{N} = \beta \mathbf{S}_{N} \mathbf{\Phi}^{\mathsf{T}} \mathbf{t}, \mathbf{S}_{N}^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi}.$$
 (3.108)

By 3.19,

$$\ln p(\mathbf{t}) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) + \frac{1}{2} \ln(\det \mathbf{S}_N) - E(\mathbf{m}_N), \quad (3.109)$$

where

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^{\mathsf{T}} \mathbf{m}_N.$$
 (3.110)

By 3.22, setting the derivatives of  $\ln p(\mathbf{t})$  with respect to  $\alpha$  and  $\beta$  to zero gives

$$\alpha = \frac{\gamma}{\mathbf{m}_{N}^{\mathsf{T}} \mathbf{m}_{N}},$$

$$\beta = \frac{N - \gamma}{\|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_{N}\|^{2}},$$
(3.111)

where

$$\gamma = \sum_{m=1}^{M} \frac{\lambda_m}{\alpha + \lambda_m} \tag{3.112}$$

and  $\lambda_1, \dots, \lambda_M$  are the eigenvalues of  $\beta \Phi^{\dagger} \Phi$ . If  $\alpha$  and  $\beta$  are set as above, then

$$E(\mathbf{m}_N) = \frac{N}{2}. (3.113)$$

#### 3.16

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}).$$
 (3.114)

where **w** and  $\phi$  are vectors in M dimensions. By the Bayes' theorem,

$$p(\mathbf{t}, \mathbf{w}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}). \tag{3.115}$$

Integrating both sides with respect to  $\mathbf{w}$  gives

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{w})p(\mathbf{w})d\mathbf{w}.$$
 (3.116)

The logarithm of the integrand of the right hand side except the terms independent of  $\mathbf{w}$  can be written as

$$-\frac{\beta}{2} \sum_{n=1}^{N} (t_n - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n))^2 - \frac{\alpha}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w} = -\frac{1}{2} \begin{bmatrix} \mathbf{w} \\ \mathbf{t} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \alpha \mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} & -\beta \mathbf{\Phi}^{\mathsf{T}} \\ -\beta \mathbf{\Phi} & \beta \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{t} \end{bmatrix}.$$

By 2.24,

$$\begin{bmatrix} \alpha \mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} & -\beta \mathbf{\Phi}^{\mathsf{T}} \\ -\beta \mathbf{\Phi} & \beta \mathbf{I} \end{bmatrix}^{-1} = \begin{bmatrix} \alpha^{-1} \mathbf{I} & \alpha^{-1} \mathbf{\Phi}^{\mathsf{T}} \\ \alpha^{-1} \mathbf{\Phi} & \alpha^{-1} \mathbf{\Phi} \mathbf{\Phi}^{\mathsf{T}} + \beta^{-1} \mathbf{I} \end{bmatrix}.$$

Therefore,

$$p(\mathbf{t}) = \mathcal{N}\left(\mathbf{t}|\mathbf{0}, \alpha^{-1}\mathbf{\Phi}\mathbf{\Phi}^{\mathsf{T}} + \beta^{-1}\mathbf{I}\right). \tag{3.117}$$

#### 3.17

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}),$$
(3.118)

where **w** and  $\phi$  are vectors in M dimensions. By the Bayes' theorem,

$$p(\mathbf{t}, \mathbf{w}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}). \tag{3.119}$$

Then,

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{w})p(\mathbf{w})d\mathbf{w}.$$
 (3.120)

The logarithm of the integrand of the right hand side can be written as

$$-\frac{N}{2}\ln\left(2\pi\beta^{-1}\right) - \frac{\beta}{2}\sum_{n=1}^{N}\left(t_n - \mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n)\right)^2 - \frac{M}{2}\ln\left(2\pi\right) - \frac{1}{2}\ln\left(\det\left(\alpha^{-1}\mathbf{I}\right)\right) - \frac{\alpha}{2}\mathbf{w}^{\mathsf{T}}\mathbf{w}.$$
(3.121)

Therefore,

$$p(\mathbf{t}) = \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \int \exp\left(-E(\mathbf{w})\right) d\mathbf{w}, \tag{3.122}$$

where

$$E(\mathbf{w}) = \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi}\mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}.$$
 (3.123)

#### 3.18

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}),$$
(3.124)

where **w** and  $\phi$  are vectors in M dimensions. By 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N), \tag{3.125}$$

where

$$\mathbf{m}_{N} = \beta \mathbf{S}_{N} \mathbf{\Phi}^{\mathsf{T}} \mathbf{t}, \mathbf{S}_{N}^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi}.$$
 (3.126)

By 3.17,

$$p(\mathbf{t}) = \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \int \exp\left(-E(\mathbf{w})\right) d\mathbf{w}, \tag{3.127}$$

where

$$E(\mathbf{w}) = \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi}\mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}.$$
 (3.128)

The first term of the definition of  $E(\mathbf{w})$  can be written as

$$\frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N - \mathbf{\Phi} (\mathbf{w} - \mathbf{m}_N)\|^2 
= \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N\|^2 - \beta (\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N)^{\mathsf{T}} \mathbf{\Phi} (\mathbf{w} - \mathbf{m}_N) + \frac{\beta}{2} (\mathbf{w} - \mathbf{m}_N)^{\mathsf{T}} \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} (\mathbf{w} - \mathbf{m}_N).$$
(3.129)

Similarly, the second term can be written as

$$\frac{\alpha}{2} (\mathbf{w} - \mathbf{m}_N + \mathbf{m}_N)^{\mathsf{T}} (\mathbf{w} - \mathbf{m}_N + \mathbf{m}_N) 
= \frac{\alpha}{2} (\mathbf{w} - \mathbf{m}_N)^{\mathsf{T}} (\mathbf{w} - \mathbf{m}_N) + \alpha \mathbf{m}_N^{\mathsf{T}} (\mathbf{w} - \mathbf{m}_N) + \frac{\alpha}{2} \mathbf{m}_N^{\mathsf{T}} \mathbf{m}_N.$$
(3.130)

Here,

$$-\beta(\mathbf{t} - \mathbf{\Phi}\mathbf{m}_N)^{\mathsf{T}}\mathbf{\Phi}(\mathbf{w} - \mathbf{m}_N) + \alpha\mathbf{m}_N^{\mathsf{T}}(\mathbf{w} - \mathbf{m}_N)$$
  
=  $(-\beta\mathbf{\Phi}^{\mathsf{T}}\mathbf{t} + \beta\mathbf{\Phi}^{\mathsf{T}}\mathbf{\Phi}\mathbf{m}_N + \alpha\mathbf{m}_N)^{\mathsf{T}}(\mathbf{w} - \mathbf{m}_N).$  (3.131)

By the definitions of  $\mathbf{m}_N$  and  $\mathbf{S}_N$  above, the right hand can be written as

$$\left(-\beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{t} + \mathbf{S}_{N}^{-1} \mathbf{m}_{N}\right)^{\mathsf{T}} (\mathbf{w} - \mathbf{m}_{N}) = 0. \tag{3.132}$$

Therefore,

$$E(\mathbf{w}) = E(\mathbf{m}_N) + \frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^{\mathsf{T}} \mathbf{S}_N^{-1} (\mathbf{w} - \mathbf{m}_N). \tag{3.133}$$

#### 3.19

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}),$$
(3.134)

where **w** and  $\phi$  are vectors in M dimensions. By 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N), \tag{3.135}$$

where

$$\mathbf{m}_{N} = \beta \mathbf{S}_{N} \mathbf{\Phi}^{\mathsf{T}} \mathbf{t}, \mathbf{S}_{N}^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi}.$$
 (3.136)

By 3.17,

$$p(\mathbf{t}) = \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \int \exp\left(-E(\mathbf{w})\right) d\mathbf{w}, \tag{3.137}$$

where

$$E(\mathbf{w}) = \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi}\mathbf{w}\|^2 + \frac{\alpha}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}.$$
 (3.138)

By 3.18,

$$E(\mathbf{w}) = E(\mathbf{m}_N) + \frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^{\mathsf{T}} \mathbf{S}_N^{-1} (\mathbf{w} - \mathbf{m}_N). \tag{3.139}$$

Therefore, the integral in the expression above of p(t) can be written as

$$\exp\left(-E(\mathbf{m}_N)\right) \int \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^{\mathsf{T}} \mathbf{S}_N^{-1}(\mathbf{w} - \mathbf{m}_N)\right) d\mathbf{w}$$

$$= (2\pi)^{\frac{M}{2}} (\det \mathbf{S}_N)^{\frac{1}{2}} \exp\left(-E(\mathbf{m}_N)\right). \tag{3.140}$$

Thus,

$$\ln p(\mathbf{t}) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) + \frac{1}{2} \ln(\det \mathbf{S}_N) - E(\mathbf{m}_N). \quad (3.141)$$

## 3.20

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}),$$
(3.142)

where **w** and  $\phi$  are vectors in M dimensions. By 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N), \tag{3.143}$$

where

$$\mathbf{m}_{N} = \beta \mathbf{S}_{N} \mathbf{\Phi}^{\mathsf{T}} \mathbf{t}, \mathbf{S}_{N}^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi}.$$
 (3.144)

By 3.19,

$$\ln p(\mathbf{t}) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) + \frac{1}{2} \ln(\det \mathbf{S}_N) - E(\mathbf{m}_N), \quad (3.145)$$

where

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^{\mathsf{T}} \mathbf{m}_N.$$
 (3.146)

Let  $\mathbf{u}_1, \cdots, \mathbf{u}_M$  be eigenvectors of  $\beta \mathbf{\Phi}^{\intercal} \mathbf{\Phi}$  such that

$$\beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} \mathbf{u}_m = \lambda_m \mathbf{u}_m. \tag{3.147}$$

Then,

$$\mathbf{S}_N^{-1}\mathbf{u}_m = (\alpha + \lambda_m)\mathbf{u}_m, \tag{3.148}$$

so that

$$\det \mathbf{S}_N = \prod_{m=1}^M \frac{1}{\alpha + \lambda_m}.$$
 (3.149)

Therefore, setting the derivative of  $\ln p(\mathbf{t}|\alpha,\beta)$  with respect to  $\alpha$  to zero gives

$$0 = \frac{M}{2\alpha} - \frac{1}{2} \sum_{m=1}^{M} \frac{1}{\alpha + \lambda_m} - \frac{1}{2} \mathbf{m}_N^{\mathsf{T}} \mathbf{m}_N.$$
 (3.150)

Multiplying both sides by  $2\alpha$  gives

$$\alpha \mathbf{m}_{N}^{\mathsf{T}} \mathbf{m}_{N} = M - \sum_{m=1}^{M} \frac{\alpha}{\alpha + \lambda_{m}}.$$
 (3.151)

The right hand side can be written as

$$\sum_{m=1}^{M} \left( 1 - \frac{\alpha}{\alpha + \lambda_m} \right) = \sum_{m=1}^{M} \frac{\lambda_i}{\alpha + \lambda_m}.$$
 (3.152)

Thus,

$$\alpha = \frac{\gamma}{\mathbf{m}_N^{\mathsf{T}} \mathbf{m}_N},\tag{3.153}$$

where

$$\gamma = \sum_{m=1}^{M} \frac{\lambda_i}{\alpha + \lambda_m}.$$
 (3.154)

# 3.21

(a)

Let  $\Sigma$  be a  $M \times M$  real symmetric matrix such that

$$\Sigma \mathbf{u}_m = \lambda_m \mathbf{u}_m, \tag{3.155}$$

where  $\mathbf{u}_1, \cdots, \mathbf{u}_M$  are unit vectors. Let

$$\Lambda = \operatorname{diag}(\lambda_1, \cdots, \lambda_M), 
\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_M].$$
(3.156)

By 2.19,

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

$$\mathbf{U}^{\mathsf{T}}\mathbf{U} = \mathbf{I}.$$
(3.157)

Therefore,

$$\det \mathbf{\Sigma} = \prod_{m=1}^{M} \lambda_m, \tag{3.158}$$

so that

$$\ln(\det \mathbf{\Sigma}) = \sum_{m=1}^{M} \ln \lambda_i. \tag{3.159}$$

Then,

$$\frac{\partial}{\partial \alpha} \ln(\det \Sigma) = \sum_{m=1}^{M} \frac{\partial \lambda_m}{\partial \alpha} \frac{1}{\lambda_m}.$$
 (3.160)

Therefore,

$$\frac{\partial}{\partial \alpha} \ln(\det \Sigma) = \operatorname{tr}\left(\Lambda^{-1} \frac{\partial \Lambda}{\partial \alpha}\right). \tag{3.161}$$

The right hand side can be written as

$$\operatorname{tr}\left(\mathbf{U}\boldsymbol{\Lambda}^{-1}\mathbf{U}^{\mathsf{T}}\frac{\partial\mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^{\mathsf{T}}}{\partial\alpha}\right) = \operatorname{tr}\left(\boldsymbol{\Sigma}^{-1}\frac{\partial\boldsymbol{\Sigma}}{\partial\alpha}\right). \tag{3.162}$$

Therefore,

$$\frac{\partial}{\partial \alpha} \ln(\det \Sigma) = \operatorname{tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial \alpha} \right). \tag{3.163}$$

(b)

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}),$$
(3.164)

where **w** and  $\phi$  are vectors in M dimensions. By 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N), \tag{3.165}$$

where

$$\mathbf{m}_{N} = \beta \mathbf{S}_{N} \mathbf{\Phi}^{\mathsf{T}} \mathbf{t}, \mathbf{S}_{N}^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi}.$$
 (3.166)

By 3.19,

$$\ln p(\mathbf{t}) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) + \frac{1}{2} \ln(\det \mathbf{S}_N) - E(\mathbf{m}_N), \quad (3.167)$$

where

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^{\mathsf{T}} \mathbf{m}_N.$$
 (3.168)

By 3.21(a),

$$\frac{\partial}{\partial \alpha} \ln \left( \det \mathbf{S}_N^{-1} \right) = \operatorname{tr} \left( \mathbf{S}_N \right). \tag{3.169}$$

The right hand side can be written as

$$\sum_{m=1}^{M} \frac{1}{\alpha + \lambda_m},\tag{3.170}$$

where  $\lambda_1, \dots, \lambda_M$  are eigenvalues of  $\beta \Phi^{\dagger} \Phi$ . Therefore, setting the derivative of  $\ln p(\mathbf{t})$  with respect to  $\alpha$  to zero gives

$$0 = \frac{M}{2\alpha} - \frac{1}{2} \sum_{m=1}^{M} \frac{1}{\alpha + \lambda_m} - \frac{1}{2} \mathbf{m}_N^{\mathsf{T}} \mathbf{m}_N, \tag{3.171}$$

Thus,

$$\alpha = \frac{\gamma}{\mathbf{m}_N^{\mathsf{T}} \mathbf{m}_N},\tag{3.172}$$

where

$$\gamma = \sum_{m=1}^{M} \frac{\lambda_m}{\alpha + \lambda_m}.$$
 (3.173)

#### 3.22

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}),$$
(3.174)

where **w** and  $\phi$  are vectors in M dimensions. By 3.7,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N), \tag{3.175}$$

where

$$\mathbf{m}_{N} = \beta \mathbf{S}_{N} \mathbf{\Phi}^{\mathsf{T}} \mathbf{t}, \mathbf{S}_{N}^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi}.$$
 (3.176)

By 3.19,

$$\ln p(\mathbf{t}) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) + \frac{1}{2} \ln(\det \mathbf{S}_N) - E(\mathbf{m}_N), \quad (3.177)$$

where

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^{\mathsf{T}} \mathbf{m}_N.$$
 (3.178)

By 3.21(a),

$$\frac{\partial}{\partial \beta} \ln \left( \det \mathbf{S}_N^{-1} \right) = \operatorname{tr} \left( \mathbf{S}_N \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} \right). \tag{3.179}$$

Since

$$\mathbf{S}_N \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} = \frac{1}{\beta} \left( \mathbf{I} - \alpha \mathbf{S}_N \right), \tag{3.180}$$

the right hand side can be written as

$$\frac{1}{\beta} \left( M - \alpha \sum_{m=1}^{M} \frac{1}{\alpha + \lambda_m} \right) = \frac{1}{\beta} \sum_{m=1}^{M} \frac{\lambda_m}{\alpha + \lambda_m}, \tag{3.181}$$

where  $\lambda_1, \dots, \lambda_M$  are eigenvalues of  $\beta \Phi^{\dagger} \Phi$ . Therefore, setting the derivative of  $\ln p(\mathbf{t})$  with respect to  $\beta$  to zero gives

$$0 = \frac{N}{2\beta} - \frac{1}{2\beta} \sum_{m=1}^{M} \frac{\lambda_i}{\alpha + \lambda_m} - \frac{1}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{m}_N\|^2.$$
 (3.182)

Thus,

$$\beta = \frac{N - \gamma}{\left\|\mathbf{t} - \mathbf{\Phi}\mathbf{m}_N\right\|^2},\tag{3.183}$$

where

$$\gamma = \sum_{m=1}^{M} \frac{\lambda_m}{\alpha + \lambda_m}.$$
 (3.184)

## 3.23

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w},\beta) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n),\beta^{-1}\right),$$
  

$$p(\mathbf{w},\beta) = \mathcal{N}\left(\mathbf{w}|\mathbf{m}_0,\beta^{-1}\mathbf{S}_0\right)\operatorname{Gam}(\beta|a_0,b_0),$$
(3.185)

where  $\mathbf{w}$  and  $\boldsymbol{\phi}$  are vectors in M dimensions. By marginalisation,

$$p(\mathbf{t}) = \int \int p(\mathbf{t}|\mathbf{w}, \beta) p(\mathbf{w}, \beta) d\mathbf{w} d\beta.$$
 (3.186)

The right hand side can be written as

$$\int \left( \int \left( \prod_{n=1}^{N} \mathcal{N} \left( t_n | \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1} \right) \right) \mathcal{N} \left( \mathbf{w} | \mathbf{m}_0, \beta^{-1} \mathbf{S}_0 \right) d\mathbf{w} \right) \operatorname{Gam}(\beta | a_0, b_0) d\beta.$$
(3.187)

The logarithm of the integrand with respect to  $\mathbf{w}$  can be written as

$$-\frac{N}{2}\ln(2\pi) - \frac{N}{2}\ln\beta^{-1} - \frac{\beta}{2}\sum_{n=1}^{N}\left(t_{n} - \mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_{n})\right)^{2}$$

$$-\frac{M}{2}\ln(2\pi) - \frac{1}{2}\ln\det(\beta^{-1}\mathbf{S}_{0}) - \frac{\beta}{2}(\mathbf{w} - \mathbf{m}_{0})^{\mathsf{T}}\mathbf{S}_{0}^{-1}(\mathbf{w} - \mathbf{m}_{0})$$

$$= -\frac{N+M}{2}\ln(2\pi) + \frac{N+M}{2}\ln\beta - \frac{1}{2}\ln(\det\mathbf{S}_{0})$$

$$-\frac{\beta}{2}\mathbf{w}^{\mathsf{T}}\left(\mathbf{S}_{0}^{-1} + \mathbf{\Phi}^{\mathsf{T}}\mathbf{\Phi}\right)\mathbf{w} + \beta\mathbf{w}^{\mathsf{T}}\left(\mathbf{S}_{0}^{-1}\mathbf{m}_{0} + \mathbf{\Phi}^{\mathsf{T}}\mathbf{t}\right) - \frac{\beta}{2}\|\mathbf{t}\|^{2} - \frac{\beta}{2}\mathbf{m}_{0}^{\mathsf{T}}\mathbf{S}_{0}^{-1}\mathbf{m}_{0}.$$
(3.188)

The right hand side can be written as

$$-\frac{N+M}{2}\ln(2\pi) + \frac{N+M}{2}\ln\beta - \frac{1}{2}\ln(\det\mathbf{S}_{0})$$

$$-\frac{\beta}{2}(\mathbf{w} - \mathbf{m}_{N})^{\mathsf{T}}\mathbf{S}_{N}^{-1}(\mathbf{w} - \mathbf{m}_{N}) + \frac{\beta}{2}\mathbf{m}_{N}^{\mathsf{T}}\mathbf{S}_{N}^{-1}\mathbf{m}_{N} - \frac{\beta}{2}\|\mathbf{t}\|^{2} - \frac{\beta}{2}\mathbf{m}_{0}^{\mathsf{T}}\mathbf{S}_{0}^{-1}\mathbf{m}_{0},$$
(3.189)

where

$$\mathbf{m}_{N} = \mathbf{S}_{N} \left( \mathbf{S}_{0}^{-1} \mathbf{m}_{0} + \mathbf{\Phi}^{\mathsf{T}} \mathbf{t} \right), \mathbf{S}_{N}^{-1} = \mathbf{S}_{0}^{-1} + \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi}.$$
 (3.190)

Therefore, the logarithm of the integral with respect to  $\mathbf{w}$  can be written as

$$-\frac{N}{2}\ln(2\pi) + \frac{N}{2}\ln\beta - \frac{1}{2}\ln(\det\mathbf{S}_0) + \frac{1}{2}\ln(\det\mathbf{S}_N) + \frac{\beta}{2}\mathbf{m}_N^{\mathsf{T}}\mathbf{S}_N^{-1}\mathbf{m}_N - \frac{\beta}{2}\|\mathbf{t}\|^2 - \frac{\beta}{2}\mathbf{m}_0^{\mathsf{T}}\mathbf{S}_0^{-1}\mathbf{m}_0.$$
(3.191)

Then, the logarithm of the integrand with respect to  $\beta$  can be written as

$$-\frac{N}{2}\ln(2\pi) + \frac{N}{2}\ln\beta - \frac{1}{2}\ln(\det\mathbf{S}_{0}) + \frac{1}{2}\ln(\det\mathbf{S}_{N}) + \frac{\beta}{2}\mathbf{m}_{N}^{\mathsf{T}}\mathbf{S}_{N}^{-1}\mathbf{m}_{N} - \frac{\beta}{2}\|\mathbf{t}\|^{2} - \frac{\beta}{2}\mathbf{m}_{0}^{\mathsf{T}}\mathbf{S}_{0}^{-1}\mathbf{m}_{0} - \ln\Gamma(a_{0}) + a_{0}\ln b_{0} + (a_{0} - 1)\ln\beta - b_{0}\beta$$

$$= -\frac{N}{2}\ln(2\pi) - \frac{1}{2}\ln(\det\mathbf{S}_{0}) + \frac{1}{2}\ln(\det\mathbf{S}_{N}) - \ln\Gamma(a_{0}) + a_{0}\ln b_{0} + (a_{N} - 1)\ln\beta - b_{N}\beta,$$
(3.192)

where

$$a_{N} = a_{0} + \frac{N}{2},$$

$$b_{N} = b_{0} + \frac{\beta}{2} \|\mathbf{t}\|^{2} + \frac{\beta}{2} \mathbf{m}_{0}^{\mathsf{T}} \mathbf{S}_{0}^{-1} \mathbf{m}_{0} - \frac{\beta}{2} \mathbf{m}_{N}^{\mathsf{T}} \mathbf{S}_{N}^{-1} \mathbf{m}_{N}.$$
(3.193)

Therefore, the logarithm of the integral with respect to  $\beta$  can be written as

$$-\frac{N}{2}\ln(2\pi) - \frac{1}{2}\ln(\det \mathbf{S}_0) + \frac{1}{2}\ln(\det \mathbf{S}_N) - \ln\Gamma(a_0) + a_0\ln b_0 + \ln\Gamma(a_N) - a_N\ln b_N.$$
(3.194)

Thus,

$$p(\mathbf{t}) = (2\pi)^{-\frac{N}{2}} \left( \frac{\det \mathbf{S}_N}{\det \mathbf{S}_0} \right)^{\frac{1}{2}} \frac{\Gamma(a_N)}{\Gamma(a_0)} \frac{b_0^{a_0}}{b_N^{a_N}}.$$
 (3.195)

# 3.24

Refer to 3.12.

# 4 Linear Models for Classification

## 4.1

Let  $x_1, \dots, x_M$  and  $y_1, \dots, y_N$  be two sets of data points. Then, the corresponding convex hulls are defined as the sets of all points  $\mathbf{x}$  and  $\mathbf{y}$  such that

$$\mathbf{x} = \sum_{m=1}^{M} \alpha_m \mathbf{x}_m,$$

$$\mathbf{y} = \sum_{n=1}^{N} \beta_n \mathbf{y}_n,$$
(4.1)

where

$$\sum_{m=1}^{M} \alpha_m = \sum_{n=1}^{N} \beta_n = 1,$$

$$\alpha_m \ge 0, \beta_n \ge 0.$$
(4.2)

Let us assume that  $\alpha_1, \dots, \alpha_M$  and  $\beta_1, \dots, \beta_N$  below are subject to the constraints above.

If the convex hulls intersect, then there exist  $\alpha_1, \dots, \alpha_M$  and  $\beta_1, \dots, \beta_N$  such that

$$\sum_{m=1}^{M} \alpha_m \mathbf{x}_m = \sum_{n=1}^{N} \beta_n \mathbf{y}_n. \tag{4.3}$$

Then,

$$\sum_{m=1}^{M} \alpha_m \left( \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{x}_m + w_0 \right) = \hat{\mathbf{w}}^{\mathsf{T}} \sum_{m=1}^{M} \alpha_m \mathbf{x}_m + w_0 \sum_m \alpha_m, \tag{4.4}$$

for any  $\hat{\mathbf{w}}$  and  $w_0$ . The right hand side can be written as

$$\hat{\mathbf{w}}^{\mathsf{T}} \sum_{n=1}^{N} \beta_n \mathbf{y}_n + w_0 \sum_{n=1}^{N} \beta_n = \sum_{n=1}^{N} \beta_n \left( \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{y}_n + w_0 \right). \tag{4.5}$$

Therefore, there do not exist  $\hat{\mathbf{w}}$  and  $w_0$  such that

$$\hat{\mathbf{w}}^{\mathsf{T}} \mathbf{x}_m + w_0 > 0, \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{y}_n + w_0 < 0.$$
 (4.6)

Conversely, if there exist  $\hat{\mathbf{w}}$  and  $w_0$  such that

$$\hat{\mathbf{w}}^{\mathsf{T}} \mathbf{x}_m + w_0 > 0, \\ \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{y}_n + w_0 < 0,$$
 (4.7)

then

$$\sum_{m=1}^{M} \alpha_m \left( \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{x}_m + w_0 \right) > 0,$$

$$\sum_{n=1}^{N} \beta_n \left( \hat{\mathbf{w}}^{\mathsf{T}} \mathbf{y}_n + w_0 \right) < 0.$$
(4.8)

The left hand sides can be written as

$$\hat{\mathbf{w}}^{\mathsf{T}} \sum_{m=1}^{M} \alpha_m \mathbf{x}_m + w_0 \sum_{m=1}^{M} \alpha_m = \hat{\mathbf{w}}^{\mathsf{T}} \sum_{m=1}^{M} \alpha_m \mathbf{x}_m + w_0,$$

$$\hat{\mathbf{w}}^{\mathsf{T}} \sum_{n=1}^{N} \beta_n \mathbf{y}_n + w_0 \sum_{n=1}^{N} \beta_n = \hat{\mathbf{w}}^{\mathsf{T}} \sum_{n=1}^{N} \beta_n \mathbf{y}_n + w_0.$$
(4.9)

Therefore, there do not exist  $\alpha_1, \dots, \alpha_M$  and  $\beta_1, \dots, \beta_N$  such that

$$\sum_{m=1}^{M} \alpha_m \mathbf{x}_m = \sum_{n=1}^{N} \beta_n \mathbf{y}_n. \tag{4.10}$$

Thus, the convex hulls do not intersect.

# 4.2 (Incomplete)

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  and  $\mathbf{w}_1, \dots, \mathbf{w}_K$  are variables in M dimensions and  $\mathbf{t}_1, \dots, \mathbf{t}_N$  are ones in K dimensions. Let

$$E(\tilde{\mathbf{W}}) = \frac{1}{2} \operatorname{tr} \left( (\tilde{\mathbf{X}} \tilde{\mathbf{W}} - \mathbf{T})^{\mathsf{T}} (\tilde{\mathbf{X}} \tilde{\mathbf{W}} - \mathbf{T}) \right), \tag{4.11}$$

where

$$\tilde{\mathbf{X}} = \begin{bmatrix} 1 & \mathbf{x}_1^{\mathsf{T}} \\ \vdots & \vdots \\ 1 & \mathbf{x}_N^{\mathsf{T}} \end{bmatrix},$$

$$\tilde{\mathbf{W}} = \begin{bmatrix} w_{10} & \cdots & w_{K0} \\ \mathbf{w}_1 & \cdots & \mathbf{w}_K \end{bmatrix}$$

and

$$\mathbf{T} = egin{bmatrix} \mathbf{t}_1^\intercal \ dots \ \mathbf{t}_N^\intercal \end{bmatrix}.$$

Setting the derivative with respect to  $\tilde{\mathbf{W}}$  to zero gives

$$\mathbf{O} = \tilde{\mathbf{X}}^{\mathsf{T}} (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}). \tag{4.12}$$

Therefore,

$$\underset{\tilde{\mathbf{W}}}{\operatorname{argmin}} E(\tilde{\mathbf{W}}) = \left(\tilde{\mathbf{X}}^{\mathsf{T}} \tilde{\mathbf{X}}\right)^{-1} \tilde{\mathbf{X}}^{\mathsf{T}} \mathbf{T}. \tag{4.13}$$

Let  $\tilde{\mathbf{W}}^*$  denote the least-square solution above. Then,

$$(\tilde{\mathbf{W}}^*)^{\mathsf{T}}\tilde{\mathbf{x}} - \mathbf{t}_n = \mathbf{T}^{\mathsf{T}}\tilde{\mathbf{X}} \left( \tilde{\mathbf{X}}^{\mathsf{T}}\tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{x}} - \mathbf{t}_n, \tag{4.14}$$

where  $\tilde{\mathbf{x}}$  is a vector in M+1 dimensions whose first element is 1. The right hand side can be written as

$$\mathbf{T}^{\mathsf{T}} \left( \tilde{\mathbf{X}} \left( \tilde{\mathbf{X}}^{\mathsf{T}} \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{x}} - \mathbf{v}_n \right) = \mathbf{0}? \tag{4.15}$$

where  $\mathbf{v}_n$  is a vector in N dimensions whose n th element is 1 and other elements are zero. Therefore,

$$(\tilde{\mathbf{W}}^*)^{\mathsf{T}}\tilde{\mathbf{x}} - \mathbf{t}_n = \mathbf{0}. \tag{4.16}$$

Thus, if

$$\mathbf{a}^{\mathsf{T}}\mathbf{t}_n + b = 0,\tag{4.17}$$

then

$$\mathbf{a}^{\mathsf{T}}(\tilde{\mathbf{W}}^*)^{\mathsf{T}}\tilde{\mathbf{x}} + b = 0. \tag{4.18}$$

# 4.3 (Incomplete)

#### 4.4

Let  $\mathbf{x}_1, \cdots, \mathbf{x}_N$  be variables and let

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n,\tag{4.19}$$

where  $N_k$  is the number of  $\mathbf{x}_n$  such that n is in  $\mathcal{C}_k$ . Setting the derivatives of

$$\mathbf{w}^{\mathsf{T}}(\mathbf{m}_2 - \mathbf{m}_1) + \lambda \left( \|\mathbf{w}\|^2 - 1 \right) \tag{4.20}$$

with respect to  $\mathbf{w}$  and  $\lambda$  to zero gives

$$\mathbf{m}_2 - \mathbf{m}_1 + 2\lambda \mathbf{w} = \mathbf{0},$$
  
$$\|\mathbf{w}\|^2 - 1 = 0.$$
 (4.21)

Therefore,  $\mathbf{w}^{\intercal}(\mathbf{m}_2 - \mathbf{m}_1)$  under the constratint

$$\|\mathbf{w}\|^2 = 1\tag{4.22}$$

is maximised if

$$\mathbf{w} \propto \mathbf{m}_2 - \mathbf{m}_1. \tag{4.23}$$

#### 4.5

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be variables and let

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n,\tag{4.24}$$

where  $N_k$  is the number of  $\mathbf{x}_n$  such that n is in  $C_k$ . Let

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2},\tag{4.25}$$

where

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2,$$

$$y_n = \mathbf{w}^{\mathsf{T}} \mathbf{x}_n,$$

$$m_k = \mathbf{w}^{\mathsf{T}} \mathbf{m}_k.$$

$$(4.26)$$

Then,  $J(\mathbf{w})$  can be written as

$$\frac{\left(\mathbf{w}^{\mathsf{T}}(\mathbf{m}_{2} - \mathbf{m}_{1})\right)^{2}}{\sum_{n \in \mathcal{C}_{1}} \left(\mathbf{w}^{\mathsf{T}}(\mathbf{x}_{n} - \mathbf{m}_{1})\right)^{2} + \sum_{n \in \mathcal{C}_{2}} \left(\mathbf{w}^{\mathsf{T}}(\mathbf{x}_{n} - \mathbf{m}_{2})\right)^{2}} = \frac{\mathbf{w}^{\mathsf{T}} \mathbf{S}_{B} \mathbf{w}}{\mathbf{w}^{\mathsf{T}} \mathbf{S}_{W} \mathbf{w}}, \quad (4.27)$$

where

$$\mathbf{S}_{\mathrm{B}} = (\mathbf{m}_{2} - \mathbf{m}_{1})(\mathbf{m}_{2} - \mathbf{m}_{1})^{\mathsf{T}},$$

$$\mathbf{S}_{\mathrm{W}} = \sum_{n \in \mathcal{C}_{1}} (\mathbf{x}_{n} - \mathbf{m}_{1})(\mathbf{x}_{n} - \mathbf{m}_{1})^{\mathsf{T}} + \sum_{n \in \mathcal{C}_{2}} (\mathbf{x}_{n} - \mathbf{m}_{2})(\mathbf{x}_{n} - \mathbf{m}_{2})^{\mathsf{T}}.$$

$$(4.28)$$

#### 4.6

Let  $\mathbf{x}_1, \dots, \mathbf{x}_N$  be variables and let

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{x}_n, \tag{4.29}$$

where  $N_k$  is the number of  $\mathbf{x}_n$  such that n is in  $\mathcal{C}_k$ . Let

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2},\tag{4.30}$$

where

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2,$$

$$y_n = \mathbf{w}^{\mathsf{T}} \mathbf{x}_n,$$

$$m_k = \mathbf{w}^{\mathsf{T}} \mathbf{m}_k.$$

$$(4.31)$$

Then, by 4.5,

$$J(\mathbf{w}) = \frac{\mathbf{w}^{\mathsf{T}} \mathbf{S}_{\mathsf{B}} \mathbf{w}}{\mathbf{w}^{\mathsf{T}} \mathbf{S}_{\mathsf{W}} \mathbf{w}},\tag{4.32}$$

where

$$\mathbf{S}_{\mathrm{B}} = (\mathbf{m}_{2} - \mathbf{m}_{1})(\mathbf{m}_{2} - \mathbf{m}_{1})^{\mathsf{T}},$$

$$\mathbf{S}_{\mathrm{W}} = \sum_{n \in \mathcal{C}_{1}} (\mathbf{x}_{n} - \mathbf{m}_{1})(\mathbf{x}_{n} - \mathbf{m}_{1})^{\mathsf{T}} + \sum_{n \in \mathcal{C}_{2}} (\mathbf{x}_{n} - \mathbf{m}_{2})(\mathbf{x}_{n} - \mathbf{m}_{2})^{\mathsf{T}}.$$

$$(4.33)$$

Let

$$E = \frac{1}{2} \sum_{n=1}^{N} (\mathbf{w}^{\mathsf{T}} \mathbf{x}_n + w_0 - t_n)^2, \qquad (4.34)$$

where

$$t_n = \begin{cases} \frac{N}{N_1}, & n \in \mathcal{C}_1, \\ -\frac{N}{N_2}, & n \in \mathcal{C}_2. \end{cases}$$
 (4.35)

Setting the derivative with respect to  $\mathbf{w}$  and  $w_0$  gives

$$0 = \sum_{n=1}^{N} (\mathbf{w}^{\mathsf{T}} \mathbf{x}_n + w_0 - t_n),$$

$$\mathbf{0} = \sum_{n=1}^{N} (\mathbf{w}^{\mathsf{T}} \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n.$$

$$(4.36)$$

The right hand side of the first equation can be written as

$$\mathbf{w}^{\mathsf{T}} \sum_{n=1}^{N} \mathbf{x}_n + N w_0 - \sum_{n=1}^{N} t_n = N \left( \mathbf{w}^{\mathsf{T}} \mathbf{m} + w_0 \right), \tag{4.37}$$

where

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n. \tag{4.38}$$

Therefore,

$$w_0 = -\mathbf{w}^{\mathsf{T}}\mathbf{m}.\tag{4.39}$$

Then, the right hand side of the second equation above can be written as

$$\sum_{n=1}^{N} (\mathbf{w}^{\mathsf{T}} (\mathbf{x}_{n} - \mathbf{m}) - t_{n}) \mathbf{x}_{n}$$

$$= \sum_{n \in C_{1}} \left( \mathbf{w}^{\mathsf{T}} (\mathbf{x}_{n} - \mathbf{m}) - \frac{N}{N_{1}} \right) \mathbf{x}_{n} + \sum_{n \in C_{2}} \left( \mathbf{w}^{\mathsf{T}} (\mathbf{x}_{n} - \mathbf{m}) + \frac{N}{N_{2}} \right) \mathbf{x}_{n}.$$
(4.40)

Since

$$\mathbf{m} = \frac{N_1}{N} \mathbf{m}_1 + \frac{N_2}{N} \mathbf{m}_2,$$

$$\sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1) = \mathbf{0},$$
(4.41)

the first term of the right hand side can be written as

$$\sum_{n \in \mathcal{C}_1} \left( \mathbf{w}^{\mathsf{T}} \left( \mathbf{x}_n - \mathbf{m}_1 + \frac{N_2}{N} (\mathbf{m}_1 - \mathbf{m}_2) \right) - \frac{N}{N_1} \right) (\mathbf{x}_n - \mathbf{m}_1 + \mathbf{m}_1)$$

$$= \left( \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1) (\mathbf{x}_n - \mathbf{m}_1)^{\mathsf{T}} \right) \mathbf{w} + \frac{N_1 N_2}{N} (\mathbf{m}_1 - \mathbf{m}_2) \mathbf{m}_1^{\mathsf{T}} \mathbf{w} - N \mathbf{m}_1.$$
(4.42)

Similarly, the second term can be written as

$$\left(\sum_{n\in\mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^{\mathsf{T}}\right) \mathbf{w} + \frac{N_1 N_2}{N} (\mathbf{m}_2 - \mathbf{m}_1) \mathbf{m}_2^{\mathsf{T}} \mathbf{w} - N \mathbf{m}_2. \quad (4.43)$$

Therefore,

$$\mathbf{0} = \left(\sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^{\mathsf{T}}\right) \mathbf{w} + \frac{N_1 N_2}{N} (\mathbf{m}_1 - \mathbf{m}_2) \mathbf{m}_1^{\mathsf{T}} \mathbf{w} - N \mathbf{m}_1$$
$$+ \left(\sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^{\mathsf{T}}\right) \mathbf{w} + \frac{N_1 N_2}{N} (\mathbf{m}_2 - \mathbf{m}_1) \mathbf{m}_2^{\mathsf{T}} \mathbf{w} - N \mathbf{m}_2.$$

$$(4.44)$$

Thus,

$$\left(\mathbf{S}_{W} + \frac{N_1 N_2}{N} \mathbf{S}_{B}\right) \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2). \tag{4.45}$$

#### 4.7

Let

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.\tag{4.46}$$

Then,

$$\sigma(-a) = \frac{1}{1 + \exp(a)}.\tag{4.47}$$

The right hand side can be written as

$$1 - \frac{\exp(a)}{1 + \exp(a)} = 1 - \frac{1}{1 + \exp(-a)}.$$
 (4.48)

Therefore,

$$\sigma(-a) = 1 - \sigma(a). \tag{4.49}$$

Additionally,

$$\exp(-a) = \frac{1}{\sigma(a)} - 1. \tag{4.50}$$

Then,

$$a = -\ln\left(\frac{1}{\sigma(a)} - 1\right). \tag{4.51}$$

Therefore,

$$\sigma^{-1}(y) = \ln\left(\frac{y}{1-y}\right). \tag{4.52}$$

#### 4.8

Let  $\mathbf{x}$  be a variable in D dimensions such that

$$p(\mathbf{x}|\mathcal{C}_k) = \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}\right), \tag{4.53}$$

where

$$p(\mathcal{C}_1) + p(\mathcal{C}_2) = 1. \tag{4.54}$$

By the Bayes' theorem,

$$p(C_1|\mathbf{x}) = \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)}.$$
(4.55)

The right hand side can be written as

$$\sigma(a) = \frac{1}{1 + \exp(-a)},\tag{4.56}$$

where

$$a = \ln \left( \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \right). \tag{4.57}$$

Substituting the expressions above of  $p(\mathbf{x}|\mathcal{C}_k)$ , we have

$$a = -\frac{D}{2}\ln(2\pi) - \frac{1}{2}\ln(\det\Sigma) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \ln p(\mathcal{C}_1) + \frac{D}{2}\ln(2\pi) + \frac{1}{2}\ln(\det\Sigma) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) - \ln p(\mathcal{C}_2).$$

$$(4.58)$$

Therefore,

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma\left(\mathbf{w}^{\mathsf{T}}\mathbf{x} + w_0\right),\tag{4.59}$$

where

$$\mathbf{w} = \mathbf{\Sigma}^{-1}(\boldsymbol{\mu}_{1} - \boldsymbol{\mu}_{2}),$$

$$w_{0} = -\frac{1}{2}\boldsymbol{\mu}_{1}^{\mathsf{T}}\mathbf{\Sigma}^{-1}\boldsymbol{\mu}_{1} + \frac{1}{2}\boldsymbol{\mu}_{2}^{\mathsf{T}}\mathbf{\Sigma}^{-1}\boldsymbol{\mu}_{2} + \ln p(\mathcal{C}_{1}) - \ln p(\mathcal{C}_{2}).$$
(4.60)

#### 4.9

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables in K dimensions such that

$$p(\mathbf{t}_n, \boldsymbol{\phi}_n) = \prod_{k=1}^K (p(\boldsymbol{\phi}_n, \mathcal{C}_k))^{t_{nk}}, \qquad (4.61)$$

where

$$\sum_{k=1}^{K} p(\mathcal{C}_k) = 1. (4.62)$$

Then,

$$p(\mathbf{T}, \mathbf{\Phi}) = \prod_{n=1}^{N} \prod_{k=1}^{K} (p(\boldsymbol{\phi}_n, \mathcal{C}_k))^{t_{nk}}.$$
 (4.63)

If

$$p(\mathcal{C}_k) = \pi_k, \tag{4.64}$$

then, by the Bayes' theorem,

$$\ln p(\mathbf{T}, \mathbf{\Phi}) = \sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \left( \ln \pi_k + \ln p(\boldsymbol{\phi}_n | \mathcal{C}_k) \right). \tag{4.65}$$

Setting the derivatives of

$$\ln p(\mathbf{T}, \mathbf{\Phi}) + \lambda \left( \sum_{k=1}^{K} \pi_k - 1 \right)$$
 (4.66)

with respect to  $\pi_k$  and  $\lambda$  to zero gives

$$0 = \frac{1}{\pi_k} \sum_{n=1}^{N} t_{nk} + \lambda,$$

$$0 = \sum_{k=1}^{K} \pi_k - 1.$$
(4.67)

Then,

$$\lambda = -\sum_{k=1}^{K} \sum_{n=1}^{N} t_{nk}.$$
(4.68)

The right hand side can be written as -N. Therefore, the maximum likelihood solution for  $\pi_k$  is given by

$$\pi_k = \frac{N_k}{N},\tag{4.69}$$

where

$$N_k = \sum_{n=1}^{N} t_{nk}. (4.70)$$

#### 4.10

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables in K dimensions such that

$$p(\mathbf{t}_n, \boldsymbol{\phi}_n) = \prod_{k=1}^K (p(\boldsymbol{\phi}_n, \mathcal{C}_k))^{t_{nk}}, \qquad (4.71)$$

where

$$\sum_{k=1}^{K} p(\mathcal{C}_k) = 1. (4.72)$$

Then,

$$p(\mathbf{T}, \mathbf{\Phi}) = \prod_{n=1}^{N} \prod_{k=1}^{K} \left( p(\boldsymbol{\phi}_n, \mathcal{C}_k) \right)^{t_{nk}}.$$
 (4.73)

If

$$p(\phi_n|\mathcal{C}_k) = \mathcal{N}(\phi_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}), \tag{4.74}$$

then, by the Bayes' theorem,

$$\ln p(\mathbf{T}, \mathbf{\Phi}) = \sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \left( \ln \mathcal{N}(\boldsymbol{\phi}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) + \ln p(\mathcal{C}_k) \right). \tag{4.75}$$

The right hand side can be written as

$$\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \left( -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln(\det \Sigma) - \frac{1}{2} (\boldsymbol{\phi}_{n} - \boldsymbol{\mu}_{k})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\boldsymbol{\phi}_{n} - \boldsymbol{\mu}_{k}) + \ln p(\mathcal{C}_{k}) \right). \tag{4.76}$$

By 3.21(a), setting the derivatives of  $\ln p(\mathbf{T}, \mathbf{\Phi})$  with respect to  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}$  to zero gives

$$\mathbf{0} = \frac{1}{2} \sum_{n=1}^{N} t_{nk} \left( \mathbf{\Sigma}^{-1} + \left( \mathbf{\Sigma}^{-1} \right)^{\mathsf{T}} \right) (\boldsymbol{\phi}_{n} - \boldsymbol{\mu}_{k}),$$

$$\mathbf{O} = -\frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \left( \left( \mathbf{\Sigma}^{-1} \right)^{\mathsf{T}} - \left( \mathbf{\Sigma}^{-1} \right)^{2} (\boldsymbol{\phi}_{n} - \boldsymbol{\mu}_{k}) (\boldsymbol{\phi}_{n} - \boldsymbol{\mu}_{k})^{\mathsf{T}} \right).$$

$$(4.77)$$

Therefore, the maximum likelihood solutions for  $\mu_k$  and  $\Sigma$  are given by

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^{N} t_{nk} \phi_n,$$

$$\Sigma = \frac{1}{N} \sum_{k=1}^{K} N_k \mathbf{S}_k,$$
(4.78)

where

$$N_k = \sum_{n=1}^{N} t_{nk},$$

$$\mathbf{S}_k = \frac{1}{N_k} \sum_{n=1}^{N} t_{nk} (\boldsymbol{\phi}_n - \boldsymbol{\mu}_k) (\boldsymbol{\phi}_n - \boldsymbol{\mu}_k)^{\mathsf{T}}.$$

$$(4.79)$$

#### 4.11

Let  $\phi_1, \dots, \phi_M$  be variables such that

$$p(\boldsymbol{\phi}_m|\mathcal{C}_k) = \prod_{l=1}^L \mu_{kml}^{\phi_{ml}}, \tag{4.80}$$

where

$$\sum_{k=1}^{K} p(\mathcal{C}_k) = 1. \tag{4.81}$$

Then,

$$p(\mathbf{\Phi}|\mathcal{C}_k) = \prod_{m=1}^{M} \prod_{l=1}^{L} \mu_{kml}^{\phi_{ml}}.$$
 (4.82)

By the Bayes' theorem,

$$p(C_k|\mathbf{\Phi}) = \frac{p(\mathbf{\Phi}|C_k)p(C_k)}{\sum_{k=1}^K p(\mathbf{\Phi}|C_k)p(C_k)}.$$
(4.83)

Therefore,

$$p(C_k|\mathbf{\Phi}) = \frac{\exp(a_k(\mathbf{\Phi}))}{\sum_{k=1}^K \exp(a_k(\mathbf{\Phi}))},$$
(4.84)

where

$$a_k(\mathbf{\Phi}) = \left(\sum_{m=1}^M \sum_{l=1}^L \phi_{ml} \ln \mu_{kml}\right) + \ln p(\mathcal{C}_k). \tag{4.85}$$

#### 4.12

Let

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.\tag{4.86}$$

Then,

$$\frac{d\sigma(a)}{da} = \frac{\exp(-a)}{\left(1 + \exp(-a)\right)^2}.$$
(4.87)

The right hand side can be written as

$$\frac{1}{1 + \exp(-a)} - \frac{1}{(1 + \exp(-a))^2} = \sigma(a) - (\sigma(a))^2.$$
 (4.88)

Therefore,

$$\frac{d\sigma(a)}{da} = \sigma(a) \left(1 - \sigma(a)\right). \tag{4.89}$$

#### 4.13

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{0, 1\},\ p(t_n | \mathbf{w}) = y_n^{t_n} (1 - y_n)^{1 - t_n},$$
(4.90)

where

$$y_n = \sigma(\mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n),$$
  
$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$
 (4.91)

Let

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}). \tag{4.92}$$

The right hand side can be written as

$$-\ln\left(\prod_{n=1}^{N} y_n^{t_n} (1 - y_n)^{1 - t_n}\right) = -\sum_{n=1}^{N} \left(t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\right). \quad (4.93)$$

Then, by 4.12,

$$\nabla E(\mathbf{w}) = -\sum_{n=1}^{N} \left( \frac{t_n}{y_n} y_n (1 - y_n) \phi_n - \frac{1 - t_n}{1 - y_n} y_n (1 - y_n) \phi_n \right). \tag{4.94}$$

The right hand side can be written as

$$-\sum_{n=1}^{N} (t_n(1-y_n)\boldsymbol{\phi}_n - (1-t_n)y_n\boldsymbol{\phi}_n) = \sum_{n=1}^{N} (y_n - t_n)\boldsymbol{\phi}_n.$$
 (4.95)

Therefore,

$$\nabla E(\mathbf{w}) = \sum_{n=1}^{N} (y_n - t_n) \phi_n. \tag{4.96}$$

# 4.14

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{0, 1\},\ p(t_n | \mathbf{w}) = y_n^{t_n} (1 - y_n)^{1 - t_n},$$
(4.97)

where

$$y_n = \sigma(\mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n),$$
  
$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$
 (4.98)

Let

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}). \tag{4.99}$$

By 4.13, setting the derivative with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = \sum_{n=1}^{N} (y_n - t_n) \, \phi_n. \tag{4.100}$$

If  $\phi_1, \cdots, \phi_N$  are linearly independent, then

$$y_n = t_n. (4.101)$$

Then,

$$\sigma\left(\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}_{n}\right) = \begin{cases} 1, & t_{n} = 1, \\ 0, & \text{otherwise.} \end{cases}$$
 (4.102)

Threrefore,

$$\mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n = \begin{cases} \infty, & t_n = 1, \\ -\infty, & \text{otherwise.} \end{cases}$$
 (4.103)

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{0, 1\},$$
  
 $p(t_n | \mathbf{w}) = y_n^{t_n} (1 - y_n)^{1 - t_n},$  (4.104)

where

$$y_n = \sigma(\mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n),$$
  
$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$
 (4.105)

Let

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}). \tag{4.106}$$

By 4.13,

$$\nabla E(\mathbf{w}) = \sum_{n=1}^{N} (y_n - t_n) \phi_n. \tag{4.107}$$

Then, by 4.12,

$$\nabla \nabla E(\mathbf{w}) = \sum_{n=1}^{N} y_n (1 - y_n) \phi_n \phi_n^{\mathsf{T}}.$$
 (4.108)

The right hand side can be written as

$$\mathbf{H} = \mathbf{\Phi}^{\mathsf{T}} \mathbf{R} \mathbf{\Phi},\tag{4.109}$$

where

$$R_{nn'} = \begin{cases} y_n(1 - y_n), & n = n', \\ 0, & \text{otherwise.} \end{cases}$$
 (4.110)

Then,

$$\mathbf{u}^{\mathsf{T}}\mathbf{H}\mathbf{u} = (\mathbf{\Phi}\mathbf{u})^{\mathsf{T}}\mathbf{R}(\mathbf{\Phi}\mathbf{u}). \tag{4.111}$$

Since

$$y_n(1-y_n) > 0, (4.112)$$

we have

$$\mathbf{u}^{\mathsf{T}}\mathbf{H}\mathbf{u} > 0. \tag{4.113}$$

Therefore,  $\mathbf{H}$  is positive definite. Thus, E is a convex function of  $\mathbf{w}$  and it has a unique minimum.

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{0, 1\},\ p(t_n = 1 | \phi_n) = \pi_n.$$
 (4.114)

Then,

$$p(t_n|\phi_n) = \pi_n^{t_n} (1 - \pi_n)^{1 - t_n}.$$
(4.115)

Therefore,

$$p(\mathbf{t}|\mathbf{\Phi}) = \prod_{n=1}^{N} \pi_n^{t_n} (1 - \pi_n)^{1 - t_n}.$$
 (4.116)

Thus,

$$-\ln p(\mathbf{t}|\mathbf{\Phi}) = -\sum_{n=1}^{N} (t_n \ln \pi_n + (1 - t_n) \ln(1 - \pi_n)). \tag{4.117}$$

# 4.17

Let

$$y_k = \frac{\exp(a_k)}{\sum_{k=1}^K \exp(a_k)}.$$
 (4.118)

Then,

$$\frac{\partial y_k}{\partial a_k} = \frac{\exp(a_k)}{\sum_{k=1}^K \exp(a_k)} - \frac{\exp(2a_k)}{\left(\sum_{k=1}^K \exp(a_k)\right)^2}.$$
 (4.119)

The right hand side can be written as  $y_k(1-y_k)$ . If  $k \neq k'$ , then

$$\frac{\partial y_k}{\partial a_{k'}} = -\frac{\exp(a_k + a_{k'})}{\left(\sum_{k=1}^K \exp(a_k)\right)^2}.$$
(4.120)

The right hand side can be written as  $-y_k y_{k'}$ . Therefore,

$$\frac{\partial y_k}{\partial a_{k'}} = y_k (I_{kk'} - y_{k'}). \tag{4.121}$$

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables such that

$$t_{nk} \in \{0, 1\},$$

$$p(\mathbf{t}_n | \mathbf{W}) = \prod_{k=1}^K y_{nk}^{t_{nk}},$$
(4.122)

where

$$y_{nk} = \frac{\exp(a_{nk})}{\sum_{k=1}^{K} \exp(a_{nk})},$$

$$a_{nk} = \mathbf{w}_{1}^{\mathsf{T}} \phi_{n}.$$
(4.123)

Then,

$$p(\mathbf{T}|\mathbf{W}) = \prod_{n=1}^{N} \prod_{k=1}^{K} y_{nk}^{t_{nk}}.$$
 (4.124)

Let

$$E(\mathbf{W}) = -\ln p(\mathbf{T}|\mathbf{W}). \tag{4.125}$$

The right hand side can be written as

$$-\sum_{n=1}^{N}\sum_{k=1}^{K}t_{nk}\ln y_{nk}.$$
(4.126)

Then, by 4.17,

$$\nabla_{\mathbf{w}_{k'}} E(\mathbf{W}) = -\sum_{n=1}^{N} \sum_{k=1}^{K} y_{nk} (I_{kk'} - y_{nk'}) \frac{t_{nk}}{y_{nk}} \phi_n.$$
 (4.127)

The right hand side can be written as

$$-\sum_{n=1}^{N} \left( \sum_{k=1}^{K} (I_{kk'} - y_{nk'}) t_{nk} \right) \phi_n = -\sum_{n=1}^{N} (t_{nk'} - y_{nk'}) \phi_n.$$
 (4.128)

Therefore,

$$\nabla_{\mathbf{w}_k} E(\mathbf{W}) = \sum_{n=1}^N (y_{nk} - t_{nk}) \, \boldsymbol{\phi}_n. \tag{4.129}$$

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{0, 1\},\ p(t_n = 1|a_n) = \Phi(a_n),$$
 (4.130)

where

$$\Phi(a) = \int_{-\infty}^{a} \mathcal{N}(\theta|0,1)d\theta,$$

$$a_n = \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n.$$
(4.131)

Then,

$$p(t_n|\phi_n) = (\Phi(a_n))^{t_n} (1 - \Phi(a_n))^{1-t_n}.$$
(4.132)

Therefore,

$$p(\mathbf{t}|\mathbf{\Phi}) = \prod_{n=1}^{N} (\Phi(a_n))^{t_n} (1 - \Phi(a_n))^{1-t_n}.$$
 (4.133)

Let

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\boldsymbol{\phi}). \tag{4.134}$$

The right hand side can be written as

$$-\sum_{n=1}^{N} (t_n \ln \Phi(a_n) + (1 - t_n) \ln (1 - \Phi(a_n))). \tag{4.135}$$

Then,

$$\nabla E(\mathbf{w}) = -\sum_{n=1}^{N} \left( t_n \frac{\mathcal{N}(a_n|0,1)}{\Phi(a_n)} - (1 - t_n) \frac{\mathcal{N}(a_n|0,1)}{1 - \Phi(a_n)} \right) \phi_n.$$
 (4.136)

The right hand side can be written as

$$-\sum_{n=1}^{N} \left( \frac{t_n}{\Phi(a_n)} - \frac{1 - t_n}{1 - \Phi(a_n)} \right) \mathcal{N}(a_n | 0, 1) \phi_n$$

$$= \sum_{n=1}^{N} \frac{\mathcal{N}(a_n | 0, 1)}{\Phi(a_n) (1 - \Phi(a_n))} \left( \Phi(a_n) - t_n \right) \phi_n.$$
(4.137)

Therefore,

$$\nabla E(\mathbf{w}) = \sum_{n=1}^{N} \frac{\mathcal{N}(a_n|0,1)}{\Phi(a_n) (1 - \Phi(a_n))} (\Phi(a_n) - t_n) \, \phi_n. \tag{4.138}$$

Then,

$$\nabla \nabla E(\mathbf{w}) = \sum_{n=1}^{N} \frac{-a_n \mathcal{N}(a_n | 0, 1)}{\Phi(a_n) (1 - \Phi(a_n))} (\Phi(a_n) - t_n) \phi_n \phi_n^{\mathsf{T}}$$

$$- \sum_{n=1}^{N} \frac{(\mathcal{N}(a_n | 0, 1))^2}{(\Phi(a_n))^2 (1 - \Phi(a_n))} (\Phi(a_n) - t_n) \phi_n \phi_n^{\mathsf{T}}$$

$$+ \sum_{n=1}^{N} \frac{(\mathcal{N}(a_n | 0, 1))^2}{\Phi(a_n) (1 - \Phi(a_n))^2} (\Phi(a_n) - t_n) \phi_n \phi_n^{\mathsf{T}}$$

$$+ \sum_{n=1}^{N} \frac{(\mathcal{N}(a_n | 0, 1))^2}{\Phi(a_n) (1 - \Phi(a_n))} \phi_n \phi_n^{\mathsf{T}}.$$
(4.139)

Therefore,

$$\nabla \nabla E(\mathbf{w}) = \sum_{n=1}^{N} b_n \phi_n \phi_n^{\mathsf{T}}, \tag{4.140}$$

where

$$b_{n} = \left(\frac{\mathcal{N}(a_{n}|0,1)}{\Phi(a_{n})(1-\Phi(a_{n}))}\right)^{2} ((\Phi(a_{n}))^{2} - 2t_{n}\Phi(a_{n}) + t_{n})$$

$$-\frac{\mathcal{N}(a_{n}|0,1)}{\Phi(a_{n})(1-\Phi(a_{n}))} a_{n} (\Phi(a_{n}) - t_{n}).$$
(4.141)

# 4.20 (Incomplete)

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables such that

$$t_{nk} \in \{0, 1\},$$

$$p(\mathbf{t}_n | \mathbf{W}) = \prod_{k=1}^{K} y_{nk}^{t_{nk}},$$
(4.142)

where

$$y_{nk} = \frac{\exp(a_{nk})}{\sum_{k=1}^{K} \exp(a_{nk})},$$

$$a_{nk} = \mathbf{w}_{k}^{\mathsf{T}} \boldsymbol{\phi}_{n}.$$

$$(4.143)$$

Then,

$$p(\mathbf{T}|\mathbf{W}) = \prod_{n=1}^{N} \prod_{k=1}^{K} y_{nk}^{t_{nk}}.$$
 (4.144)

Let

$$E(\mathbf{W}) = -\ln p(\mathbf{T}|\mathbf{W}). \tag{4.145}$$

By 4.18,

$$\nabla_{\mathbf{w}_k} E(\mathbf{W}) = \sum_{n=1}^N (y_{nk} - t_{nk}) \, \boldsymbol{\phi}_n. \tag{4.146}$$

Additionally, by 4.17,

$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_{k'}} E(\mathbf{W}) = \sum_{n=1}^N y_{nk} (I_{kk'} - y_{nk'}) \phi_n \phi_n^{\mathsf{T}}. \tag{4.147}$$

The right hand side can be written as

$$\mathbf{H}_{kk'} = \mathbf{\Phi}^{\mathsf{T}} \mathbf{R}_{kk'} \mathbf{\Phi}, \tag{4.148}$$

where

$$R_{kk'nn'} = \begin{cases} y_{nk}(I_{kk'} - y_{nk'}), & n = n', \\ 0, & \text{otherwise.} \end{cases}$$
 (4.149)

Let

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_{11} & \cdots & \mathbf{H}_{1K} \\ \vdots & \ddots & \vdots \\ \mathbf{H}_{K1} & \cdots & \mathbf{H}_{KK} \end{bmatrix},$$

and

$$\mathbf{u} = \begin{bmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_K \end{bmatrix},$$

where  $\mathbf{u}_1, \cdots, \mathbf{u}_K$  are vectors in the same dimension as  $\mathbf{w}$ . Then,

$$\mathbf{u}^{\mathsf{T}}\mathbf{H}\mathbf{u} = \sum_{k=1}^{K} \sum_{k'=1}^{K} \mathbf{u}_{k}^{\mathsf{T}}\mathbf{H}_{kk'}\mathbf{u}_{k'}, \tag{4.150}$$

Then, the right hand side can be written as

$$\sum_{k=1}^{K} \sum_{k'=1}^{K} (\mathbf{\Phi} \mathbf{u}_k)^{\mathsf{T}} \mathbf{R}_{kk'} (\mathbf{\Phi} \mathbf{u}_{k'}). \tag{4.151}$$

#### 4.21

Let

$$\Phi(a) = \int_{-\infty}^{a} \mathcal{N}(\theta|0,1)d\theta. \tag{4.152}$$

The right hand side can be written as

$$\int_{-\infty}^{0} \mathcal{N}(\theta|0,1)d\theta + \int_{0}^{a} \mathcal{N}(\theta|0,1)d\theta = \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \int_{0}^{a} \exp\left(-\frac{\theta^{2}}{2}\right) d\theta. \quad (4.153)$$

The second term of the right hand side can be written as

$$\frac{1}{\sqrt{2\pi}} \int_0^{\frac{a}{\sqrt{2}}} \exp\left(-t^2\right) \sqrt{2} dt = \frac{1}{2} \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right), \tag{4.154}$$

where

$$\operatorname{erf}(a) = \frac{2}{\sqrt{\pi}} \int_0^a \exp(-t^2) dt. \tag{4.155}$$

Therefore,

$$\Phi(a) = \frac{1}{2} \left( 1 + \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \right). \tag{4.156}$$

# 4.22

Let  $\theta$  be a variable in M dimensions. By a Taylor expansion,

$$\ln (p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})) \simeq \ln (p(\mathcal{D}|\boldsymbol{\theta}_0)p(\boldsymbol{\theta}_0)) + \mathbf{v}(\boldsymbol{\theta}_0)^{\mathsf{T}}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) - \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^{\mathsf{T}}\mathbf{A}(\boldsymbol{\theta}_0)(\boldsymbol{\theta} - \boldsymbol{\theta}_0),$$
(4.157)

where

$$\mathbf{v}(\boldsymbol{\theta}) = \nabla \ln \left( p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) \right), \mathbf{A}(\boldsymbol{\theta}) = -\nabla \nabla \ln \left( p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) \right).$$
(4.158)

Let  $\boldsymbol{\theta}_{\text{MAP}}$  be a stationary point of  $p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$ . Then,

$$\ln (p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})) \simeq \ln (p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}})p(\boldsymbol{\theta}_{\text{MAP}})) - \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{MAP}})^{\mathsf{T}} \mathbf{A}(\boldsymbol{\theta}_{\text{MAP}})(\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{MAP}}),$$
(4.159)

so that

$$p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \simeq p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}})p(\boldsymbol{\theta}_{\text{MAP}})\exp\left(-\frac{1}{2}(\boldsymbol{\theta}-\boldsymbol{\theta}_{\text{MAP}})^{\mathsf{T}}\mathbf{A}(\boldsymbol{\theta}_{\text{MAP}})(\boldsymbol{\theta}-\boldsymbol{\theta}_{\text{MAP}})\right). \tag{4.160}$$

By marginalisation, integrating both sides with respect to  $\theta$  gives

$$p(\mathcal{D}) \simeq p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}})p(\boldsymbol{\theta}_{\text{MAP}}) \int \exp\left(-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{MAP}})^{\mathsf{T}} \mathbf{A}(\boldsymbol{\theta}_{\text{MAP}})(\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{MAP}})\right) d\boldsymbol{\theta}.$$
(4.161)

The integral of the right hand side can be written as

$$(2\pi)^{\frac{M}{2}} \left( \det \mathbf{A}(\boldsymbol{\theta}_{MAP})^{-1} \right)^{\frac{1}{2}} = (2\pi)^{\frac{M}{2}} \left( \det \mathbf{A}(\boldsymbol{\theta}_{MAP}) \right)^{-\frac{1}{2}}.$$
 (4.162)

Therefore,

$$p(\mathcal{D}) \simeq p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}})p(\boldsymbol{\theta}_{\text{MAP}})(2\pi)^{\frac{M}{2}} \left(\det \mathbf{A}(\boldsymbol{\theta}_{\text{MAP}})\right)^{-\frac{1}{2}},$$
 (4.163)

so that

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) + \ln p(\boldsymbol{\theta}_{\text{MAP}}) + \frac{M}{2}\ln(2\pi) - \frac{1}{2}\ln\left(\det \mathbf{A}(\boldsymbol{\theta}_{\text{MAP}})\right). \tag{4.164}$$

# 4.23

Let  $\boldsymbol{\theta}$  be a variable in M dimensions. By 4.22,

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) + \ln p(\boldsymbol{\theta}_{\text{MAP}}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln \left( \det \mathbf{A}(\boldsymbol{\theta}_{\text{MAP}}) \right), \tag{4.165}$$

where  $\boldsymbol{\theta}_{\text{MAP}}$  is a stationary point of  $p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$  and

$$\mathbf{A}(\boldsymbol{\theta}) = -\nabla\nabla \ln\left(p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})\right). \tag{4.166}$$

If

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \mathbf{V}_0), \tag{4.167}$$

then

$$\nabla\nabla \ln p(\theta) = -\mathbf{V}_0^{-1},\tag{4.168}$$

so that

$$\mathbf{A}(\boldsymbol{\theta}) = \mathbf{H}(\boldsymbol{\theta}) + \mathbf{V}_0^{-1},\tag{4.169}$$

where

$$\mathbf{H}(\boldsymbol{\theta}) = -\nabla \nabla \ln p(\mathcal{D}|\boldsymbol{\theta}). \tag{4.170}$$

Then, the right hand side of the approximation above can be written as

$$\ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) - \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln(\det \mathbf{V}_{0}) - \frac{1}{2} (\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m}_{0})^{\mathsf{T}} \mathbf{V}_{0}^{-1} (\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m}_{0})$$

$$+ \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln\left(\det\left(\mathbf{H}(\boldsymbol{\theta}_{\text{MAP}}) + \mathbf{V}_{0}^{-1}\right)\right)$$

$$= \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) + \frac{1}{2} \ln\left(\det\mathbf{V}_{0}^{-1}\right) - \frac{1}{2} (\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m}_{0})^{\mathsf{T}} \mathbf{V}_{0}^{-1} (\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m}_{0})$$

$$- \frac{1}{2} \ln\left(\det\left(\mathbf{H}(\boldsymbol{\theta}_{\text{MAP}}) + \mathbf{V}_{0}^{-1}\right)\right).$$
(4.171)

If  $\mathbf{V}_0^{-1}$  can be neglected, the right hand side can be written as

$$\ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) - \frac{1}{2} \ln \left( \det \mathbf{H}(\boldsymbol{\theta}_{\text{MAP}}) \right). \tag{4.172}$$

If each data point is independent and identically distributed, then

$$\mathbf{H}(\boldsymbol{\theta}) = N\bar{\mathbf{H}}(\boldsymbol{\theta}),\tag{4.173}$$

where

$$\bar{\mathbf{H}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^{N} \mathbf{H}_n(\boldsymbol{\theta}), \tag{4.174}$$

and  $\mathbf{H}_n(\boldsymbol{\theta})$  is the one for each data point. Then,

$$\det \mathbf{H}(\boldsymbol{\theta}_{\text{MAP}}) = N^M \det \bar{\mathbf{H}}(\boldsymbol{\theta}_{\text{MAP}}). \tag{4.175}$$

Threrefore,

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) - \frac{M}{2} \ln N.$$
 (4.176)

# 4.24 (Incomplete)

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{0, 1\},\ p(t_n | \mathbf{w}) = y_n^{t_n} (1 - y_n)^{1 - t_n},$$
(4.177)

where

$$y_n = \sigma(\mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n),$$
  
$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$
 (4.178)

By the Bayes' theorem,

$$p(\mathbf{w}|\mathbf{t})p(\mathbf{t}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}). \tag{4.179}$$

If

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0), \tag{4.180}$$

then the lograrithm of the right hand side except the terms independent of  ${\bf w}$  and  ${\bf t}$  can be written as

$$\sum_{n=1}^{N} (t_n \ln y_n + (1 - t_n) \ln(1 - y_n)) - \frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^{\mathsf{T}} \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0).$$
 (4.181)

Then, by 4.22,

$$p(\mathbf{w}|\mathbf{t}) \simeq \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{S}_N)$$
? (4.182)

where  $\mathbf{w}_{\text{MAP}}$  is the maximum likelihood solution for  $p(\mathbf{w})$  and

$$\mathbf{S}_N = -\nabla \nabla \ln p(\mathbf{w}|\mathbf{t}). \tag{4.183}$$

By marginalisation,

$$p(\mathcal{C}_1|\mathbf{t}) = \int p(\mathcal{C}_1|\mathbf{w})p(\mathbf{w}|\mathbf{t})d\mathbf{w}.$$
 (4.184)

The logarithm of the integrand of the right hand side except the terms independent of  $\mathbf{w}$  can be approximated as

$$-\ln\left(1 + \exp\left(-\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}\right)\right) - \frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\mathsf{T}}\mathbf{S}_{N}^{-1}(\mathbf{w} - \mathbf{w}_{\text{MAP}}) = (4.185)$$

Let

$$\sigma(a) = \frac{1}{1 + \exp(-a)},$$

$$\Phi(a) = \int_{-\infty}^{a} \mathcal{N}(\theta|0, 1)d\theta.$$
(4.186)

By 4.12,

$$\frac{d\sigma(a)}{da} = \sigma(a) \left(1 - \sigma(a)\right). \tag{4.187}$$

On the other hand, the right hand side can be written as

$$\frac{d\Phi(\lambda a)}{da} = \lambda \mathcal{N}(\lambda a|0,1). \tag{4.188}$$

Let us assume that

$$\left. \frac{d\sigma(a)}{da} \right|_{a=0} = \left. \frac{d\Phi(\lambda a)}{da} \right|_{a=0}. \tag{4.189}$$

Then,

$$\frac{1}{4} = \lambda (2\pi)^{-\frac{1}{2}}. (4.190)$$

Therefore,

$$\lambda^2 = \frac{\pi}{8}.\tag{4.191}$$

# 4.26

Let

$$I(\mu) = \int \Phi(\lambda a) \mathcal{N}(a|\mu, \sigma^2) da, \qquad (4.192)$$

where

$$\Phi(a) = \int_{-\infty}^{a} \mathcal{N}(\theta|0,1)d\theta. \tag{4.193}$$

By the transformation

$$z = \frac{a - \mu}{\sigma},\tag{4.194}$$

the right hand side can be written as

$$\int \Phi(\lambda(\mu + \sigma z)) \mathcal{N}(\mu + \sigma z | \mu, \sigma^2) \sigma dz = \int \Phi(\lambda(\mu + \sigma z)) \mathcal{N}(z | 0, 1) dz.$$
(4.195)

Then,

$$\frac{\partial}{\partial \mu}I(\mu) = \lambda \int \mathcal{N}\left(\lambda(\mu + \sigma z)|0,1\right) \mathcal{N}(z|0,1)dz. \tag{4.196}$$

The logarithm of the integrand of the right hand side can be written as

$$-\frac{1}{2}\ln(2\pi) - \frac{\lambda^2(\mu + \sigma z)^2}{2} - \frac{1}{2}\ln(2\pi) - \frac{z^2}{2}$$

$$= -\ln(2\pi) - \frac{1 + \sigma^2\lambda^2}{2} \left(z + \frac{\mu\sigma\lambda^2}{1 + \sigma^2\lambda^2}\right)^2 + \frac{\mu^2\sigma^2\lambda^4}{2(1 + \sigma^2\lambda^2)} - \frac{\mu^2\lambda^2}{2}.$$
(4.197)

The right hand side can be written as

$$-\ln(2\pi) - \frac{1+\sigma^2\lambda^2}{2} \left(z + \frac{\mu\sigma\lambda^2}{1+\sigma^2\lambda^2}\right)^2 - \frac{\mu^2\lambda^2}{2(1+\sigma^2\lambda^2)}$$

$$= -\frac{1}{2}\ln(2\pi) - \frac{1}{2}\ln(1+\sigma^2\lambda^2)^{-1} - \frac{1+\sigma^2\lambda^2}{2} \left(z + \frac{\mu\sigma\lambda^2}{1+\sigma^2\lambda^2}\right)^2 - \ln\lambda - \frac{1}{2}\ln(2\pi) - \frac{1}{2}\ln(\lambda^{-2} + \sigma^2) - \frac{\mu^2}{2(\lambda^{-2} + \sigma^2)}.$$
(4.198)

Then, the integral can be written as

$$\int \mathcal{N}\left(z| - \frac{\mu\sigma\lambda^2}{1 + \sigma^2\lambda^2}, \left(1 + \sigma^2\lambda^2\right)^{-1}\right) \frac{1}{\lambda} \mathcal{N}\left(\mu|0, \lambda^{-2} + \sigma^2\right) dz$$

$$= \frac{1}{\lambda} \mathcal{N}\left(\mu|0, \lambda^{-2} + \sigma^2\right).$$
(4.199)

Therefore,

$$\frac{\partial}{\partial \mu} I(\mu) = \mathcal{N} \left( \mu | 0, \lambda^{-2} + \sigma^2 \right). \tag{4.200}$$

Integrating both sides with respect to  $\mu$  gives

$$I(\mu) = \int_{-\infty}^{\mu} \mathcal{N}(m|0, \lambda^{-2} + \sigma^2) dm.$$
 (4.201)

By the transformation

$$m' = \frac{m}{(\lambda^{-2} + \sigma^2)^{\frac{1}{2}}},\tag{4.202}$$

the right hand side can be written as

$$\int_{-\infty}^{\frac{\mu}{(\lambda^{-2} + \sigma^2)^{\frac{1}{2}}}} \left(\lambda^{-2} + \sigma^2\right)^{-\frac{1}{2}} \mathcal{N}\left(m'|0,1\right) \left(\lambda^{-2} + \sigma^2\right)^{\frac{1}{2}} dm' = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{\frac{1}{2}}}\right). \tag{4.203}$$

Therefore,

$$I(\mu) = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{\frac{1}{2}}}\right).$$
 (4.204)

# 5 Neural Networks

#### 5.1

Let

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left( \sum_{m=1}^{M} w_{km}^{(2)} \sigma \left( \sum_{d=1}^{D} w_{md}^{(1)} x_d + w_{m0}^{(1)} \right) + w_{k0}^{(2)} \right), \tag{5.1}$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.\tag{5.2}$$

Here,

$$\sigma(a) = \frac{\exp\left(\frac{a}{2}\right)}{\exp\left(\frac{a}{2}\right) + \exp\left(-\frac{a}{2}\right)}.$$
 (5.3)

The right hand side can be written as

$$\tanh\left(\frac{a}{2}\right) + \sigma(-a) = \tanh\left(\frac{a}{2}\right) + 1 - \sigma(a). \tag{5.4}$$

Therefore,

$$\sigma(a) = \frac{1}{2} \left( 1 + \tanh\left(\frac{a}{2}\right) \right). \tag{5.5}$$

Then, the argument of the right hand side can be written as

$$\sum_{m=1}^{M} w_{km}^{(2)} \left( \frac{1}{2} \left( 1 + \tanh \left( \frac{1}{2} \left( \sum_{d=1}^{D} w_{md}^{(1)} x_d + w_{m0}^{(1)} \right) \right) \right) \right) + w_{k0}^{(2)} \\
= \frac{1}{2} \sum_{m=1}^{M} w_{km}^{(2)} \tanh \left( \frac{1}{2} \sum_{d=1}^{D} w_{md}^{(1)} x_d + \frac{1}{2} w_{m0}^{(1)} \right) + \frac{1}{2} \sum_{m=1}^{M} w_{km}^{(2)} + w_{k0}^{(2)}.$$
(5.6)

Therefore,

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left( \frac{1}{2} \sum_{m=1}^{M} w_{km}^{(2)} \tanh \left( \frac{1}{2} \sum_{d=1}^{D} w_{md}^{(1)} x_d + \frac{1}{2} w_{m0}^{(1)} \right) + \frac{1}{2} \sum_{m=1}^{M} w_{km}^{(2)} + w_{k0}^{(2)} \right).$$
(5.7)

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables such that

$$p(\mathbf{t}_n|\mathbf{x}_n, \mathbf{w}) = \mathcal{N}\left(\mathbf{t}_n|\mathbf{y}(\mathbf{x}_n, \mathbf{w}), \beta^{-1}\mathbf{I}\right). \tag{5.8}$$

Then, the logarithm of the likelihood except the terms independent of  $\mathbf{w}$  can be written as

$$-\frac{1}{2}\sum_{n=1}^{N}\left(\mathbf{t}_{n}-\mathbf{y}(\mathbf{x}_{n},\mathbf{w})\right)^{\intercal}\left(\beta^{-1}\mathbf{I}\right)^{-1}\left(\mathbf{t}_{n}-\mathbf{y}(\mathbf{x}_{n},\mathbf{w})\right)=-\frac{\beta}{2}\sum_{n=1}^{N}\|\mathbf{y}(\mathbf{x}_{n},\mathbf{w})-\mathbf{t}_{n}\|^{2}.$$
(5.9)

Setting the derivative with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = -\beta \sum_{n=1}^{N} \frac{\partial \mathbf{y}(\mathbf{x}_{n}, \mathbf{w})}{\partial \mathbf{w}} (\mathbf{y}(\mathbf{x}_{n}, \mathbf{w}) - \mathbf{t}_{n}).$$
 (5.10)

Let

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \|\mathbf{y}(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2.$$
 (5.11)

Setting the derivative with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = \sum_{n=1}^{N} \frac{\partial \mathbf{y}(\mathbf{x}_{n}, \mathbf{w})}{\partial \mathbf{w}} (\mathbf{y}(\mathbf{x}_{n}, \mathbf{w}) - \mathbf{t}_{n}).$$
 (5.12)

Therefore, maximising the likelihood is equivalent to minimising  $E(\mathbf{w})$ .

#### 5.3

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables such that

$$p(\mathbf{t}_n|\mathbf{x}_n,\mathbf{w}) = \mathcal{N}\left(\mathbf{t}_n|\mathbf{y}(\mathbf{x}_n,\mathbf{w}),\boldsymbol{\Sigma}\right). \tag{5.13}$$

Then, the logarithm of the likelihood except the terms independent of  ${\bf w}$  and  ${\bf \Sigma}$  can be written as

$$-\frac{1}{2}\ln(\det \mathbf{\Sigma}) - \frac{1}{2}\sum_{n=1}^{N} (\mathbf{t}_n - \mathbf{y}(\mathbf{x}_n, \mathbf{w}))^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{y}(\mathbf{x}_n, \mathbf{w})).$$
 (5.14)

Setting the derivatives with respect to  $\mathbf{w}$  and  $\Sigma$  to zero gives

$$\mathbf{0} = -\sum_{n=1}^{N} \frac{\partial \mathbf{y}(\mathbf{x}_{n}, \mathbf{w})}{\partial \mathbf{w}} \mathbf{\Sigma}^{-1} \left( \mathbf{t}_{n} - \mathbf{y}(\mathbf{x}_{n}, \mathbf{w}) \right),$$

$$\mathbf{O} = -\frac{1}{2} \mathbf{\Sigma}^{-1} + \frac{1}{2} \left( \mathbf{\Sigma}^{-1} \right)^{2} \sum_{n=1}^{N} \left( \mathbf{t}_{n} - \mathbf{y}(\mathbf{x}_{n}, \mathbf{w}) \right) \left( \mathbf{t}_{n} - \mathbf{y}(\mathbf{x}_{n}, \mathbf{w}) \right)^{\mathsf{T}}.$$
(5.15)

Therefore, the maximum likelihood solution for  $\Sigma$  is given by

$$\Sigma = \sum_{n=1}^{N} (\mathbf{t}_n - \mathbf{y}(\mathbf{x}_n, \mathbf{w})) (\mathbf{t}_n - \mathbf{y}(\mathbf{x}_n, \mathbf{w}))^{\mathsf{T}}.$$
 (5.16)

On the other hand, if  $\Sigma$  is fixed and known, then the maximum likelihood solution for  $\mathbf{w}$  is given by minimising

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} (\mathbf{t}_n - \mathbf{y}(\mathbf{x}_n, \mathbf{w}))^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{y}(\mathbf{x}_n, \mathbf{w})).$$
 (5.17)

#### 5.4

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{0, 1\},\ p(t_n = 1 | \mathbf{x}_n) = (1 - \epsilon)y_n + \epsilon (1 - y_n),$$
 (5.18)

where

$$y_n = y(\mathbf{x}_n, \mathbf{w}). \tag{5.19}$$

Then,

$$p(t_n|\mathbf{x}_n) = ((1 - \epsilon)y_n + \epsilon (1 - y_n))^{t_n} (\epsilon y_n + (1 - \epsilon) (1 - y_n))^{1 - t_n}.$$
 (5.20)

Therefore,

$$-\ln\left(\prod_{n=1}^{N} p(t_n|\mathbf{x}_n)\right)$$

$$= -\sum_{n=1}^{N} \left(t_n \ln\left((1-\epsilon)y_n + \epsilon(1-y_n)\right) + (1-t_n) \ln\left(\epsilon y_n + (1-\epsilon)(1-y_n)\right)\right).$$
(5.21)

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables in K dimensions such that

$$t_{nk} \in \{0, 1\},\ p(t_{nk} = 1 | \mathbf{x}_n) = y_{nk},$$
 (5.22)

where

$$y_{nk} = y_k(\mathbf{x}_n, \mathbf{w}). \tag{5.23}$$

Then,

$$p(t_{nk}|\mathbf{x}_n) = y_{nk}^{t_{nk}},\tag{5.24}$$

so that

$$p(\mathbf{t}_n|\mathbf{x}_n) = \prod_{k=1}^K y_{nk}^{t_{nk}}.$$
 (5.25)

Therefore,

$$\ln\left(\prod_{n=1}^{N} p(\mathbf{t}_{n}|\mathbf{x}_{n})\right) = \sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln y_{nk}.$$
 (5.26)

# 5.6

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{0, 1\},\ p(t_n = 1 | \mathbf{x}_n) = y_n,$$
 (5.27)

where

$$y_n = y(\mathbf{x}_n, \mathbf{w}). \tag{5.28}$$

Then,

$$p(t_n|\mathbf{x}_n) = y_n^{t_n} (1 - y_n)^{1 - t_n}.$$
 (5.29)

Let

$$E(\mathbf{w}) = -\ln\left(\prod_{n=1}^{N} p(t_n|\mathbf{x}_n)\right).$$
 (5.30)

Then,

$$E(\mathbf{w}) = -\sum_{n=1}^{N} (t_n \ln y_n + (1 - t_n) \ln(1 - y_n)).$$
 (5.31)

If

$$y_n = \sigma(a_n), \tag{5.32}$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)},\tag{5.33}$$

then, by 4.12,

$$\frac{\partial E(\mathbf{w})}{\partial a_n} = -y_n (1 - y_n) \left( \frac{t_n}{y_n} - \frac{1 - t_n}{1 - y_n} \right). \tag{5.34}$$

Therefore,

$$\frac{\partial E(\mathbf{w})}{\partial a_n} = y_n - t_n. \tag{5.35}$$

# 5.7

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables such that

$$t_{nk} \in \{0, 1\},\$$

$$p(t_{nk} = 1 | \mathbf{x}_n) = y_{nk},$$
(5.36)

where

$$y_{nk} = y_k(\mathbf{x}_n, \mathbf{w}),$$

$$y_{nk} = y_k(\mathbf{x}_n, \mathbf{w}),$$

$$\sum_{k=1}^K y_{nk} = 1.$$
(5.37)

Then,

$$p(\mathbf{t}_n|\mathbf{x}_n) = \prod_{k=1}^K y_{nk}^{t_{nk}}.$$
 (5.38)

Let

$$E(\mathbf{w}) = -\ln\left(\prod_{n=1}^{N} p(\mathbf{t}_n | \mathbf{x}_n)\right). \tag{5.39}$$

Then,

$$E(\mathbf{w}) = -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln y_{nk}.$$
 (5.40)

If

$$y_{nk} = \frac{\exp\left(a_k(\mathbf{x}_n, \mathbf{w})\right)}{\sum_{k=1}^K \exp\left(a_k(\mathbf{x}_n, \mathbf{w})\right)},$$
(5.41)

then, by 4.17,

$$\frac{\partial E(\mathbf{w})}{\partial a_{k'}} = -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} y_{nk} (I_{kk'} - y_{nk}) \frac{1}{y_{nk}}.$$
 (5.42)

The right hand side can be written as

$$-\sum_{n=1}^{N}\sum_{k=1}^{K}t_{nk}(I_{kk'}-y_{nk}) = -\sum_{n=1}^{N}\left(\sum_{k=1}^{K}t_{nk}y_{nk}-t_{nk'}\right).$$
 (5.43)

Therefore,

$$\frac{\partial E(\mathbf{w})}{\partial a_k} = \sum_{n=1}^{N} (y_{nk} - t_{nk}). \tag{5.44}$$

# 5.8

Setting the derivative of

$$tanh a = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}$$
(5.45)

gives

$$\frac{d}{da}\tanh a = 1 - \left(\frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}\right)^2. \tag{5.46}$$

Therefore,

$$\frac{d}{da}\tanh a = 1 - (\tanh a)^2. \tag{5.47}$$

#### 5.9

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{-1, 1\},\$$

$$p(t_n = 1 | \mathbf{x}_n) = \frac{1 + y_n}{2},$$
(5.48)

where

$$y_n = y(\mathbf{x}_n, \mathbf{w}). \tag{5.49}$$

Then,

$$p(t_n|\mathbf{x}_n) = \left(\frac{1+y_n}{2}\right)^{\frac{1+t_n}{2}} \left(\frac{1-y_n}{2}\right)^{\frac{1-t_n}{2}}.$$
 (5.50)

Let

$$E(\mathbf{w}) = -\ln\left(\prod_{n=1}^{N} p(t_n|\mathbf{x}_n)\right). \tag{5.51}$$

Then,

$$E(\mathbf{w}) = -\sum_{n=1}^{N} \left( \frac{1+t_n}{2} \ln \frac{1+y_n}{2} + \frac{1-t_n}{2} \ln \frac{1-y_n}{2} \right).$$
 (5.52)

The appropriate choice of y is tanh.

#### 5.10

Let

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}),\tag{5.53}$$

where E is a real function of real vectors. Then,  $\mathbf{H}$  is a real symmetric matrix. Therefore, by 2.20,  $\mathbf{H}$  is positive if and only if its eigenvalues are positive.

#### 5.11

Let  $\mathbf{w}$  be a real vector in M dimensions. Let E be a real function of  $\mathbf{w}$ . Let  $\mathbf{w}^*$  be a vector such that

$$\nabla E\left(\mathbf{w}^*\right) = \mathbf{0}.\tag{5.54}$$

Then, by a Taylor expansion,

$$E(\mathbf{w}) \simeq E(\mathbf{w}^*) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^{\mathsf{T}} \mathbf{H} (\mathbf{w} - \mathbf{w}^*),$$
 (5.55)

where

$$\mathbf{H} = \left. \nabla \nabla E(\mathbf{w}) \right|_{\mathbf{w} = \mathbf{w}^*}. \tag{5.56}$$

Let  $\mathbf{u}_1, \dots, \mathbf{u}_M$  be eigenvectors such that

$$\mathbf{H}\mathbf{u}_m = \lambda_m \mathbf{u}_m. \tag{5.57}$$

Note that **H** is a real symmetric matrix. Then, by ??, we have

$$\mathbf{u}_{m}^{\mathsf{T}}\mathbf{u}_{m'} = I_{mm'}.\tag{5.58}$$

Therefore, there exists  $\alpha_1, \dots, \alpha_M$  such that

$$\mathbf{w} - \mathbf{w}^* = \sum_{m=1}^{M} \alpha_m \mathbf{u}_m. \tag{5.59}$$

Then, the approximation can be written as

$$E(\mathbf{w}) \simeq E(\mathbf{w}^*) + \frac{1}{2} \sum_{m=1}^{M} \lambda_m \alpha_m^2.$$
 (5.60)

Therefore, the contours of constant E are ellipses whose axes are aligned with  $\mathbf{u}_1, \dots, \mathbf{u}_M$  with lengths which are proportional to  $\lambda_1^{-\frac{1}{2}}, \dots, \lambda_M^{-\frac{1}{2}}$ .

#### 5.12

Let **w** be a real vector. Let E be a real function of **w**. Let  $\mathbf{w}^*$  be a vector such that

$$\nabla E\left(\mathbf{w}^*\right) = \mathbf{0}.\tag{5.61}$$

Then, by a Taylor expansion,

$$E(\mathbf{w}) \simeq E(\mathbf{w}^*) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^{\mathsf{T}} \mathbf{H} (\mathbf{w} - \mathbf{w}^*),$$
 (5.62)

where

$$\mathbf{H} = \left. \nabla \nabla E(\mathbf{w}) \right|_{\mathbf{w} = \mathbf{w}^*}. \tag{5.63}$$

If  $\mathbf{H}$  is positive definite, then the second term of the right hand side is positive unless

$$\mathbf{w} = \mathbf{w}^*. \tag{5.64}$$

Therefore,  $\mathbf{w}^*$  is a local minimum of the right hand side. On the other hand, if  $\mathbf{w}^*$  is a local minimum of the right hand side, then the second term of the right hand side is positive unless

$$\mathbf{w} = \mathbf{w}^*. \tag{5.65}$$

Therefore,  $\mathbf{H}$  is positive definite. Thus, the necessary and sufficient condition for  $\mathbf{w}^*$  to be a local minimum is that  $\mathbf{H}$  be positive definite.

Let  $\mathbf{w}$  be a vector in M dimensions. Let E be a function of  $\mathbf{w}$ . Then, by a Taylor expansion,

$$E(\mathbf{w}) \simeq E(\hat{\mathbf{w}}) + (\mathbf{w} - \hat{\mathbf{w}})^{\mathsf{T}} \mathbf{b} + \frac{1}{2} (\mathbf{w} - \hat{\mathbf{w}})^{\mathsf{T}} \mathbf{H} (\mathbf{w} - \hat{\mathbf{w}}),$$
 (5.66)

where

$$\mathbf{b} = \nabla E(\mathbf{w})|_{\mathbf{w} = \hat{\mathbf{w}}}.$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w})|_{\mathbf{w} = \hat{\mathbf{w}}}.$$
(5.67)

Since **b** is a vector in M dimensions and **H** is a  $M \times M$  symmetric matrix, the number of independent elements of the right hand side is

$$M + \frac{M(M+1)}{2} = \frac{M(M+3)}{2}. (5.68)$$

#### 5.14

Let w be a variable. Let  $E_n$  be a function of w. Then, by a Taylor expansion,

$$E_{n}(w_{mm'} + \epsilon) = E_{n}(w_{mm'}) + \frac{\partial E_{n}}{\partial w}\Big|_{w=w_{mm'}} \epsilon + O(\epsilon^{2}),$$

$$E_{n}(w_{mm'} - \epsilon) = E_{n}(w_{mm'}) - \frac{\partial E_{n}}{\partial w}\Big|_{w=w_{mm'}} \epsilon + O(\epsilon^{2}).$$
(5.69)

Therefore,

$$\left. \frac{\partial E_n}{\partial w} \right|_{w=w} = \frac{E_n(w_{mm'} + \epsilon) - E_n(w_{mm'} - \epsilon)}{2\epsilon} + O\left(\epsilon^2\right). \tag{5.70}$$

# 5.15 (Incomplete)

#### 5.16

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be vectors. Let  $\mathbf{y}_1, \dots, \mathbf{y}_N$  be vectors which are dependent on a vector  $\mathbf{w}$ . Let

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \|\mathbf{y}_n - \mathbf{t}_n\|^2.$$
 (5.71)

Then,

$$\nabla E(\mathbf{w}) = \sum_{n=1}^{N} (\nabla \mathbf{y}_n)^{\mathsf{T}} (\mathbf{y}_n - \mathbf{t}_n), \qquad (5.72)$$

so that

$$\nabla \nabla E(\mathbf{w}) = \sum_{n=1}^{N} (\nabla \operatorname{vec} (\nabla \mathbf{y}_n)^{\mathsf{T}})^{\mathsf{T}} (\mathbf{y}_n - \mathbf{t}_n) + \sum_{n=1}^{N} (\nabla \mathbf{y}_n)^{\mathsf{T}} (\nabla \mathbf{y}_n). \quad (5.73)$$

# 5.17

Let t be a variable. Let y be a function of a vector  $\mathbf{x}$  and a vector  $\mathbf{w}$ . Let

$$E(\mathbf{w}) = \frac{1}{2} \int \int (y-t)^2 p(\mathbf{x}, t) d\mathbf{x} dt.$$
 (5.74)

Then,

$$\nabla E(\mathbf{w}) = \int \int (y - t)p(\mathbf{x}, t)\nabla y d\mathbf{x} dt.$$
 (5.75)

The right hand side can be written as

$$\int y \nabla y \left( \int p(\mathbf{x}, t) dt \right) d\mathbf{x} - \int \nabla y \left( \int t p(t|\mathbf{x}) dt \right) p(\mathbf{x}) d\mathbf{x}$$

$$= \int y \nabla y p(\mathbf{x}) d\mathbf{x} - \int \nabla y \, \mathbf{E}(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$
(5.76)

Then,

$$\nabla \nabla E(\mathbf{w}) = \int \nabla y (\nabla y)^{\mathsf{T}} p(\mathbf{x}) d\mathbf{x} + \int y \nabla \nabla y p(\mathbf{x}) d\mathbf{x} - \int \nabla \nabla y E(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$
(5.77)

The second and the third terms of the right hand side can be written as

$$E(y\nabla\nabla y) - E(\nabla\nabla y E(t|\mathbf{x})) = E((y - E(t|\mathbf{x})) \nabla\nabla y).$$
 (5.78)

Therefore, if

$$y = \mathcal{E}(t|\mathbf{x}),\tag{5.79}$$

then

$$\nabla \nabla E(\mathbf{w}) = \int \nabla y (\nabla y)^{\mathsf{T}} p(\mathbf{x}) d\mathbf{x}. \tag{5.80}$$

Let  $t_1, \dots, t_N$  be variables. Let  $y_1, \dots, y_N$  be variables such that

$$y_n = \mathbf{w}_n^{(2)^{\mathsf{T}}} \mathbf{z} + \mathbf{w}_n^{(0)^{\mathsf{T}}} \mathbf{x},$$
  
$$z_m = \tanh\left(\mathbf{w}_m^{(1)^{\mathsf{T}}} \mathbf{x}\right).$$
 (5.81)

Let

$$E = \frac{1}{2} \sum_{n=1}^{N} (y_n - t_n)^2.$$
 (5.82)

Then,

$$\frac{\partial E}{\partial \mathbf{w}_{n}^{(0)}} = \frac{\partial E}{\partial y_{n}} \frac{\partial y_{n}}{\partial \mathbf{w}_{n}^{(0)}},$$

$$\frac{\partial E}{\partial \mathbf{w}_{m}^{(1)}} = \frac{\partial E}{\partial y_{n}} \frac{\partial y_{n}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{w}_{m}^{(1)}},$$

$$\frac{\partial E}{\partial \mathbf{w}_{n}^{(2)}} = \frac{\partial E}{\partial y_{n}} \frac{\partial y_{n}}{\partial \mathbf{w}_{n}^{(2)}}.$$
(5.83)

Therefore,

$$\frac{\partial E}{\partial \mathbf{w}_{n}^{(0)}} = (y_{n} - t_{n})\mathbf{x},$$

$$\frac{\partial E}{\partial \mathbf{w}_{m}^{(1)}} = (y_{n} - t_{n})\mathbf{A}\mathbf{w}_{n}^{(2)},$$

$$\frac{\partial E}{\partial \mathbf{w}_{n}^{(2)}} = (y_{n} - t_{n})\mathbf{z},$$
(5.84)

where

$$A_{mm'} = (1 - z_m^2) x_{m'}. (5.85)$$

#### 5.19

Let  $t_1, \dots, t_N$  be variables such that

$$t_n \in \{0, 1\},\ p(t_n | \mathbf{w}) = y_n^{t_n} (1 - y_n)^{1 - t_n},$$
(5.86)

where

$$y_n = \sigma(a_n(\mathbf{w})),$$
  

$$\sigma(a) = \frac{1}{1 + \exp(-a)},$$
(5.87)

Let

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}). \tag{5.88}$$

The right hand side can be written as

$$-\ln\left(\prod_{n=1}^{N} y_n^{t_n} (1-y_n)^{1-t_n}\right) = -\sum_{n=1}^{N} \left(t_n \ln y_n + (1-t_n) \ln(1-y_n)\right). \quad (5.89)$$

Then, by 4.12,

$$\nabla E(\mathbf{w}) = -\sum_{n=1}^{N} \left( \frac{t_n}{y_n} y_n (1 - y_n) \nabla a_n - \frac{1 - t_n}{1 - y_n} y_n (1 - y_n) \nabla a_n \right).$$
 (5.90)

The right hand side can be written as

$$-\sum_{n=1}^{N} (t_n(1-y_n)\nabla a_n - (1-t_n)y_n\nabla a_n) = \sum_{n=1}^{N} (y_n - t_n)\nabla a_n.$$
 (5.91)

Then, by 4.13,

$$\nabla E(\mathbf{w}) = \sum_{n=1}^{N} (y_n - t_n) \nabla a_n.$$
 (5.92)

Therefore, by 4.12,

$$\nabla \nabla E(\mathbf{w}) = \sum_{n=1}^{N} y_n (1 - y_n) \nabla a_n (\nabla a_n)^{\mathsf{T}} + \sum_{n=1}^{N} (y_n - t_n) \nabla \nabla a_n.$$
 (5.93)

#### 5.20

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables such that

$$t_{nk} \in \{0, 1\},\ p(t_{nk} = 1 | \mathbf{x}_n) = y_{nk},$$
 (5.94)

where

$$y_{nk} = y_k(\mathbf{x}_n, \mathbf{w}),$$

$$\sum_{k=1}^K y_{nk} = 1.$$
(5.95)

Then,

$$p(\mathbf{t}_n|\mathbf{x}_n) = \prod_{k=1}^K y_{nk}^{t_{nk}}.$$
 (5.96)

Let

$$E(\mathbf{w}) = -\ln\left(\prod_{n=1}^{N} p(\mathbf{t}_n | \mathbf{x}_n)\right). \tag{5.97}$$

Then,

$$E(\mathbf{w}) = -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln y_{nk}.$$
 (5.98)

If

$$y_{nk} = \frac{\exp(a_k(\mathbf{x}_n, \mathbf{w}))}{\sum_{k=1}^K \exp(a_k(\mathbf{x}_n, \mathbf{w}))},$$
 (5.99)

then, by 5.7,

$$\frac{\partial E(\mathbf{w})}{\partial a_k} = \sum_{n=1}^{N} (y_{nk} - t_{nk}). \tag{5.100}$$

Then,

$$\nabla E(\mathbf{w}) = \sum_{n=1}^{N} \sum_{k=1}^{K} (y_{nk} - t_{nk}) \nabla a_k.$$
 (5.101)

Then,

$$\nabla \nabla E(\mathbf{w}) = \sum_{n=1}^{N} \sum_{k=1}^{K} y_{nk} (1 - y_{nk}) \nabla a_k (\nabla a_k)^{\mathsf{T}} + \sum_{n=1}^{N} \sum_{k=1}^{K} (y_{nk} - t_{nk}) \nabla \nabla a_k.$$
(5.102)

# 5.21 (Incomplete)

# 5.22

Let  $y_1, \dots, y_N$  be variables such that

$$y_n = \mathbf{w}_n^{(2)^{\mathsf{T}}} \mathbf{z},$$

$$z_m = h(a_m),$$

$$a_m = \mathbf{w}_m^{(1)^{\mathsf{T}}} \mathbf{x}.$$
(5.103)

Let E be a function of  $y_1, \dots, y_N$ . Then,

$$\frac{\partial E}{\partial \mathbf{w}_{m}^{(1)}} = \frac{\partial E}{\partial y_{n}} \frac{\partial y_{n}}{\partial z_{m}} \frac{\partial z_{m}}{\partial a_{m}} \frac{\partial a_{m}}{\partial \mathbf{w}_{m}^{(1)}},$$

$$\frac{\partial E}{\partial \mathbf{w}_{n}^{(2)}} = \frac{\partial E}{\partial y_{n}} \frac{\partial y_{n}}{\partial \mathbf{w}_{n}^{(2)}}.$$
(5.104)

Therefore,

$$\frac{\partial E}{\partial \mathbf{w}_{m}^{(1)}} = \frac{\partial E}{\partial y_{n}} w_{nm}^{(2)} h'(a_{m}) \mathbf{x}, 
\frac{\partial E}{\partial \mathbf{w}_{n}^{(2)}} = \frac{\partial E}{\partial y_{n}} \mathbf{z}.$$
(5.105)

Thus,

$$\frac{\partial^{2} E}{\partial \mathbf{w}_{m}^{(1)} \partial \mathbf{w}_{m'}^{(1)}} = w_{nm}^{(2)} \mathbf{x} \left( \frac{\partial^{2} E}{\partial y_{n}^{2}} w_{nm'}^{(2)} h'(a_{m'}) h'(a_{m}) \mathbf{x} + \frac{\partial E}{\partial y_{n}} h''(a_{m}) I_{mm'} \mathbf{x} \right)^{\mathsf{T}},$$

$$\frac{\partial^{2} E}{\partial \mathbf{w}_{m}^{(1)} \partial \mathbf{w}_{n}^{(2)}} = h'(a_{m}) \mathbf{x} \left( \frac{\partial^{2} E}{\partial y_{n}^{2}} w_{nm}^{(2)} \mathbf{z} + \frac{\partial E}{\partial y_{n}} \mathbf{v} \right)^{\mathsf{T}},$$

$$\frac{\partial^{2} E}{\partial \mathbf{w}_{n}^{(2)} \partial \mathbf{w}_{n'}^{(2)}} = \frac{\partial^{2} E}{\partial y_{n} \partial y_{n'}} \mathbf{z} \mathbf{z}^{\mathsf{T}},$$
(5.106)

where

$$v_m = \begin{cases} 1, & m = n, \\ 0 & \text{otherwise.} \end{cases}$$
 (5.107)

# 5.23

Let  $y_1, \dots, y_N$  be variables such that

$$y_n = \mathbf{w}_n^{(2)^{\mathsf{T}}} \mathbf{z} + \mathbf{w}_n^{(0)^{\mathsf{T}}} \mathbf{x},$$

$$z_m = h(a_m),$$

$$a_m = \mathbf{w}_m^{(1)^{\mathsf{T}}} \mathbf{x}.$$
(5.108)

Let E be a function of  $y_1, \dots, y_N$ . Then,

$$\frac{\partial E}{\partial \mathbf{w}_{n}^{(0)}} = \frac{\partial E}{\partial y_{n}} \frac{\partial y_{n}}{\partial \mathbf{w}_{n}^{(0)}},$$

$$\frac{\partial E}{\partial \mathbf{w}_{m}^{(1)}} = \frac{\partial E}{\partial y_{n}} \frac{\partial y_{n}}{\partial z_{m}} \frac{\partial z_{m}}{\partial a_{m}} \frac{\partial a_{m}}{\partial \mathbf{w}_{m}^{(1)}},$$

$$\frac{\partial E}{\partial \mathbf{w}_{n}^{(2)}} = \frac{\partial E}{\partial y_{n}} \frac{\partial y_{n}}{\partial \mathbf{w}_{n}^{(2)}}.$$
(5.109)

Therefore,

$$\frac{\partial E}{\partial \mathbf{w}_{n}^{(0)}} = \frac{\partial E}{\partial y_{n}} \mathbf{x},$$

$$\frac{\partial E}{\partial \mathbf{w}_{m}^{(1)}} = \frac{\partial E}{\partial y_{n}} w_{nm}^{(2)} h'(a_{m}) \mathbf{x},$$

$$\frac{\partial E}{\partial \mathbf{w}_{n}^{(2)}} = \frac{\partial E}{\partial y_{n}} \mathbf{z}.$$
(5.110)

Thus,

$$\frac{\partial^{2} E}{\partial \mathbf{w}_{n}^{(0)} \partial \mathbf{w}_{n'}^{(0)}} = \frac{\partial^{2} E}{\partial y_{n} \partial y_{n'}} \mathbf{x} \mathbf{x}^{\mathsf{T}}, 
\frac{\partial^{2} E}{\partial \mathbf{w}_{n}^{(0)} \partial \mathbf{w}_{m}^{(1)}} = \frac{\partial^{2} E}{\partial y_{n}^{2}} w_{nm}^{(2)} h'(a_{m}) \mathbf{x} \mathbf{x}^{\mathsf{T}}, 
\frac{\partial^{2} E}{\partial \mathbf{w}_{n}^{(0)} \partial \mathbf{w}_{n'}^{(2)}} = \frac{\partial^{2} E}{\partial y_{n} \partial y_{n'}} \mathbf{x} \mathbf{z}^{\mathsf{T}}, 
\frac{\partial^{2} E}{\partial \mathbf{w}_{m}^{(1)} \partial \mathbf{w}_{m'}^{(1)}} = w_{nm}^{(2)} \mathbf{x} \left( \frac{\partial^{2} E}{\partial y_{n}^{2}} w_{nm'}^{(2)} h'(a_{m'}) h'(a_{m}) \mathbf{x} + \frac{\partial E}{\partial y_{n}} h''(a_{m}) I_{mm'} \mathbf{x} \right)^{\mathsf{T}}, 
\frac{\partial^{2} E}{\partial \mathbf{w}_{m}^{(1)} \partial \mathbf{w}_{n'}^{(2)}} = h'(a_{m}) \mathbf{x} \left( \frac{\partial^{2} E}{\partial y_{n}^{2}} w_{nm}^{(2)} \mathbf{z} + \frac{\partial E}{\partial y_{n}} \mathbf{v} \right)^{\mathsf{T}}, 
\frac{\partial^{2} E}{\partial \mathbf{w}_{n}^{(2)} \partial \mathbf{w}_{n'}^{(2)}} = \frac{\partial^{2} E}{\partial y_{n} \partial y_{n'}} \mathbf{z} \mathbf{z}^{\mathsf{T}},$$

$$(5.111)$$

where

$$v_m = \begin{cases} 1, & m = n, \\ 0 & \text{otherwise.} \end{cases}$$
 (5.112)

Let  $y_1, \dots, y_n$  be variables such that

$$y_n = \mathbf{w}_n^{\mathsf{T}} \mathbf{z} + w_{n0},$$
  

$$z_m = h \left( \mathbf{w}_m^{\mathsf{T}} \mathbf{x} + w_{m0} \right).$$
(5.113)

(a)

Let

$$\tilde{\mathbf{x}} = a\mathbf{x} + b\mathbf{v},\tag{5.114}$$

where

$$\mathbf{v} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}.$$

Then,

$$z_m = h\left(\frac{1}{a}\mathbf{w}_m^{\mathsf{T}}\left(\tilde{\mathbf{x}} - b\mathbf{v}\right) + w_{m0}\right). \tag{5.115}$$

Therefore,

$$\tilde{z}_m = h\left(\tilde{\mathbf{w}}_m^{\dagger} \mathbf{x} + \tilde{w}_{m0}\right). \tag{5.116}$$

where

$$\tilde{\mathbf{w}}_{m} = \frac{1}{a} \mathbf{w}_{m},$$

$$\tilde{w}_{m0} = w_{m0} - \frac{b}{a} \mathbf{w}_{m}^{\mathsf{T}} \mathbf{v}.$$
(5.117)

(b)

Let

$$\tilde{y}_n = cy_n + d. (5.118)$$

Then,

$$\frac{\tilde{y}_n - d}{c} = \mathbf{w}_n^{\mathsf{T}} \mathbf{z} + w_{n0}. \tag{5.119}$$

Therefore,

$$\tilde{y}_n = \tilde{\mathbf{w}}_n^{\mathsf{T}} \mathbf{z} + \tilde{w}_{n0}, \tag{5.120}$$

where

$$\tilde{\mathbf{w}}_n = c\mathbf{w}_n, 
\tilde{w}_{n0} = cw_{n0} + d.$$
(5.121)

# 5.25 (Incomplete)

Let E be a quadratic error function such that

$$E(\mathbf{w}) = E_0 + \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^{\mathsf{T}} \mathbf{H} (\mathbf{w} - \mathbf{w}^*), \qquad (5.122)$$

where  $\mathbf{w}^*$  represents the minimum and  $\mathbf{H}$  is a positive definite and constant matrix whose eignevectors and eigenvalues are  $\mathbf{u}_1, \dots, \mathbf{u}_M$  and  $\eta_1, \dots, \eta_M$ . Let

$$\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \rho \nabla E,$$
  

$$\mathbf{w}^{(0)} = \mathbf{0}.$$
(5.123)

(a)

We have

$$\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \rho \mathbf{H} \left( \mathbf{w}^{(\tau-1)} - \mathbf{w}^* \right), \tag{5.124}$$

so that

$$\mathbf{w}^{(\tau)} = (\mathbf{I} - \rho \mathbf{H}) \,\mathbf{w}^{(\tau - 1)} + \rho \mathbf{H} \mathbf{w}^*. \tag{5.125}$$

Then,

$$\mathbf{u}_{m}^{\mathsf{T}}\mathbf{w}^{(\tau)} = \mathbf{u}_{m}^{\mathsf{T}}\mathbf{w}^{(\tau-1)} - \rho\mathbf{u}_{m}^{\mathsf{T}}\mathbf{H}\mathbf{w}^{(\tau-1)} + \rho\mathbf{u}_{m}^{\mathsf{T}}\mathbf{H}\mathbf{w}^{*}. \tag{5.126}$$

Since  $\mathbf{H}$  is symmetric,

$$\mathbf{u}_{m}^{\mathsf{T}}\mathbf{H}\mathbf{w} = \mathbf{w}^{\mathsf{T}}\mathbf{H}\mathbf{u}_{m}.\tag{5.127}$$

The right hand side can be written as

$$\eta_m \mathbf{w}^{\mathsf{T}} \mathbf{u}_m = \eta_m w_m, \tag{5.128}$$

where

$$w_m = \mathbf{w}^\mathsf{T} \mathbf{u}_m. \tag{5.129}$$

Then,

$$w_m^{(\tau)} = (1 - \rho \eta_m) w_m^{(\tau - 1)} + \rho \eta_m w_m^*, \tag{5.130}$$

so that

$$w_m^{(\tau)} - w_m^* = (1 - \rho \eta_m) \left( w_m^{(\tau - 1)} - w_m^* \right). \tag{5.131}$$

Therefore,

$$w_m^{(\tau)} = w_m^* - (1 - \rho \eta_m)^{\tau} \left( w_m^{(0)} - w_m^* \right), \tag{5.132}$$

so that

$$w_m^{(\tau)} = (1 - (1 - \rho \eta_m)^{\tau}) w_m^*. \tag{5.133}$$

(b)

$$|1 - \rho \eta_m| < 1, (5.134)$$

then

$$\lim_{\tau \to \infty} w_m^{(\tau)} = w_m^*. \tag{5.135}$$

(c)

By (a) and a Taylor series

$$(1 - \rho \eta_m)^{\tau} = 1 - \rho \tau \eta_m + \rho^2 \tau^2 O(\eta_m^2), \qquad (5.136)$$

we have

$$w_m^{(\tau)} = \left(\rho \tau \eta_m - \rho^2 \tau^2 O\left(\eta_m^2\right)\right) w_m^*, \tag{5.137}$$

so that

$$w_m^{(\tau)} = \rho \tau \eta_m \left( 1 - \rho \tau \eta_m O(\eta_m) \right) w_m^*. \tag{5.138}$$

Therefore,

$$w_m^{(\tau)} \simeq w_m^*, \quad \text{if} \quad \eta_m \gg (\rho \tau)^{-1}?$$

$$\left| w_m^{(\tau)} \right| \ll w_m^*, \quad \text{if} \quad \eta_m \ll (\rho \tau)^{-1}.$$

$$(5.139)$$

# 5.26 (Incomplete)

Let

$$\tilde{E} = E + \lambda \Omega, \tag{5.140}$$

where

$$\Omega = \frac{1}{2} \sum_{n=1}^{N} \sum_{k=1}^{K} \left( \frac{\partial y_{nk}}{\partial \xi} \right)^{2}.$$
 (5.141)

(a)

We have

$$\frac{\partial y_{nk}}{\partial \xi} = \mathcal{G}y_{nk},\tag{5.142}$$

where

$$\mathcal{G} = \sum_{d=1}^{D} \tau_d \frac{\partial}{\partial x_d},$$

$$\tau_d = \frac{\partial x_d}{\partial \xi}.$$
(5.143)

Therefore,

$$\Omega = \sum_{n=1}^{N} \Omega_n, \tag{5.144}$$

where

$$\Omega_n = \frac{1}{2} \sum_{k=1}^K (\mathcal{G} y_{nk})^2.$$
 (5.145)

(b)

$$\alpha_j = \mathcal{G}z_j = h'(a_j)\beta_j?$$

$$\beta_j = \mathcal{G}a_j = \sum_{d=1}^D w_{jd}\alpha_d?$$
(5.146)

#### 5.27

Let

$$\tilde{E} = \frac{1}{2} \iiint (y(\mathbf{x} + \boldsymbol{\xi}) - t)^2 p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) d\mathbf{x} dt d\boldsymbol{\xi},$$
 (5.147)

where

$$p(\boldsymbol{\xi}) = \mathcal{N}(\boldsymbol{\xi}|\mathbf{0}, \mathbf{I}). \tag{5.148}$$

By a Taylor series

$$y(\mathbf{x} + \boldsymbol{\xi}) = y(\mathbf{x}) + \boldsymbol{\xi}^{\mathsf{T}} \nabla y(\mathbf{x}) + \frac{1}{2} \boldsymbol{\xi}^{\mathsf{T}} \nabla \nabla y(\mathbf{x}) \boldsymbol{\xi} + O\left(\|\boldsymbol{\xi}\|^{3}\right), \tag{5.149}$$

the integrand can be written as

$$(y(\mathbf{x}) - t)^{2} + 2(y(\mathbf{x}) - t) \left(\boldsymbol{\xi}^{\mathsf{T}} \nabla y(\mathbf{x}) + \frac{1}{2} \boldsymbol{\xi}^{\mathsf{T}} \nabla \nabla y(\mathbf{x}) \boldsymbol{\xi} + O(\|\boldsymbol{\xi}\|^{3})\right) + \left(\boldsymbol{\xi}^{\mathsf{T}} \nabla y(\mathbf{x}) + \frac{1}{2} \boldsymbol{\xi}^{\mathsf{T}} \nabla \nabla y(\mathbf{x}) \boldsymbol{\xi} + O(\|\boldsymbol{\xi}\|^{3})\right)^{2}.$$
(5.150)

Omitting the terms of  $O(\|\boldsymbol{\xi}\|^3)$ , the integral can be approximated as

$$\iiint (y(\mathbf{x}) - t)^{2} p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) d\mathbf{x} dt d\boldsymbol{\xi} 
+ 2 \iiint (y(\mathbf{x}) - t) \left(\boldsymbol{\xi}^{\mathsf{T}} \nabla y(\mathbf{x}) + \frac{1}{2} \boldsymbol{\xi}^{\mathsf{T}} \nabla \nabla y(\mathbf{x}) \boldsymbol{\xi}\right) p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) d\mathbf{x} dt d\boldsymbol{\xi} 
+ \iiint \left(\boldsymbol{\xi}^{\mathsf{T}} \nabla y(\mathbf{x}) + \frac{1}{2} \boldsymbol{\xi}^{\mathsf{T}} \nabla \nabla y(\mathbf{x}) \boldsymbol{\xi}\right)^{2} p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) d\mathbf{x} dt d\boldsymbol{\xi}.$$
(5.151)

Omitting the terms of  $O(\|\boldsymbol{\xi}\|^3)$ , the second and third terms can be approximated as

$$2 \iiint (y(\mathbf{x}) - t) \boldsymbol{\xi}^{\mathsf{T}} \nabla y(\mathbf{x}) p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) d\mathbf{x} dt d\boldsymbol{\xi}$$

$$+ \iiint (y(\mathbf{x}) - t) \boldsymbol{\xi}^{\mathsf{T}} \nabla \nabla y(\mathbf{x}) \boldsymbol{\xi} p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) d\mathbf{x} dt d\boldsymbol{\xi}$$

$$+ \iiint \boldsymbol{\xi}^{\mathsf{T}} (\nabla y(\mathbf{x})) (\nabla y(\mathbf{x}))^{\mathsf{T}} \boldsymbol{\xi} p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) d\mathbf{x} dt d\boldsymbol{\xi}.$$

$$(5.152)$$

Since

$$\mathbf{E}\,\boldsymbol{\xi} = \mathbf{0},\tag{5.153}$$

the integral of the first term can be written as

$$\left(\int \boldsymbol{\xi} p(\boldsymbol{\xi}) d\boldsymbol{\xi}\right)^{\mathsf{T}} \int \int (y(\mathbf{x}) - t) \nabla y(\mathbf{x}) p(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} dt = \mathbf{0}.$$
 (5.154)

Since

$$cov \boldsymbol{\xi} = \mathbf{I}, \tag{5.155}$$

The integrals of the second term can be written as

$$\int \boldsymbol{\xi}^{\mathsf{T}} \left( \int \int (y(\mathbf{x}) - t) \, \nabla \nabla y(\mathbf{x}) p(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} dt \right) \boldsymbol{\xi} p(\boldsymbol{\xi}) d\boldsymbol{\xi} 
= \operatorname{tr} \left( \int \int (y(\mathbf{x}) - t) \, \nabla \nabla y(\mathbf{x}) p(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} dt \right).$$
(5.156)

The integral of the right hand side can be written as

$$\int \left( y(\mathbf{x}) - \int t p(t|\mathbf{x}) dt \right) \nabla \nabla y(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

$$= \int \left( y(\mathbf{x}) - \mathbf{E}(t|\mathbf{x}) \right) \nabla \nabla y(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$
(5.157)

Similarly, the integral of the third term can be written as

$$\int \boldsymbol{\xi}^{\mathsf{T}} \left( \int \int (\nabla y(\mathbf{x})) (\nabla y(\mathbf{x}))^{\mathsf{T}} p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) d\mathbf{x} dt \right) \boldsymbol{\xi} d\boldsymbol{\xi} 
= \operatorname{tr} \left( \int \int (\nabla y(\mathbf{x})) (\nabla y(\mathbf{x}))^{\mathsf{T}} p(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} dt \right).$$
(5.158)

The integral of the right hand side can be written as

$$\int \left( \int p(t|\mathbf{x})dt \right) (\nabla y(\mathbf{x})) (\nabla y(\mathbf{x}))^{\mathsf{T}} p(\mathbf{x})d\mathbf{x} 
= \int (\nabla y(\mathbf{x})) (\nabla y(\mathbf{x}))^{\mathsf{T}} p(\mathbf{x})d\mathbf{x}.$$
(5.159)

If

$$y(\mathbf{x}) = \mathbf{E}(t|\mathbf{x}),\tag{5.160}$$

then

$$\tilde{E} \simeq E + \Omega,$$
 (5.161)

where

$$E = \frac{1}{2} \iiint (y(\mathbf{x}) - t)^2 p(t|\mathbf{x}) p(\mathbf{x}) p(\boldsymbol{\xi}) d\mathbf{x} dt d\boldsymbol{\xi},$$

$$\Omega = \frac{1}{2} \operatorname{tr} \left( \int \nabla y(\mathbf{x}) (\nabla y(\mathbf{x}))^{\mathsf{T}} p(\mathbf{x}) d\mathbf{x}. \right).$$
(5.162)

# 5.28 (Incomplete)

#### 5.29

Let  $w_1, \dots, w_M$  be variables such that

$$p(w_m) = \sum_{k=1}^{K} \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right). \tag{5.163}$$

Let

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w}),$$
 (5.164)

where

$$\Omega(\mathbf{w}) = -\sum_{m=1}^{M} \ln \left( \sum_{k=1}^{K} \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right) \right).$$
 (5.165)

We have

$$\frac{\partial \Omega}{\partial w_m} = -\frac{\sum_{k=1}^K \pi_k \frac{w_m - \mu_k}{\sigma_k^2} \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right)}{\sum_{k=1}^K \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right)}.$$
 (5.166)

Therefore,

$$\frac{\partial \tilde{E}}{\partial w_m} = \frac{\partial E}{\partial w_m} - \lambda \sum_{k=1}^K \gamma_k(w_m) \frac{w_m - \mu_k}{\sigma_k^2}, \tag{5.167}$$

where

$$\gamma_k(w_m) = \frac{\pi_k \mathcal{N}(w_m | \mu_k, \sigma_k^2)}{\sum_{k=1}^K \pi_k \mathcal{N}(w_m | \mu_k, \sigma_k^2)}.$$
 (5.168)

#### 5.30

Let  $w_1, \dots, w_M$  be variables such that

$$p(w_m) = \sum_{k=1}^{K} \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right). \tag{5.169}$$

Let

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w}),$$
 (5.170)

where

$$\Omega(\mathbf{w}) = -\sum_{m=1}^{M} \ln \left( \sum_{k=1}^{K} \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right) \right).$$
 (5.171)

We have

$$\frac{\partial \Omega}{\partial \mu_k} = -\frac{\sum_{k=1}^K \pi_k \frac{w_m - \mu_k}{\sigma_k^2} \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right)}{\sum_{k=1}^K \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right)}.$$
 (5.172)

Therefore,

$$\frac{\partial \tilde{E}}{\partial \mu_k} = -\lambda \sum_{k=1}^K \gamma_k(w_m) \frac{w_m - \mu_k}{\sigma_k^2},\tag{5.173}$$

where

$$\gamma_k(w_m) = \frac{\pi_k \mathcal{N}(w_m | \mu_k, \sigma_k^2)}{\sum_{k=1}^K \pi_k \mathcal{N}(w_m | \mu_k, \sigma_k^2)}.$$
 (5.174)

#### 5.31

Let  $w_1, \dots, w_M$  be variables such that

$$p(w_m) = \sum_{k=1}^{K} \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right). \tag{5.175}$$

Let

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w}),$$
 (5.176)

where

$$\Omega(\mathbf{w}) = -\sum_{m=1}^{M} \ln \left( \sum_{k=1}^{K} \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right) \right).$$
 (5.177)

We have

$$\frac{\partial \Omega}{\partial \sigma_k} = -\frac{\sum_{k=1}^K \pi_k \left( -\frac{1}{\sigma_k} + \frac{(w_m - \mu_k)^2}{\sigma_k^3} \right) \mathcal{N}\left( w_m | \mu_k, \sigma_k^2 \right)}{\sum_{k=1}^K \pi_k \mathcal{N}\left( w_m | \mu_k, \sigma_k^2 \right)}.$$
 (5.178)

Therefore,

$$\frac{\partial \tilde{E}}{\partial \sigma_k} = -\lambda \sum_{k=1}^K \gamma_k(w_m) \left( -\frac{1}{\sigma_k} + \frac{(w_m - \mu_k)^2}{\sigma_k^3} \right), \tag{5.179}$$

where

$$\gamma_k(w_m) = \frac{\pi_k \mathcal{N}(w_m | \mu_k, \sigma_k^2)}{\sum_{k=1}^K \pi_k \mathcal{N}(w_m | \mu_k, \sigma_k^2)}.$$
 (5.180)

#### 5.32

Let  $w_1, \dots, w_M$  be variables such that

$$p(w_m) = \sum_{k=1}^{K} \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right), \qquad (5.181)$$

where

$$\pi_k = \frac{\exp(\eta_k)}{\sum_{k=1}^K \exp(\eta_k)}.$$
 (5.182)

Let

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w}),$$
 (5.183)

$$\Omega(\mathbf{w}) = -\sum_{m=1}^{M} \ln \left( \sum_{k=1}^{K} \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right) \right).$$
 (5.184)

(a)

If  $k \neq k'$ , then

$$\frac{\partial \pi_k}{\partial \eta_{k'}} = -\frac{\exp(\eta_{k'}) \exp(\eta_k)}{\left(\sum_{k=1}^K \exp(\eta_k)\right)^2}.$$
 (5.185)

We have

$$\frac{\partial \pi_k}{\partial \eta_k} = \frac{\exp(\eta_k)}{\sum_{k=1}^K \exp(\eta_k)} - \left(\frac{\exp(\eta_k)}{\sum_{k=1}^K \exp(\eta_k)}\right)^2.$$
 (5.186)

Therefore,

$$\frac{\partial \pi_k}{\partial \eta_{k'}} = I_{kk'} \pi_k - \pi_k \pi_{k'}. \tag{5.187}$$

(b)

We have

$$\frac{\partial \Omega}{\partial \eta_k} = \sum_{k'=1}^K \frac{\partial \Omega}{\partial \pi_{k'}} \frac{\partial \pi_{k'}}{\partial \eta_k}.$$
 (5.188)

We have

$$\frac{\partial \Omega}{\partial \pi_{k'}} = -\frac{\sum_{m=1}^{M} \mathcal{N}\left(w_m | \mu_{k'}, \sigma_{k'}^2\right)}{\sum_{k=1}^{K} \pi_k \mathcal{N}\left(w_m | \mu_k, \sigma_k^2\right)}.$$
 (5.189)

By (a),

$$\frac{\partial \Omega}{\partial \eta_k} = -\sum_{k'=1}^K \frac{\sum_{m=1}^M \mathcal{N}(w_m | \mu_{k'}, \sigma_{k'}^2)}{\sum_{k=1}^K \pi_k \mathcal{N}(w_m | \mu_k, \sigma_k^2)} (I_{kk'} \pi_k - \pi_k \pi_{k'}).$$
 (5.190)

The right hand side can be written as

$$-\sum_{m=1}^{M} \frac{\pi_{k} \mathcal{N}\left(w_{m} | \mu_{k}, \sigma_{k}^{2}\right)}{\sum_{k=1}^{K} \pi_{k} \mathcal{N}\left(w_{m} | \mu_{k}, \sigma_{k}^{2}\right)} + \pi_{k} \sum_{m=1}^{M} \frac{\sum_{k'=1}^{K} \pi_{k'} \mathcal{N}\left(w_{m} | \mu_{k}, \sigma_{k}^{2}\right)}{\sum_{k=1}^{K} \pi_{k} \mathcal{N}\left(w_{m} | \mu_{k}, \sigma_{k}^{2}\right)}$$

$$= \sum_{m=1}^{M} \left(\pi_{k} - \frac{\pi_{k} \mathcal{N}\left(w_{m} | \mu_{k}, \sigma_{k}^{2}\right)}{\sum_{k=1}^{K} \pi_{k} \mathcal{N}\left(w_{m} | \mu_{k}, \sigma_{k}^{2}\right)}\right).$$
(5.191)

Therefore,

$$\frac{\partial \tilde{E}}{\partial \eta_k} = \lambda \sum_{m=1}^{M} (\pi_k - \gamma_k(w_m)), \qquad (5.192)$$

where

$$\gamma_k(w_m) = \frac{\pi_k \mathcal{N}(w_m | \mu_k, \sigma_k^2)}{\sum_{k=1}^K \pi_k \mathcal{N}(w_m | \mu_k, \sigma_k^2)}.$$
 (5.193)

#### 5.33

Let  $(x_1, x_2)$  be the Cartesian coordinates of the end-effector of a two-link robot arm whose joint angles are  $\theta_1$  and  $\theta_2$  and whose arm lengths are  $L_1$ and  $L_2$ . Let us assume that the origin of the coordinate system is given by the attachment point of the lower arm. We have

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} L_1 \cos \theta_1 \\ L_1 \sin \theta_1 \end{bmatrix} + \begin{bmatrix} L_2 \cos (\theta_2 - (\pi - \theta_1)) \\ L_2 \sin (\theta_2 - (\pi - \theta_1)) \end{bmatrix}.$$
 (5.194)

The second term of the right hand side can be written as

$$\begin{bmatrix} L_2 \cos (\theta_1 + \theta_2 - \pi) \\ L_2 \sin (\theta_1 + \theta_2 - \pi) \end{bmatrix} = \begin{bmatrix} -L_2 \cos(\theta_1 + \theta_2) \\ -L_2 \sin(\theta_1 + \theta_2) \end{bmatrix}.$$
 (5.195)

Therefore,

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} L_1 \cos \theta_1 - L_2 \cos(\theta_1 + \theta_2) \\ L_1 \sin \theta_1 - L_2 \sin(\theta_1 + \theta_2) \end{bmatrix}.$$
 (5.196)

#### 5.34

Let

$$E_n(\mathbf{w}) = -\ln\left(\sum_{k=1}^K \pi_k \mathcal{N}_{nk}\right),\tag{5.197}$$

where

$$\pi_k = \frac{\exp(a_k)}{\sum_{k=1}^K \exp(a_k)},$$

$$\mathcal{N}_{nk} = \mathcal{N}\left(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})\mathbf{I}\right).$$
(5.198)

Then,

$$\frac{\partial E_n}{\partial a_k^{\pi}} = \sum_{k'=1}^K \frac{\partial E_n}{\partial \pi_{k'}} \frac{\partial \pi_{k'}}{\partial a_k^{\pi}}.$$
 (5.199)

We have

$$\frac{\partial E_n}{\partial \pi_{k'}} = -\frac{\mathcal{N}_{nk'}}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}}.$$
 (5.200)

By 5.32(a),

$$\frac{\partial \pi_{k'}}{\partial a_k^{\pi}} = I_{kk'} \pi_{k'} - \pi_k \pi_{k'}. \tag{5.201}$$

Then,

$$\frac{\partial E_n}{\partial a_k^{\pi}} = -\sum_{k'=1}^K \frac{\mathcal{N}_{nk'}}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}} (I_{kk'} \pi_{k'} - \pi_k \pi_{k'}). \tag{5.202}$$

The right hand side can be written as

$$\pi_k \frac{\sum_{k'=1}^K \pi_{k'} \mathcal{N}_{nk'}}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}} - \frac{\pi_k \mathcal{N}_{nk}}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}} = \pi_k - \frac{\pi_k \mathcal{N}_{nk}}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}}.$$
 (5.203)

Therefore,

$$\frac{\partial E_n}{\partial a_k^{\pi}} = \pi_k - \gamma_{nk},\tag{5.204}$$

where

$$\gamma_{nk} = \frac{\pi_k \mathcal{N}_{nk}}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}}.$$
 (5.205)

## 5.35

Let

$$E_n(\mathbf{w}) = -\ln\left(\sum_{k=1}^K \pi_k \mathcal{N}_{nk}\right),\tag{5.206}$$

where

$$\mathcal{N}_{nk} = \mathcal{N}\left(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})\mathbf{I}\right),$$
  

$$\mu_{km} = a_{km}^{\mu}.$$
(5.207)

Then,

$$\frac{\partial E_n}{\partial a_{km}^{\mu}} = \frac{\partial E_n}{\partial \mathcal{N}_{nk}} \frac{\partial \mathcal{N}_{nk}}{\partial a_{km}^{\mu}}.$$
 (5.208)

We have

$$\frac{\partial E_n}{\partial \mathcal{N}_{nk}} = -\frac{\pi_k}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}}.$$
 (5.209)

We have

$$\frac{\partial \mathcal{N}_{nk}}{\partial a_{km}^{\mu}} = \frac{t_{nm} - \mu_{km}}{\sigma_k^2} \mathcal{N}_{nk}.$$
 (5.210)

Then,

$$\frac{\partial E_n}{\partial a_{km}^{\mu}} = -\frac{\pi_k}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}} \frac{t_{nm} - \mu_{km}}{\sigma_k^2} \mathcal{N}_{nk}.$$
 (5.211)

Therefore,

$$\frac{\partial E_n}{\partial a_{km}^{\mu}} = \gamma_{nk} \frac{\mu_{km} - t_{nm}}{\sigma_k^2},\tag{5.212}$$

where

$$\gamma_{nk} = \frac{\pi_k \mathcal{N}_{nk}}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}}.$$
 (5.213)

#### 5.36

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables in M dimensions. Let

$$E_n(\mathbf{w}) = -\ln\left(\sum_{k=1}^K \pi_k \mathcal{N}_{nk}\right),\tag{5.214}$$

where

$$\mathcal{N}_{nk} = \mathcal{N} \left( \mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}) \mathbf{I} \right),$$

$$\sigma_k = \exp\left(a_k^{\sigma}\right).$$
(5.215)

Then,

$$\frac{\partial E_n}{\partial a_k^\sigma} = \frac{\partial E_n}{\partial \mathcal{N}_{nk}} \frac{\partial \mathcal{N}_{nk}}{\partial a_k^\sigma}.$$
 (5.216)

We have

$$\frac{\partial E_n}{\partial \mathcal{N}_{nk}} = -\frac{\pi_k}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}}.$$
 (5.217)

We have

$$\frac{\partial \mathcal{N}_{nk}}{\partial a_k^{\sigma}} = \frac{\partial \mathcal{N}_{nk}}{\partial \sigma_k} \frac{\partial \sigma_k}{\partial a_k^{\sigma}}.$$
 (5.218)

The right hand side can be written as

$$\left(-M\sigma_k^{-1} + \frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{\sigma_k^3}\right) \mathcal{N}_{nk} \sigma_k = \left(-M + \frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{\sigma_k^2}\right) \mathcal{N}_{nk}.$$
(5.219)

Then,

$$\frac{\partial E_n}{\partial a_k^{\sigma}} = -\frac{\pi_k}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}} \left( -M + \frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{\sigma_k^2} \right) \mathcal{N}_{nk}. \tag{5.220}$$

Therefore,

$$\frac{\partial E_n}{\partial a_k^{\sigma}} = \gamma_{nk} \left( M - \frac{\|\mathbf{t}_n - \boldsymbol{\mu}_k\|^2}{\sigma_k^2} \right), \tag{5.221}$$

where

$$\gamma_{nk} = \frac{\pi_k \mathcal{N}_{nk}}{\sum_{k=1}^K \pi_k \mathcal{N}_{nk}}.$$
 (5.222)

#### 5.37

Let  $\mathbf{t}$  be a variable such that

$$p(\mathbf{t}|\mathbf{x}) = \sum_{k=1}^{K} \pi_k(\mathbf{x}) \mathcal{N}\left(\mathbf{t}|\boldsymbol{\mu}_k(\mathbf{x}), \sigma_k^2(\mathbf{x})\mathbf{I}\right).$$
 (5.223)

(a)

We have

$$E(\mathbf{t}|\mathbf{x}) = \int \mathbf{t}p(\mathbf{t}|\mathbf{x})d\mathbf{t}.$$
 (5.224)

The right hand side can be written as

$$\sum_{k=1}^{K} \pi_k(\mathbf{x}) \int \mathbf{t} \mathcal{N} \left( \mathbf{t} | \boldsymbol{\mu}_k(\mathbf{x}), \sigma_k^2(\mathbf{x}) \mathbf{I} \right) d\mathbf{t} = \sum_{k=1}^{K} \pi_k(\mathbf{x}) \boldsymbol{\mu}_k(\mathbf{x}).$$
 (5.225)

Therefore,

$$E(\mathbf{t}|\mathbf{x}) = \sum_{k=1}^{K} \pi_k(\mathbf{x}) \boldsymbol{\mu}_k(\mathbf{x}). \tag{5.226}$$

(b)

We have

$$cov(\mathbf{t}|\mathbf{x}) = \int (\mathbf{t} - E(\mathbf{t}|\mathbf{x})) (\mathbf{t} - E(\mathbf{t}|\mathbf{x}))^{\mathsf{T}} p(\mathbf{t}|\mathbf{x}) d\mathbf{t}.$$
 (5.227)

The right hand side can be written as

$$\int \mathbf{t} \mathbf{t}^{\mathsf{T}} p(\mathbf{t}|\mathbf{x}) d\mathbf{t} - \mathbf{E}(\mathbf{t}|\mathbf{x}) \left( \int \mathbf{t} p(\mathbf{t}|\mathbf{x}) d\mathbf{t} \right)^{\mathsf{T}} - \left( \int \mathbf{t} p(\mathbf{t}|\mathbf{x}) d\mathbf{t} \right) \mathbf{E}(\mathbf{t}|\mathbf{x})^{\mathsf{T}} 
+ \mathbf{E}(\mathbf{t}|\mathbf{x}) \mathbf{E}(\mathbf{t}|\mathbf{x})^{\mathsf{T}} \int p(\mathbf{t}|\mathbf{x}) d\mathbf{t} 
= \int \mathbf{t} \mathbf{t}^{\mathsf{T}} p(\mathbf{t}|\mathbf{x}) d\mathbf{t} - \mathbf{E}(\mathbf{t}|\mathbf{x}) \mathbf{E}(\mathbf{t}|\mathbf{x})^{\mathsf{T}}.$$
(5.228)

The first term of the right hand side can be written as

$$\sum_{k=1}^{K} \pi_k(\mathbf{x}) \int \mathbf{t} \mathbf{t}^{\mathsf{T}} \mathcal{N}\left(\mathbf{t} | \boldsymbol{\mu}_k(\mathbf{x}), \sigma_k^2(\mathbf{x}) \mathbf{I}\right) d\mathbf{t}. \tag{5.229}$$

The integral of the right hand side can be written as

$$\int (\mathbf{t} - \boldsymbol{\mu}_{k}(\mathbf{x}) + \boldsymbol{\mu}_{k}(\mathbf{x})) (\mathbf{t} - \boldsymbol{\mu}_{k}(\mathbf{x}) + \boldsymbol{\mu}_{k}(\mathbf{x}))^{\mathsf{T}} \mathcal{N} (\mathbf{t} | \boldsymbol{\mu}_{k}(\mathbf{x}), \sigma_{k}^{2}(\mathbf{x}) \mathbf{I}) d\mathbf{t}$$

$$= \int (\mathbf{t} - \boldsymbol{\mu}_{k}(\mathbf{x})) (\mathbf{t} - \boldsymbol{\mu}_{k}(\mathbf{x}))^{\mathsf{T}} \mathcal{N} (\mathbf{t} | \boldsymbol{\mu}_{k}(\mathbf{x}), \sigma_{k}^{2}(\mathbf{x}) \mathbf{I}) d\mathbf{t}$$

$$+ \boldsymbol{\mu}_{k}(\mathbf{x}) \left( \int (\mathbf{t} - \boldsymbol{\mu}_{k}(\mathbf{x})) \mathcal{N} (\mathbf{t} | \boldsymbol{\mu}_{k}(\mathbf{x}), \sigma_{k}^{2}(\mathbf{x}) \mathbf{I}) d\mathbf{t} \right)^{\mathsf{T}}$$

$$+ \left( \int (\mathbf{t} - \boldsymbol{\mu}_{k}(\mathbf{x})) \mathcal{N} (\mathbf{t} | \boldsymbol{\mu}_{k}(\mathbf{x}), \sigma_{k}^{2}(\mathbf{x}) \mathbf{I}) d\mathbf{t} \right) \boldsymbol{\mu}_{k}(\mathbf{x})^{\mathsf{T}}$$

$$+ \boldsymbol{\mu}_{k}(\mathbf{x}) \boldsymbol{\mu}_{k}(\mathbf{x})^{\mathsf{T}} \int \mathcal{N} (\mathbf{t} | \boldsymbol{\mu}_{k}(\mathbf{x}), \sigma_{k}^{2}(\mathbf{x}) \mathbf{I}) d\mathbf{t}.$$
(5.230)

The right hand side can be written as

$$\sigma_k^2(\mathbf{x})\mathbf{I} + \boldsymbol{\mu}_k(\mathbf{x})\boldsymbol{\mu}_k(\mathbf{x})^{\mathsf{T}}.$$
 (5.231)

Therefore, by (a),

$$cov(\mathbf{t}|\mathbf{x}) = \sum_{k=1}^{K} \pi_k(\mathbf{x}) \left( \sigma_k^2(\mathbf{x}) \mathbf{I} + \boldsymbol{\mu}_k(\mathbf{x}) \boldsymbol{\mu}_k(\mathbf{x})^{\mathsf{T}} \right) - \left( \sum_{k=1}^{K} \pi_k(\mathbf{x}) \boldsymbol{\mu}_k(\mathbf{x}) \right) \left( \sum_{k=1}^{K} \pi_k(\mathbf{x}) \boldsymbol{\mu}_k(\mathbf{x}) \right)^{\mathsf{T}}.$$
 (5.232)

#### 5.38

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|y(\mathbf{x}, \mathbf{w}), \beta^{-1}\right),$$
  

$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}\right).$$
(5.233)

By marginalisation,

$$p(t|\mathbf{t}) = \int p(t|\mathbf{w})p(\mathbf{w}|\mathbf{t})d\mathbf{w}.$$
 (5.234)

By the Bayes' theorem,

$$p(\mathbf{w}|\mathbf{t})p(\mathbf{t}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}). \tag{5.235}$$

The logarithm of the right hand side except the terms independent of  $\mathbf{w}$  can be written as  $-E(\mathbf{w})$  where

$$E(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^{N} (t_n - y(\mathbf{x}_n, \mathbf{w}))^2 + \frac{\alpha}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}.$$
 (5.236)

Let  $\mathbf{w}_{\text{MAP}}$  be a stationary point of  $p(\mathbf{w}|\mathbf{t})$ . Then, we have a Taylor series

$$E(\mathbf{w}) = E(\mathbf{w}_{\text{MAP}}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\mathsf{T}} \mathbf{A} (\mathbf{w} - \mathbf{w}_{\text{MAP}}) + O(\|\mathbf{w} - \mathbf{w}_{\text{MAP}}\|^{3}),$$
(5.237)

where

$$\mathbf{A} = \left. \nabla \nabla E(\mathbf{w}) \right|_{\mathbf{w} = \mathbf{w}_{\text{MAP}}}. \tag{5.238}$$

Then,

$$p(\mathbf{w}|\mathbf{t}) \simeq \mathcal{N}\left(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{A}^{-1}\right).$$
 (5.239)

By a Taylor series

$$y(\mathbf{x}, \mathbf{w}) = y(\mathbf{x}, \mathbf{w}_{MAP}) + \mathbf{g}^{\mathsf{T}}(\mathbf{w} - \mathbf{w}_{MAP}) + O(\|\mathbf{w} - \mathbf{w}_{MAP}\|^{2}),$$
 (5.240)

where

$$\mathbf{g} = \left. \nabla_{\mathbf{w}} y(\mathbf{x}, \mathbf{w}) \right|_{\mathbf{w} = \mathbf{w}_{\text{MAP}}},\tag{5.241}$$

we have

$$p(t|\mathbf{w}) \simeq \mathcal{N}\left(t|y(\mathbf{x}, \mathbf{w}_{\text{MAP}}) + \mathbf{g}^{\mathsf{T}}(\mathbf{w} - \mathbf{w}_{\text{MAP}}), \beta^{-1}\right).$$
 (5.242)

Then, the logarithm of the integrand except the terms independent of  $\mathbf{w}$  can be approximated as

$$-\frac{\beta}{2} (t - y(\mathbf{x}, \mathbf{w}_{\text{MAP}}) - \mathbf{g}^{\mathsf{T}} (\mathbf{w} - \mathbf{w}_{\text{MAP}}))^{2} - \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\mathsf{T}} \mathbf{A} (\mathbf{w} - \mathbf{w}_{\text{MAP}})$$

$$= -\frac{1}{2} \begin{bmatrix} \mathbf{w} - \mathbf{w}_{\text{MAP}} \\ t - y(\mathbf{x}, \mathbf{w}_{\text{MAP}}) \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \mathbf{A} + \beta \mathbf{g} \mathbf{g}^{\mathsf{T}} & -\beta \mathbf{g} \\ -\beta \mathbf{g}^{\mathsf{T}} & \beta \end{bmatrix} \begin{bmatrix} \mathbf{w} - \mathbf{w}_{\text{MAP}} \\ t - y(\mathbf{x}, \mathbf{w}_{\text{MAP}}) \end{bmatrix}.$$
(5.243)

By 2.24,

$$\begin{bmatrix} \mathbf{A} + \beta \mathbf{g} \mathbf{g}^{\mathsf{T}} & -\beta \mathbf{g} \\ -\beta \mathbf{g}^{\mathsf{T}} & \beta \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{A}^{-1} \mathbf{g} \\ \mathbf{g}^{\mathsf{T}} \mathbf{A}^{-1} & \beta^{-1} + \mathbf{g}^{\mathsf{T}} \mathbf{A}^{-1} \mathbf{g} \end{bmatrix}. \tag{5.244}$$

Therefore,

$$p(t|\mathbf{t}) \simeq \mathcal{N}\left(t|y(\mathbf{x}, \mathbf{w}_{\text{MAP}}), \beta^{-1} + \mathbf{g}^{\mathsf{T}}\mathbf{A}^{-1}\mathbf{g}\right).$$
 (5.245)

#### 5.39

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|y(\mathbf{x}, \mathbf{w}), \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}\right),$$
 (5.246)

where  $\mathbf{w}$  is a vector in M dimensions. By marginalisation,

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{w})p(\mathbf{w})d\mathbf{w}.$$
 (5.247)

The logarithm of the integrand of the right hand side can be written as

$$-\frac{N}{2}\ln\left(2\pi\beta^{-1}\right) - \frac{\beta}{2}\sum_{n=1}^{N}\left(t_{n} - y(\mathbf{x}_{n}, \mathbf{w})\right)^{2}$$

$$-\frac{M}{2}\ln\left(2\pi\right) - \frac{1}{2}\ln\left(\det\left(\alpha^{-1}\mathbf{I}\right)\right) - \frac{\alpha}{2}\mathbf{w}^{\mathsf{T}}\mathbf{w}$$

$$= -E(\mathbf{w}) - \frac{N+M}{2}\ln(2\pi) + \frac{N}{2}\ln\beta + \frac{M}{2}\ln\alpha,$$

$$(5.248)$$

$$E(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^{N} (t_n - y(\mathbf{x}_n, \mathbf{w}))^2 + \frac{\alpha}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}.$$
 (5.249)

Let  $\mathbf{w}_{\text{MAP}}$  be a stationary point of E. Then, we have a Taylor series

$$E(\mathbf{w}) = E(\mathbf{w}_{\text{MAP}}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\mathsf{T}} \mathbf{A} (\mathbf{w} - \mathbf{w}_{\text{MAP}}) + O(\|\mathbf{w} - \mathbf{w}_{\text{MAP}}\|^{3}),$$
(5.250)

where

$$\mathbf{A} = \left. \nabla \nabla E(\mathbf{w}) \right|_{\mathbf{w} = \mathbf{w}_{\text{MAR}}}. \tag{5.251}$$

Then, the logarithm of the integrand can be approximated as

$$-E(\mathbf{w}_{\text{MAP}}) - \frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\mathsf{T}} \mathbf{A}(\mathbf{w} - \mathbf{w}_{\text{MAP}})$$
$$-\frac{N+M}{2} \ln(2\pi) + \frac{N}{2} \ln\beta + \frac{M}{2} \ln\alpha.$$
(5.252)

Therefore,

$$\ln p(\mathbf{t}) \simeq -E(\mathbf{w}_{\text{MAP}}) - \frac{1}{2} \ln |\det \mathbf{A}| - \frac{N}{2} \ln(2\pi) + \frac{N}{2} \ln \beta + \frac{M}{2} \ln \alpha.$$
 (5.253)

## 5.40 (Incomplete)

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables such that

$$p(\mathbf{t}_n|\mathbf{w}) = \prod_{k=1}^K y_{nk}^{t_{nk}},$$

$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}\right),$$
(5.254)

where

$$t_{nk} \in \{0, 1\},$$

$$y_{nk} = y_k(\mathbf{x}_n, \mathbf{w}),$$

$$\sum_{k=1}^{K} y_{nk} = 1.$$

$$(5.255)$$

By marginalisation,

$$p(\mathbf{t}|\mathbf{T}) = \int p(\mathbf{t}|\mathbf{w})p(\mathbf{w}|\mathbf{T})d\mathbf{w}.$$
 (5.256)

By the Bayes' theorem,

$$p(\mathbf{w}|\mathbf{T})p(\mathbf{T}) = p(\mathbf{T}|\mathbf{w})p(\mathbf{w}). \tag{5.257}$$

The logarithm of the right hand side except the terms independent of  $\mathbf{w}$  can be written as  $-E(\mathbf{w})$  where

$$E(\mathbf{w}) = -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln y_{nk} + \frac{\alpha}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}.$$
 (5.258)

Let  $\mathbf{w}_{\text{MAP}}$  be a stationary point of E. Then, we have a Taylor series

$$E(\mathbf{w}) = E(\mathbf{w}_{\text{MAP}}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\mathsf{T}} \mathbf{A} (\mathbf{w} - \mathbf{w}_{\text{MAP}}) + O(\|(\mathbf{w} - \mathbf{w}_{\text{MAP}})\|^{3}),$$
(5.259)

where

$$\mathbf{A} = \left. \nabla \nabla E(\mathbf{w}) \right|_{\mathbf{w} = \mathbf{w}_{\text{MAP}}}. \tag{5.260}$$

Then,

$$p(\mathbf{w}|\mathbf{T}) \simeq \mathcal{N}\left(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{A}^{-1}\right).$$
 (5.261)

#### 5.41

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables such that

$$p(\mathbf{t}_n|\mathbf{w}) = \prod_{k=1}^K y_{nk}^{t_{nk}},$$

$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}\right),$$
(5.262)

where

$$t_{nk} \in \{0, 1\},$$

$$y_{nk} = y_k(\mathbf{x}_n, \mathbf{w}),$$

$$\sum_{k=1}^{K} y_{nk} = 1.$$

$$(5.263)$$

By marginalisation,

$$p(\mathbf{T}) = \int p(\mathbf{T}|\mathbf{w})p(\mathbf{w})d\mathbf{w}.$$
 (5.264)

The logarithm of the integrand of the right hand side can be written as

$$\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln y_{nk} - \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln\left(\det\left(\alpha^{-1}\mathbf{I}\right)\right) - \frac{\alpha}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}$$

$$= -E(\mathbf{w}) - \frac{M}{2} \ln(2\pi) + \frac{M}{2} \ln \alpha,$$
(5.265)

where

$$E(\mathbf{w}) = -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln y_{nk} + \frac{\alpha}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}.$$
 (5.266)

Let  $\mathbf{w}_{\text{MAP}}$  be a stationary point of E. Then, we have a Taylor series

$$E(\mathbf{w}) = E(\mathbf{w}_{\text{MAP}}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\mathsf{T}} \mathbf{A} (\mathbf{w} - \mathbf{w}_{\text{MAP}}) + O(\|(\mathbf{w} - \mathbf{w}_{\text{MAP}})\|^{3}),$$
(5.267)

where

$$\mathbf{A} = \left. \nabla \nabla E(\mathbf{w}) \right|_{\mathbf{w} = \mathbf{w}_{\text{MAP}}}. \tag{5.268}$$

Then, the logarithm of the integrand can be approximated as

$$-E(\mathbf{w}_{\text{MAP}}) - \frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MAP}})^{\mathsf{T}}\mathbf{A}(\mathbf{w} - \mathbf{w}_{\text{MAP}}) - \frac{M}{2}\ln(2\pi) + \frac{M}{2}\ln\alpha. \quad (5.269)$$

Therefore,

$$\ln p(\mathbf{T}) \simeq -E(\mathbf{w}_{\text{MAP}}) - \frac{1}{2} \ln |\det \mathbf{A}| + \frac{M}{2} \ln \alpha.$$
 (5.270)

# 6 Kernel Methods

## 6.1

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n), 1),$$
  

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}).$$
(6.1)

By the Bayes' theorem,

$$p(\mathbf{w}|\mathbf{t})p(\mathbf{t}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}). \tag{6.2}$$

The logarithm of the right hand side except the terms independent of  $\mathbf{w}$  can be written as  $-E_{\mathbf{w}}(\mathbf{w})$  where

$$E_{\mathbf{w}}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} (t_n - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n))^2 + \frac{\lambda}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w},$$
(6.3)

so that

$$E_{\mathbf{w}}(\mathbf{w}) = \frac{1}{2} \|\mathbf{t} - \mathbf{\Phi}\mathbf{w}\|^2 + \frac{\lambda}{2} \mathbf{w}^{\mathsf{T}} \mathbf{w}, \tag{6.4}$$

where

$$\mathbf{\Phi} = \begin{bmatrix} \boldsymbol{\phi}(\mathbf{x}_1)^{\mathsf{T}} \\ \vdots \\ \boldsymbol{\phi}(\mathbf{x}_N)^{\mathsf{T}} \end{bmatrix}. \tag{6.5}$$

(a)

Setting the derivative of  $E_{\mathbf{w}}$  with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = -\mathbf{\Phi}^{\mathsf{T}}(\mathbf{t} - \mathbf{\Phi}\mathbf{w}) + \lambda\mathbf{w}. \tag{6.6}$$

Let the stationary point of  $E_{\mathbf{w}}$  be  $\mathbf{w}_{\text{MAP}}$ . Then,

$$\mathbf{w}_{\text{MAP}} = \mathbf{\Phi}^{\mathsf{T}} \mathbf{a}_{\text{MAP}},\tag{6.7}$$

$$\mathbf{a}_{\text{MAP}} = \frac{1}{\lambda} \left( \mathbf{t} - \mathbf{\Phi} \mathbf{w}_{\text{MAP}} \right). \tag{6.8}$$

(b)

Let

$$\mathbf{w} = \mathbf{\Phi}^{\mathsf{T}} \mathbf{a}.\tag{6.9}$$

Then,

$$E_{\mathbf{w}}(\mathbf{w}) = E_{\mathbf{a}}(\mathbf{a}),\tag{6.10}$$

where

$$E_{\mathbf{a}}(\mathbf{a}) = \frac{1}{2} \|\mathbf{t} - \mathbf{\Phi} \mathbf{\Phi}^{\mathsf{T}} \mathbf{a}\|^{2} + \frac{\lambda}{2} \mathbf{a}^{\mathsf{T}} \mathbf{\Phi} \mathbf{\Phi}^{\mathsf{T}} \mathbf{a}. \tag{6.11}$$

The right hand side can be written as

$$\frac{1}{2} \|\mathbf{t}\|^2 - \mathbf{t}^{\mathsf{T}} \mathbf{K} \mathbf{a} + \frac{1}{2} \mathbf{a}^{\mathsf{T}} \mathbf{K} \mathbf{K} \mathbf{a} + \frac{\lambda}{2} \mathbf{a}^{\mathsf{T}} \mathbf{K} \mathbf{a}, \tag{6.12}$$

where

$$\mathbf{K} = \mathbf{\Phi} \mathbf{\Phi}^{\mathsf{T}}.\tag{6.13}$$

Setting the derivative of  $E_{\mathbf{a}}$  with respect to  $\mathbf{a}$  to zero gives

$$\mathbf{0} = -\mathbf{K}\mathbf{t} + \mathbf{K}\mathbf{K}\mathbf{a} + \lambda \mathbf{K}\mathbf{a}. \tag{6.14}$$

Then, the stationary point of  $E_{\mathbf{a}}$  is given by

$$\mathbf{a}_{\text{MAP}} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t}. \tag{6.15}$$

(c)

By (a) and (b),

$$\frac{1}{\lambda} \left( \mathbf{t} - \mathbf{\Phi} \mathbf{w}_{\text{MAP}} \right) = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t}, \tag{6.16}$$

so that

$$\mathbf{t} - \mathbf{\Phi} \mathbf{w}_{\text{MAP}} = \lambda (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t}. \tag{6.17}$$

The right hand side can be written as

$$(\mathbf{K} + \lambda \mathbf{I} - \mathbf{K})(\mathbf{K} + \lambda \mathbf{I})^{-1}\mathbf{t} = \mathbf{t} - \mathbf{K}(\mathbf{K} + \lambda \mathbf{I})^{-1}\mathbf{t}.$$
 (6.18)

Then,

$$\mathbf{\Phi w}_{\text{MAP}} = \mathbf{K}(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t}, \tag{6.19}$$

so that

$$\mathbf{\Phi}^{\mathsf{T}}\mathbf{\Phi}\mathbf{w}_{\mathrm{MAP}} = \mathbf{\Phi}^{\mathsf{T}}\mathbf{K}(\mathbf{K} + \lambda \mathbf{I})^{-1}\mathbf{t}.$$
 (6.20)

$$\mathbf{w}_{\text{MAP}} = \mathbf{\Phi}^{\mathsf{T}} (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t}, \tag{6.21}$$

so that

$$\lambda \mathbf{w}_{\text{MAP}} = \lambda \mathbf{\Phi}^{\mathsf{T}} (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t}. \tag{6.22}$$

Then,

$$(\mathbf{\Phi}^{\mathsf{T}}\mathbf{\Phi} + \lambda \mathbf{I})\mathbf{w}_{\mathrm{MAP}} = \mathbf{\Phi}^{\mathsf{T}}(\mathbf{K} + \lambda \mathbf{I})(\mathbf{K} + \lambda \mathbf{I})^{-1}\mathbf{t}.$$
 (6.23)

Therefore,

$$\mathbf{w}_{\text{MAP}} = (\mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} + \lambda \mathbf{I})^{-1} \mathbf{\Phi}^{\mathsf{T}} \mathbf{t}. \tag{6.24}$$

# 6.2 (Incomplete)

Let  $t_1, \dots, t_N$  be vaariables such that

$$t_n \in \{-1, 1\},\$$

$$p(t_n | \mathbf{w}) = \left(\frac{1 + y_n}{2}\right)^{\frac{1 + t_n}{2}} \left(\frac{1 - y_n}{2}\right)^{\frac{1 - t_n}{2}},$$
(6.25)

where

$$y_n = f(\mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n),$$

$$f(a) = \begin{cases} 1, & a \ge 0, \\ -1, & \text{otherwise.} \end{cases}$$
(6.26)

Let

$$E(\mathbf{w}) = -\sum_{n \in \mathcal{M}} \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n t_n, \tag{6.27}$$

where  $\mathcal{M}$  is the set of all misclassification. Setting the derivative of E with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = -\sum_{n \in \mathcal{M}} t_n \boldsymbol{\phi}_n. \tag{6.28}$$

#### 6.3

We have

$$\|\mathbf{x} - \mathbf{x}_n\|^2 = \mathbf{x}^\mathsf{T} \mathbf{x} - 2\mathbf{x}^\mathsf{T} \mathbf{x}_n + \mathbf{x}_n^\mathsf{T} \mathbf{x}_n. \tag{6.29}$$

The right hand side can be written as

$$k(\mathbf{x}, \mathbf{x}) - 2k(\mathbf{x}, \mathbf{x}_n) + k(\mathbf{x}_n, \mathbf{x}_n), \tag{6.30}$$

where

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^{\mathsf{T}} \mathbf{x}'. \tag{6.31}$$

#### 6.4

Let **A** be a  $2 \times 2$  matrix with positive eigenvalues and with at least one negative element. We have

$$\det(\lambda \mathbf{I} - \mathbf{A}) = (\lambda - A_{11})(\lambda - A_{22}) - A_{12}A_{21}. \tag{6.32}$$

The right hand side can be written as

$$\lambda^2 - (A_{11} + A_{22})\lambda + A_{11}A_{22} - A_{12}A_{21}. \tag{6.33}$$

Then,

$$(A_{11} + A_{22})^{2} - 4(A_{11}A_{22} - A_{12}A_{21}) > 0,$$

$$A_{11} + A_{22} > 0,$$

$$A_{11}A_{22} - A_{12}A_{21} > 0.$$
(6.34)

Therefore, an example is

$$\mathbf{A} = \begin{bmatrix} 2 & 0 \\ -1 & 1 \end{bmatrix}. \tag{6.35}$$

#### 6.5

Let

$$k_1(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}_1(\mathbf{x})^{\mathsf{T}} \boldsymbol{\phi}_1(\mathbf{x}'). \tag{6.36}$$

(a)

Let

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}'), \tag{6.37}$$

where c is a positive constant. The right hand side can be written as

$$c\phi_1(\mathbf{x})^{\mathsf{T}}\phi_1(\mathbf{x}') = \phi(\mathbf{x})^{\mathsf{T}}\phi(\mathbf{x}'), \tag{6.38}$$

where

$$\phi = \sqrt{c}\phi_1. \tag{6.39}$$

Therefore, k is a valid kernel.

(b)

Let

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}'). \tag{6.40}$$

The right hand side can be written as

$$f(\mathbf{x})\phi_1(\mathbf{x})^{\mathsf{T}}\phi_1(\mathbf{x}')f(\mathbf{x}') = \phi(\mathbf{x})^{\mathsf{T}}\phi(\mathbf{x}'), \tag{6.41}$$

where

$$\phi = f\phi_1. \tag{6.42}$$

Therefore, k is a valid kernel.

6.6

Let

$$k_1(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}_1(\mathbf{x})^{\mathsf{T}} \boldsymbol{\phi}_1(\mathbf{x}'). \tag{6.43}$$

(a)

Let

$$k(\mathbf{x}, \mathbf{x}') = \sum_{j=0}^{J} c_j q_j(\mathbf{x}, \mathbf{x}'), \tag{6.44}$$

where  $c_0, \dots, c_J$  are nonnegative constants and

$$q_j(\mathbf{x}, \mathbf{x}') = (k_1(\mathbf{x}, \mathbf{x}'))^j. \tag{6.45}$$

We have

$$q_0(\mathbf{x}, \mathbf{x}') = 1. \tag{6.46}$$

Let us assume that  $q_j$  is a valid kernel. Then,

$$q_{j+1}(\mathbf{x}, \mathbf{x}') = q_j(\mathbf{x}, \mathbf{x}')\phi_1(\mathbf{x})^{\mathsf{T}}\phi_1(\mathbf{x}'). \tag{6.47}$$

The right hand side can be written as

$$q_j(\mathbf{x}, \mathbf{x}') \sum_{m=1}^M \phi_{1m}(\mathbf{x}) \phi_{1m}(\mathbf{x}') = \sum_{m=1}^M \phi_{1m}(\mathbf{x}) q_j(\mathbf{x}, \mathbf{x}') \phi_{1m}(\mathbf{x}').$$
 (6.48)

By 6.5(b) and 6.7(a), the right hand side is a valid kernel. Then,  $q_{j+1}$  is a valid kernel. Then, the assumption is proved by induction on j. Therefore, by 6.5(a) and 6.7(a), k is a valid kernel.

(b)

Let

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(k_1(\mathbf{x}, \mathbf{x}')\right). \tag{6.49}$$

The right hand side can be written as

$$\sum_{j=0}^{\infty} \frac{1}{j!} \left( k_1(\mathbf{x}, \mathbf{x}') \right)^j. \tag{6.50}$$

Therefore, by (a), k is a valid kernel.

6.7

Let

$$k_1(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}_1(\mathbf{x})^{\mathsf{T}} \boldsymbol{\phi}_1(\mathbf{x}'), k_2(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}_2(\mathbf{x})^{\mathsf{T}} \boldsymbol{\phi}_2(\mathbf{x}').$$
(6.51)

(a)

Let

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}'). \tag{6.52}$$

The right hand side can be written as

$$\phi_1(\mathbf{x})^{\mathsf{T}}\phi_1(\mathbf{x}') + \phi_2(\mathbf{x})^{\mathsf{T}}\phi_2(\mathbf{x}') = \phi(\mathbf{x})^{\mathsf{T}}\phi(\mathbf{x}'), \tag{6.53}$$

where

$$\phi = \begin{bmatrix} \phi_1 \\ \phi_2 \end{bmatrix}. \tag{6.54}$$

Therefore, k is a valid kernel.

(b)

Let

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}'). \tag{6.55}$$

Then,

$$\ln k(\mathbf{x}, \mathbf{x}') = \ln k_1(\mathbf{x}, \mathbf{x}') + \ln k_2(\mathbf{x}, \mathbf{x}'). \tag{6.56}$$

By (a) and 6.6(b), the right hand side is a valid kernel. Then,  $\ln k$  is a valid kernel. Therefore, by 6.6(b), k is a valid kernel.

6.8

(a)

Let

$$k(\mathbf{x}, \mathbf{x}') = k_1(\phi(\mathbf{x}), \phi(\mathbf{x}')), \qquad (6.57)$$

where  $k_1$  is a valid kernel. Then, k is a valid kernel.

(b)

Let

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^{\mathsf{T}} \mathbf{A} \mathbf{x}',\tag{6.58}$$

where  $\mathbf{A}$  is a symmetric positive semidefinite matrix. There exists  $\mathbf{M}$  such that

$$\mathbf{A} = \mathbf{M}^{\mathsf{T}} \mathbf{M}.\tag{6.59}$$

Then,

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^{\mathsf{T}} \phi(\mathbf{x}'), \tag{6.60}$$

where

$$\phi(\mathbf{x}) = \mathbf{M}\mathbf{x}.\tag{6.61}$$

Therefore, k is a valid kernel.

6.9

Let

$$k_a(\mathbf{x}_a, \mathbf{x}_a') = \boldsymbol{\phi}_a(\mathbf{x}_a)^{\mathsf{T}} \boldsymbol{\phi}_a(\mathbf{x}_a'), k_b(\mathbf{x}_b, \mathbf{x}_b') = \boldsymbol{\phi}_b(\mathbf{x}_b)^{\mathsf{T}} \boldsymbol{\phi}_b(\mathbf{x}_b').$$
(6.62)

(a)

Let

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b), \tag{6.63}$$

where

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}. \tag{6.64}$$

The right hand side can be written as

$$\phi_a(\mathbf{x}_a)^{\mathsf{T}}\phi_a(\mathbf{x}_a') + \phi_b(\mathbf{x}_b)^{\mathsf{T}}\phi_b(\mathbf{x}_b') = \phi(\mathbf{x})^{\mathsf{T}}\phi(\mathbf{x}'), \tag{6.65}$$

where

$$\phi(\mathbf{x}) = \begin{bmatrix} \phi_a(\mathbf{x}_a) \\ \phi_b(\mathbf{x}_b) \end{bmatrix}. \tag{6.66}$$

Therefore, k is a valid kernel.

(b)

Let

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b), \tag{6.67}$$

where

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}. \tag{6.68}$$

Then,

$$\ln k(\mathbf{x}, \mathbf{x}') = \ln k_a(\mathbf{x}_a, \mathbf{x}'_a) + \ln k_b(\mathbf{x}_b, \mathbf{x}'_b). \tag{6.69}$$

By (a) and 6.6(b), the right hand side is a valid kernel. Then,  $\ln k$  is a valid kernel. Therefore, by 6.6(b), k is a valid kernel.

# 6.10 (Incomplete)

Let

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})f(\mathbf{x}'). \tag{6.70}$$

#### 6.11

Let

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right). \tag{6.71}$$

The right hand side can be written as

$$\exp\left(-\frac{\mathbf{x}^{\mathsf{T}}\mathbf{x}}{2\sigma^{2}}\right)\exp\left(\frac{\mathbf{x}^{\mathsf{T}}\mathbf{x}'}{\sigma^{2}}\right)\exp\left(-\frac{\mathbf{x}'^{\mathsf{T}}\mathbf{x}'}{2\sigma^{2}}\right). \tag{6.72}$$

We have

$$\exp\left(\frac{\mathbf{x}^{\mathsf{T}}\mathbf{x}'}{\sigma^2}\right) = \sum_{j=0}^{\infty} q_j(\mathbf{x}, \mathbf{x}'), \tag{6.73}$$

where

$$q_j(\mathbf{x}, \mathbf{x}') = \frac{1}{j!} \left( \frac{\mathbf{x}^\mathsf{T} \mathbf{x}'}{\sigma^2} \right)^j. \tag{6.74}$$

By 6.6(a),  $q_j$  is a valid kernel. Then, there exists  $\psi_j$  such that

$$q_j(\mathbf{x}, \mathbf{x}') = \boldsymbol{\psi}_j(\mathbf{x})^{\mathsf{T}} \boldsymbol{\psi}_j(\mathbf{x}'). \tag{6.75}$$

Then,

$$\exp\left(\frac{\mathbf{x}^{\mathsf{T}}\mathbf{x}'}{\sigma^2}\right) = \sum_{j=0}^{\infty} \boldsymbol{\psi}_j(\mathbf{x})^{\mathsf{T}} \boldsymbol{\psi}_j(\mathbf{x}'). \tag{6.76}$$

The right hand side can be written as  $\phi(\mathbf{x})^{\mathsf{T}}\phi(\mathbf{x}')$ , where

$$\boldsymbol{\phi} = \begin{bmatrix} \boldsymbol{\psi}_1 \\ \vdots \\ \boldsymbol{\psi}_j \\ \vdots \end{bmatrix} . \tag{6.77}$$

Therefore, k can be written as the inner product of an infinite-dimensional feature vector.

#### 6.12

Let  $A_1, \dots, A_{2^{|D|}}$  are all the subsets of D ordered in the alphabetical order where |.| is the number of elements in (.). Let

$$k(A_m, A_{m'}) = 2^{|A_m \cap A_{m'}|}. (6.78)$$

Let  $\phi$  be a vector such that

$$\phi_m(A) = \begin{cases} 1, & \text{if } A_m \subseteq A, \\ 0, & \text{otherwise.} \end{cases}$$
 (6.79)

If

$$|A_m \cap A_{m'}| = 0, (6.80)$$

then

$$k(A_m, A_{m'}) = 1. (6.81)$$

The right hand side can be written as

$$\boldsymbol{\phi}(A_m)^{\mathsf{T}} \boldsymbol{\phi}(A_{m'}). \tag{6.82}$$

Let us assume that if

$$|A_m \cap A_{m'}| = n, \tag{6.83}$$

then

$$k(A_m, A_{m'}) = \boldsymbol{\phi}(A_m)^{\mathsf{T}} \boldsymbol{\phi}(A_{m'}). \tag{6.84}$$

Let

$$A_{m''} = A_{m'} \cup \{x\},\tag{6.85}$$

where x is an element such that

$$x \in A_m - A_{m'}. \tag{6.86}$$

Then,

$$|A_m \cap A_{m''}| = n + 1, (6.87)$$

so that

$$k(A_m, A_{m''}) = 2^{n+1}. (6.88)$$

We have

$$\phi(A_m)^{\mathsf{T}}\phi(A_{m''}) = \phi(A_m)^{\mathsf{T}}\phi(A_{m'}) + \phi(A_m)^{\mathsf{T}}(\phi(A_{m''}) - \phi(A_{m'})). \quad (6.89)$$

By the assumption, the first term of the right hand side is  $2^n$ . By the assumption,  $\phi(A_m)$  and  $\phi(A_{m'})$  have a set of  $2^n$  1s in common. Then,  $\phi(A_m)$  and  $\phi(A_{m''}) - \phi(A_{m'})$  have another set of  $2^n$  1s in common. Then, the second term of the right hand side is  $2^n$ . Then,

$$\phi(A_m)^{\mathsf{T}}\phi(A_{m''}) = 2^{n+1},$$
 (6.90)

so that

$$k(A_m, A_{m''}) = \phi(A_m)^{\mathsf{T}} \phi(A_{m''}).$$
 (6.91)

Therefore, the assumption is proved by induction on n.

6.13

Let

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^{\mathsf{T}} \mathbf{F}^{-1} \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}'),$$
 (6.92)

where

$$\mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) = \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x}|\boldsymbol{\theta}),$$

$$\mathbf{F} = \mathbf{E}_{\mathbf{x}} \left( \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^{\mathsf{T}} \right).$$
(6.93)

Let

$$k_{\psi}(\mathbf{x}, \mathbf{x}') = \mathbf{g}_{\psi}(\boldsymbol{\theta}, \mathbf{x})^{\mathsf{T}} \mathbf{F}_{\psi}^{-1} \mathbf{g}_{\psi}(\boldsymbol{\theta}, \mathbf{x}'),$$
 (6.94)

where  $\psi$  is invertible and differentiable and

$$\mathbf{g}_{\psi}(\boldsymbol{\theta}, \mathbf{x}) = \nabla_{\psi(\boldsymbol{\theta})} \ln p(\mathbf{x}|\boldsymbol{\theta}),$$
  

$$\mathbf{F}_{\psi} = \mathbf{E}_{\mathbf{x}} \left( \mathbf{g}_{\psi}(\boldsymbol{\theta}, \mathbf{x}) \mathbf{g}_{\psi}(\boldsymbol{\theta}, \mathbf{x})^{\mathsf{T}} \right).$$
(6.95)

We have

$$\mathbf{g}_{\psi}(\boldsymbol{\theta}, \mathbf{x}) = \frac{\partial \boldsymbol{\theta}}{\partial \psi(\boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x}|\boldsymbol{\theta}). \tag{6.96}$$

Then,

$$\mathbf{g}_{\psi}(\boldsymbol{\theta}, \mathbf{x}) = \left(\frac{\partial \psi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)^{-1} \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}). \tag{6.97}$$

We have

$$\mathbf{F}_{\psi} = \int \mathbf{g}_{\psi}(\boldsymbol{\theta}, \mathbf{x}) \mathbf{g}_{\psi}(\boldsymbol{\theta}, \mathbf{x})^{\mathsf{T}} p(\mathbf{x}) d\mathbf{x}. \tag{6.98}$$

The right hand side can be written as

$$\int \left(\frac{\partial \boldsymbol{\psi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)^{-1} \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^{\mathsf{T}} \left(\left(\frac{\partial \boldsymbol{\psi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)^{-1}\right)^{\mathsf{T}} p(\mathbf{x}) d\mathbf{x}$$

$$= \left(\frac{\partial \boldsymbol{\psi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)^{-1} \left(\int \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^{\mathsf{T}} p(\mathbf{x}) d\mathbf{x}\right) \left(\left(\frac{\partial \boldsymbol{\psi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)^{-1}\right)^{\mathsf{T}}.$$
(6.99)

Then,

$$\mathbf{F}_{\psi} = \left(\frac{\partial \psi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)^{-1} \mathbf{F} \left(\left(\frac{\partial \psi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)^{-1}\right)^{\mathsf{T}}.$$
 (6.100)

Then,

$$k_{\psi}(\mathbf{x}, \mathbf{x}')$$

$$= \mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^{\mathsf{T}} \left( \left( \frac{\partial \boldsymbol{\psi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^{-1} \right)^{\mathsf{T}} \left( \frac{\partial \boldsymbol{\psi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^{\mathsf{T}} \mathbf{F}^{-1} \left( \frac{\partial \boldsymbol{\psi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) \left( \frac{\partial \boldsymbol{\psi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^{-1} \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}').$$
(6.101)

The right hand side can be written as

$$\mathbf{g}(\boldsymbol{\theta}, \mathbf{x})^{\mathsf{T}} \mathbf{F}^{-1} \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}'). \tag{6.102}$$

Therefore,

$$k_{\psi}(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}'). \tag{6.103}$$

#### 6.14

Let  $\mathbf{x}$  be a variable such that

$$p(\mathbf{x}|\boldsymbol{\mu}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{S}). \tag{6.104}$$

Let

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{g}(\boldsymbol{\mu}, \mathbf{x})^{\mathsf{T}} \mathbf{F}^{-1} \mathbf{g}(\boldsymbol{\mu}, \mathbf{x}'), \tag{6.105}$$

where

$$\mathbf{g}(\boldsymbol{\mu}, \mathbf{x}) = \nabla_{\boldsymbol{\mu}} \ln p(\mathbf{x}|\boldsymbol{\mu}),$$
  
$$\mathbf{F} = \mathbf{E}_{\mathbf{x}} \left( \mathbf{g}(\boldsymbol{\mu}, \mathbf{x}) \mathbf{g}(\boldsymbol{\mu}, \mathbf{x})^{\mathsf{T}} \right).$$
 (6.106)

We have

$$\mathbf{g}(\boldsymbol{\mu}, \mathbf{x}) = -\mathbf{S}^{-1}(\mathbf{x} - \boldsymbol{\mu}). \tag{6.107}$$

Then,

$$\mathbf{F} = \int \mathbf{S}^{-1} (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{S}^{-1} \mathcal{N} (\mathbf{x} | \boldsymbol{\mu}, \mathbf{S}) d\mathbf{x}. \tag{6.108}$$

The right hand side can be written as

$$\mathbf{S}^{-1} \left( \int (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{S}) d\mathbf{x} \right) \mathbf{S}^{-1} = \mathbf{S}^{-1}.$$
 (6.109)

Then,

$$\mathbf{F} = \mathbf{S}^{-1},\tag{6.110}$$

so that

$$\mathbf{g}(\boldsymbol{\mu}, \mathbf{x})^{\mathsf{T}} \mathbf{F}^{-1} \mathbf{g}(\boldsymbol{\mu}, \mathbf{x}') = (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{S}^{-1} (\mathbf{x} - \boldsymbol{\mu}). \tag{6.111}$$

Therefore,

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{S}^{-1} (\mathbf{x} - \boldsymbol{\mu}). \tag{6.112}$$

#### 6.15

Let k be a positive semidefinite kernel function. Then,

$$k(x,x)k(x',x') - k(x,x')^2 = \det \mathbf{K},$$
 (6.113)

where

$$\mathbf{K} = \begin{bmatrix} k(x,x) & k(x,x') \\ k(x,x') & k(x',x') \end{bmatrix}. \tag{6.114}$$

Since  $\mathbf{K}$  is positive semidefinite,

$$\det \mathbf{K} \ge 0. \tag{6.115}$$

Therefore,

$$k(x, x')^2 \le k(x, x)k(x', x').$$
 (6.116)

# 6.16 (Incomplete)

Let

$$J(\mathbf{w}) = f(\mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_1), \cdots, \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_N)) + g(\mathbf{w}^{\mathsf{T}} \mathbf{w}), \qquad (6.117)$$

where g is a monotonically increasing function. Setting the derivative with respect to  $\mathbf{w}$  to zero gives

$$\mathbf{0} = \sum_{n=1}^{N} \frac{\partial f}{\partial \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n)} \boldsymbol{\phi}(\mathbf{x}_n) + 2g'(\mathbf{w}^{\mathsf{T}} \mathbf{w}) \mathbf{w}. \tag{6.118}$$

Therefore,

$$\underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w}) = -\frac{1}{2g'(\mathbf{w}^{\mathsf{T}}\mathbf{w})} \sum_{n=1}^{N} \frac{\partial f}{\partial \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n)} \boldsymbol{\phi}(\mathbf{x}_n). \tag{6.119}$$

#### 6.17

Let

$$E = \frac{1}{2} \sum_{n=1}^{N} \int (y(\mathbf{x}_n + \boldsymbol{\xi}) - t_n)^2 p(\boldsymbol{\xi}) d\boldsymbol{\xi}.$$
 (6.120)

By the transformation

$$\mathbf{z} = \mathbf{x}_n + \boldsymbol{\xi},\tag{6.121}$$

we have

$$E = \frac{1}{2} \sum_{n=1}^{N} \int (y(\mathbf{z}) - t_n)^2 p(\mathbf{z} - \mathbf{x}_n) d\mathbf{z}.$$
 (6.122)

Setting the variation with respect to y to zero gives

$$0 = \sum_{n=1}^{N} (y(\mathbf{z}) - t_n) p(\mathbf{z} - \mathbf{x}_n).$$

$$(6.123)$$

The right hand side can be written as

$$y(\mathbf{z}) \sum_{n=1}^{N} p(\mathbf{z} - \mathbf{x}_n) - \sum_{n=1}^{N} t_n p(\mathbf{z} - \mathbf{x}_n).$$
 (6.124)

Therefore,

$$y(\mathbf{x}) = \sum_{n=1}^{N} t_n k(\mathbf{x}, \mathbf{x}_n), \tag{6.125}$$

where

$$k(\mathbf{x}, \mathbf{x}_n) = \frac{p(\mathbf{x} - \mathbf{x}_n)}{\sum_{n=1}^{N} p(\mathbf{x} - \mathbf{x}_n)}.$$
 (6.126)

### 6.18

Let x be a variable such that

$$p(x,t) = \frac{1}{N} \sum_{n=1}^{N} f(x - x_n, t - t_n), \tag{6.127}$$

where

$$f(x,t) = \mathcal{N}(x|0,\sigma^2) \mathcal{N}(t|0,\sigma^2).$$
 (6.128)

(a)

By the Bayes' theorem,

$$p(t|x) = \frac{p(x,t)}{p(x)}. (6.129)$$

By marginalisation,

$$p(x) = \int p(x,t)dt. \tag{6.130}$$

The right hand side can be written as

$$\frac{1}{N} \sum_{n=1}^{N} \mathcal{N}\left(x - x_n | 0, \sigma^2\right) \int \mathcal{N}\left(t - t_n | 0, \sigma^2\right) dt$$

$$= \frac{1}{N} \sum_{n=1}^{N} \mathcal{N}\left(x - x_n | 0, \sigma^2\right).$$
(6.131)

Then,

$$p(t|x) = \frac{\frac{1}{N} \sum_{n=1}^{N} \mathcal{N}(x - x_n | 0, \sigma^2) \mathcal{N}(t - t_n | 0, \sigma^2)}{\frac{1}{N} \sum_{n=1}^{N} \mathcal{N}(x - x_n | 0, \sigma^2)}.$$
 (6.132)

Therefore,

$$p(t|x) = \sum_{n=1}^{N} k(x, x_n) \mathcal{N} (t - t_n | 0, \sigma^2), \qquad (6.133)$$

where

$$k(x, x_n) = \frac{\mathcal{N}(x - x_n | 0, \sigma^2)}{\sum_{n=1}^{N} \mathcal{N}(x - x_n | 0, \sigma^2)}.$$
 (6.134)

(b)

We have

$$E(t|x) = \int tp(t|x)dt.$$
 (6.135)

By (a), the right hand side can be written as

$$\sum_{n=1}^{N} k(x, x_n) \int t \mathcal{N}\left(t - t_n | 0, \sigma^2\right) dt. \tag{6.136}$$

By the transformation

$$t' = t - t_n, \tag{6.137}$$

the integral can be written as

$$\int (t'+t_n)\mathcal{N}\left(t'|0,\sigma^2\right)dt' = \int t'\mathcal{N}\left(t'|0,\sigma^2\right)dt' + t_n \int \mathcal{N}\left(t'|0,\sigma^2\right)dt'. \quad (6.138)$$

The right hand side can be written as  $t_n$ . Therefore,

$$E(t|x) = \sum_{n=1}^{N} k(x, x_n) t_n.$$
 (6.139)

(c)

We have

$$var(t|x) = \int (t - E(t|x))^2 p(t|x)dt.$$
 (6.140)

By (a) and (b), the right hand side can be written as

$$\sum_{n=1}^{N} k(x, x_n) \int \left( t - \sum_{n=1}^{N} k(x, x_n) t_n \right)^2 \mathcal{N} \left( t - t_n | 0, \sigma^2 \right) dt.$$
 (6.141)

By the transformation

$$t' = t - t_n, \tag{6.142}$$

the integral can be written as

$$\int \left( t' + t_n - \sum_{n=1}^{N} k(x, x_n) t_n \right)^2 \mathcal{N} \left( t' | 0, \sigma^2 \right) dt' 
= \int t'^2 \mathcal{N} \left( t' | 0, \sigma^2 \right) dt' + 2 \left( t_n - \sum_{n=1}^{N} k(x, x_n) t_n \right) \int t' \mathcal{N} \left( t' | 0, \sigma^2 \right) dt' \quad (6.143) 
+ \left( t_n - \sum_{n=1}^{N} k(x, x_n) t_n \right)^2 \int \mathcal{N} \left( t' | 0, \sigma^2 \right) dt'.$$

The right hand side can be written as

$$\sigma^{2} + \left(t_{n} - \sum_{n=1}^{N} k(x, x_{n})t_{n}\right)^{2}.$$
 (6.144)

Then,

$$var(t|x) = \sum_{n=1}^{N} k(x, x_n) \left( \sigma^2 + \left( t_n - \sum_{n'=1}^{N} k(x, x_{n'}) t_{n'} \right)^2 \right).$$
 (6.145)

Therefore,

$$\operatorname{var}(t|x) = \sigma^2 + \sum_{n=1}^{N} k(x, x_n) \left( t_n - \sum_{n'=1}^{N} k(x, x_{n'}) t_{n'} \right)^2.$$
 (6.146)

6.19

Let

$$E = \frac{1}{2} \sum_{n=1}^{N} \int (y(\mathbf{x}_n - \boldsymbol{\xi}_n) - t_n)^2 p(\boldsymbol{\xi}_n) d\boldsymbol{\xi}_n.$$
 (6.147)

By the transformation

$$\mathbf{z} = \mathbf{x}_n - \boldsymbol{\xi}_n,\tag{6.148}$$

we have

$$E = \frac{1}{2} \sum_{n=1}^{N} \int (y(\mathbf{z}) - t_n)^2 p(\mathbf{x}_n - \mathbf{z}) d\mathbf{z}.$$
 (6.149)

Setting the variation with respect to y to zero gives

$$0 = \sum_{n=1}^{N} (y(\mathbf{z}) - t_n) p(\mathbf{x}_n - \mathbf{z}).$$

$$(6.150)$$

The right hand side can be written as

$$y(\mathbf{z})\sum_{n=1}^{N}p(\mathbf{x}_{n}-\mathbf{z})-\sum_{n=1}^{N}t_{n}p(\mathbf{x}_{n}-\mathbf{z}).$$
(6.151)

Therefore,

$$y(\mathbf{x}) = \sum_{n=1}^{N} t_n k(\mathbf{x}, \mathbf{x}_n), \qquad (6.152)$$

where

$$k(\mathbf{x}, \mathbf{x}_n) = \frac{p(\mathbf{x}_n - \mathbf{x})}{\sum_{n=1}^{N} p(\mathbf{x}_n - \mathbf{x})}.$$
(6.153)

#### 6.20

Let  $t_1, \dots, t_N$  be variables such that

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}\left(\mathbf{t}|\mathbf{y}, \beta^{-1}\mathbf{I}\right),$$
  

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}),$$
(6.154)

$$K_{nn'} = k(\mathbf{x}_n, \mathbf{x}_{n'}). \tag{6.155}$$

By the Bayes' theorem,

$$p(t_{N+1}|\mathbf{t})p(\mathbf{t}) = p(\mathbf{t}'), \tag{6.156}$$

where

$$\mathbf{t}' = \begin{bmatrix} \mathbf{t} \\ t_{N+1} \end{bmatrix}. \tag{6.157}$$

By marginalisation,

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y})d\mathbf{y}.$$
 (6.158)

The logarighm of the integrand of the right hand side except the terms independent of  $\mathbf{y}_N$  can be written as

$$-\frac{\beta}{2}(\mathbf{t} - \mathbf{y})^{\mathsf{T}}(\mathbf{t} - \mathbf{y}) - \frac{1}{2}\mathbf{y}^{\mathsf{T}}\mathbf{K}^{-1}\mathbf{y}$$

$$= -\frac{1}{2}\begin{bmatrix} \mathbf{y} \\ \mathbf{t} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \beta \mathbf{I} + \mathbf{K}^{-1} & -\beta \mathbf{I} \\ -\beta \mathbf{I} & \beta \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{t} \end{bmatrix}.$$
(6.159)

By 2.24,

$$\begin{bmatrix} \beta \mathbf{I} + \mathbf{K}^{-1} & -\beta \mathbf{I} \\ -\beta \mathbf{I} & \beta \mathbf{I} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{K} & \mathbf{K} \\ \mathbf{K} & \mathbf{K} + \beta^{-1} \mathbf{I} \end{bmatrix}.$$
 (6.160)

Then,

$$p(\mathbf{t}) = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}), \tag{6.161}$$

where

$$\mathbf{C} = \mathbf{K} + \beta^{-1} \mathbf{I}. \tag{6.162}$$

Then,

$$p(\mathbf{t}') = \mathcal{N}\left(\mathbf{t}'|\mathbf{0}, \mathbf{C}'\right), \tag{6.163}$$

where

$$\mathbf{C}' = \begin{bmatrix} \mathbf{C} & \mathbf{k} \\ \mathbf{k}^{\mathsf{T}} & c \end{bmatrix}, \tag{6.164}$$

where

$$k_n = k(\mathbf{x}_n, \mathbf{x}_{N+1}),$$
  
 $c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}.$  (6.165)

By 2.24,

$$\begin{bmatrix} c & \mathbf{k}^{\mathsf{T}} \\ \mathbf{k} & \mathbf{C} \end{bmatrix}^{-1} = \begin{bmatrix} s^{-1} & -s^{-1}\mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1} \\ -s^{-1}\mathbf{C}^{-1}\mathbf{k} & \mathbf{C}^{-1} + s^{-1}\mathbf{C}^{-1}\mathbf{k}\mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1} \end{bmatrix}, \tag{6.166}$$

where

$$s = c - \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{k}. \tag{6.167}$$

Therefore,

$$p(t_{N+1}|\mathbf{t}) = \mathcal{N}(t_{N+1}|m,s),$$
 (6.168)

where

$$m = \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{t}. \tag{6.169}$$

#### 6.21

Let  $t_1, \dots, t_N$  be variables such that

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}\left(\mathbf{t}|\mathbf{y}, \beta^{-1}\mathbf{I}\right),$$

$$p(\mathbf{y}) = \mathcal{N}\left(\mathbf{y}|\mathbf{0}, \mathbf{\Phi}\mathbf{\Phi}^{\mathsf{T}}\right).$$
(6.170)

By 3.10,

$$p(t_{N+1}|\mathbf{t}) = \mathcal{N}(t_{N+1}|m,s),$$
 (6.171)

where

$$m = \beta \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \left( \mathbf{I} + \beta \boldsymbol{\Phi}^{\mathsf{T}} \boldsymbol{\Phi} \right)^{-1} \boldsymbol{\Phi}^{\mathsf{T}} \mathbf{t},$$
  

$$s = \beta^{-1} + \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \left( \mathbf{I} + \beta \boldsymbol{\Phi}^{\mathsf{T}} \boldsymbol{\Phi} \right)^{-1} \boldsymbol{\phi}_{N+1}.$$
(6.172)

By 2.26,

$$(\mathbf{I} + \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi})^{-1} = \mathbf{I} - \mathbf{\Phi}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{\Phi}, \tag{6.173}$$

where

$$\mathbf{C} = \beta^{-1} \mathbf{I} + \mathbf{\Phi} \mathbf{\Phi}^{\mathsf{T}}.\tag{6.174}$$

Then,

$$m = \beta \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \left( \mathbf{I} - \boldsymbol{\Phi}^{\mathsf{T}} \mathbf{C}^{-1} \boldsymbol{\Phi} \right) \boldsymbol{\Phi}^{\mathsf{T}} \mathbf{t}. \tag{6.175}$$

The right hand side can be written as

$$\beta \left( \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \boldsymbol{\Phi}^{\mathsf{T}} - \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \boldsymbol{\Phi}^{\mathsf{T}} \mathbf{C}^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^{\mathsf{T}} \right) \mathbf{t} = \beta \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \boldsymbol{\Phi}^{\mathsf{T}} \mathbf{C}^{-1} \left( \mathbf{C} - \boldsymbol{\Phi} \boldsymbol{\Phi}^{\mathsf{T}} \right) \mathbf{t}. \quad (6.176)$$

The right hand side can be written as

$$\beta \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \beta^{-1} \mathbf{I} \mathbf{t} = \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{t}, \tag{6.177}$$

$$\mathbf{k} = \mathbf{\Phi} \phi_{N+1}. \tag{6.178}$$

Similarly,

$$s = \beta^{-1} + \boldsymbol{\phi}_{N+1}^{\mathsf{T}} \left( \mathbf{I} - \boldsymbol{\Phi}^{\mathsf{T}} \mathbf{C}^{-1} \boldsymbol{\Phi} \right) \boldsymbol{\phi}_{N+1}. \tag{6.179}$$

The right hand side can be written as

$$\beta^{-1} + \phi_{N+1}^{\mathsf{T}} \phi_{N+1} - \phi_{N+1}^{\mathsf{T}} \Phi^{\mathsf{T}} \mathbf{C}^{-1} \Phi \phi_{N+1} = c - \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{k}, \tag{6.180}$$

where

$$c = \beta^{-1} + \phi_{N+1}^{\mathsf{T}} \phi_{N+1}. \tag{6.181}$$

Therefore,

$$m = \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{t},$$
  

$$s = c - \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{k}.$$
(6.182)

## 6.22

Let  $t_1, \dots, t_N$  be variables such that

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}\left(\mathbf{t}|\mathbf{y}, \beta^{-1}\mathbf{I}\right),$$
  

$$p(\mathbf{y}) = \mathcal{N}\left(\mathbf{y}|\mathbf{0}, \mathbf{K}\right),$$
(6.183)

where

$$K_{nn'} = k(\mathbf{x}_n, \mathbf{x}_{n'}). \tag{6.184}$$

(a)

By the Bayes' theorem,

$$p(t_{N+1}, \cdots, t_{N+L}|\mathbf{t})p(\mathbf{t}) = p(\mathbf{t}'), \tag{6.185}$$

where

$$\mathbf{t}' = \begin{bmatrix} \mathbf{t} \\ t_{N+1} \\ \vdots \\ t_{N+L} \end{bmatrix}. \tag{6.186}$$

By 6.20,

$$p(\mathbf{t}) = \mathcal{N}\left(\mathbf{t}|\mathbf{0}, \mathbf{C}\right),\tag{6.187}$$

$$\mathbf{C} = \beta^{-1} \mathbf{I} + \mathbf{K}. \tag{6.188}$$

Let

$$K'_{nl} = k(\mathbf{x}_n, \mathbf{x}_{N+l}),$$
  

$$\Gamma_{ll'} = \beta^{-1} I_{ll'} + k(\mathbf{x}_{N+l}, \mathbf{x}_{N+l'}).$$
(6.189)

Then,

$$p(\mathbf{t}') = \mathcal{N}\left(\mathbf{t}'|\mathbf{0}, \mathbf{C}'\right),\tag{6.190}$$

where

$$\mathbf{C}' = \begin{bmatrix} \mathbf{C} & \mathbf{K}' \\ \mathbf{K}'^{\mathsf{T}} & \mathbf{\Gamma} \end{bmatrix}. \tag{6.191}$$

By 2.24,

$$\begin{bmatrix} \mathbf{\Gamma} & \mathbf{K}'^{\mathsf{T}} \\ \mathbf{K}' & \mathbf{C} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{S}^{-1} & -\mathbf{S}^{-1}\mathbf{K}'^{\mathsf{T}}\mathbf{C}^{-1} \\ -\mathbf{C}^{-1}\mathbf{K}'\mathbf{S}^{-1} & \mathbf{C}^{-1} + \mathbf{C}^{-1}\mathbf{K}'\mathbf{S}^{-1}\mathbf{K}'^{\mathsf{T}}\mathbf{C}^{-1} \end{bmatrix}, \quad (6.192)$$

where

$$\mathbf{S} = \mathbf{\Gamma} - \mathbf{K}'^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{K}'. \tag{6.193}$$

Therefore,

$$p(t_{N+1}, \dots, t_{N+L} | \mathbf{t}_N) = \mathcal{N}(t_{N+1}, \dots, t_{N+L} | \mathbf{m}, \mathbf{S}),$$
 (6.194)

where

$$\mathbf{m} = \mathbf{K}^{\prime \mathsf{T}} \mathbf{C}^{-1} \mathbf{t}. \tag{6.195}$$

(b)

Let

$$k_n = k(\mathbf{x}_n, \mathbf{x}_{N+1}),$$
  
 $c = \beta^{-1} + k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}).$  (6.196)

By (a),

$$p(t_{N+1}|\mathbf{t}) = \mathcal{N}(t_{N+1}|m,s),$$
 (6.197)

$$m = \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{t},$$
  

$$s = c - \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{k}.$$
(6.198)

Let  $\mathbf{t}_1, \dots, \mathbf{t}_N$  be variables such that

$$p(\mathbf{t}_n|\mathbf{y}_n) = \mathcal{N}\left(\mathbf{t}_n|\mathbf{y}_n, \beta^{-1}\mathbf{I}\right),$$
  

$$p(\mathbf{Y}) = \mathcal{N}\left(\mathbf{Y}|\mathbf{O}, \mathbf{K}\right),$$
(6.199)

where

$$K_{nn'} = k(\mathbf{x}_n, \mathbf{x}_{n'}). \tag{6.200}$$

By the Bayes' theorem,

$$p(\mathbf{t}_{N+1}|\mathbf{T})p(\mathbf{T}) = p(\mathbf{T}'). \tag{6.201}$$

By marginalisation,

$$p(\mathbf{T}) = \int p(\mathbf{T}|\mathbf{Y})p(\mathbf{Y})d\mathbf{Y}.$$
 (6.202)

The logarithm of the integrannd of the right hand side except the terms independent of T and Y can be written as

$$-\frac{\beta}{2} \sum_{n=1}^{N} \|\mathbf{t}_n - \mathbf{y}_n\|^2 - \frac{1}{2} \operatorname{tr} \left( \mathbf{Y}^{\mathsf{T}} \mathbf{K}^{-1} \mathbf{Y} \right) = -\frac{1}{2} \operatorname{tr} \mathbf{M}, \tag{6.203}$$

where

$$\mathbf{M} = \begin{bmatrix} \mathbf{Y} \\ \mathbf{T} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \beta \mathbf{I} + \mathbf{K}^{-1} & -\beta \mathbf{I} \\ -\beta \mathbf{I} & \beta \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{Y} \\ \mathbf{T} \end{bmatrix}. \tag{6.204}$$

By 2.24,

$$\begin{bmatrix} \beta \mathbf{I} + \mathbf{K}^{-1} & -\beta \mathbf{I} \\ -\beta \mathbf{I} & \beta \mathbf{I} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{K} & \mathbf{K} \\ \mathbf{K} & \beta^{-1} \mathbf{I} + \mathbf{K} \end{bmatrix}.$$
 (6.205)

Then,

$$p(\mathbf{T}) = \mathcal{N}\left(\mathbf{T}|\mathbf{O},\mathbf{C}\right),\tag{6.206}$$

where

$$\mathbf{C} = \beta^{-1}\mathbf{I} + \mathbf{K}.\tag{6.207}$$

Then,

$$p(\mathbf{T}') = \mathcal{N}(\mathbf{T}'|\mathbf{O}, \mathbf{C}'), \qquad (6.208)$$

$$\mathbf{C}' = \begin{bmatrix} \mathbf{C} & \mathbf{k} \\ \mathbf{k}^{\mathsf{T}} & c \end{bmatrix}, \tag{6.209}$$

where

$$k_n = k(\mathbf{x}_n, \mathbf{x}_{N+1}),$$
  
 $c = \beta^{-1} + k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}).$  (6.210)

By 2.24,

$$\begin{bmatrix} c & \mathbf{k}^{\mathsf{T}} \\ \mathbf{k} & \mathbf{C} \end{bmatrix}^{-1} = \begin{bmatrix} s^{-1} & -s^{-1}\mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1} \\ -s^{-1}\mathbf{C}^{-1}\mathbf{k} & \mathbf{C}^{-1} + s^{-1}\mathbf{C}^{-1}\mathbf{k}\mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1} \end{bmatrix}, \tag{6.211}$$

where

$$s = c - \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{k}. \tag{6.212}$$

Therefore,

$$p(\mathbf{t}_{N+1}|\mathbf{T}) = \mathcal{N}\left(\mathbf{t}_{N+1}|\mathbf{m}, s\mathbf{I}\right), \tag{6.213}$$

where

$$\mathbf{m} = \mathbf{T}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{k}. \tag{6.214}$$

## 6.24

(a)

Let **M** be a  $D \times D$  diagonal matrix such that  $0 < M_{dd} < 1$ . For any vector **v** in D dimensions,

$$\mathbf{v}^{\mathsf{T}} \mathbf{M} \mathbf{v} = \sum_{d=1}^{D} \sum_{d'=1}^{D} v_d M_{dd'} v_{d'}.$$
 (6.215)

The right hand side can be written as

$$\sum_{d=1}^{D} M_{dd} v_d^2 > 0. (6.216)$$

Therefore, M is positive definite.

(b)

Let  $\mathbf{M}_1$  and  $\mathbf{M}_2$  be  $D \times D$  positive definite matrices. Then, for any vector  $\mathbf{v}$  in D dimensions,

$$\mathbf{v}^{\mathsf{T}}\mathbf{M}_{1}\mathbf{v} > 0,$$
  
$$\mathbf{v}^{\mathsf{T}}\mathbf{M}_{2}\mathbf{v} > 0.$$
 (6.217)

Then,

$$\mathbf{v}^{\mathsf{T}}(\mathbf{M}_1 + \mathbf{M}_2)\mathbf{v} > 0. \tag{6.218}$$

Therefore,  $\mathbf{M}_1 + \mathbf{M}_2$  is positive definite.

#### 6.25

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|a_n) = \sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1 - t_n},$$
  

$$p(\mathbf{a}) = \mathcal{N}(\mathbf{a}|\mathbf{0}, \mathbf{C}),$$
(6.219)

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)},$$

$$C_{nn'} = k(\mathbf{x}_n, \mathbf{x}_{n'}) + \nu I_{nn'}.$$
(6.220)

(a)

By the Bayes' theorem,

$$p(\mathbf{a}|\mathbf{t})p(\mathbf{t}) = p(\mathbf{t}|\mathbf{a})p(\mathbf{a}). \tag{6.221}$$

The logarithm of the right hand side except the terms independent of  ${\bf t}$  and  ${\bf a}$  can be written as

$$\Psi(\mathbf{a}) = \sum_{n=1}^{N} (t_n \ln \sigma(a_n) + (1 - t_n) \ln (1 - \sigma(a_n))) - \frac{1}{2} \mathbf{a}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{a}.$$
 (6.222)

In order to maximise  $p(\mathbf{a}|\mathbf{t})$  with respect to  $\mathbf{a}$ ,  $\Psi$  needs to be maximised. We have

$$\frac{d}{da}\sigma(a) = \sigma(a)\left(1 - \sigma(a)\right). \tag{6.223}$$

Then,

$$(\nabla \Psi(\mathbf{a}))_n = (t_n (1 - \sigma(a_n)) - (1 - t_n)\sigma(a_n)) - (\mathbf{C}^{-1}\mathbf{a})_n,$$
 (6.224)

so that

$$\nabla \Psi(\mathbf{a}) = \mathbf{t} - \boldsymbol{\sigma} - \mathbf{C}^{-1}\mathbf{a},\tag{6.225}$$

where

$$\sigma_n = \sigma(a_n). \tag{6.226}$$

Since  $\sigma$  non-linearly depends on  $\mathbf{a}$ , the Taylor series

$$\nabla \Psi(\mathbf{a}^{\text{new}}) = \nabla \Psi(\mathbf{a}) + \nabla \nabla \Psi(\mathbf{a}) \left(\mathbf{a}^{\text{new}} - \mathbf{a}\right) + O\left(\|\mathbf{a} - \mathbf{a}^{\text{new}}\|^{2}\right). \quad (6.227)$$

is used to iteratively find  $\mathbf{a}$  which maximises  $\Psi$ . Setting the left hand side to zero and neglecting the second order term of the right hand side gives

$$\mathbf{0} = \nabla \Psi(\mathbf{a}) + \nabla \nabla \Psi(\mathbf{a}) \left( \mathbf{a}^{\text{new}} - \mathbf{a} \right). \tag{6.228}$$

Then,

$$\mathbf{a}^{\text{new}} = \mathbf{a} - (\nabla \nabla \Psi(\mathbf{a}))^{-1} \nabla \Psi(\mathbf{a}). \tag{6.229}$$

We have

$$\nabla \nabla \Psi(\mathbf{a}) = -\mathbf{W} - \mathbf{C}^{-1},\tag{6.230}$$

where

$$W_{nn'} = \sigma(a_n) \left(1 - \sigma(a_n)\right) I_{nn'}. \tag{6.231}$$

Then, the right hand side can be written as

$$\mathbf{a}^{\text{new}} = \mathbf{a} + \left(\mathbf{W} + \mathbf{C}^{-1}\right)^{-1} \left(\mathbf{t} - \boldsymbol{\sigma} - \mathbf{C}^{-1}\mathbf{a}\right). \tag{6.232}$$

The right hand side can be written as

$$(\mathbf{W} + \mathbf{C}^{-1})^{-1} ((\mathbf{W} + \mathbf{C}^{-1}) \mathbf{a} + \mathbf{t} - \boldsymbol{\sigma} - \mathbf{C}^{-1} \mathbf{a})$$

$$= \mathbf{C} (\mathbf{W} \mathbf{C} + \mathbf{I})^{-1} (\mathbf{W} \mathbf{a} + \mathbf{t} - \boldsymbol{\sigma}).$$
(6.233)

Therefore,

$$\mathbf{a}^{\text{new}} = \mathbf{C} \left( \mathbf{W} \mathbf{C} + \mathbf{I} \right)^{-1} \left( \mathbf{W} \mathbf{a} + \mathbf{t} - \boldsymbol{\sigma} \right). \tag{6.234}$$

(b)

Let  $\mathbf{a}_{\text{MAP}}$  be a stationary point of  $p(\mathbf{a}|\mathbf{t})$  which is found by the iterative process given in (a). We have the Taylor series

$$\begin{split} \Psi(\mathbf{a}) = & \Psi(\mathbf{a}_{MAP}) + \nabla \Psi(\mathbf{a}_{MAP})(\mathbf{a} - \mathbf{a}_{MAP}) \\ & + \frac{1}{2}(\mathbf{a} - \mathbf{a}_{MAP})^{\intercal} \nabla \nabla \Psi(\mathbf{a}_{MAP})(\mathbf{a} - \mathbf{a}_{MAP}) + O\left(\|\mathbf{a} - \mathbf{a}_{MAP}\|^{3}\right). \end{split} \tag{6.235}$$

Then,

$$\Psi(\mathbf{a}) \simeq \Psi(\mathbf{a}_{MAP}) - \frac{1}{2}(\mathbf{a} - \mathbf{a}_{MAP})^{\mathsf{T}}\mathbf{H}(\mathbf{a} - \mathbf{a}_{MAP}),$$
 (6.236)

where

$$\mathbf{H} = \mathbf{W} \mid_{\mathbf{a} = \mathbf{a}_{\text{MAP}}} + \mathbf{C}^{-1}. \tag{6.237}$$

Therefore,

$$p(\mathbf{a}|\mathbf{t}) \simeq \mathcal{N}\left(\mathbf{a}|\mathbf{a}_{\text{MAP}}, \mathbf{H}^{-1}\right).$$
 (6.238)

(c)

By (a),

$$\mathbf{a}_{\text{MAP}} = \mathbf{C} \left( \mathbf{WC} + \mathbf{I} \right)^{-1} \left( \mathbf{W} \mathbf{a}_{\text{MAP}} + \mathbf{t} - \boldsymbol{\sigma} \right). \tag{6.239}$$

Then,

$$(\mathbf{WC} + \mathbf{I}) \mathbf{C}^{-1} \mathbf{a}_{MAP} = \mathbf{W} \mathbf{a}_{MAP} + \mathbf{t} - \boldsymbol{\sigma}. \tag{6.240}$$

Therefore,

$$\mathbf{a}_{\text{MAP}} = \mathbf{C}(\mathbf{t} - \boldsymbol{\sigma}). \tag{6.241}$$

## 6.26

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|a_n) = \sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1 - t_n},$$
  

$$p(\mathbf{a}) = \mathcal{N}(\mathbf{a}|\mathbf{0}, \mathbf{C}),$$
(6.242)

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)},$$

$$C_{nn'} = k(\mathbf{x}_n, \mathbf{x}_{n'}) + \nu I_{nn'}.$$
(6.243)

By marginalisation,

$$p(a_{N+1}|\mathbf{t}) = \int p(a_{N+1}|\mathbf{a})p(\mathbf{a}|\mathbf{t})d\mathbf{a}.$$
 (6.244)

Let

$$\mathbf{a}' = \begin{bmatrix} \mathbf{a} \\ a_{N+1} \end{bmatrix},$$

$$k_n = k(\mathbf{x}_n, \mathbf{x}_{N+1}),$$

$$c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \nu.$$
(6.245)

Then,

$$p(\mathbf{a}') = \mathcal{N}(\mathbf{a}'|\mathbf{0}, \mathbf{C}'),$$
 (6.246)

where

$$\mathbf{C}' = \begin{bmatrix} \mathbf{C} & \mathbf{k} \\ \mathbf{k}^{\mathsf{T}} & c \end{bmatrix}. \tag{6.247}$$

By 2.24,

$$\begin{bmatrix} c & \mathbf{k}^{\mathsf{T}} \\ \mathbf{k} & \mathbf{C} \end{bmatrix}^{-1} = \begin{bmatrix} s^{-1} & -s^{-1}\mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1} \\ -s^{-1}\mathbf{C}^{-1}\mathbf{k} & \mathbf{C}^{-1} + s^{-1}\mathbf{C}^{-1}\mathbf{k}\mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1} \end{bmatrix}, \tag{6.248}$$

where

$$s = c - \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{k}. \tag{6.249}$$

Then,

$$p(a_{N+1}|\mathbf{a}) = \mathcal{N}(a_{N+1}|m,s),$$
 (6.250)

where

$$m = \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{a}. \tag{6.251}$$

By 6.25(b),

$$p(\mathbf{a}|\mathbf{t}) \simeq \mathcal{N}\left(\mathbf{a}|\mathbf{a}_{\text{MAP}}, \mathbf{H}^{-1}\right),$$
 (6.252)

where  $\mathbf{a}_{\text{MAP}}$  is a stationary point of  $p(\mathbf{a}|\mathbf{t})$  and

$$\mathbf{H} = \mathbf{W} \mid_{\mathbf{a} = \mathbf{a}_{\text{MAP}}} + \mathbf{C}^{-1},$$

$$W_{nn'} = \sigma(a_n) \left(1 - \sigma(a_n)\right) I_{nn'}.$$
(6.253)

Then, the logarithm of the integrand except the terms independent of  ${\bf t}$  and  ${\bf a}$  can be approximated as

$$-\frac{1}{2s} (a_{N+1} - m)^{2} - \frac{1}{2} (\mathbf{a} - \mathbf{a}_{MAP})^{\mathsf{T}} \mathbf{H} (\mathbf{a} - \mathbf{a}_{MAP})$$

$$= -\frac{1}{2} \mathbf{a}'^{\mathsf{T}} \mathbf{M} \mathbf{a}' + \mathbf{a}'^{\mathsf{T}} \mathbf{v} - \frac{1}{2} \mathbf{a}_{MAP}^{\mathsf{T}} \mathbf{H} \mathbf{a}_{MAP}.$$
(6.254)

where

$$\mathbf{M} = \begin{bmatrix} \mathbf{H} + s^{-1}\mathbf{C}^{-1}\mathbf{k}\mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1} & -s^{-1}\mathbf{C}^{-1}\mathbf{k} \\ -s^{-1}\mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1} & s^{-1} \end{bmatrix},$$

$$\mathbf{v} = \begin{bmatrix} \mathbf{H}\mathbf{a}_{\text{MAP}} \\ 0 \end{bmatrix}.$$
(6.255)

By 2.24,

$$\mathbf{M}^{-1} = \begin{bmatrix} \mathbf{H}^{-1} & \mathbf{H}^{-1}\mathbf{C}^{-1}\mathbf{k} \\ \mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1}\mathbf{H}^{-1} & s + \mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1}\mathbf{H}^{-1}\mathbf{C}^{-1}\mathbf{k} \end{bmatrix}, \tag{6.256}$$

so that

$$\mathbf{M}^{-1}\mathbf{v} = \begin{bmatrix} \mathbf{a}_{\text{MAP}} \\ \mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1}\mathbf{a}_{\text{MAP}} \end{bmatrix}. \tag{6.257}$$

Then,

$$p(a_{N+1}|\mathbf{t}) \simeq \mathcal{N}\left(a_{N+1}|\mu,\sigma^2\right),$$
 (6.258)

$$\mu = \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{a}_{\text{MAP}},$$
  

$$\sigma^{2} = s + \mathbf{k}^{\mathsf{T}} \mathbf{C}^{-1} \mathbf{H}^{-1} \mathbf{C}^{-1} \mathbf{k}.$$
(6.259)

By 
$$6.25(c)$$
,

$$\mathbf{a}_{\text{MAP}} = \mathbf{C}(\mathbf{t} - \boldsymbol{\sigma}),\tag{6.260}$$

where

$$\sigma_n = \sigma(a_n). \tag{6.261}$$

Therefore,

$$\mu = \mathbf{k}^{\mathsf{T}}(\mathbf{t} - \boldsymbol{\sigma}),$$
  

$$\sigma^{2} = c - \mathbf{k}^{\mathsf{T}}\mathbf{C}^{-1}(\mathbf{C} - \mathbf{H}^{-1})\mathbf{C}^{-1}\mathbf{k}.$$
(6.262)

## 6.27 (Incomplete)

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|a_n) = \sigma(a_n)^{t_n} (1 - \sigma(a_n))^{1 - t_n},$$
  

$$p(\mathbf{a}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{a}|\mathbf{0}, \mathbf{C}),$$
(6.263)

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)} \tag{6.264}$$

and  ${\bf C}$  is dependent on  ${\boldsymbol \theta}$  with M dimensions.

(a)

By marginalisation,

$$p(\mathbf{t}|\boldsymbol{\theta}) = \int p(\mathbf{t}|\mathbf{a})p(\mathbf{a}|\boldsymbol{\theta})d\mathbf{a}.$$
 (6.265)

Let  $\mathbf{a}_{\text{MAP}}$  be a stationary point of  $p(\mathbf{a}|\mathbf{t})$  and

$$\Psi(\mathbf{a}) = \ln p(\mathbf{t}|\mathbf{a}) + \ln p(\mathbf{a}|\boldsymbol{\theta}). \tag{6.266}$$

By 6.25(b),

$$\Psi(\mathbf{a}) \simeq \Psi(\mathbf{a}_{\text{MAP}}) - \frac{1}{2} (\mathbf{a} - \mathbf{a}_{\text{MAP}})^{\mathsf{T}} \mathbf{H} (\mathbf{a} - \mathbf{a}_{\text{MAP}}),$$
 (6.267)

where

$$\mathbf{H} = \mathbf{W} \mid_{\mathbf{a} = \mathbf{a}_{\text{MAP}}} + \mathbf{C}^{-1},$$

$$W_{nn'} = \sigma(a_n) (1 - \sigma(a_n)) I_{nn'}.$$
(6.268)

Therefore,

$$\ln p(\mathbf{t}|\boldsymbol{\theta}) \simeq \Psi(\mathbf{a}_{\text{MAP}}) + \frac{N}{2}\ln(2\pi) - \frac{1}{2}\ln(\det \mathbf{H}). \tag{6.269}$$

(b)

We have

$$\frac{\partial \Psi(\mathbf{a}_{\text{MAP}})}{\partial \theta_m} = \sum_{n=1}^{N} \sum_{n'=1}^{N} \frac{\partial \Psi(\mathbf{a}_{\text{MAP}})}{\partial C_{nn'}} \frac{\partial C_{nn'}}{\partial \theta_m}.$$
 (6.270)

The right hand side can be written as

$$\frac{1}{2} \operatorname{tr} \left( \mathbf{a}_{\text{MAP}} \mathbf{a}_{\text{MAP}}^{\mathsf{T}} \left( \mathbf{C}^{-1} \right)^{2} \frac{\partial \mathbf{C}}{\partial \theta_{m}} \right). \tag{6.271}$$

We have

$$\frac{\partial \ln(\det \mathbf{H})}{\partial \theta_m} = \sum_{n=1}^{N} \sum_{n'=1}^{N} \frac{\partial \ln(\det \mathbf{H})}{\partial H_{nn'}} \frac{\partial H_{nn'}}{\partial \theta_m}.$$
 (6.272)

By 3.21(a), the right hand side can be written as

$$\operatorname{tr}\left(\mathbf{H}^{-1}\frac{\partial\mathbf{C}^{-1}}{\partial\theta_{m}}\right).\tag{6.273}$$

Therefore,

$$\frac{\partial \ln p(\mathbf{t}|\boldsymbol{\theta})}{\partial \theta_m} \simeq \frac{1}{2} \operatorname{tr} \left( \mathbf{a}_{\text{MAP}} \mathbf{a}_{\text{MAP}}^{\mathsf{T}} \left( \mathbf{C}^{-1} \right)^2 \frac{\partial \mathbf{C}}{\partial \theta_m} - \mathbf{H}^{-1} \frac{\partial \mathbf{C}^{-1}}{\partial \theta_m} \right). \tag{6.274}$$

# 7 Sparse Kernel Machines

# 7.1 (Incomplete)

# 7.2

Let  $t_1, \dots, t_N$  be variables. In order to minimise  $\|\mathbf{w}\|^2$  under the constraint

$$t_n\left(\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}(\mathbf{x}_n) + b\right) \ge \gamma,$$
 (7.1)

for  $\gamma > 0$ , let

$$L = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^{N} a_n \left( t_n \left( \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n) + b \right) - \gamma \right). \tag{7.2}$$

Setting the derivatives of L with respect to  $\mathbf{w}$  and b to zero gives

$$\mathbf{0} = \mathbf{w} - \sum_{n=1}^{N} a_n t_n \phi(\mathbf{x}_n),$$

$$0 = -\sum_{n=1}^{N} a_n t_n.$$
(7.3)

Therefore,

$$\mathbf{w} = \sum_{n=1}^{N} a_n t_n \phi(\mathbf{x}_n),$$

$$0 = \sum_{n=1}^{N} a_n t_n,$$
(7.4)

which is independent of  $\gamma$ .

- 7.3 (Incomplete)
- 7.4 (Incomplete)
- 7.5 (Incomplete)

## 7.6

Let  $t_1, \dots, t_N$  be a variable such that

$$t_n \in \{-1, 1\},\ p(t_n|y_n) = \sigma(y_n)^{\frac{1+t_n}{2}} \left(1 - \sigma(y_n)\right)^{\frac{1-t_n}{2}},$$
(7.5)

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)},$$

$$y_n = \mathbf{w}^{\mathsf{T}} \phi(\mathbf{x}_n) + b.$$
(7.6)

We have

$$1 - \sigma(a) = \sigma(-a). \tag{7.7}$$

Then,

$$p(t_n|y_n) = \sigma(y_n t_n). \tag{7.8}$$

Then,

$$p(\mathbf{t}|\mathbf{y}) = \prod_{n=1}^{N} \sigma(y_n t_n), \tag{7.9}$$

so that

$$-\ln p(\mathbf{t}|\mathbf{y}) = -\sum_{n=1}^{N} \ln \sigma(y_n t_n). \tag{7.10}$$

Therefore,

$$-\ln p(\mathbf{t}|\mathbf{y}) = \sum_{n=1}^{N} \ln \left(1 + \exp(-y_n t_n)\right). \tag{7.11}$$

## 7.7

Let  $t_1, \dots, t_N$  be vaariables. Let

$$E = C \sum_{n=1}^{N} E_{\epsilon}(y_n - t_n) + \frac{1}{2} \|\mathbf{w}\|^2,$$
 (7.12)

where  $\epsilon > 0$  and

$$E_{\epsilon}(a) = \begin{cases} 0, & |a| < \epsilon, \\ |a| - \epsilon, & \text{otherwise,} \end{cases}$$

$$y_n = \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n) + b.$$
(7.13)

In order to minimise E under the constraints

$$t_n \le y_n + \epsilon + \xi_n,$$
  

$$t_n > y_n - \epsilon - \hat{\xi}_n,$$
(7.14)

where  $\xi_n \geq 0$  and  $\hat{\xi}_n \geq 0$ , let us minimise

$$L = C \sum_{n=1}^{N} (\xi_n + \hat{\xi}_n) + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^{N} \left( \mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n \right) - \sum_{n=1}^{N} a_n (\epsilon + \xi_n + y_n - t_n) - \sum_{n=1}^{N} \hat{a}_n (\epsilon + \hat{\xi}_n - y_n + t_n).$$
(7.15)

Setting the derivatives of L with respect to  $\mathbf{w}$ , b,  $\xi_n$  and  $\hat{\xi}_n$  to zero gives

$$\mathbf{0} = \mathbf{w} - \sum_{n=1}^{N} a_n \phi(\mathbf{x}_n) + \sum_{n=1}^{N} \hat{a}_n \phi(\mathbf{x}_n),$$

$$0 = -\sum_{n=1}^{N} a_n + \sum_{n=1}^{N} \hat{a}_n,$$

$$0 = C - \mu_n - a_n,$$

$$0 = C - \hat{\mu}_n - \hat{a}_n.$$
(7.16)

Then, L can be written as

$$\sum_{n=1}^{N} \left( (\mu_n + a_n) \xi_n + (\hat{\mu}_n + \hat{a}_n) \hat{\xi}_n \right) + \frac{1}{2} \left\| \sum_{n=1}^{N} (\hat{a}_n - a_n) \phi(\mathbf{x}_n) \right\|^2 - \sum_{n=1}^{N} \left( \mu_n \xi_n + \hat{\mu}_n \hat{\xi}_n \right) - \sum_{n=1}^{N} a_n (\epsilon + \xi_n + y_n - t_n)$$

$$- \sum_{n=1}^{N} \hat{a}_n (\epsilon + \hat{\xi}_n - y_n + t_n).$$
(7.17)

Therefore, minimising L is equilvalent to maximising

$$\tilde{L} = -\frac{1}{2} \sum_{n=1}^{N} \sum_{n'=1}^{N} (\hat{a}_n - a_n)(\hat{a}_{n'} - a_{n'}) k(\mathbf{x}_n, \mathbf{x}_{n'})$$

$$+ \epsilon \sum_{n=1}^{N} (a_n + \hat{a}_n) + \sum_{n=1}^{N} (a_n - \hat{a}_n) (y_n - t_n),$$
(7.18)

where

$$k(\mathbf{x}_n, \mathbf{x}_{n'}) = \phi(\mathbf{x}_n)^{\mathsf{T}} \phi(\mathbf{x}_{n'}). \tag{7.19}$$

# 7.8 (Incomplete)

Let  $t_1, \dots, t_N$  be vaariables. Let

$$E = C \sum_{n=1}^{N} E_{\epsilon}(y_n - t_n) + \frac{1}{2} \|\mathbf{w}\|^2,$$
 (7.20)

where  $\epsilon > 0$  and

$$E_{\epsilon}(a) = \begin{cases} 0, & |a| < \epsilon, \\ |a| - \epsilon, & \text{otherwise,} \end{cases}$$

$$y_n = \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_n) + b.$$
(7.21)

(a)

In order to minimise E under the constraints

$$t_n \le y_n + \epsilon + \xi_n, \tag{7.22}$$

where  $\xi_n \geq 0$ , let us minimise

$$L = C \sum_{n=1}^{N} \xi_n + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^{N} \mu_n \xi_n - \sum_{n=1}^{N} a_n (\epsilon + \xi_n + y_n - t_n).$$
 (7.23)

Setting the derivatives of L with respect to  $\mathbf{w}$ , b and  $\xi_n$  to zero gives

$$\mathbf{0} = \mathbf{w} - \sum_{n=1}^{N} a_n \boldsymbol{\phi}(\mathbf{x}_n),$$

$$0 = -\sum_{n=1}^{N} a_n,$$

$$0 = C - \mu_n - a_n.$$

$$(7.24)$$

Then, L can be written as

$$\sum_{n=1}^{N} (\mu_n + a_n) \xi_n + \frac{1}{2} \left\| \sum_{n=1}^{N} a_n \phi(\mathbf{x}_n) \right\|^2 - \sum_{n=1}^{N} \mu_n \xi_n - \sum_{n=1}^{N} a_n (\epsilon + \xi_n + y_n - t_n).$$
(7.25)

Therefore, minimising L is equilvalent to maximising

$$\tilde{L} = -\frac{1}{2} \sum_{n=1}^{N} \sum_{n'=1}^{N} a_n a_{n'} k(\mathbf{x}_n, \mathbf{x}_{n'}) + \sum_{n=1}^{N} a_n (y_n - t_n),$$
 (7.26)

where

$$k(\mathbf{x}_n, \mathbf{x}_{n'}) = \boldsymbol{\phi}(\mathbf{x}_n)^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}_{n'}). \tag{7.27}$$

#### 7.9

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}_n, \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w}|\mathbf{0}, \mathbf{A}^{-1}\right).$$
 (7.28)

By the Bayes' theorem,

$$p(\mathbf{w}|\mathbf{t})p(\mathbf{t}) = p(\mathbf{t}|\mathbf{w})p(\mathbf{w}). \tag{7.29}$$

The logarithm of the right hand side except the terms independent of  ${\bf t}$  and  ${\bf w}$  can be written as

$$-\frac{\beta}{2} \sum_{n=1}^{N} \|t_n - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n\|^2 - \frac{1}{2} \mathbf{w}^{\mathsf{T}} \mathbf{A} \mathbf{w}$$

$$= -\frac{1}{2} \begin{bmatrix} \mathbf{w} \\ \mathbf{t} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} + \mathbf{A} & -\beta \mathbf{\Phi}^{\mathsf{T}} \\ -\beta \mathbf{\Phi} & \beta \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{t} \end{bmatrix}.$$
(7.30)

Therefore,

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S}), \qquad (7.31)$$

$$\mathbf{m} = \beta \mathbf{S} \mathbf{\Phi}^{\mathsf{T}} \mathbf{t}, \mathbf{S} = (\beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} + \mathbf{A})^{-1}.$$
 (7.32)

#### 7.10

Let  $t_1, \dots, t_N$  be variables such that

$$p(t_n|\mathbf{w}) = \mathcal{N}\left(t_n|\mathbf{w}^{\mathsf{T}}\boldsymbol{\phi}_n, \beta^{-1}\right),$$
  
$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w}|\mathbf{0}, \mathbf{A}^{-1}\right).$$
 (7.33)

By marginalisation,

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{w})p(\mathbf{w})d\mathbf{w}.$$
 (7.34)

The logarithm of the integrand of the right hand side except the terms independent of  $\mathbf{t}$  and  $\mathbf{w}$  can be written as

$$-\frac{\beta}{2} \sum_{n=1}^{N} \|t_n - \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}_n\|^2 - \frac{1}{2} \mathbf{w}^{\mathsf{T}} \mathbf{A} \mathbf{w}$$

$$= -\frac{1}{2} \begin{bmatrix} \mathbf{w} \\ \mathbf{t} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} \beta \boldsymbol{\Phi}^{\mathsf{T}} \boldsymbol{\Phi} + \mathbf{A} & -\beta \boldsymbol{\Phi}^{\mathsf{T}} \\ -\beta \boldsymbol{\Phi} & \beta \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{t} \end{bmatrix}.$$
(7.35)

By 2.24,

$$\begin{bmatrix} \beta \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi} + \mathbf{A} & -\beta \mathbf{\Phi}^{\mathsf{T}} \\ -\beta \mathbf{\Phi} & \beta \mathbf{I} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{A}^{-1} \mathbf{\Phi}^{\mathsf{T}} \\ \mathbf{\Phi} \mathbf{A}^{-1} & \beta^{-1} \mathbf{I} + \mathbf{\Phi} \mathbf{A}^{-1} \mathbf{\Phi}^{\mathsf{T}} \end{bmatrix}. \tag{7.36}$$

Therefore,

$$p(\mathbf{t}) = \mathcal{N}\left(\mathbf{t}|\mathbf{0}, \beta^{-1}\mathbf{I} + \mathbf{\Phi}\mathbf{A}^{-1}\mathbf{\Phi}^{\mathsf{T}}\right). \tag{7.37}$$