

More PRML Errata

Yousuke Takada
yousuketakada@gmail.com

February 24, 2018

Preface

This report is an unofficial list of errata for *Pattern Recognition and Machine Learning* or PRML by Bishop (2006). In this report, I have compiled only errata that are not (yet) listed in the official errata document (Svensén and Bishop, 2011) at the time of this writing. Currently, there are three versions of the official errata document, corresponding to the first (2006), second (2007), and third (2009) printings of PRML; consult the [support page](#) for how to identify which printing your copy of PRML is from.¹

I have tried to follow the terminology and the notation used in PRML and the official errata as closely as possible. In particular, when specifying the location of an error, I follow the notational conventions (such as “Paragraph 2, Line –1”) adopted by Svensén and Bishop (2011). As the official errata document “is intended to be complete,” this report also tries to correct even trivial typographical errors as well.

PRML is arguably such a great textbook in the field of machine learning that it is extremely helpful and easier to understand than any other similar account. That said, there are a few subtleties that some readers might have hard time to appreciate. In hopes to help such readers get out of struggle or obtain a better grasp on some important concepts, I have also included in this report some comments and suggestions for improving the readability to which I would have liked to refer when I first read PRML.

It should be noted that the readers of the Japanese edition of PRML will find its [support page](#) (in Japanese) useful. Along with other information such as the contents, it lists errata specific to the Japanese edition as well as some additional errata for the English edition, which have also been included in this report for the reader’s convenience.

I welcome all comments and suggestions regarding this report; please send me any such feedback via email or, preferably, by creating an “issue” or a “pull request” at the following GitHub repository

https://github.com/yousuketakada/prml_errata

where you can find the source code of this report as well as other supporting material.

¹The last line but one of the bibliographic information page (the page immediately preceding the dedication page) of my copy of PRML reads “9 8 7 (corrected at 6th printing 2007).” Note that, although it says it is from the “6th printing,” it is actually from the *second* printing according to the official errata document (Svensén and Bishop, 2011) so that I refer to Version 2 of the official errata.

Acknowledgements

I would like to thank those who have informed me of yet more errata and clarifications incorporated in this report. In particular, I am grateful to Christopher Sahnwaldt and Mark-Jan Nederhof for their helpful comments and discussions.

Corrections and Comments

Page xi

Paragraph –2, Line 1: $|f(x)/g(x)|$ should read $|g(x)/f(x)|$ (with the functions swapped). Moreover, the limit we take is *not* necessarily the one specified in the text, i.e., $x \rightarrow \infty$, but is often implied by the context as follows.

Big O notation The big O notation $g(x) = O(f(x))$ generally denotes that $|g(x)/f(x)|$ is bounded as $x \rightarrow c$ where, if c is not given explicitly, $c = 0$ for a Taylor series such as (2.299) or (D.1); or $c = \infty$ for an asymptotic series such as (10.241) or for computational complexity (see, e.g., Section 5.2.3), for example. See [Olver et al. \(2017\)](#) for other asymptotic and order notations.

Page 5

Equation (1.1): The lower ellipsis (\dots) should be centered (\cdots).² Specifically, (1.1) should read

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \cdots + w_Mx^M = \sum_{j=0}^M w_jx^j. \quad (1)$$

Page 10

The text after (1.4): The lower ellipsis (\dots) should be centered (\cdots).

Page 18

Equation (1.27): It should be noted that the transformation $g : \mathcal{Y} \rightarrow \mathcal{X}$ must be *bijective* or, equivalently, *invertible* in order for the change of variables (1.27) to be meaningful where \mathcal{X} and \mathcal{Y} are the domains of the distributions $p_x(\cdot)$ and $p_y(\cdot)$, respectively.³ This can be easily understood by noting that, if, for any *measurable*⁴ subset $\mathcal{X}_0 \subset \mathcal{X}$ of \mathcal{X} , the

²The \LaTeX command `\cdots` or, with the `amsmath` or `mathtools` package, `\dots` (in most cases) will do.

³A function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is said to be *bijective* (or *one-to-one correspondence*) if, for any $y \in \mathcal{Y}$, there exists a unique $x \in \mathcal{X}$ such that $y = f(x)$. Note also that bijectivity is equivalent to invertibility: A function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is bijective if and only if f is *invertible*, i.e., there exists an *inverse function* $f^{-1} : \mathcal{Y} \rightarrow \mathcal{X}$ (which is, of course, also bijective) such that $f^{-1} \circ f$ is the identity function on \mathcal{X} and $f \circ f^{-1}$ is the identity function on \mathcal{Y} (an identity function is a function that maps every input to itself).

⁴A *measurable* set is such that we can consider its “size” (or *measure*) in some sense so that the integration over it is meaningful; this is a concept formally defined in a branch of mathematics called *measure theory* (see, e.g., [Tao \(2011\)](#)), which however “lies outside the scope of [PRML]” (Page 19, Paragraph 3, Line –5). The reason why we restrict ourselves to measurable subsets here is, of course, that we indeed have “pathological” ones that are not measurable. However, since it is safe to say that all the sets we meet in practice are measurable (for example, measurable subsets of \mathbb{R} include all the open sets $(a_1, b_1), (a_2, b_2), \dots$ and their countable

*preimage*⁵ $\mathcal{Y}_0 = g^{-1}(\mathcal{X}_0)$ of \mathcal{X}_0 under g is again measurable and we can make the change of variable $x = g(y)$ as

$$\int_{\mathcal{X}_0} p_x(x) dx = \int_{\mathcal{Y}_0} p_x(x) \left| \frac{dx}{dy} \right| dy \quad (3)$$

$$= \int_{\mathcal{Y}_0} p_x(g(y)) |g'(y)| dy \quad (4)$$

then we can identify the integrand of the right hand side as $p_y(y)$.

Multivariate change of variables Similarly, the multivariate version of the change of variables formula is given by

$$p(\mathbf{y}) = p(\mathbf{x}) \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right) \right| \quad (5)$$

where the transformation between \mathbf{x} and \mathbf{y} is again assumed to be bijective.

Page 42

Paragraph 1, Line 1: “the class j ” should read “class \mathcal{C}_j .”

Page 44

Paragraph 2, Line –3: Insert a space before the sentence starting “There has been. . .” (*the third printing only*).⁶

Page 46

Equation (1.85): A period (.) should terminate (1.85).

Page 47

Paragraph 1, Line –1: $\mathbb{E}_t[t|\mathbf{x}]$ should read $\mathbb{E}_t[t|\mathbf{x}]$ (the subscript should be the vector \mathbf{t}).

Page 47

Equation (1.90): The quantity $\text{var}[t|\mathbf{x}]$ is the *conditional variance*, which is defined similarly to the *conditional expectation* (1.37) so that

$$\text{var}[t|\mathbf{x}] = \mathbb{E}[(t - \mathbb{E}[t|\mathbf{x}])^2|\mathbf{x}] \quad (6)$$

where we have omitted the subscript t in what should be $\mathbb{E}_t[\cdot|\mathbf{x}]$

unions $(a_1, b_1) \cup (a_2, b_2) \cup \dots$), we omit the “measurable” qualifier for brevity in the rest of this report.

⁵Let $f: \mathcal{X} \rightarrow \mathcal{Y}$ be some function from a set \mathcal{X} (the *domain*) to another set \mathcal{Y} (the *codomain*). The *preimage* (or *inverse image*) of a subset $\mathcal{Y}_0 \subset \mathcal{Y}$ of the codomain \mathcal{Y} under f is defined by

$$f^{-1}(\mathcal{Y}_0) \equiv \{x \in \mathcal{X} \mid f(x) \in \mathcal{Y}_0\} \quad (2)$$

so that $f^{-1}(\mathcal{Y}_0) \subset \mathcal{X}$ is a subset of the domain \mathcal{X} .

⁶Something strange must have happened in the third (2009) printing, leading to some spacing issues where a sentence ends with a reference number such as “Figure 1.27.” We can also find other “regression” errors in the third printing. Such errors are marked “*the third printing only*” in this report.

Page 51

Equation (1.98): Following the notation (1.93) for the entropy, we should write the left hand side of (1.98) as $H[X]$ instead of $H[p]$ so that

$$H[X] = - \sum_i p(x_i) \ln p(x_i). \quad (7)$$

As suggested in Appendix D, if we regard the (differential) entropy $H[\cdot]$ as a *functional*, then we see that “the entropy could equally well have been written as $H[p]$ ” (Page 703, Paragraph 1, Lines –2 and –1). However, it is probably better to maintain the notational consistency here.

Pages 55 and 56

The text around (1.114): There are some inaccuracies in the definitions and the properties of *convex* and *strictly convex* functions. First, a convex function is not necessarily differentiable (consider, e.g., the absolute value function $f(x) = |x|$, which is convex but not differentiable at $x = 0$). Second, even for twice differentiable functions, strict positivity of the second derivative is not necessary for convexity nor for strict convexity. Third, the condition for strict convexity that “the equality [in (1.114)] is satisfied only for $\lambda = 0$ and $\lambda = 1$ ” (Page 56, Paragraph 1, Line –4) is meaningless because the equality holds for any λ when $a = b$.

In the following, instead of correcting these errors one by one, I would like to present slightly more general definitions for convex and strictly convex functions where we let the parameter λ vary only on the open set $(0, 1)$, rather than on the closed set $[0, 1]$ as in PRML, in order to avoid edge cases. I also give some well-known properties regarding convex and strictly convex functions that are twice differentiable (which I think are intended to be addressed in PRML).

Convexity and strict convexity Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a real-valued, continuous⁷ function defined on some *convex set* \mathcal{X} such that, if $x_0, x_1 \in \mathcal{X}$, then $(1 - \lambda)x_0 + \lambda x_1 \in \mathcal{X}$ for any $\lambda \in (0, 1)$. The function $f(x)$ is said to be *convex* if

$$f((1 - \lambda)x_0 + \lambda x_1) \leq (1 - \lambda)f(x_0) + \lambda f(x_1) \quad (8)$$

for any $x_0, x_1 \in \mathcal{X}$ and $\lambda \in (0, 1)$. If the equality in (8) holds only when $x_0 = x_1$, then $f(x)$ is said to be *strictly convex*.

If $f(x)$ is twice differentiable, the following properties hold:

- (i) The function $f(x)$ is convex if and only if the second derivative $f''(x)$ is nonnegative for all $x \in \mathcal{X}$.

⁷Strictly speaking, we do not need to assume continuity here because convexity implies continuity. More specifically, a convex function $f : \mathcal{X} \rightarrow \mathbb{R}$ is continuous on the entire domain \mathcal{X} (except the boundary $\partial\mathcal{X}$). To see this, consider three points A, B, C on the convex curve $y = f(x)$ whose x values are $x_0, x_1, x_2 \in \mathcal{X}$, respectively, where we assume that $x_0 < x_1 < x_2$. First, we see from convexity that B lies below the line AC . Again from convexity, we have (i) that, on the interval $[x_0, x_1]$, the curve must be below AB and above BC ; and (ii) that, on the interval $[x_1, x_2]$, the curve must be below BC and above AB . These requirements (i) and (ii) imply continuity of $f(x)$ at $x = x_1$. Since x_1 can be any element in \mathcal{X} (except $\partial\mathcal{X}$), we see that $f(x)$ is continuous (on $\mathcal{X} \setminus \partial\mathcal{X}$).

- (ii) If $f''(x)$ is strictly positive for all $x \in \mathcal{X}$, then $f(x)$ is strictly convex. Note however that the converse of this does not hold (consider, e.g., $f(x) = x^4$, which is strictly convex but $f''(0) = 0$).

It is easy to see that convexity implies $f''(x) \geq 0$. In fact, from Taylor's theorem,⁸ we can write $f''(x)$ in the form

$$f''(x) = \lim_{h \rightarrow 0} \frac{f(x-h) - 2f(x) + f(x+h)}{h^2} \quad (12)$$

where we see that the right hand side is nonnegative from the inequality condition (8) in which we let $\lambda = 1/2$ and $x_i = x + (2i-1)h$ where $i = 0, 1$. To show the converse, we again make use of Taylor's theorem and expand $f(x)$ around $x_\lambda = (1-\lambda)x_0 + \lambda x_1$ so that we have

$$f(x_i) = f(x_\lambda) + (x_i - x_\lambda)f'(x_\lambda) + \frac{(x_i - x_\lambda)^2}{2}f''(\xi_i) \quad (13)$$

for some ξ_i between x_λ and x_i . With this expansion, we can write the right hand side of (8) in the form

$$(1-\lambda)f(x_0) + \lambda f(x_1) = f(x_\lambda) + \lambda(1-\lambda)\frac{(x_1 - x_0)^2}{2} \{\lambda f''(\xi_0) + (1-\lambda)f''(\xi_1)\} \quad (14)$$

from which we see that $f''(x) \geq 0$ implies convexity and also that $f''(x) > 0$ implies strict convexity.

Page 56

Equation (1.116): In general, we cannot interpret λ_i in *Jensen's inequality* (1.115) as the probability distribution over a discrete random variable x such that $\lambda_i \equiv p(x = x_i)$ because, since (1.115) holds for any $\{x_i\}$, we can take, say, $x_i = x_j$ and $\lambda_i \neq \lambda_j$ where $i \neq j$, assigning different probabilities to the same value of x . Actually, (1.116) is a special case of (1.115). An equivalent of (1.115) in terms of random variables can be derived as follows.

⁸*Taylor's theorem* (Abramowitz and Stegun, 1964) states that an n times differentiable function $f(x)$ can be expanded around a given point x_0 in the form

$$f(x) = f(x_0) + (x - x_0)f'(x_0) + \frac{(x - x_0)^2}{2}f''(x_0) + \dots + \frac{(x - x_0)^{n-1}}{(n-1)!}f^{(n-1)}(x_0) + R_n(x) \quad (9)$$

where the remainder term $R_n(x)$ (known as the *Lagrange remainder*) satisfies

$$R_n(x) = \frac{(x - x_0)^n}{n!}f^{(n)}(\xi) \quad (10)$$

for some ξ between x_0 and x . If we expand $f(x)$ around x with displacement $\pm h$, then the Taylor series expansion (up to the third order term) is given by

$$f(x \pm h) = f(x) \pm hf'(x) + \frac{h^2}{2}f''(x) \pm \frac{h^3}{3!}f^{(3)}(x) + O(h^4). \quad (11)$$

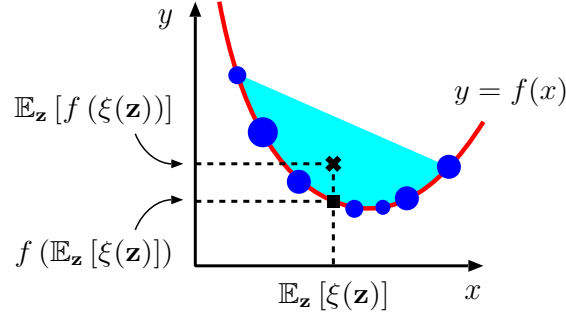


Figure 1 A physical “proof” of Jensen’s inequality (MacKay, 2003). Let us suppose that we have a set of point masses $m_i = p(\mathbf{z} = \mathbf{z}_i)$, denoted by filled blue circles (●) with areas proportional to m_i , and place them at the corresponding locations $(x, y) = (\xi(\mathbf{z}_i), f(\xi(\mathbf{z}_i)))$ on a convex curve $y = f(x)$. The center of gravity of those masses, which is $(\mathbb{E}_{\mathbf{z}}[\xi(\mathbf{z})], \mathbb{E}_{\mathbf{z}}[f(\xi(\mathbf{z}))])$, denoted by a cross sign (×), must lie above the convex curve and thus right above the point $(\mathbb{E}_{\mathbf{z}}[\xi(\mathbf{z})], f(\mathbb{E}_{\mathbf{z}}[\xi(\mathbf{z})]))$ on the curve, denoted by a filled square (■), showing Jensen’s inequality (15). One can also see that, if $f(\cdot)$ is strictly convex, the equality in (15) implies that $\xi(\mathbf{z})$ is almost surely constant (it is trivial to show that the converse is true).

Jensen’s inequality in terms of random variables In order to interpret (1.115) probabilistically, we instead introduce another set of underlying random variables \mathbf{z} such that $\lambda_i \equiv p(\mathbf{z} = \mathbf{z}_i)$ and a function $\xi(\cdot)$ such that $x_i \equiv \xi(\mathbf{z}_i)$, giving a result slightly more general than (1.116)

$$f(\mathbb{E}_{\mathbf{z}}[\xi(\mathbf{z})]) \leq \mathbb{E}_{\mathbf{z}}[f(\xi(\mathbf{z}))] \quad (15)$$

where $f(\cdot)$ is a convex function but $\xi(\cdot)$ can be any. Moreover, if $f(\cdot)$ is strictly convex, the equality in (15) holds if and only if $\xi(\mathbf{z})$ is constant with probability one or *almost surely*,⁹ meaning that there exists some constant ξ_0 such that $\xi(\mathbf{z}) = \xi_0$ on the range of \mathbf{z} *almost everywhere*, in which case we have $\mathbb{E}_{\mathbf{z}}[\xi(\mathbf{z})] = \xi_0$ and the both sides of (15) equal $f(\xi_0)$. See Figure 1 for an intuitive, physical “proof” of the inequality (15).

Since the random variables \mathbf{z} as well as their probability $p(\mathbf{z})$ can be chosen arbitrarily, it makes sense to write \mathbf{z} implicit in (15), giving a simpler form of Jensen’s inequality

$$f(\mathbb{E}[\xi]) \leq \mathbb{E}[f(\xi)]. \quad (16)$$

For continuous random variables, we have

$$f\left(\int \xi(\mathbf{x})p(\mathbf{x})d\mathbf{x}\right) \leq \int f(\xi(\mathbf{x}))p(\mathbf{x})d\mathbf{x} \quad (17)$$

where we have used \mathbf{x} to denote the underlying random variables for which we take the expectations. By making use of (17), one can show that the Kullback-Leibler divergence $\text{KL}(p\|q)$ given by (1.113) satisfies *Gibbs’s inequality*

$$\text{KL}(p\|q) \geq 0 \quad (18)$$

with equality if and only if $p(\mathbf{x}) = q(\mathbf{x})$ almost everywhere. See the following erratum for more details.

⁹Here, the proviso *almost surely* (often abbreviated as a.s.) or *almost everywhere* (a.e.) means that there may be some exceptions but they can occur only with probability zero (or *measure zero*) so that we can safely ignore them; this is a concept formally defined in a branch of mathematics called *measure theory* (see, e.g., Tao (2011)). As in PRML, we omit such “almost” provisos for brevity in the rest of this report.

Equation (1.118): There are some difficulties in the derivation (1.118) of Gibbs's inequality (18). First, the quantity $\xi(\mathbf{x}) = q(\mathbf{x})/p(\mathbf{x})$ is undefined for \mathbf{x} such that $p(\mathbf{x}) = 0$. Second, the convex function $f(\xi) = -\ln \xi$ is undefined for $\xi = 0$, which occurs where $q(\mathbf{x}) = 0$ and $p(\mathbf{x}) > 0$. In order to avoid these difficulties, we shall take a different approach (MacKay, 2003; Kullback and Leibler, 1951) in which we make use of Jensen's inequality (17) with respect to $q(\mathbf{x})$ where we identify $f(\xi) = \xi \ln \xi$ and $\xi(\mathbf{x}) = p(\mathbf{x})/q(\mathbf{x})$. Note that we can safely proceed with this approach because we can assume $q(\mathbf{x}) > 0$ without loss of generality.

In the following, we first show Gibbs's inequality (18) along this line, after which we also see that the Kullback-Leibler divergence $\text{KL}(p\|q)$ is convex in the sense that it satisfies the inequality (8) with respect to the pair of the distributions (p, q) . Finally, we give an alternative proof of Gibbs's inequality (18) in terms of a generalized version of the Kullback-Leibler divergence such that it is extended for unnormalized distributions.

A proof of Gibbs's inequality in terms of Jensen's inequality Let us first examine the behavior of the integrand of the Kullback-Leibler divergence $\text{KL}(p\|q)$

$$p(\mathbf{x}) \ln p(\mathbf{x}) - p(\mathbf{x}) \ln q(\mathbf{x}) \quad (19)$$

where $q(\mathbf{x})$ or $p(\mathbf{x})$ vanishes. We notice that, if $q(\mathbf{x}) \rightarrow 0$ for \mathbf{x} such that $p(\mathbf{x}) > 0$, the integrand (19) diverges so that $\text{KL}(p\|q) \rightarrow \infty$. On the other hand, the integrand (19) always vanishes for \mathbf{x} such that $p(\mathbf{x}) = 0$ regardless of the values of $q(\mathbf{x})$.¹⁰ Therefore, in order for $\text{KL}(p\|q)$ to be well-defined, we must have $p(\mathbf{x}) = 0$ for all \mathbf{x} such that $q(\mathbf{x}) = 0$ or, stated differently, the *support* of $p(\mathbf{x})$ must be contained in that of $q(\mathbf{x})$, i.e.,¹¹

$$\text{supp}(p) \subset \text{supp}(q) \quad (20)$$

where $\text{supp}(p) = \{\mathbf{x} \mid p(\mathbf{x}) > 0\}$ and so on.¹² Note that, for two sets A and B , we write $A \subset B$ or $B \supset A$ if $a \in B$ for all $a \in A$ so that $\text{supp}(p)$ may equal $\text{supp}(q)$ in (20).¹³

Assuming the condition (20) under which the Kullback-Leibler divergence $\text{KL}(p\|q)$ is well-defined, we can restrict the integration in $\text{KL}(p\|q)$ only over the support $\Omega \equiv \text{supp}(q)$

¹⁰Recall that we have defined $0 \log_2 0 \equiv 0$ or, equivalently, $0 \ln 0 \equiv 0$ (Page 49, Paragraph 2, Line -2) so that the entropy in "bits" (1.93) or "nats" (7) is well-defined.

¹¹One can understand (20) intuitively from the perspective of information theory as follows. As we have seen in Section 1.6.1, the Kullback-Leibler divergence $\text{KL}(p\|q)$ can be interpreted as the average amount of information (in nats) wasted to encode samples generated from the source p with an encoder optimized for q . In order for this relative entropy $\text{KL}(p\|q)$ to be well-defined (i.e., in order that we can encode every sample), the support of the source p must be contained in that of the encoder q . Note also that, in the context of variational inference in which we minimize $\text{KL}(p\|q)$ by optimizing p given q (*variational Bayes*) or q given p (*expectation propagation*), the property (20) is referred to as *zero-forcing* or *zero-avoiding* because $q = 0$ implies $p = 0$ or, equivalently, $p > 0$ implies $q > 0$ (see Section 10.1.2).

¹²Here, we define the *support* $\text{supp}(f)$ of a real-valued function $f : \mathcal{X} \rightarrow \mathbb{R}$ by

$$\text{supp}(f) \equiv \{\mathbf{x} \in \mathcal{X} \mid f(\mathbf{x}) \neq 0\} \quad (21)$$

so that, for a probability density function $p : \mathcal{X} \rightarrow [0, \infty)$, we have $\text{supp}(p) = \{\mathbf{x} \in \mathcal{X} \mid p(\mathbf{x}) > 0\}$.

¹³In fact, a set A equals another set B or $A = B$ if and only if $A \subset B$ and $A \supset B$.

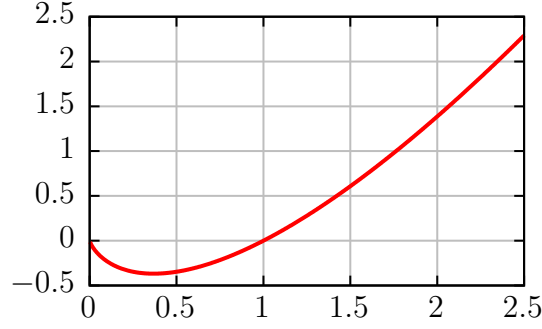


Figure 2 Plot of $f(x) = x \ln x$. The function $f(x)$ is a strictly convex function defined over $[0, \infty)$ where we have defined $f(0) = 0 \ln 0 \equiv 0$. The curve $y = f(x)$ takes the minimum at $(x, y) = (e^{-1}, -e^{-1})$. The roots (the values of x such that $f(x) = 0$) are $x = 0$ and $x = 1$.

of $q(\mathbf{x})$. Identifying $f(\xi) = \xi \ln \xi$ (see Figure 2) and $\xi(\mathbf{x}) = p(\mathbf{x})/q(\mathbf{x})$, we have

$$\text{KL}(p\|q) = \int_{\Omega} q(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} \ln \left\{ \frac{p(\mathbf{x})}{q(\mathbf{x})} \right\} d\mathbf{x} \quad (22)$$

$$= \int_{\Omega} q(\mathbf{x}) f(\xi(\mathbf{x})) d\mathbf{x} \quad (23)$$

$$\geq f\left(\int_{\Omega} q(\mathbf{x}) \xi(\mathbf{x}) d\mathbf{x}\right) \quad (24)$$

$$= f\left(\int_{\Omega} p(\mathbf{x}) d\mathbf{x}\right) \quad (25)$$

$$= f(1) = 0 \quad (26)$$

where we have used Jensen's inequality (17) with respect to $q(\mathbf{x})$ (instead of $p(\mathbf{x})$). Note that, since $q(\mathbf{x}) > 0$ for all $\mathbf{x} \in \Omega$, we see that $\xi(\mathbf{x}) = p(\mathbf{x})/q(\mathbf{x}) \geq 0$ is well-defined for all $\mathbf{x} \in \Omega$ and so is $f(\xi(\mathbf{x})) = \xi(\mathbf{x}) \ln \xi(\mathbf{x})$. Since $f(\xi)$ is strictly convex, the equality $\text{KL}(p\|q) = 0$ holds if and only if $\xi(\mathbf{x}) = p(\mathbf{x})/q(\mathbf{x})$ is constant for all $\mathbf{x} \in \Omega$, which, together with (20), yields the equality condition that $p(\mathbf{x}) = q(\mathbf{x})$ for all \mathbf{x} .

Convexity of Kullback-Leibler divergence Let us next consider the weighted sum of two well-defined Kullback-Leibler divergences $\text{KL}(p_1\|q_1), \text{KL}(p_2\|q_2) < \infty$

$$S \equiv \lambda \text{KL}(p_1\|q_1) + (1 - \lambda) \text{KL}(p_2\|q_2) \quad (27)$$

$$= \lambda \int_{\Omega_1} p_1(\mathbf{x}) \ln \left\{ \frac{p_1(\mathbf{x})}{q_1(\mathbf{x})} \right\} d\mathbf{x} + (1 - \lambda) \int_{\Omega_2} p_2(\mathbf{x}) \ln \left\{ \frac{p_2(\mathbf{x})}{q_2(\mathbf{x})} \right\} d\mathbf{x} \quad (28)$$

where $\lambda \in (0, 1)$. Here, we again assume (20) for each pair of the distributions (p_i, q_i) , i.e., $p_i(\mathbf{x}) = 0$ for all $\mathbf{x} \notin \Omega_i \equiv \text{supp}(q_i)$ where $i = 1, 2$. Writing

$$\begin{aligned} a_1 &\equiv \lambda p_1(\mathbf{x}), & a_2 &\equiv (1 - \lambda) p_2(\mathbf{x}) \\ b_1 &\equiv \lambda q_1(\mathbf{x}), & b_2 &\equiv (1 - \lambda) q_2(\mathbf{x}) \end{aligned} \quad (29)$$

and noting that $a_i = b_i = 0$ outside Ω_i and $b_i > 0$ otherwise, we have

$$S = \int_{\Omega_1} a_1 \ln \frac{a_1}{b_1} d\mathbf{x} + \int_{\Omega_2} a_2 \ln \frac{a_2}{b_2} d\mathbf{x} \quad (30)$$

$$= \int_{\Omega_1 \setminus \Omega_2} a_1 \ln \frac{a_1}{b_1} d\mathbf{x} + \int_{\Omega_2 \setminus \Omega_1} a_2 \ln \frac{a_2}{b_2} d\mathbf{x} + \int_{\Omega_1 \cap \Omega_2} \left[a_1 \ln \frac{a_1}{b_1} + a_2 \ln \frac{a_2}{b_2} \right] d\mathbf{x} \quad (31)$$

$$\geq \int_{\Omega_1 \cup \Omega_2} (a_1 + a_2) \ln \frac{a_1 + a_2}{b_1 + b_2} d\mathbf{x} \quad (32)$$

$$= \text{KL}(\lambda p_1 + (1 - \lambda)p_2 \| \lambda q_1 + (1 - \lambda)q_2) \quad (33)$$

from which we see that $\text{KL}(p \| q)$ is convex with respect to (p, q) . Here, we have used the inequality¹⁴

$$a_1 \ln \frac{a_1}{b_1} + a_2 \ln \frac{a_2}{b_2} = (b_1 + b_2) \left[\frac{b_1}{b_1 + b_2} f\left(\frac{a_1}{b_1}\right) + \frac{b_2}{b_1 + b_2} f\left(\frac{a_2}{b_2}\right) \right] \quad (35)$$

$$\geq (b_1 + b_2) f\left(\frac{a_1 + a_2}{b_1 + b_2}\right) \quad (36)$$

$$= (a_1 + a_2) \ln \frac{a_1 + a_2}{b_1 + b_2} \quad (37)$$

where we have again used Jensen's inequality (15) with $f(\xi) = \xi \ln \xi$.

We also see that (i) convexity of $\text{KL}(p \| q)$ with respect to (p, q) implies (ii) convexity of $\text{KL}(p \| q)$ with respect to p . Although the former convexity (i) is not strict in general (consider, e.g., the case where $\Omega_1 \cap \Omega_2 = \emptyset$), the latter convexity (ii) can be shown to be strict with a similar discussion as above (where we let $q_1 = q_2$ and thus $\Omega_1 = \Omega_2$), i.e., $\text{KL}(p \| q)$ is strictly convex with respect to p so that

$$\lambda \text{KL}(p_1 \| q) + (1 - \lambda) \text{KL}(p_2 \| q) \geq \text{KL}(\lambda p_1 + (1 - \lambda)p_2 \| q), \quad \lambda \in (0, 1) \quad (38)$$

with equality if and only if $p_1 = p_2$.

Extended Kullback-Leibler divergence Let us now define what we call the *extended Kullback-Leibler divergence* (Minka, 2005; Zhu and Rohwer, 1995) as

$$\text{EKL}(\tilde{p} \| \tilde{q}) \equiv \int \tilde{p}(\mathbf{x}) \ln \left\{ \frac{\tilde{p}(\mathbf{x})}{\tilde{q}(\mathbf{x})} \right\} d\mathbf{x} + \int \{\tilde{q}(\mathbf{x}) - \tilde{p}(\mathbf{x})\} d\mathbf{x}. \quad (39)$$

Note that the definition (39) of the extended Kullback-Leibler divergence includes a correction term of the form $\int \{\tilde{q}(\mathbf{x}) - \tilde{p}(\mathbf{x})\} d\mathbf{x}$ so that it applies to unnormalized distributions $\tilde{p}(\mathbf{x})$ and $\tilde{q}(\mathbf{x})$. One can easily see that, for correctly normalized distributions $p(\mathbf{x})$ and $q(\mathbf{x})$, the correction term vanishes so that

$$\text{KL}(p \| q) = \text{EKL}(p \| q). \quad (40)$$

¹⁴The inequality used here is a special case of the *log sum inequality*, which states that, for nonnegative a_i, b_i ,

$$\sum_i a_i \ln \frac{a_i}{b_i} \geq a \ln \frac{a}{b} \quad (34)$$

with equality if and only if there exists some constant c such that $a_i = cb_i$ for all i where $a = \sum_i a_i$ and $b = \sum_i b_i$.

The extended Kullback-Leibler divergence (39) can also be written in the form

$$\text{EKL}(\tilde{p} \parallel \tilde{q}) = \int \tilde{p}(\mathbf{x}) F\left(\ln \left\{ \frac{\tilde{p}(\mathbf{x})}{\tilde{q}(\mathbf{x})} \right\}\right) d\mathbf{x} \quad (41)$$

where

$$F(t) = t + e^{-t} - 1, \quad t \in (-\infty, \infty). \quad (42)$$

Since $F(t) \geq 0$ with equality if and only if $t = 0$,¹⁵ we see that an analogue of Gibbs's inequality (18) holds also for the extended Kullback-Leibler divergence (39), i.e., we have

$$\text{EKL}(\tilde{p} \parallel \tilde{q}) \geq 0 \quad (43)$$

with equality if and only if $\tilde{p}(\mathbf{x}) = \tilde{q}(\mathbf{x})$ for all \mathbf{x} . Thus, we see that Gibbs's inequality (18) for the ordinary Kullback-Leibler divergence (1.113) can also be shown from the Gibbs's inequality (43) for the extended Kullback-Leibler divergence (39) via the relationship (40).

Moreover, one can easily see from the linearity of the correction term that the extended Kullback-Leibler divergence (39) also enjoys convexity similar to that of the ordinary Kullback-Leibler divergence (1.113). Specifically, $\text{EKL}(\tilde{p} \parallel \tilde{q})$ is (i) convex with respect to (\tilde{p}, \tilde{q}) and (ii) strictly convex with respect to \tilde{p} .

It is interesting to note that, if the two distributions are close so that $p(\mathbf{x}) \approx q(\mathbf{x})$, then we can approximate the Kullback-Leibler divergence as

$$\text{KL}(p \parallel q) \approx \frac{1}{2} \int p(\mathbf{x}) \{\ln p(\mathbf{x}) - \ln q(\mathbf{x})\}^2 d\mathbf{x} \quad (44)$$

where we have again used (40) and the Taylor expansion $F(t) \simeq t^2/2$ of $F(t)$ around $t = 0$.¹⁶

Page 59

Exercise 1.7, Line 2: "To do this consider, the..." should be "To do this, consider the..."

Page 61

Exercise 1.15, Line -1: Add a period (.) at the end of the last sentence.

Page 62

Exercise 1.18, the text after (1.142): "Gamma" should read "gamma" (without capitalization).

Page 64

Exercise 1.28, Line 1: In Section 1.6, the quantity $h(x)$ is introduced as a measure of the information gained on observing the random variable x , whereas the *entropy* is the average of $h(x)$ over x . The first sentence should thus read, e.g., "In Section 1.6, we introduced the idea of entropy as the average of the information $h(x)$ gained..."

¹⁵Note that $F(t)$ is a strictly convex function (because $F''(t) > 0$) and takes the minimum $F(t) = 0$ at $t = 0$.

¹⁶Although the sign \simeq is used for any type of approximate equality in PRML, I make some (soft) distinction between \simeq and \approx in this report: I use \simeq for series expansions where approximation can be exact at some point; and \approx for general (e.g., numerical) approximation including the Laplace approximation (see Section 4.4).

Page 64

Exercise 1.28, Lines 3 and 4: “the entropy functions are additive, so that. . .” should read, e.g., “the information $h(\cdot)$ is additive so that. . .” (see also the previous erratum).

Page 69

Equation (2.2): Although the *Bernoulli* distribution $\text{Bern}(x|\mu)$ is a valid, correctly normalized probability distribution for any value of the parameter $0 \leq \mu \leq 1$ (Page 69, Paragraph 1, Line 1),¹⁷ it becomes *degenerate*, i.e., x is fixed to a single value so that $x = 0$ or $x = 1$ if $\mu = 0$ or $\mu = 1$, respectively.¹⁸ Such degenerate distributions are often difficult to treat with some generality so that they actually seem to have been excluded (implicitly or explicitly) in most of the discussions in PRML. For example, we should be unable to take the logarithm of the Bernoulli likelihood (2.5) to give (2.6) if the distribution can be degenerate (because the logarithm diverges if $\mu = 0$ or $\mu = 1$). More generally, one cannot identify a degenerate distribution with any member of the *exponential family* (see Section 2.4). For instance, the degenerate Bernoulli cannot be expressed as the exponential of the logarithm as in (2.197) because its *natural parameter* (2.198) is again not well-defined for $\mu = 0$ or $\mu = 1$.

Restriction on probability of success In order for the Bernoulli distribution $\text{Bern}(x|\mu)$ *not* to be degenerate, we assume in this report that the parameter μ (called the *probability of success*) is restricted on the open set $(0, 1)$, rather than on the closed set $[0, 1]$ as in PRML, so that we can write

$$\text{Bern}(x|\mu) = \mu^x(1 - \mu)^{1-x} = \exp\{x \ln \mu + (1 - x) \ln(1 - \mu)\}, \quad \mu \in (0, 1) \quad (45)$$

while we can still consider the degenerate case as the limit $\mu \rightarrow 0$ or $\mu \rightarrow 1$. Similar discussions also apply to other discrete distributions, i.e., the *binomial* distribution (2.9) and the *multinomial* distribution (2.34), for which we shall, therefore, assume the restrictions $\mu \in (0, 1)$ and $\mu_k \in (0, 1)$ for all k , respectively. See also our definition (193) of the *multinoulli* distribution.

Accordingly, the same restrictions $\mu \in (0, 1)$ and $\mu_k \in (0, 1)$ for all k are assumed also on the domains of the (conjugate) prior distributions, i.e., the *beta* distribution (2.13) and the *Dirichlet* distribution (2.38). Note that these restrictions do not affect correct normalization of the distributions; and also that, with these restrictions, we can avoid potential difficulty that the density function can diverge at the boundary of the domain (see Figures 2.2 and 2.5).

Page 69

Equation (2.6): In order to take the logarithm of the likelihood (2.5), we should assume $\mu \in (0, 1)$ so that the Bernoulli distribution is not degenerate. See (45).

Page 70

Paragraph –1, Line –1: The lower ellipsis (. . .) should be centered (\cdots).

¹⁷Similarly to $0 \ln 0 \equiv 0$, we have defined $0^0 \equiv 1$ here.

¹⁸Note however that, in the context of Bayesian inference in which we regard the parameter μ as a random variable, the Bernoulli distribution $\text{Bern}(x|\mu)$ cannot be degenerate because $\mu \in (0, 1)$ almost surely so that the edge cases, i.e., $\mu = 0$ and $\mu = 1$, do not matter after all.

Page 75

The text following (2.26): We assume in this report that $\mu_k \in (0, 1)$ for all k . See (45).

Page 76

Equations (2.34) and (2.35): The *multinomial coefficient* (2.35) is better written as

$$\binom{N}{m_1, m_2, \dots, m_K} \equiv \frac{N!}{m_1! m_2! \dots m_K!} \quad (46)$$

where m_1, m_2, \dots, m_K are comma separated in the left hand side so as not to be confused with a function of the product of m_1, m_2, \dots, m_K . Also, the lower ellipsis (...) in the right hand side of (2.35) should be centered (...) as shown in (46).

Page 76

The text following (2.37): We assume in this report that $\mu_k \in (0, 1)$ for all k . See (45).

Page 77

The caption of Figure 2.4: We assume in this report that $\mu_k \in (0, 1)$ for all k . See (45).

Page 78

The caption of Figure 2.5: “ $\{\alpha_k\} = 0.1$ ” should read “ $\alpha_k = 0.1$ for all k ” and so on.

Page 80

Equation (2.52): We usually take eigenvectors \mathbf{u}_i to be the columns of \mathbf{U} as in (C.37). If we follow this convention, (2.52) and the following text should read

$$\mathbf{y} = \mathbf{U}^T(\mathbf{x} - \boldsymbol{\mu}) \quad (47)$$

where \mathbf{U} is a matrix whose columns are given by \mathbf{u}_i so that $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_D)$. From (2.46) it follows that \mathbf{U} is an *orthogonal* matrix, i.e., it satisfies $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ and hence also $\mathbf{U} \mathbf{U}^T = \mathbf{I}$ where \mathbf{I} is the identity matrix.

Page 81

Equations (2.53) and (2.54): If we write the change of variable from \mathbf{x} to \mathbf{y} as (47) instead of (2.52), the Jacobian matrix $\mathbf{J} = (J_{ij})$ is simply given by \mathbf{U} . Equation (2.53) should read

$$J_{ij} = \frac{\partial x_i}{\partial y_j} = U_{ij} \quad (48)$$

where U_{ij} is the (ij) -th element of \mathbf{U} . The square of the determinant of the Jacobian matrix (2.54) can then be evaluated as

$$|\mathbf{J}|^2 = |\mathbf{U}|^2 = |\mathbf{U}^T| |\mathbf{U}| = |\mathbf{U}^T \mathbf{U}| = |\mathbf{I}| = 1. \quad (49)$$

Page 81

The text after (2.54): Since the Jacobian matrix \mathbf{J} is only assumed to be orthogonal here, the determinant of \mathbf{J} can be either positive or negative so that we should write $|\mathbf{J}| = \pm 1$ instead of $|\mathbf{J}| = 1$.

Page 82

Equation (2.56): We should take the absolute value of the determinant for the same reason given in the above erratum so that the factor $|\mathbf{J}|$ should read $|\det(\mathbf{J})|$. Note however that it is *not* recommended to write $\|\mathbf{J}\|$ to mean $|\det(\mathbf{J})|$ because $\|\mathbf{J}\|$ is confusingly similar to the *matrix norm* $\|\mathbf{J}\|$, which usually refers to the largest singular value of \mathbf{J} (Golub and Van Loan, 2013). This notational inconsistency is caused by the abuse of the notation $|\cdot|$ for both the absolute value and the matrix determinant; if we always use $\det(\cdot)$ for the determinant, confusion will not arise and the notation be consistent.

Notation for absolute determinant An alternative solution to the problem of notational inconsistency mentioned above would be to explicitly define $|\mathbf{A}|$ as the absolute value of the determinant of a square matrix \mathbf{A} , i.e.,

$$|\mathbf{A}| \equiv |\det(\mathbf{A})| \quad (50)$$

so that we have $|\mathbf{J}| = 1$ and (2.56) holds as is. Note also that this notation (50) is mostly consistent in other part of PRML because we have $|\mathbf{A}| = \det(\mathbf{A})$ for any positive-semidefinite matrix $\mathbf{A} \succeq 0$ (see Appendix C) and most matrices for which we take determinants are in fact positive definite.¹⁹ Such positive-definite matrices include the covariance Σ or the precision Λ of the multivariate Gaussian distribution and the scale matrix \mathbf{W} of the Wishart distribution (see Appendix B).

Page 82

Two lines above (2.59): “the term in \mathbf{z} in the factor $(\mathbf{z} + \boldsymbol{\mu})$ ” should read “the term \mathbf{z} in the factor $(\mathbf{z} + \boldsymbol{\mu})$ ” (Remove the first occurrence of “in”).

Page 90

Paragraph –1, Line 2: The partitioned vector should read (2.65) or $\mathbf{x} = (\mathbf{x}_a^T, \mathbf{x}_b^T)^T$.

Page 91

Paragraph 1, Line 1: “linear Gaussian” should read “linear-Gaussian” (with hyphenation) for consistency with other part of PRML (see, e.g., Section 8.1.4).

¹⁹In this report, we assume as customary that the concept of positive/negative (semi)definiteness is restricted to symmetric matrices. For example, when we say “ \mathbf{A} is positive definite” or $\mathbf{A} \succ 0$, we implicitly assume that \mathbf{A} is also symmetric so that $\mathbf{A}^T = \mathbf{A}$, though we still sometimes say “ \mathbf{A} is symmetric positive definite” to avoid confusion.

Equation (2.120): The vector derivative operator $\frac{\partial}{\partial \boldsymbol{\mu}}$ should read the *gradient* $\nabla_{\boldsymbol{\mu}}$ if we use the notation (294) we adopt in this report.

Equations (2.121) and (2.122): We obtain the maximum likelihood solutions $\boldsymbol{\mu}_{\text{ML}}$ and $\boldsymbol{\Sigma}_{\text{ML}}$ for the Gaussian by setting the derivatives of the log likelihood function $\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ given by (2.118) with respect to $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ equal to zero, which, however, only implies that $\boldsymbol{\mu}_{\text{ML}}$ and $\boldsymbol{\Sigma}_{\text{ML}}$ are stationary points. We should also show that $\boldsymbol{\mu}_{\text{ML}}$ and $\boldsymbol{\Sigma}_{\text{ML}}$ indeed maximize the likelihood as discussed in the following.

Maximum likelihood for Gaussian Let us first maximize the likelihood function with respect to the mean $\boldsymbol{\mu}$. This can be easily done by noting that the log likelihood (2.118) is quadratic in $\boldsymbol{\mu}$ so that

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{N}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_{\text{ML}})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_{\text{ML}}) + \text{const} \quad (51)$$

where $\boldsymbol{\mu}_{\text{ML}}$ is given by (2.121) and the terms independent of $\boldsymbol{\mu}$ have been absorbed into “const.” Since the covariance $\boldsymbol{\Sigma}$ is positive definite and so is its inverse $\boldsymbol{\Sigma}^{-1}$, we see that the log likelihood (51) is concave with respect to $\boldsymbol{\mu}$ and that $\boldsymbol{\mu}_{\text{ML}}$ indeed maximizes the likelihood.

Next, we consider maximization with respect to the covariance $\boldsymbol{\Sigma}$. The maximum likelihood solution $\boldsymbol{\Sigma}_{\text{ML}}$ given by (2.122) can be obtained by solving

$$\nabla_{\boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}_{\text{ML}}, \boldsymbol{\Sigma}) = \mathbf{O} \quad (52)$$

where $\nabla_{\mathbf{A}}$ is the gradient operator with respect to a matrix \mathbf{A} defined by (319) and \mathbf{O} is a zero matrix. Making use of the eigenvalue expansion (2.48) of $\boldsymbol{\Sigma}$, we can write the log likelihood (2.118) in terms of the eigenvalues $\{\lambda_i\}$ so that

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{N}{2} \sum_{i=1}^D \left\{ \ln \lambda_i + \frac{S_i}{\lambda_i} \right\} + \text{const} \quad (53)$$

where

$$S_i = \frac{1}{N} \sum_{n=1}^N y_{ni}^2, \quad y_{ni} = \mathbf{u}_i^T (\mathbf{x}_n - \boldsymbol{\mu}). \quad (54)$$

Although the log likelihood (53) is not a concave function of $\boldsymbol{\Sigma}$ (one can easily see this by considering the univariate case), one can observe that $\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \rightarrow -\infty$ if $\boldsymbol{\Sigma}$ approaches the boundary of the space of symmetric positive-definite matrices, i.e., if $\lambda_i \rightarrow 0$ or $\lambda_i \rightarrow \infty$ for any i . Therefore, if (52) has a unique solution $\boldsymbol{\Sigma}_{\text{ML}} \succ 0$, then $\boldsymbol{\mu}_{\text{ML}}$ and $\boldsymbol{\Sigma}_{\text{ML}}$ jointly maximize the likelihood.

Note that a similar observation holds when we maximize the log likelihood in terms of the precision $\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$, in which case the corresponding log likelihood for $\boldsymbol{\Lambda}$ is given by

$$\ln p(\mathbf{X}|\boldsymbol{\mu}_{\text{ML}}, \boldsymbol{\Lambda}) = \frac{N}{2} \ln |\boldsymbol{\Lambda}| - \frac{N}{2} \text{Tr}(\boldsymbol{\Sigma}_{\text{ML}} \boldsymbol{\Lambda}) + \text{const}. \quad (55)$$

Setting the derivative of (55) with respect to Λ equal to zero, we indeed obtain $\Lambda_{\text{ML}} = \Sigma_{\text{ML}}^{-1}$. One can also see that (55) is actually a strictly concave function of Λ due to the strict concavity of $\ln |\Lambda|$ for $\Lambda \succ 0$ (Magnus and Neudecker, 2007) together with the linearity of $\text{Tr}(\Sigma_{\text{ML}}\Lambda)$. See Anderson and Olkin (1985) for further discussions.

Page 100

Equations (2.147) and (2.148): In addition to the mean $\mathbb{E}[\lambda]$ and the variance $\text{var}[\lambda]$, given by (2.147) and (2.148), respectively, we are also interested in the *log expectation* $\mathbb{E}[\ln \lambda]$, given by (B.30), of the gamma distribution (2.146), which is necessary to evaluate the entropy $H[\lambda]$, given by (B.31). Note that the log expectation of the Dirichlet distribution (2.38) is derived in Exercise 2.11 by differentiating its probability with respect to the parameters (the mean and the covariance are concerned in Exercise 2.10). Applying this technique of differentiation, we can calculate the log expectation of the gamma distribution. Here, I would like to state the technique in more general terms (see Section 2.4 for even more general exposition in terms of the *exponential family*), after which we show (B.30). We also find an alternative form of the log expectation $\mathbb{E}[\ln \lambda]$ in terms of the logarithm of the mean $\ln \mathbb{E}[\lambda]$ and an interesting function related to the digamma function, namely, the *log minus digamma function*.

Score function For a correctly normalized probability distribution $p(\mathbf{x}|\boldsymbol{\theta})$ over some random variables \mathbf{x} parameterized by parameters $\boldsymbol{\theta}$ and differentiable with respect to $\boldsymbol{\theta}$, let us consider how $p(\mathbf{x}|\boldsymbol{\theta})$ changes under perturbations in $\boldsymbol{\theta}$. Specifically, the first-order relative difference in the direction $\boldsymbol{\eta}$ is given by

$$\frac{1}{p(\mathbf{x}|\boldsymbol{\theta})} \lim_{\epsilon \rightarrow 0} \left\{ \frac{p(\mathbf{x}|\boldsymbol{\theta} + \epsilon \boldsymbol{\eta}) - p(\mathbf{x}|\boldsymbol{\theta})}{\epsilon} \right\} = \boldsymbol{\eta}^T \mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) \quad (56)$$

where we have assumed that $p(\mathbf{x}|\boldsymbol{\theta})$ remains correctly normalized under sufficiently small perturbations in $\boldsymbol{\theta}$; and defined the *score function*, denoted by $\mathbf{g}(\boldsymbol{\theta}, \mathbf{x})$, as the derivative of the log probability with respect to $\boldsymbol{\theta}$ so that

$$\mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) \equiv \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x}|\boldsymbol{\theta}) = \frac{\nabla_{\boldsymbol{\theta}} p(\mathbf{x}|\boldsymbol{\theta})}{p(\mathbf{x}|\boldsymbol{\theta})}. \quad (57)$$

Note that the score function (57) is called the Fisher score (6.32) in PRML. In fact, the first-order relative difference (56) is zero on average in whatever the direction $\boldsymbol{\eta}$ because the expectation of the score function (57) vanishes so that

$$\mathbb{E}_{\mathbf{x}} [\mathbf{g}(\boldsymbol{\theta}, \mathbf{x}) | \boldsymbol{\theta}] = \mathbf{0} \quad (58)$$

where $\mathbb{E}_{\mathbf{x}} [\cdot | \boldsymbol{\theta}]$ denotes the conditional expectation (1.37) so that the above expectation is taken with respect to $p(\mathbf{x}|\boldsymbol{\theta})$.

We can show the general identity (58) by differentiating the both sides of the integral identity

$$\int p(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} = 1 \quad (59)$$

with respect to θ , giving

$$\nabla_{\theta} \int p(\mathbf{x}|\theta) d\mathbf{x} = \mathbf{0} \quad (60)$$

$$\int \nabla_{\theta} p(\mathbf{x}|\theta) d\mathbf{x} = \mathbf{0} \quad (61)$$

$$\int p(\mathbf{x}|\theta) \nabla_{\theta} \ln p(\mathbf{x}|\theta) d\mathbf{x} = \mathbf{0} \quad (62)$$

where we have assumed that we can interchange the order of the derivative and the integral; and used the *log derivative identity*

$$\nabla f = f \nabla \ln f. \quad (63)$$

Although we have assumed here that the variables \mathbf{x} are continuous, the same discussion holds if some or all of \mathbf{x} are discrete by replacing the integrations with summations as required.

At this moment, I would like to point out a subtlety in the identity (58). Recall that, when we introduce the score function (57), we have assumed that sufficiently small perturbations in θ do not affect the correct normalization of $p(\mathbf{x}|\theta)$. This assumption is required to show (58) because otherwise the right hand side of (60) would not vanish. Let us take the *multinoulli* distribution $\text{Mult}(\mathbf{x}|\boldsymbol{\mu})$ defined by (193) as an example. Since we cannot change a single parameter μ_k (the normalized probability of observing $x_k = 1$) independently of the others μ_j where $j \neq k$ due to the sum-to-one constraint $\sum_k \mu_k = 1$, it is *not* valid to substitute $\nabla_{\mu_k} \ln \text{Mult}(\mathbf{x}|\boldsymbol{\mu})$ into (58). Instead, we should consider the derivatives with respect to *independent* parameters. The unnormalized probabilities $\tilde{\mu}_k$ related to μ_k through (195) are among such parameters; the corresponding score function is given by

$$\nabla_{\tilde{\mu}_k} \ln \text{Mult}(\mathbf{x}|\boldsymbol{\mu}) = \frac{x_k}{\tilde{\mu}_k} - \frac{1}{\sum_j \tilde{\mu}_j}. \quad (64)$$

Substituting (64) into (58), we indeed obtain a valid result that $\mathbb{E}[x_k] = \mu_k$.

Log expectation of gamma distribution Now, let us return to the gamma distribution (2.146). The derivative of the log probability with respect to a is given by

$$\frac{\partial}{\partial a} \ln \text{Gam}(\lambda|a, b) = \ln \lambda - \psi(a) + \ln b \quad (65)$$

where $\psi(\cdot)$ is the *digamma function* given by (101). Substituting (65) into (58), we obtain

$$\mathbb{E}[\ln \lambda] = \psi(a) - \ln b \quad (66)$$

showing (B.30). Similarly, one can reproduce the result (2.147) for the mean $\mathbb{E}[\lambda]$ by substituting the derivative of the log probability with respect to b into (58).

Log minus digamma function It follows from Jensen's inequality (15) that the log expectation $\mathbb{E}[\ln \lambda]$ is less than the logarithm of the mean $\ln \mathbb{E}[\lambda]$ because $\ln \xi$ is strictly concave where $\xi > 0$ so that

$$\mathbb{E}[\ln \lambda] < \ln \mathbb{E}[\lambda] = \ln \frac{a}{b} \quad (67)$$

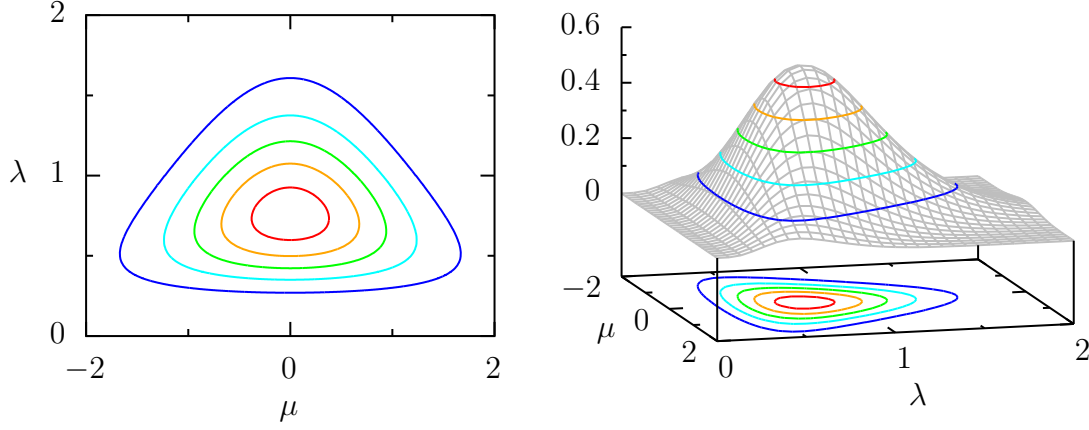


Figure 3 Contour and surface plots of the Gaussian-gamma (normal-gamma) distribution (2.154) where $\mu_0 = 0$, $\beta = 2$, $a = 5$, and $b = 6$. Contours are plotted at densities from 0.1 (the outermost contour, shown in blue) to 0.5 (innermost, red) with an equal step of 0.1.

where we have used (2.147). The difference between $\ln \mathbb{E}[\lambda]$ and $\mathbb{E}[\ln \lambda]$ can be evaluated analytically in this case by noting (66) or (B.30) so that

$$\mathbb{E}[\ln \lambda] = \ln \mathbb{E}[\lambda] - \varphi(a) \quad (68)$$

where $\varphi(a) > 0$ is what we call the *log minus digamma function* defined by

$$\varphi(a) \equiv \ln a - \psi(a), \quad a \in (0, \infty). \quad (69)$$

The log minus digamma function (69) naturally arises also in deriving the maximum likelihood solution for the gamma distribution as we shall see shortly.

Page 102

Figure 2.14: Since no contour labels are given, this contour plot alone does not convey very useful information regarding the shape of the distribution. For a better grasp, we can use the contour and the surface plots in combination as shown in Figure 3.

Page 102

Equation (2.155): Although an interpretation for the parameters of the gamma distribution (2.146) has been given, no such an interpretation for the parameters of the Wishart distribution (2.155) is given here nor in Exercise 2.45. Generally speaking, when we construct a probabilistic model with priors, we must choose some reasonable (initial) values for their parameters, known as *hyperparameters*; this calls for an intuitive interpretation for the parameters of such priors. We can give an interpretation for the parameters of the Wishart distribution as follows.

Interpreting parameters of Wishart Let us consider a simple Bayesian inference problem in which, given a set of N observations $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ for a zero-mean Gaussian random variable, we infer the covariance matrix Σ or, equivalently, the precision matrix $\Lambda \equiv \Sigma^{-1}$.

The likelihood $p(\mathbf{X}|\mathbf{\Lambda})$ in terms of the precision $\mathbf{\Lambda}$ is given by

$$p(\mathbf{X}|\mathbf{\Lambda}) = \prod_{n=1}^N p(\mathbf{x}_n|\mathbf{\Lambda}) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n|\mathbf{0}, \mathbf{\Lambda}^{-1}). \quad (70)$$

If we choose the prior $p(\mathbf{\Lambda})$ over $\mathbf{\Lambda}$ to be a Wishart distribution so that

$$p(\mathbf{\Lambda}) = \mathcal{W}(\mathbf{\Lambda}|\mathbf{W}_0, \nu_0) \quad (71)$$

our analysis can be simplified because it is the conjugate prior. In fact, the posterior $p(\mathbf{\Lambda}|\mathbf{X})$ is given by

$$p(\mathbf{\Lambda}|\mathbf{X}) \propto p(\mathbf{X}|\mathbf{\Lambda}) p(\mathbf{\Lambda}) \quad (72)$$

$$\propto |\mathbf{\Lambda}|^{N/2} \exp\left\{-\frac{1}{2} \sum_{n=1}^N \mathbf{x}_n^T \mathbf{\Lambda} \mathbf{x}_n\right\} |\mathbf{\Lambda}|^{(\nu_0-D-1)/2} \exp\left\{-\frac{1}{2} \text{Tr}(\mathbf{W}_0^{-1} \mathbf{\Lambda})\right\} \quad (73)$$

$$= |\mathbf{\Lambda}|^{(\nu_N-D-1)/2} \exp\left\{-\frac{1}{2} \text{Tr}(\mathbf{W}_N^{-1} \mathbf{\Lambda})\right\} \quad (74)$$

where

$$\nu_N = \nu_0 + N \quad (75)$$

$$\mathbf{W}_N^{-1} = \mathbf{W}_0^{-1} + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T. \quad (76)$$

Reinstating the normalization constant, we indeed see that the posterior becomes again a Wishart distribution of the form

$$p(\mathbf{\Lambda}|\mathbf{X}) = \mathcal{W}(\mathbf{\Lambda}|\mathbf{W}_N, \nu_N). \quad (77)$$

This result suggests us how we can interpret the parameters of the Wishart distribution (2.155), namely the scale matrix \mathbf{W} and the number of degrees of freedom ν . Since observing N data points increases the number of degrees of freedom ν by N , we can interpret ν_0 in the prior (71) as the number of “effective” prior observations. The N observations also contribute $N\Sigma_{\text{ML}}$ to the inverse of the scale matrix \mathbf{W} where Σ_{ML} is the maximum likelihood estimate for the covariance of the observations given by

$$\Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \quad (78)$$

suggesting an interpretation of \mathbf{W} in terms of the “covariance” parameter

$$\Sigma \equiv (\nu \mathbf{W})^{-1}. \quad (79)$$

More specifically, we can interpret $\Sigma_0 = (\nu_0 \mathbf{W}_0)^{-1}$ as the covariance of the ν_0 “effective” prior observations. Note that this interpretation is in accordance with another observation that the expectation of $\mathbf{\Lambda}$ with respect to the prior (71) is indeed given by $\mathbb{E}[\mathbf{\Lambda}] = \nu_0 \mathbf{W}_0 = \Sigma_0^{-1}$ where we have used (B.80).

Equation (2.157): Again, no interpretation is given for the parameters of the Gaussian-Wishart distribution (2.157) nor for those of the Gaussian-gamma distribution (2.154). Since the Gaussian-gamma can be obtained as a special case of the Gaussian-Wishart where the dimension is one so that $D = 1$, we shall make an interpretation only for the parameters of the Gaussian-Wishart here.

Interpreting parameters of Gaussian-Wishart Let us consider a problem of inferring the mean $\boldsymbol{\mu}$ and the precision $\boldsymbol{\Lambda}$ given the Gaussian likelihood

$$p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (80)$$

and the Gaussian-Wishart prior

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) = \mathcal{N}(\boldsymbol{\mu}|\boldsymbol{\mu}_0, (\beta_0 \boldsymbol{\Lambda})^{-1}) \mathcal{W}(\boldsymbol{\Lambda}|\mathbf{W}_0, \nu_0). \quad (81)$$

At this moment, we introduce notations for the maximum likelihood estimates for the mean and the covariance given the N observations $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, i.e.,

$$\boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad \boldsymbol{\Sigma}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T \quad (82)$$

respectively. Evaluating the posterior, we have

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda}|\mathbf{X}) \propto p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Lambda}) p(\boldsymbol{\mu}, \boldsymbol{\Lambda}) \quad (83)$$

$$\begin{aligned} & \propto |\boldsymbol{\Lambda}|^{N/2} \exp \left\{ -\frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{x}_n - \boldsymbol{\mu}) \right\} \\ & \times |\boldsymbol{\Lambda}|^{(\nu_0 - D)/2} \exp \left\{ -\frac{1}{2} \text{Tr} \left(\left\{ \mathbf{W}_0^{-1} + \beta_0 (\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \right\} \boldsymbol{\Lambda} \right) \right\} \end{aligned} \quad (84)$$

$$= |\boldsymbol{\Lambda}|^{(\nu_N - D)/2} \exp \left\{ -\frac{1}{2} \text{Tr} \left(\left\{ \mathbf{W}_N^{-1} + \beta_N (\boldsymbol{\mu} - \boldsymbol{\mu}_N)(\boldsymbol{\mu} - \boldsymbol{\mu}_N)^T \right\} \boldsymbol{\Lambda} \right) \right\} \quad (85)$$

where²⁰

$$\beta_N = \beta_0 + N \quad (90)$$

$$\beta_N \boldsymbol{\mu}_N = \beta_0 \boldsymbol{\mu}_0 + N \boldsymbol{\mu}_{\text{ML}} \quad (91)$$

$$\nu_N = \nu_0 + N \quad (92)$$

$$\mathbf{W}_N^{-1} = \mathbf{W}_0^{-1} + N \left[\boldsymbol{\Sigma}_{\text{ML}} + \frac{\beta_0}{\beta_N} (\boldsymbol{\mu}_{\text{ML}} - \boldsymbol{\mu}_0) (\boldsymbol{\mu}_{\text{ML}} - \boldsymbol{\mu}_0)^T \right]. \quad (93)$$

Thus, we find that the posterior is again a Gaussian-Wishart of the form

$$p(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{X}) = \mathcal{N}(\boldsymbol{\mu} | \boldsymbol{\mu}_N, (\beta_N \boldsymbol{\Lambda})^{-1}) \mathcal{W}(\boldsymbol{\Lambda} | \mathbf{W}_N, \nu_N). \quad (94)$$

Note that a similar result is obtained in Section 10.2.1 for a Bayesian mixture of Gaussians model in which we assume a Gaussian-Wishart prior for each Gaussian component.

From the above result, we see that the parameters β_0 and $\boldsymbol{\mu}_0$ for the mean $\boldsymbol{\mu}$ can be interpreted somewhat independently of those ν_0 and \mathbf{W}_0 for the precision $\boldsymbol{\Lambda}$. We can interpret β_0 as the number of “effective” prior observations for $\boldsymbol{\mu}$ and $\boldsymbol{\mu}_0$ as the mean of the β_0 prior observations. The interpretation of ν_0 and \mathbf{W}_0 is similar to the one we have made in the previous erratum except that we have in (93) a term due to the uncertainty in $\boldsymbol{\mu}$, that is, a term involving the outer product of the difference between the maximum likelihood mean $\boldsymbol{\mu}_{\text{ML}}$ and the prior mean $\boldsymbol{\mu}_0$, scaled by β_0/β_N .

Page 102

Paragraph –1, Line –2: “Gamma” should read “gamma” (without capitalization).

Page 103

Figure 2.15: The tails of Student’s t-distributions are too high; one can easily see that, if compared to the corresponding Gaussian distribution labeled $\nu \rightarrow \infty$, the t-distributions are not correctly normalized. Figure 4 gives the correct plot.

Page 103

Paragraph –1, Line –3: As pointed out in the text, the maximum likelihood solution for Student’s t-distribution can be most easily found by the *expectation maximization* (EM) algorithm, which we study for discrete and continuous latent variables in Chapters 9 and 12, respectively; it is not until Exercise 12.24 that we apply the EM algorithm to the problem of

²⁰The form (93) of \mathbf{W}_N^{-1} is a little tricky to obtain so that I would like to show a more detailed derivation here. Collecting and evaluating the coefficients of $\boldsymbol{\Lambda}$ inside $\text{Tr}(\cdot)$ in the posterior (84), we have

$$\mathbf{W}_0^{-1} + \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^T + \beta_0 (\boldsymbol{\mu} - \boldsymbol{\mu}_0) (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \quad (86)$$

$$= \mathbf{W}_0^{-1} + N \boldsymbol{\Sigma}_{\text{ML}} + N (\boldsymbol{\mu} - \boldsymbol{\mu}_{\text{ML}}) (\boldsymbol{\mu} - \boldsymbol{\mu}_{\text{ML}})^T + \beta_0 (\boldsymbol{\mu} - \boldsymbol{\mu}_0) (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \quad (87)$$

$$= \mathbf{W}_0^{-1} + N \boldsymbol{\Sigma}_{\text{ML}} + \beta_N (\boldsymbol{\mu} - \boldsymbol{\mu}_N) (\boldsymbol{\mu} - \boldsymbol{\mu}_N)^T - \beta_N \boldsymbol{\mu}_N \boldsymbol{\mu}_N^T + N \boldsymbol{\mu}_{\text{ML}} \boldsymbol{\mu}_{\text{ML}}^T + \beta_0 \boldsymbol{\mu}_0 \boldsymbol{\mu}_0^T \quad (88)$$

$$= \mathbf{W}_0^{-1} + N \boldsymbol{\Sigma}_{\text{ML}} + \beta_N (\boldsymbol{\mu} - \boldsymbol{\mu}_N) (\boldsymbol{\mu} - \boldsymbol{\mu}_N)^T + \frac{\beta_0 N}{\beta_0 + N} (\boldsymbol{\mu}_{\text{ML}} - \boldsymbol{\mu}_0) (\boldsymbol{\mu}_{\text{ML}} - \boldsymbol{\mu}_0)^T. \quad (89)$$

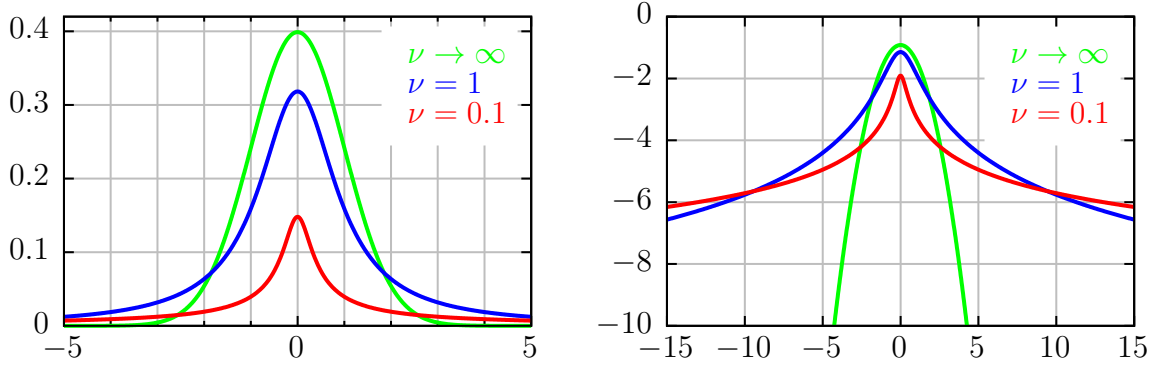


Figure 4 Plot of Student's t-distribution density functions $\text{St}(x|\mu, \lambda, \nu)$ (left) and corresponding log density functions $\ln \text{St}(x|\mu, \lambda, \nu)$ (right) for various values of ν where we have fixed $\mu = 0$ and $\lambda = 1$.

maximum likelihood for the (multivariate) Student's t-distribution (2.162). Although we have to defer the derivation of the above mentioned EM algorithm for some time, it is useful to have considered a related problem of maximum likelihood for the gamma distribution (2.146) here in advance, because, since the t-distribution is obtained by marginalizing over a gamma distributed precision as we have seen in (2.158), we need to estimate the gamma distribution as a subproblem of the EM for the t-distribution.

Maximum likelihood for gamma distribution Given a data set $\mathbf{x} = \{x_1, \dots, x_N\}$, we consider a likelihood function of the form

$$p(\mathbf{x}|a, b) = \prod_{n=1}^N \text{Gam}(x_n|a, b) \quad (95)$$

where the gamma distribution $\text{Gam}(\lambda|a, b)$ is given by (2.146). The log likelihood is given by

$$\ln p(\mathbf{x}|a, b) = N \{-\ln \Gamma(a) + a \ln b + (a-1) \ln \hat{x} - b\bar{x}\} \quad (96)$$

where \bar{x} and \hat{x} denote the *arithmetic mean* and the *geometric mean*, respectively, so that

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n, \quad \ln \hat{x} = \frac{1}{N} \sum_{n=1}^N \ln x_n. \quad (97)$$

We see from (96) that \bar{x} and $\ln \hat{x}$ are the sufficient statistics of the gamma distribution. Here, we assume that $x_n > 0$ for all n (which holds with probability one if x_n has been drawn from a gamma distribution) so that we have $\bar{x} > 0$ and $\hat{x} > 0$.

Let us first assume that $a > 0$ is known. It is easy to see that the log likelihood (96) is a strictly concave function of $b > 0$. We can maximize the likelihood by setting the derivative of (96) with respect to b equal to zero, which gives $b = a/\bar{x}$. Back substituting this into (96), we have

$$\ln p(\mathbf{x}|a, b)|_{b=a/\bar{x}} = N \{-\ln \Gamma(a) + a \ln a - a \ln \bar{x} + (a-1) \ln \hat{x} - a\}. \quad (98)$$

Next we maximize (98) with respect to $a > 0$. This can be done by setting the derivative of (98) with respect to a equal to zero, which gives a nonlinear equation of the form

$$\varphi(a) = \ln \bar{x} - \ln \hat{x} \quad (99)$$

where $\varphi(\cdot)$ is the *log minus digamma function* given by (69). One can see that (98) is again a strictly concave function of $a > 0$ because $\varphi(a)$ is a strictly monotonically decreasing function so that $\varphi'(a) < 0$.

It follows from Jensen's inequality (15) that $\ln \bar{x} \geq \ln \hat{x}$ (which implies $\bar{x} \geq \hat{x}$ because of the monotonicity of the logarithm). Here, we further assume that the strict inequality $\ln \bar{x} > \ln \hat{x}$ holds so that the right hand side of (99) lies among $(0, \infty)$. Since the log minus digamma function $\varphi : (0, \infty) \rightarrow (0, \infty)$ is bijective and thus has the inverse function $\varphi^{-1} : (0, \infty) \rightarrow (0, \infty)$, we can solve (99) uniquely for $a > 0$. Substituting this into $b = a/\bar{x}$, we finally obtain the maximum likelihood solution for the gamma distribution

$$a_{\text{ML}} = \varphi^{-1}(\ln \bar{x} - \ln \hat{x}), \quad b_{\text{ML}} = \frac{a_{\text{ML}}}{\bar{x}}. \quad (100)$$

Page 104

Paragraph 1, Line 4: In practical applications, the importance of *robustness to outliers* cannot be overemphasized. Here, I would like to point out that, particularly in the context of robust regression, there exist historically a number of heuristic approaches to robustness such as *M-estimators* (Press et al., 1992; Szeliski, 2010), in which the standard least squares method is modified so as to use a more “robust” cost function. In this respect, the robust regression in terms of Student's t-distribution can be regarded as an M-estimator where the cost function is derived from its negative log likelihood.

An M-estimator can be solved iteratively by approximating the cost functions successively in terms of quadratic bounds. Called *iteratively reweighted least squares* or IRLS, this algorithm closely resembles the EM algorithm. In fact, one can identify the IRLS and the EM for the robust regression in terms of Student's t-distribution. Note also that the successive quadratic approximation in IRLS can be regarded as a *local variational method* discussed in Chapter 10.

Although we are free to choose from a broad class of cost functions in M-estimators, such a heuristic choice makes our Bayesian analysis difficult. For instance, M-estimators need a separate evaluation data set for selecting hyperparameters. On the other hand, probabilistic models such as the robust regression model in terms of Student's t-distribution allow us to perform model selection in a consistent way without the need for an evaluation data set, while they, in principle, do not suffer from overfitting.

Page 104

The text after (2.160): The Gaussian $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda})$ should read $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$.

Page 110

Paragraph –1, Line 3: Insert a space before the sentence starting “This is known. . .” (*the third printing only*).

Page 112

Paragraph –1, Line 3: Insert a comma (,) after the ellipsis (. . .).

Page 114

Equation (2.210): We assume in this report that $\mu_k \in (0, 1)$ for all k and thus $\sum_{k=1}^{K-1} \mu_k \in (0, 1)$. See (45).

Page 116

Equation (2.224): The right hand side should be a zero vector $\mathbf{0}$ (instead of a scalar zero 0). Furthermore, I would like to point out that the gradient operator ∇ , which is first used here in PRML, has not been properly defined. Actually, although Appendix C introduces the “vector derivative” $\frac{\partial}{\partial \mathbf{x}}$, which is (perhaps confusingly) used interchangeably with the gradient $\nabla_{\mathbf{x}}$ throughout PRML, the gradient itself is not defined anywhere in PRML. Moreover, the vector derivatives given in Appendix C are not well-defined (we shall come back to this issue later in this report). See (294) for the definition of the gradient ∇ we adopt in this report.

Page 119

The line before (2.239): “for choices” should be “for all choices.”

Page 126

The caption of Figure 2.28: The red, green, and blue points correspond to the “homogeneous,” “annular,” and “laminar” (or “stratified”) classes, respectively.

Page 127

Paragraph 2, Lines 1 and 2: Remove the two commas before and after the phrase “and the kernel density estimator.”

Page 128

Exercise 2.4, Line 3: “the mean of n ” should be “the mean of m .”

Page 129

Exercise 2.9, Line 1: Remove the period (.) after [www](#).

Page 129

Equation (2.275): Insert a comma (,) between μ_j and μ_l so that the left hand side of (2.275) reads $\text{cov}[\mu_j, \mu_l]$.

Page 130

Equation (2.277): In order to be consistent with the mathematical notation in PRML, the differential operator d should be an upright d . Specifically, the *digamma function* is given by

$$\psi(a) \equiv \frac{d}{da} \ln \Gamma(a) = \frac{\Gamma'(a)}{\Gamma(a)}. \quad (101)$$

Note that the digamma function (101) is also known as the *psi function* (Abramowitz and Stegun, 1964; Olver et al., 2017).

Page 132

Exercise 2.28, Line –2: “given (2.99)” should be “given by (2.99).”

Page 133

Exercise 2.35, Line 1: Remove the comma in “the results (2.59), and (2.62).”

Page 138

Equation (3.1): The lower ellipsis (\dots) should be centered (\cdots).

Page 141

Equation (3.13): The use of the gradient operator ∇ is not consistent here. As in, e.g., (2.224), the gradient of a scalar function is usually defined as a column vector of derivatives so that (3.13) should read²¹

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w}, \beta) = \beta \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n) \quad (102)$$

where we have written the variable \mathbf{w} with respect to which we take the gradient in the subscript of ∇ explicitly. See (294) for the definition of the gradient operator ∇ adopted in this report.

Page 142

Equation (3.14): The left hand side should be a zero vector $\mathbf{0}$ instead of a scalar zero 0 so that (3.14) should read

$$\mathbf{0} = \sum_{n=1}^N t_n \phi(\mathbf{x}_n) - \left(\sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right) \mathbf{w} \quad (103)$$

where we have used the gradient of the form (102) instead of (3.13).

Page 142

Equation (3.16): Note that the *design matrix* Φ can be expressed more concisely as

$$\Phi = (\phi_1, \phi_2, \dots, \phi_N)^T \quad (104)$$

where we have written $\phi_n = \phi(\mathbf{x}_n)$. Using this representation (104), one can more easily see that the likelihood function (3.10) can also be written as a multivariate Gaussian so that

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{t}|\Phi\mathbf{w}, \beta^{-1}\mathbf{I}) \quad (105)$$

where the target variables $\{t_n\}$ have been grouped into a column vector

$$\mathbf{t} = (t_1, t_2, \dots, t_N)^T. \quad (106)$$

²¹Note that we use a different typeface (from a D -dimensional target variable \mathbf{t}) for the N -dimensional vector $\mathbf{t} = (t_1, \dots, t_N)^T$ consisting of one-dimensional target variables $\{t_n\}$ where $n = 1, \dots, N$. See also “Mathematical Notation” for PRML on Pages xi–xii.

Similarly, the sum-of-squares error function (3.12) can be written in the form

$$E_D(\mathbf{w}) = \frac{1}{2} \|\mathbf{t} - \Phi \mathbf{w}\|^2 \quad (107)$$

where $\|\varepsilon\| = \sqrt{\varepsilon^T \varepsilon}$ is the *norm* of a vector ε . Taking the gradient of (107) with respect to \mathbf{w} , we have

$$\nabla_{\mathbf{w}} E_D(\mathbf{w}) = -\Phi^T (\mathbf{t} - \Phi \mathbf{w}) \quad (108)$$

where we have used the identity (305) together with the chain rule (297) and the identity (302). Setting the gradient (108) equal to zero, we directly obtain the *normal equations* (3.15).

Page 146

Equation (3.31): The left hand side should be $\mathbf{y}(\mathbf{x}, \mathbf{W})$ instead of $\mathbf{y}(\mathbf{x}, \mathbf{w})$.

Page 147

Paragraph –2: The argument that “the phenomenon of [overfitting] does not arise when we marginalize over parameters in a Bayesian setting” is simply an overstatement. Bayesian methods, like any other machine learning methods, can overfit because the *true* model from which the data set has been generated is unknown in general so that one could possibly assume an inappropriate model. For instance, if too broad the prior distribution (3.52) is used in the Bayesian regression model of Section 3.3, this effectively leads to insufficient regularization and thus overfitting.

Moreover, one should also be aware of a subtlety here, that is, (i) the *generalization error*, which can be measured by cross-validation, and (ii) the *model evidence* (or the *marginal likelihood*) are closely related but different criteria, although, in practice, a higher model evidence often tends to imply a lower generalization error and vice versa. For more (advanced) discussions, see Watanabe (2010, 2013).

Page 166

Paragraph 2, Line 1: “Gamma” should read “gamma” (without capitalization).

Pages 168–169, and 177

Equations (3.88), (3.93), and (3.117) as well as the text before (3.93): The derivative operators should be partial differentials. For example, (3.117) should read

$$\frac{\partial}{\partial \alpha} \ln |\mathbf{A}| = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial}{\partial \alpha} \mathbf{A} \right). \quad (109)$$

Page 170

Figure 3.15: The eigenvectors \mathbf{u}_1 and \mathbf{u}_2 are unit vectors so that their orientations should be shown as in Figure 2.7 on Page 81. Or, the scaled vectors \mathbf{u}_1 and \mathbf{u}_2 should be labeled as $\lambda_1^{-1/2} \mathbf{u}_1$ and $\lambda_2^{-1/2} \mathbf{u}_2$, respectively.

Page 174

Exercise 3.4, Line –4: The *Kronecker delta* δ_{ij} should read I_{ij} for consistency with other part of PRML where I_{ij} is the (i, j) -th element of the identity matrix \mathbf{I} (see “Mathematical Notation” on Pages xi–xii).

Page 179

Paragraph 1, Line –4: The decision surfaces are defined by linear *equations* of the input vector \mathbf{x} and thus are $(D - 1)$ -dimensional hyperplanes within the D -dimensional input space.

Page 186

Paragraph 2, Line 2: Insert a space before the sentence starting “This shows a . . .” (*the third printing only*).

Page 190

Equation (4.33): The right hand side should be a zero vector $\mathbf{0}$ instead of a scalar zero 0.

Page 205

Equation (4.88): The differential operator d should be an upright d.

Page 207

Equation (4.92): The gradient and the Hessian in the right hand side, which are in general functions of the parameter \mathbf{w} , must be evaluated at the previous estimate \mathbf{w}^{old} for the parameter. Thus, (4.92) should read

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - [\mathbf{H}(\mathbf{w}^{\text{old}})]^{-1} \nabla E(\mathbf{w}^{\text{old}}) \quad (110)$$

where $\mathbf{H}(\mathbf{w}) \equiv \nabla \nabla E(\mathbf{w})$ is the Hessian matrix whose elements comprise the second derivatives of $E(\mathbf{w})$ with respect to the components of \mathbf{w} .

Page 210

Equation (4.110) and the preceding text: The left hand side of (4.110) is obtained by taking the gradient of $\nabla_{\mathbf{w}_j} E$ given in (4.109) with respect to \mathbf{w}_k and corresponds to the (k, j) -th block of the Hessian, *not* the (j, k) -th. Thus, (4.110) should read

$$\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N y_{nj} (I_{kj} - y_{nk}) \phi_n \phi_n^T. \quad (111)$$

To be clear, we have used the following notation. If we group all the parameters $\mathbf{w}_1, \dots, \mathbf{w}_K$ into a column vector

$$\mathbf{w} = \begin{pmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_K \end{pmatrix} \quad (112)$$

the gradient and the Hessian of the error function $E(\mathbf{w})$ with respect to \mathbf{w} are given by

$$\nabla_{\mathbf{w}} E = \begin{pmatrix} \nabla_{\mathbf{w}_1} E \\ \vdots \\ \nabla_{\mathbf{w}_K} E \end{pmatrix}, \quad \nabla_{\mathbf{w}} \nabla_{\mathbf{w}} E = \begin{pmatrix} \nabla_{\mathbf{w}_1} \nabla_{\mathbf{w}_1} E & \cdots & \nabla_{\mathbf{w}_1} \nabla_{\mathbf{w}_K} E \\ \vdots & \ddots & \vdots \\ \nabla_{\mathbf{w}_K} \nabla_{\mathbf{w}_1} E & \cdots & \nabla_{\mathbf{w}_K} \nabla_{\mathbf{w}_K} E \end{pmatrix} \quad (113)$$

respectively.

Pages 212–214

Equations (4.119), (4.122), (4.126), and (4.128): The differential operator d should be an upright d .

Page 213

Paragraph 1, Line 1: “must related” should be “must be related.”

Page 218

Equation (4.144): The covariance should be the one \mathbf{S}_N evaluated at \mathbf{w}_{MAP} . To make this point clear, we can write the approximate posterior in the form

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{\text{MAP}}, \mathbf{S}_{\text{MAP}}) \quad (114)$$

where

$$\mathbf{S}_{\text{MAP}} = \mathbf{S}_N |_{\mathbf{w}=\mathbf{w}_{\text{MAP}}} \quad (115)$$

and \mathbf{S}_N is given by (4.143).

Page 219

Equation (4.150): \mathbf{m}_N should read \mathbf{w}_{MAP} . Note that the notation \mathbf{m}_N is the one used for the mean (3.50) of the posterior (3.49) for the Bayesian linear regression whereas \mathbf{w}_{MAP} , which however cannot be represented analytically, is the mean of the approximate posterior (114) for the Bayesian logistic regression. Furthermore, if we adopt the notation (115) for the covariance of the approximate posterior (114), then we have the variance σ_a^2 in the form

$$\sigma_a^2 = \phi^T \mathbf{S}_{\text{MAP}} \phi. \quad (116)$$

Page 237

Equation (5.26): The right hand side should be a zero vector $\mathbf{0}$ instead of a scalar zero 0.

Page 237

Paragraph 4, Line 2: $\nabla E(\mathbf{w}) = 0$ should read $\nabla E(\mathbf{w}) = \mathbf{0}$ (the right hand side should be a zero vector $\mathbf{0}$).

Page 238

Paragraph 2, Line 3: $\nabla E(\mathbf{w}) = 0$ should read $\nabla E(\mathbf{w}) = \mathbf{0}$ (the right hand side should be a zero vector $\mathbf{0}$).

Page 238

Equation (5.32): Since we refer to the right hand side of (5.32) later, let us write it as, say, $\tilde{E}(\mathbf{w})$ so that (5.32) reads

$$E(\mathbf{w}) \simeq \tilde{E}(\mathbf{w}) \equiv E(\mathbf{w}^*) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^T \mathbf{H} (\mathbf{w} - \mathbf{w}^*). \quad (117)$$

Perhaps we should have used different notation because we use \tilde{E} differently in Section 5.5.5, though no big confusion will occur even if we use $\tilde{E}(\mathbf{w})$ here.

Page 238

Equation (5.34): The Kronecker delta δ_{ij} should read I_{ij} for consistency with other part of PRML.

Page 238

Equation (5.36): The left hand side should read $\tilde{E}(\mathbf{w})$ where $\tilde{E}(\mathbf{w})$ is given by (117).

Page 239

Figure 5.6: The eigenvectors \mathbf{u}_1 and \mathbf{u}_2 are unit vectors so that their orientations should be shown as in Figure 2.7 on Page 81. Or, the scaled vectors \mathbf{u}_1 and \mathbf{u}_2 should be labeled as $\lambda_1^{-1/2} \mathbf{u}_1$ and $\lambda_2^{-1/2} \mathbf{u}_2$, respectively.

Page 246

The line following (5.65): “ δs ” should read “ $\delta's$ ” for consistency with the line above (5.65).

Page 248

Equations (5.75) and (5.76): The Kronecker delta δ_{kl} should read I_{kl} for consistency with other part of PRML. See also (4.106).

Page 251

Paragraph 2, Line 1: The outer product approximation to the Hessian of the form (5.84) is usually referred to as the *Gauss-Newton* approximation (Press et al., 1992), which not only eliminates the computation of second derivatives but also guarantees that the Hessian thus approximated is positive (semi)definite, whereas the *Levenberg-Marquardt* method (Press et al., 1992) is a method that improves the numerical stability of (Gauss-)Newton type iterations by correcting the Hessian matrix so as to be more diagonal dominant. Let us now compare the two types of approximation to the Hessian, i.e., Gauss-Newton and Levenberg-Marquardt, more specifically in the following. We first observe that the Gauss-Newton approximation to the Hessian given in the right hand side of (5.84) can be written succinctly in terms of matrix product as

$$\mathbf{H}_{\text{GN}} = \mathbf{J}^T \mathbf{J} \quad (118)$$

where $\mathbf{J} = (\nabla a_1, \dots, \nabla a_N)^T$ is the Jacobian of the activations a_1, \dots, a_N with respect to the parameters (weights and biases). The Levenberg-Marquardt approximation to the above Hessian typically takes the form

$$\mathbf{H}_{\text{LM}} = \mathbf{J}^T \mathbf{J} + \lambda \mathbf{I} \quad (119)$$

or

$$\mathbf{H}_{\text{LM}} = \mathbf{J}^T \mathbf{J} + \lambda \text{diag}(\mathbf{J}^T \mathbf{J}) \quad (120)$$

where we have introduced an adjustable damping factor $\lambda \geq 0$ (which will be adjusted through the iterations) and defined that, for a square matrix $\mathbf{A} = (A_{ij})$, $\text{diag}(\mathbf{A})$ is a diagonal matrix obtained by setting the off-diagonal elements equal to zero so that $\text{diag}(\mathbf{A}) = \text{diag}(A_{ii})$.

Page 259

Paragraph 1, Line -1: The parameters rescaling should be $\lambda_1 \rightarrow a^2 \lambda_1$ and $\lambda_2 \rightarrow c^{-2} \lambda_2$.

Page 266

The unlabeled equation following the line starting “Substituting into the mean error function (5.130). . . ”:²² Add the term $O(\xi^3)$ to the right hand side (*the third printing only*).

Page 266

Equation (5.132): The third occurrence of the superscript T for matrix transpose should be an upright T .

Page 267

Equation (5.134): The superscript T should be an upright T .

Page 275

The text after (5.154): The identity matrix \mathbf{I} should multiply $\sigma_k^2(\mathbf{x}_n)$.

Page 277

Equation (5.160): The factor L should multiply $\sigma_k^2(\mathbf{x})$ because we have

$$s^2(\mathbf{x}) = \mathbb{E} \left[\text{Tr} \left\{ (\mathbf{t} - \mathbb{E}[\mathbf{t}|\mathbf{x}]) (\mathbf{t} - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T \right\} \middle| \mathbf{x} \right] \quad (121)$$

$$= \sum_{k=1}^K \pi_k(\mathbf{x}) \text{Tr} \left\{ \sigma_k^2(\mathbf{x}) \mathbf{I} + (\boldsymbol{\mu}_k(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}]) (\boldsymbol{\mu}_k(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T \right\} \quad (122)$$

$$= \sum_{k=1}^K \pi_k(\mathbf{x}) \{ L \sigma_k^2(\mathbf{x}) + \|\boldsymbol{\mu}_k(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}]\|^2 \} \quad (123)$$

where L is the dimensionality of \mathbf{t} .

²²Such an unlabeled equation makes me upset because it is simply difficult to make reference. Called *Fisher's rule* (Mermin, 1989), it is a good practice to *number all displayed equations* (including those not referenced therein).

Page 279

Equation (5.165): We should add conditioning on α and β so that the left hand side reads $\ln p(\mathbf{w}|\mathcal{D}, \alpha, \beta)$.

Page 279

Equation (5.167): The conditioning on \mathcal{D} (or on any other variable) does not make sense for an approximate distribution $q(\cdot)$ unless properly defined. Hence, the conditioning for $q(\cdot)$ should be removed or the variables on which the posterior (5.164) is conditioned should instead be specified as parameters for $q(\cdot)$ so that the left hand side of (5.167) reads $q(\mathbf{w})$ or $q(\mathbf{w}; \mathcal{D}, \alpha, \beta)$, respectively.

Page 279

Equation (5.168): The equality ($=$) should be approximate (\approx). Also, we should again add conditioning on α and β . Thus, (5.168) should read

$$p(t|\mathbf{x}, \mathcal{D}, \alpha, \beta) \approx \int p(t|\mathbf{x}, \mathbf{w}, \beta) q(\mathbf{w}) d\mathbf{w} \quad (124)$$

where we have written the approximate posterior as $q(\mathbf{w})$.

Page 279

Equations (5.169) and (5.171): The superscripts **T** (in a bold typeface) should read T (in a roman typeface).

Page 279

Equation (5.172): The equality ($=$) should be approximate (\approx).

Page 289

Equation (5.208): The Kronecker delta δ_{jk} should read I_{jk} for consistency with other part of PRML; see, e.g., (5.95).

Page 295

Paragraph 1, Line 1: The vector \mathbf{x} should be a column vector so that $\mathbf{x} = (x_1, x_2)^T$.

Page 300

Paragraph –1, Line 4: Remove the comma (,) before the clause “which retain the...”

Pages 307 and 314

Equations (6.62) and (6.75): The Kronecker delta δ_{nm} should read I_{nm} for consistency with other part of PRML.

Page 318

Equations (6.93) and (6.94) as well as the text before (6.93): The text and the equations should read: We can evaluate the derivative of a_n^* with respect to θ_j by differentiating the relation (6.84) with respect to θ_j to give

$$\frac{\partial \mathbf{a}_N^*}{\partial \theta_j} = \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N) - \mathbf{C}_N \mathbf{W}_N \frac{\partial \mathbf{a}_N^*}{\partial \theta_j} \quad (125)$$

where the derivatives are Jacobians defined by (C.16) for a vector and analogously by (318) for a matrix. Rearranging (125) then gives

$$\frac{\partial \mathbf{a}_N^*}{\partial \theta_j} = (\mathbf{I} + \mathbf{C}_N \mathbf{W}_N)^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_j} (\mathbf{t}_N - \boldsymbol{\sigma}_N). \quad (126)$$

Page 319

Paragraph –1, Line –1: Insert a comma (,) before the clause “which breaks translation. . .”

Page 326

Paragraph –1, Line –3: Insert a comma (,) before the clause “who consider a. . .”

Page 333

Equation (7.29): $\frac{\partial L}{\partial \mathbf{w}} = 0$ should read $\nabla_{\mathbf{w}} L = 0$.

Page 335

Paragraph 1, Line 12: The term “protected conjugate gradients” should read “*projected* conjugate gradients.”

Page 341

Equation (7.57): $\frac{\partial L}{\partial \mathbf{w}} = 0$ should read $\nabla_{\mathbf{w}} L = 0$.

Page 349

Paragraph 1, Line 2: Remove the article “a” before “Gaussian processes.”

Page 354

Equation (7.112): The mean \mathbf{w}^* of the Laplace approximation to the posterior $p(\mathbf{w}|\mathbf{t}, \boldsymbol{\alpha})$ can only be obtained iteratively by, say, IRLS as described in the text so that (7.112) does not represent an explicit solution and is thus best removed.

As we shall see shortly, it is however useful to note that, at the convergence of IRLS, we have the following implicit equations for \mathbf{w}^*

$$\nabla_{\mathbf{w}} \ln p(\mathbf{w}^*|\mathbf{t}, \boldsymbol{\alpha}) = \Phi^T(\mathbf{t} - \mathbf{y}^*) - \mathbf{A}\mathbf{w}^* = 0 \quad (127)$$

where \mathbf{y}^* is \mathbf{y} evaluated at $\mathbf{w} = \mathbf{w}^*$ so that

$$\mathbf{y}^* = \mathbf{y}|_{\mathbf{w}=\mathbf{w}^*}. \quad (128)$$

Equation (7.113): Let us note that the precision (the inverse of the covariance Σ) of the Laplace approximation to the posterior $p(\mathbf{w}|\mathbf{t}, \alpha)$ is given by the Hessian of the negative log posterior evaluated at \mathbf{w}^* so that

$$\Sigma^{-1} = -\nabla_{\mathbf{w}} \nabla_{\mathbf{w}} \ln p(\mathbf{w}^*|\mathbf{t}, \alpha). \quad (129)$$

The covariance Σ should thus be given by

$$\Sigma = (\mathbf{A} + \Phi^T \mathbf{B}^* \Phi)^{-1} \quad (130)$$

where \mathbf{B}^* is \mathbf{B} evaluated at \mathbf{w}^* so that

$$\mathbf{B}^* = \mathbf{B}|_{\mathbf{w}=\mathbf{w}^*}. \quad (131)$$

Equation (7.117): The typeface of the vector \mathbf{y} in (7.117) should be that in (7.110), i.e., \mathbf{y} . Moreover, \mathbf{B} and \mathbf{y} should be those evaluated at $\mathbf{w} = \mathbf{w}^*$ so that $\hat{\mathbf{t}}$ is given by

$$\hat{\mathbf{t}} = \Phi \mathbf{w}^* + (\mathbf{B}^*)^{-1} (\mathbf{t} - \mathbf{y}^*). \quad (132)$$

It should also be noted here that we define $\hat{\mathbf{t}}$ as (132) in order that we can write the posterior mean \mathbf{w}^* in terms of $\hat{\mathbf{t}}$ so that

$$\mathbf{w}^* = \Sigma \Phi^T \mathbf{B}^* \hat{\mathbf{t}} \quad (133)$$

because we have

$$\Sigma \Phi^T \mathbf{B}^* \hat{\mathbf{t}} = \Sigma \Phi^T \mathbf{B}^* \Phi \mathbf{w}^* + \Sigma \Phi^T (\mathbf{t} - \mathbf{y}^*) \quad (134)$$

$$= \Sigma \Phi^T \mathbf{B}^* \Phi \mathbf{w}^* + \Sigma \mathbf{A} \mathbf{w}^* \quad (135)$$

$$= \Sigma (\mathbf{A} + \Phi^T \mathbf{B}^* \Phi) \mathbf{w}^* \quad (136)$$

$$= \mathbf{w}^* \quad (137)$$

where we have made use of (132), (127), and (130). We shall make use of (133) when we analyze the RVM classification problem (see below).

Equation (7.118): Although this marginal distribution cannot be obtained directly from the Laplace approximation (7.114) to the marginal $p(\mathbf{t}|\alpha)$ of the RVM classification problem, it can be shown to be an (approximate) marginal for a “linearized” version of the classification problem where $\hat{\mathbf{t}}$, given by (132), serves as the target (Tipping and Faul, 2003). Since the linearized problem can be regarded as an RVM regression problem having data-dependent precisions, we first review such a regression problem, after which we derive the marginal for the linearized classification problem.

RVM regression The likelihood for the RVM regression problem with data-dependent precisions $\beta = \{\beta_1, \dots, \beta_N\}$ is given by

$$p(\mathbf{t}|\mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi_n, \beta_n^{-1}) \quad (138)$$

$$= \mathcal{N}(\mathbf{t} | \Phi \mathbf{w}, \mathbf{B}^{-1}) \quad (139)$$

where we have omitted the conditioning on $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ to keep the notation uncluttered; and written $\phi_n = \phi(\mathbf{x}_n)$ and

$$\mathbf{t} = (t_1, \dots, t_N)^T \quad (140)$$

$$\Phi = (\phi_1, \dots, \phi_N)^T \quad (141)$$

$$\mathbf{B} = \text{diag}(\beta_1, \dots, \beta_N). \quad (142)$$

The prior is the same as (7.80) so that

$$p(\mathbf{w}|\alpha) = \prod_{i=1}^M \mathcal{N}(w_i | 0, \alpha_i^{-1}) \quad (143)$$

$$= \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{A}^{-1}) \quad (144)$$

where

$$\mathbf{w} = (w_1, \dots, w_M)^T \quad (145)$$

$$\mathbf{A} = \text{diag}(\alpha_1, \dots, \alpha_M). \quad (146)$$

The joint distribution is given by a linear-Gaussian model of the form

$$p(\mathbf{t}, \mathbf{w} | \alpha, \beta) = p(\mathbf{t} | \mathbf{w}, \beta) p(\mathbf{w} | \alpha) \quad (147)$$

$$= \mathcal{N}(\mathbf{t} | \Phi \mathbf{w}, \mathbf{B}^{-1}) \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{A}^{-1}). \quad (148)$$

Making use of the general results (2.115) and (2.116) for the marginal and the conditional Gaussians, we can readily evaluate the marginal and the posterior distributions again as Gaussians. The posterior is given by

$$p(\mathbf{w} | \mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w} | \mathbf{w}^*, \Sigma) \quad (149)$$

where

$$\mathbf{w}^* = \Sigma \Phi^T \mathbf{B} \mathbf{t} \quad (150)$$

$$\Sigma = (\mathbf{A} + \Phi^T \mathbf{B} \Phi)^{-1}. \quad (151)$$

The marginal is given by

$$p(\mathbf{t} | \alpha, \beta) = \mathcal{N}(\mathbf{t} | \mathbf{0}, \mathbf{C}) \quad (152)$$

where

$$\mathbf{C} = \mathbf{B}^{-1} + \Phi \mathbf{A}^{-1} \Phi^T. \quad (153)$$

RVM classification Let us now return to the RVM classification problem. We have already seen that the posterior can be approximated by the Laplace approximation so that

$$p(\mathbf{w}|\mathbf{t}, \boldsymbol{\alpha}) \approx \mathcal{N}(\mathbf{w}|\mathbf{w}^*, \boldsymbol{\Sigma}) \quad (154)$$

where the mean \mathbf{w}^* can be obtained by the IRLS algorithm as we have discussed; and the covariance $\boldsymbol{\Sigma}$ is given by (130).

Here, we note that \mathbf{w}^* can be written in the form (133). Comparing (133) with (150) and (130) with (151), we see that the Laplace approximation locally maps the classification problem to a regression problem with the data-dependent precision matrix \mathbf{B}^* , given by (131), where the target vector \mathbf{t} is replaced by the “linearized” target $\hat{\mathbf{t}}$, given by (132).

Assuming that the distribution over $\hat{\mathbf{t}}$ can be approximated by the Laplace approximation (as we have done in (154) for \mathbf{w}); and making use of the linear-Gaussian relation, we can obtain the corresponding marginal for the linearized problem in the form

$$p(\hat{\mathbf{t}}|\boldsymbol{\alpha}) \approx \mathcal{N}(\hat{\mathbf{t}}|\mathbf{0}, \mathbf{C}) \quad (155)$$

where

$$\mathbf{C} = (\mathbf{B}^*)^{-1} + \boldsymbol{\Phi} \mathbf{A}^{-1} \boldsymbol{\Phi}^T. \quad (156)$$

The right hand side of (155) takes the same form as the marginal (152) of the regression problem so that “we can apply the same analysis of sparsity and obtain the same fast learning algorithm” (Page 355, Paragraph –4, Line –3).

Page 355

Equation (7.119): Both \mathbf{A} and \mathbf{B} should be inverted and, moreover, \mathbf{B} should be that evaluated at $\mathbf{w} = \mathbf{w}^*$ so that (7.119) should read (156).

Page 357

Exercise 7.1, Line –2: Remove the second occurrence of “that.”

Page 386

Paragraph –2, Line –2: “involves” should be “involve.”

Page 390

Paragraph –1, Line 5: Insert a space before the sentence starting “Here the joint. . .” (*the third printing only*).

Page 399

Paragraph 3, Line 3: Remove the first occurrence of an indefinite article “a.”

Page 403

Paragraph 1, Line 2: The conjunction “and” in the phrase “. . . , and is equivalent to. . .” should be replaced by a relative pronoun “which.”

Page 411

Paragraph –2, Line –7: Insert “to” after “corresponding.”

Page 414

The caption of Figure 8.53, Line 6: The term “max-product” should be “max-sum.”

Page 417

Paragraph 2, Line 3: The clause “that can broadly be called *variational* methods” is restrictive so that the enclosing commas (,) should be removed.

Page 418

Paragraph 1, Line 5: “give” should be “gives.”

Page 424

Paragraph –2, Line –2: Remove the comma (,) after μ_k .

Page 425

Equation (9.3): The right hand side should be a zero vector $\mathbf{0}$ instead of a scalar zero 0.

Page 432

The text after (9.13): I would like to point out for clarity that the prior $p(\mathbf{z})$ given by (9.10) is a multinomial distribution or, more precisely, a *multinoulli* distribution (193) so that

$$p(\mathbf{z}) = \text{Mult}(\mathbf{z}|\boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_k}. \quad (157)$$

Moreover, we see that the posterior $p(\mathbf{z}|\mathbf{x})$ again becomes a multinoulli distribution of the form

$$p(\mathbf{z}|\mathbf{x}) = \text{Mult}(\mathbf{z}|\boldsymbol{\gamma}) = \prod_{k=1}^K \gamma_k^{z_k} \quad (158)$$

where we have written $\gamma_k \equiv \gamma(z_k)$. Called the *responsibility*, γ_k is given by (9.13), which can also be found directly by inspecting the functional form of the joint distribution

$$p(\mathbf{z}) p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\}^{z_k} \quad (159)$$

and noting that the multinoulli distribution (193) can be expressed in terms of unnormalized probabilities as shown in (194) where the normalized probabilities are given by (195). This observation helps the reader understand that evaluating the responsibilities γ_k indeed corresponds to the E step of the EM algorithm.

Page 433

The caption of Figure 9.5, Line –1: Add a period (.) at the end of the last sentence.

Page 434

Equation (9.15): Although the official errata document (Svensén and Bishop, 2011) states that σ_j in the right hand side should be raised to a power of D , the whole right hand side should be raised to D so that (9.15) should read

$$\mathcal{N}(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{(2\pi\sigma_j^2)^{D/2}}. \quad (160)$$

Page 435

Equation (9.16): The left hand side should be a zero vector $\mathbf{0}$ instead of a scalar zero 0.

Page 440

Paragraph 2, Lines 1 and 2: In order to be consistent with the expression “for each observation in \mathbf{X} ,” the phrase “the corresponding value of the latent variable \mathbf{Z} ” should read, e.g., “the value of the corresponding latent variable in \mathbf{Z} .”

Page 440

Paragraph 3, Line 5: Throughout the discussion, \mathbf{Z} denotes a set of latent variables. Hence, “the latent variable” should be “the latent variables.”

Page 450

Paragraph 1, Line –1: “maxmization” should be “maximization.”

Page 453

Paragraph 1: The old and new parameters should read θ^{old} and θ^{new} (without parentheses), respectively, as in (9.74) and the text.

Page 453

Paragraph 2, Line 4: “. . . , by marginalizing over the $\{\mathbf{z}_n\}$ we have. . . ” should read, e.g., “. . . , by marginalizing over the latent variables $\{\mathbf{z}_n\}$, we have. . . ”

Page 465

Equations (10.6) and (10.7): The lower bound of the form (10.6) for variational Bayes will be later recognized as “a negative Kullback-Leibler divergence between $q_j(\mathbf{Z}_j)$ and $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ ” (Page 465, Paragraph –1, Line –2). However, there is no point in taking the Kullback-Leibler divergence between two probability distributions over different sets of random variables; such a quantity is undefined. Moreover, the discussion here seems to be somewhat redundant. Without introducing an intermediate quantity like $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$, we can rewrite (10.6) and (10.7)

directly in terms of $q_j^*(\mathbf{Z}_j)$. Specifically, writing down the terms dependent on one of the factors $q_j(\mathbf{Z}_j)$, we obtain the lower bound $\mathcal{L}(q)$ in the form

$$\mathcal{L}(q) = \int q_j(\mathbf{Z}_j) \mathbb{E}_{\mathbf{Z} \setminus \mathbf{Z}_j} [\ln p(\mathbf{X}, \mathbf{Z})] d\mathbf{Z}_j - \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \text{const} \quad (161)$$

$$= -\text{KL}(q_j \| q_j^*) + \text{const} \quad (162)$$

where we have assumed that the expectation $\mathbb{E}_{\mathbf{Z} \setminus \mathbf{Z}_j}[\cdot]$ is taken with respect to \mathbf{Z} but \mathbf{Z}_j so that

$$\mathbb{E}_{\mathbf{Z} \setminus \mathbf{Z}_j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int \cdots \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (163)$$

and defined a new distribution $q_j^*(\mathbf{Z}_j)$ over \mathbf{Z}_j by the relation

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{\mathbf{Z} \setminus \mathbf{Z}_j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{const}. \quad (164)$$

It directly follows from (162) that, since the lower bound $\mathcal{L}(q)$ is the negative Kullback-Leibler divergence between $q_j(\mathbf{Z}_j)$ and $q_j^*(\mathbf{Z}_j)$ up to some additive constant, the maximum of $\mathcal{L}(q)$ occurs when $q_j(\mathbf{Z}_j) = q_j^*(\mathbf{Z}_j)$.

Page 465

The text before (10.8): The latent variable \mathbf{z}_i should read \mathbf{Z}_i .

Page 465

Paragraph –1, Line –1: If we adopt the representation (162), the factor $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ should read $q_j^*(\mathbf{Z}_j)$.

Page 466

Paragraph 1, Line 1: Again, $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ should read $q_j^*(\mathbf{Z}_j)$. The sentence “Thus we obtain. . .” should read, e.g., “Thus we see that we have already obtained a general expression for the optimal solution in (164).”

Page 466

Paragraph 3, Line –2: “bound” should be “the bound” (with the definite article). Also, the (lower) bound $\mathcal{L}(q)$ is concave (*not* convex) with respect to each of the factors $q_i(\mathbf{Z}_i)$. The concavity of the bound $\mathcal{L}(q)$ follows from the convexity (38) of the Kullback-Leibler divergence. To see this, recall that $\mathcal{L}(q)$ can be written in the form (162).

Page 467

The text before (10.12) and after (10.15): As Svensén and Bishop (2011) correct the left hand side of (10.12), we should write $q_1^*(z_1)$ instead of $q^*(z_1)$ and so on also in the text. Or, we should clarify that we simply write $q(\mathbf{z}_i)$ to denote the variational distribution over the latent variables \mathbf{z}_i in the same manner as the notation for the probability $p(\cdot)$ is “overloaded” by its argument(s).

Page 468

The text after (10.16): The constant term in (10.16) is the *negative* entropy of $p(\mathbf{Z})$.

Page 470

The text after (10.19): “zero forcing” should be “zero-forcing” (with hyphenation).

Page 470

The text after (10.23): “Gaussian-Gamma” should read “Gaussian-gamma” (without capitalization for “gamma”).

Page 478

Equation (10.63): The additive constant $+1$ on the right hand side should be omitted so that (10.63) should read

$$\nu_k = \nu_0 + N_k. \quad (165)$$

A quick check for the correctness of the re-estimation equations would be to consider the limit of $N \rightarrow 0$, in which the effective number of observations N_k also goes to zero and the re-estimation equations should reduce to identities. Equation (10.63) does not reduce to $\nu_k = \nu_0$, failing the test. Note that the solution for Exercise 10.13 given by [Svensén and Bishop \(2009\)](#) correctly derives the result (165).

Page 489

Equations (10.107) through (10.112): Some of the notations for the expectation are inconsistent with the one (1.36) employed in PRML; they should read $\mathbb{E}_{\mathbf{Z}}[\cdot]$ where \mathbf{Z} is replaced with the corresponding latent variables. For example, $\mathbb{E}_{\alpha} [\ln q(\mathbf{w})]_{\mathbf{w}}$ in the last line of (10.107) and $\mathbb{E} [\ln p(\mathbf{t}|\mathbf{w})]_{\mathbf{w}}$ in the left hand side of (10.108) should read $\mathbb{E}_{\mathbf{w}} [\ln q(\mathbf{w})]$ and $\mathbb{E}_{\mathbf{w}} [\ln p(\mathbf{t}|\mathbf{w})]$, respectively, where we have assumed that the expectation $\mathbb{E}_{\mathbf{Z}}[\cdot]$ is taken with respect to the variational distribution $q(\mathbf{Z})$.

Note however that we can safely omit the subscripts \mathbf{Z} of the expectations $\mathbb{E}_{\mathbf{Z}}[\cdot]$ here, as we have done in, e.g., (10.70), because the variables over which we take the expectations are clear; we take the expectations over all the latent variables when we calculate the lower bound. We only need to make the subscripts explicit when we find an optimal factor $q^*(\mathbf{Z}_i)$, in which case we take expectation selectively, that is, over all the latent variables but \mathbf{Z}_i ; see, e.g., (10.92) and (10.96).

Page 490

Paragraph –1, Line 2: Insert a comma (,) after the ellipsis (...).

Page 496

Equation (10.140): The differential operator d should be an upright d . Moreover, the derivative of x with respect to x^2 should be written with parentheses as $\frac{\mathrm{d}x}{\mathrm{d}(x^2)}$, instead of $\frac{\mathrm{d}x}{\mathrm{d}x^2}$, to avoid ambiguity.

Paragraph 1, Line –2: The sentence reads “Once [the right hand side of (10.152)] is normalized to give a variational posterior distribution $q(\mathbf{w})$, however, it no longer represents a bound.” The statement does not make sense because the right hand side of (10.152) is a lower bound in terms of the variational parameters ξ and thus not directly dependent on the variational distribution $q(\mathbf{w})$. Moreover, as we shall see shortly, we obtain the optimal solution for $q(\mathbf{w})$ by making use of the general result (168) for local variational Bayes, but *not* by normalizing the right hand side of (10.152). Therefore, this sentence is irrelevant and can be safely removed.

Pages 500 and 501

Equations (10.156) and (10.160): It is not very clear why the variational posterior is obtained in the form (10.156) and the variational parameters can be optimized by maximizing (10.160). This EM-like algorithm is not the same as *the* EM algorithm we have seen in Chapter 9; it can be derived by maximizing the lower bound (10.3) as follows. Note that the discussion here is similar to, but more general than, that of Section 10.6.3.

In a more general setting, we consider a local variational approximation to the joint distribution of the form

$$p(\mathbf{X}, \mathbf{Z}) \geq \tilde{p}(\mathbf{X}, \mathbf{Z}; \xi) \quad (166)$$

where ξ denotes the set of variational parameters, assuming that we can bound the likelihood $p(\mathbf{X}|\mathbf{Z}) \geq \tilde{p}(\mathbf{X}|\mathbf{Z}; \xi)$ or the prior $p(\mathbf{Z}) \geq \tilde{p}(\mathbf{Z}; \xi)$, or both. Then, we can again bound the lower bound (10.3) as

$$\mathcal{L}(q) \geq \tilde{\mathcal{L}}(q, \xi) \equiv \mathbb{E}_{\mathbf{Z}} [\ln \tilde{p}(\mathbf{X}, \mathbf{Z}; \xi)] - \mathbb{E}_{\mathbf{Z}} [\ln q(\mathbf{Z})] \quad (167)$$

where the expectation $\mathbb{E}_{\mathbf{Z}}[\cdot]$ is taken with respect to the variational distribution $q(\mathbf{Z})$. With much the same discussion as the derivation of the optimal solution (164) for the standard variational Bayesian method where we assume some appropriate factorization (10.5) for $q(\mathbf{Z})$, the optimal solution for the factor $q_j(\mathbf{Z}_j)$ that maximizes the lower bound $\tilde{\mathcal{L}}(q, \xi)$ can be obtained by the relation

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{\mathbf{Z} \setminus \mathbf{Z}_j} [\ln \tilde{p}(\mathbf{X}, \mathbf{Z}; \xi)] + \text{const} \quad (168)$$

which leads to the variational approximation to the posterior given by (10.156).

The optimization of the variational parameters ξ can be done by maximizing the first term of the lower bound $\tilde{\mathcal{L}}(q, \xi)$, i.e.,

$$\mathcal{Q}(\xi) = \mathbb{E}_{\mathbf{Z}} [\ln \tilde{p}(\mathbf{X}, \mathbf{Z}; \xi)] \quad (169)$$

which leads to the \mathcal{Q} function given by (10.160).

Page 501

The text after (10.162): We have that the variational parameter $\lambda(\xi)$ is a monotonic function of ξ for $\xi \geq 0$, but not that its derivative $\lambda'(\xi)$ is.

Page 503

The text after (10.168): A period (.) should be appended at the end of the sentence that follows (10.168).

Page 504

Paragraph 1, Line 1: In order to obtain the optimized variational distribution (10.174), we should use the optimal solution (168) for *local* variational Bayes. Note that the result (168) is different from the result (164), or (10.9), for standard variational Bayes in that (168) is given in terms of the lower bound $\tilde{p}(\mathbf{X}, \mathbf{Z}; \xi)$ to the joint distribution $p(\mathbf{X}, \mathbf{Z})$.

Page 504

Equation (10.177): The factor $a_N^{b_N}$ in the right hand side should be $b_N^{a_N}$ (it is probably safe to omit the right hand side at all because it is nothing but a gamma distribution with which the reader is fairly familiar).

Pages 511 and 512

Equations (10.212) and (10.213): It is helpful to note here that, if we employ the factors $\tilde{f}_n(\boldsymbol{\theta})$ of the form (10.213) where $n = 1, \dots, N$ together with $\tilde{f}_0(\boldsymbol{\theta}) = f_0(\boldsymbol{\theta}) = p(\boldsymbol{\theta})$ where $p(\boldsymbol{\theta})$ is given by (10.210), then we indeed obtain the approximate posterior $q(\boldsymbol{\theta})$ in the form (10.212). To see this, let us evaluate the product of all the factors $\tilde{f}_n(\boldsymbol{\theta})$ where $n = 0, 1, \dots, N$, giving

$$\prod_{n=0}^N \tilde{f}_n(\boldsymbol{\theta}) = f_0(\boldsymbol{\theta}) \prod_{n=1}^N \tilde{f}_n(\boldsymbol{\theta}) \quad (170)$$

$$= \mathcal{N}(\boldsymbol{\theta} | \mathbf{0}, b\mathbf{I}) \prod_{n=1}^N s_n \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_n, v_n \mathbf{I}) \quad (171)$$

$$= \frac{1}{(2\pi b)^{D/2}} \left[\prod_{n=1}^N \frac{s_n}{(2\pi v_n)^{D/2}} \right] \exp \left\{ -\frac{1}{2v} \|\boldsymbol{\theta} - \mathbf{m}\|^2 \right\} \exp \left(\frac{B}{2} \right) \quad (172)$$

where we have used the fact that an EP update leaves the factor $\tilde{f}_0(\boldsymbol{\theta})$ unchanged so that $\tilde{f}_0(\boldsymbol{\theta}) = f_0(\boldsymbol{\theta}) = p(\boldsymbol{\theta})$ holds (see Exercise 10.37); and defined v , \mathbf{m} , and B by

$$\frac{1}{v} = \frac{1}{b} + \sum_{n=1}^N \frac{1}{v_n} \quad (173)$$

$$\frac{\mathbf{m}}{v} = \sum_{n=1}^N \frac{\mathbf{m}_n}{v_n} \quad (174)$$

$$B = \frac{\mathbf{m}^T \mathbf{m}}{v} - \sum_{n=1}^N \frac{\mathbf{m}_n^T \mathbf{m}_n}{v_n}. \quad (175)$$

From (172), we see that the approximate posterior (10.203) is given by (10.212) where v and \mathbf{m} are given by (173) and (174); and also that the approximate model evidence (10.208) is

given by

$$p(\mathcal{D}) \approx \left(\frac{v}{b}\right)^{D/2} \exp\left(\frac{B}{2}\right) \left[\prod_{n=1}^N \frac{s_n}{(2\pi v_n)^{D/2}}\right] \quad (176)$$

where B is given by (175).

Page 512

Paragraph 2, Lines 1 and 3: $f_n(\boldsymbol{\theta})$ and $f_0(\boldsymbol{\theta})$ should read $\tilde{f}_n(\boldsymbol{\theta})$ and $\tilde{f}_0(\boldsymbol{\theta})$, respectively.

Page 512

Equation (10.222): The factor $(2\pi v_n)^{D/2}$ in the denominator of the right hand side of (10.222) should be removed because it has been already included in the normalization constant of the Gaussian in the approximate factors (10.213).²³

Page 513

Equation (10.223): The right hand side should read that of (176) where v and \mathbf{m} are replaced by v^{new} and \mathbf{m}^{new} , respectively.²⁴

Page 513

Equation (10.224): v should read v^{new} for consistency with (10.223).

Page 515

Equations (10.228) and (10.229): Although [Svensén and Bishop \(2011\)](#) correct (10.228) so that $q^{\setminus b}(\mathbf{x})$ is a normalized distribution, we do not need the normalization of $q^{\setminus b}(\mathbf{x})$ here and, even with this normalization, we cannot ensure that $\hat{p}(\mathbf{x})$ given by (10.229) is normalized. Similarly to (10.195), we can proceed with the unnormalized $q^{\setminus b}(\mathbf{x})$ given by the original (10.228) and, rather than correcting (10.228), we should correct (10.229) so that

$$\hat{p}(\mathbf{x}) \propto q^{\setminus b}(\mathbf{x}) f_b(x_2, x_3) = \dots \quad (177)$$

implying that $\hat{p}(\mathbf{x})$ is a normalized distribution.

Page 515

The text after (10.229): The new distribution $q^{\text{new}}(\mathbf{z})$ should read $q^{\text{new}}(\mathbf{x})$.

Page 516

Equation (10.240): The subscript k of the product $\prod_k \dots$ should read $k \neq j$ because we have already removed the term $\tilde{f}_j(\boldsymbol{\theta}_j)$.

²³Note that, in PRML, we use the approximate factors (10.213) slightly different from those used by [Minka \(2001\)](#).

²⁴One might notice that the approximate evidence derived in [Minka \(2001\)](#) looks more like the original (10.223); this is however due to the different definition for the factors (10.213) and the fact that the product begins from $n = 0$ (not $n = 1$) in [Minka \(2001\)](#).

Page 526

Equation (11.6): The transformation $f : (0, 1) \rightarrow (-\infty, \infty)$ between the random variables z and y , which is, of course, bijective as is assumed in (5), is also assumed to be monotonically increasing in (11.6) so that $p(y \in (-\infty, y_0)) = p(z \in (0, z_0))$ where $y_0 = f(z_0)$.

Page 528

Paragraph –2, Lines 1, 2, and 4: z should be z for consistency with (11.13).

Page 539

The caption of Figure 11.9: Insert “the” before “Metropolis algorithm.”

Page 542

Paragraph 1, Line –6: “steps sizes” should be “step sizes.”

Page 542

Paragraph 1, Line –1: “Metropolis Hastings” should read “Metropolis-Hastings” (with hyphenation) for consistency.

Pages 554 and 555

Equation (11.72), Line –2 and the text after (11.72): The expectation in the last line but one of (11.72) is taken with respect to the probability $p_G(\mathbf{z})$. This is probably better expressed in words, rather than the unclear notation like $\mathbb{E}_{G(\mathbf{z})}[\cdot]$. Specifically, the expectation should read

$$\mathbb{E}_{\mathbf{z}} [\exp (-E(\mathbf{z}) + G(\mathbf{z}))] \quad (178)$$

where we have written the argument \mathbf{z} for $E(\mathbf{z})$ and $G(\mathbf{z})$ for clarity; and the text following (11.72) should read “where $\mathbb{E}_{\mathbf{z}}[\cdot]$ is taken with respect to $p_G(\mathbf{z})$ and $\{\mathbf{z}^{(l)}\}$ are samples drawn from the distribution defined by $p_G(\mathbf{z})$.”

Page 555

Paragraph 3, Line 2: The term “importance-sampling distribution” (with hyphenation) is inconsistent with “importance sampling distribution” (without hyphenation) found in other part of PRML (e.g., Paragraph 2 on the same page).

Page 556

Exercise 11.7, Line 1: The interval should be $[-\pi/2, \pi/2]$ instead of $[0, 1]$.

Page 557

Exercise 11.14, Line 2: The variance should be σ_i^2 instead of σ_i .

Page 563

Equation (12.7): The Kronecker delta δ_{ij} should read I_{ij} for consistency with other part of PRML.

Page 564

The text after (12.12): The derivative we consider here is that with respect to b_j (*not* that with respect to b_i).

Page 564

Paragraph –1, Line 2: $\mathbf{u}_i = 0$ should read $\mathbf{u}_i = \mathbf{0}$ (the right hand side should be a zero vector $\mathbf{0}$).

Page 575

Paragraph –2, Line 5: The zero vector should be a row vector instead of a column vector so that we have $\mathbf{v}^T \mathbf{U} = \mathbf{0}^T$. Or, the both sides are transposed to give $\mathbf{U}^T \mathbf{v} = \mathbf{0}$.

Page 578

Equation (12.53): As stated in the text preceding (12.53), we should substitute $\boldsymbol{\mu} = \bar{\mathbf{x}}$ into (12.53).

Page 578

The text before (12.56): For the maximization with respect to \mathbf{W} , we use (C.25) and (C.27) instead of (C.24).

Page 579

Paragraph 1, Line 5: The eigendecomposition requires $O(D^3)$ computations (in the plural form).

Page 587

The text before (12.75): The phrase “the eigenvector equations tells us. . .” should read, e.g., “the eigenvector equation (12.74) tells us. . .” if we see (12.74) as one equation as a whole.

Page 588

Paragraph –2, Lines 1 and 2: The term “principal component projections” is inconsistent with “principal components projection” (with “components” in the plural form) found in Paragraph 2, Line –1. Probably, we should write “principal components projection(s).”

Page 599

Exercise 12.1, Line –1: The quantity λ_{M+1} is an eigenvalue (not an eigenvector).

Page 599

The first line before (12.93): Either remove the comma or add another one after “notation.”

Page 602

Exercise 12.25, Line 2: The latent space distribution should read $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$.

Page 610

Paragraph 1, Line –4: The text “our predictions for \mathbf{x}_{n+1} depends on. . .” should read: “our predictions for \mathbf{x}_{n+1} depend on. . .” (Remove the trailing ‘s’ from the verb).

Page 616

Equation (13.15): The summation should run over \mathbf{z}_n in the rightmost expression so that it reads

$$\sum_{\mathbf{z}_n} \gamma(\mathbf{z}_n) z_{nk} \quad (179)$$

(the third printing only).²⁵

Page 617

Paragraph 1, Line 1: Remove the space preceding the comma.

Page 619

Paragraph 3, Line –1: “. . . because these fixed throughout” should read “. . . because these are fixed throughout.”

Page 620

Paragraph –1, Line 4 and the following (unlabeled) equation: The last sentence before the equation and the equation should each be terminated with a period (.).

Page 621

Paragraph 1, Line –2: “scaled” should read “scales.”

Pages 621 and 622

Figures 13.12 and 13.13: It should be clarified that, similarly to the notations $\alpha(z_{nk})$ and $\beta(z_{nk})$ defined in Section 13.2.2, the notation $p(\mathbf{x}_n|z_{nk})$ denotes the value of $p(\mathbf{x}_n|\mathbf{z}_n)$ when $z_{nk} = 1$ so that

$$p(\mathbf{x}_n|z_{nk}) \equiv p(\mathbf{x}_n|z_{nk} = 1). \quad (180)$$

²⁵Note that this erratum is taken from the official errata (Svensén and Bishop, 2011) for the first and the second printings, which is however missing in the errata for the third printing.

Page 622

Paragraph –1, Line 1: “M step equations” should read “M-step equations” (with hyphenation for the adjectival term “M-step”) for consistency with the following line as well as other part of PRML.

Page 622

Equation (13.40): The summations should read $\sum_{n=1}^N$.

Page 623

Paragraph 1, Line –2: z_{nk} should read $z_{n-1,k}$.

Page 627

Paragraph 1, Line –2: “send” should read “sent.”

Page 627

Paragraph 3, Line 2: “forward backward algorithm” should read “forward-backward algorithm” for consistency (see Section 13.2.2).

Page 630

The caption of Figure 13.16: $p(\mathbf{x}_n|k)$ should read $p(\mathbf{x}_n|z_{nk})$ where we have used (180).

Page 631

Equation (13.73): The equation should read

$$\sum_{r=1}^R \ln \left\{ \frac{p(\mathbf{X}_r|\boldsymbol{\theta}_{m_r}) p(m_r)}{\sum_{l=1}^M p(\mathbf{X}_r|\boldsymbol{\theta}_l) p(l)} \right\}. \quad (181)$$

Page 635

Paragraph 1, Line –3: “algorithms” should read “algorithm.”

Page 637

Equations (13.81), (13.82), and (13.83): The distribution (13.81) over \mathbf{w} should read

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \boldsymbol{\Gamma}) \quad (182)$$

and so on.

Page 638

Paragraph 1, Line 2: “conditional on” should read “conditioned on.”

Page 641

Equation (13.104) and the preceding text: The form of the Gaussian is unclear. Since a multivariate Gaussian is usually defined over a column vector, we should construct a column vector from the concerned random variables to clearly define the mean and the covariance. Specifically, (13.104) and the preceding text should read: ... we see that $\xi(\mathbf{z}_{n-1}, \mathbf{z}_n)$ is a Gaussian of the form

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = \mathcal{N} \left(\begin{pmatrix} \mathbf{z}_{n-1} \\ \mathbf{z}_n \end{pmatrix} \middle| \begin{pmatrix} \hat{\boldsymbol{\mu}}_{n-1} \\ \hat{\boldsymbol{\mu}}_n \end{pmatrix}, \begin{pmatrix} \hat{\mathbf{V}}_{n-1} & \hat{\mathbf{V}}_{n-1,n} \\ \hat{\mathbf{V}}_{n-1,n}^T & \hat{\mathbf{V}}_n \end{pmatrix} \right) \quad (183)$$

where the mean $\hat{\boldsymbol{\mu}}_n$ and the covariance $\hat{\mathbf{V}}_n$ of \mathbf{z}_n are given by (13.100) and (13.101), respectively; and the covariance $\hat{\mathbf{V}}_{n-1,n}$ between \mathbf{z}_{n-1} and \mathbf{z}_n is given by

$$\hat{\mathbf{V}}_{n-1,n} = \text{cov} [\mathbf{z}_{n-1}, \mathbf{z}_n] = \mathbf{J}_{n-1} \hat{\mathbf{V}}_n. \quad (184)$$

Pages 642 and 643

Equation (13.109) and the following equations: If we follow the notation in Chapter 9, the typeface of the Q function should be \mathcal{Q} .

Page 642

Equation (13.109): If we follow the notation for the conditional expectation (1.37), the Q function (13.109) should read

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}] \quad (185)$$

$$= \int d\mathbf{Z} p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) \quad (186)$$

which corresponds to (9.30).

Page 643

Equation (13.111): $\mathbf{V}_0^{\text{new}}$ should read $\mathbf{P}_0^{\text{new}}$. [Svensén and Bishop \(2011\)](#) have failed to mention (13.111).

Page 643

Equation (13.114): The size of the opening curly brace “{” should match that of the closing curly brace “}.”

Page 647

The caption of Figure 13.23, Line -1: $p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}^{(l)})$ should read $p(\mathbf{x}_{n+1} | z_{n+1}^{(l)})$.

Page 649

Exercise 13.14, Line 1: (8.67) should be (8.64).

Page 650

Equations (13.127) and (13.128): The equal signs should be aligned.

Page 651

Exercises 13.25 through 13.28: A zero matrix is denoted by \mathbf{O} (*not* by $\mathbf{0}$ nor 0) so that we should write $\mathbf{A} = \mathbf{O}$ and so on.

Page 651

Exercises 13.29: “backwards recursion” should read “backward recursion” for consistency.

Page 657

Paragraph –2, Line 1: Insert a comma before “such as” or remove the comma that follows.

Page 658

The equation at the bottom of Figure 14.1: The subscript of the summation in the right hand side should read $m = 1$.

Page 659

Paragraph 2, Line 3: “learners” should read “learner.”

Page 659

Paragraph 2, Line 4: “stumps” should read “stump.”

Page 666

Paragraph –1, Line 2: “mixtures” should read “mixture.”

Page 668

Equation (14.37): The arguments of the probability are notationally inconsistent with those of (14.34), (14.35), and (14.36). Specifically, the conditioning on ϕ_n should read that on t_n and the probability $p(k|\dots)$ be the value of $p(\mathbf{z}_n|\dots)$ when $z_{nk} = 1$, which we write $p(z_{nk} = 1|\dots)$. Moreover, strictly speaking, the old parameters $\pi_k, \mathbf{w}_k, \beta$ should read $\pi_k^{\text{old}}, \mathbf{w}_k^{\text{old}}, \beta^{\text{old}} \in \boldsymbol{\theta}^{\text{old}}$. In order to solve these problems, we should rewrite (14.37) as, for example,

$$\gamma_{nk} = \mathbb{E} [z_{nk} | t_n, \boldsymbol{\theta}^{\text{old}}] \quad (187)$$

where we have written the conditioning in the expectation explicitly and the expectation is given by

$$\mathbb{E} [z_{nk} | t_n, \boldsymbol{\theta}] = p(z_{nk} = 1 | t_n, \boldsymbol{\theta}) = \frac{\pi_k \mathcal{N}(t_n | \mathbf{w}_k^T \phi_n, \beta^{-1})}{\sum_j \pi_j \mathcal{N}(t_n | \mathbf{w}_j^T \phi_n, \beta^{-1})}. \quad (188)$$

Page 668

The unlabeled equation between (14.37) and (14.38): If we write the implicit conditioning in the expectation explicitly (similarly to the above equations), the unlabeled equation should read

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{t}, \mathbf{Z} | \boldsymbol{\theta}) | \mathbf{t}, \boldsymbol{\theta}^{\text{old}}] \quad (189)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left\{ \ln \pi_k + \ln \mathcal{N}(t_n | \mathbf{w}_k^T \boldsymbol{\phi}_n, \beta^{-1}) \right\}. \quad (190)$$

Page 669

Equations (14.40) and (14.41): The left hand sides should both read a zero vector $\mathbf{0}$ instead of a scalar zero 0.

Page 669

Equation (14.41): Φ is undefined. The text following (14.41) should read for example: where $\mathbf{R}_k = \text{diag}(\gamma_{nk})$ is a diagonal matrix of size $N \times N$ and $\Phi = (\phi_1, \dots, \phi_N)^T$ is an $N \times M$ matrix. Here, N is the size of the data set and M is the dimensionality of the feature vectors ϕ_n .

Page 669

Equation (14.43): “+const” should be added to the right hand side.

Page 671

The text after (14.46): The text should read: “where we have omitted the dependence on $\{\phi_n\}$ and defined $y_{nk} = \dots$ ” Or, ϕ should have been omitted from the left hand side of (14.45) in the first place.

Page 671

Equation (14.48): The notation should be corrected similarly to (187) and (188).

Page 671

Equation (14.49): The notation should be corrected similarly to (189).

Page 672

Equation (14.52): The negation should be removed so that the Hessian is given by $\mathbf{H}_k \equiv \nabla_k \nabla_k \mathcal{Q}$ where

$$\nabla_k \nabla_k \mathcal{Q} = - \sum_{n=1}^N \gamma_{nk} y_{nk} (1 - y_{nk}) \phi_n \phi_n^T. \quad (191)$$

Page 674

Exercise 14.1, Line 1: “of” should be inserted after “set.”

Page 685

Paragraph –1, Line 3: We assume in this report that $\mu \in (0, 1)$. See (45).

Page 686

Paragraph 1, Line 1: We assume in this report that $\mu \in (0, 1)$. See (45).

Page 686

Equation (B.9): The mode (B.9) of the beta distribution exists “if $a > 1$ and $b > 1$.”

Page 686

Paragraph 1, Line –3: Since we assume in this report that $\mu \in (0, 1)$, we have no singularity. See (45).

Page 686

Paragraph –1, Line 3: We assume in this report that $\mu \in (0, 1)$. See (45).

Page 686

Paragraph –1, Line –3: The comma in the first inline math should be removed so that the product reads: $m \times (m - 1) \times \cdots \times 2 \times 1$.

Page 687

Equation (B.15) and the preceding text: We assume in this report that $\mu_k \in (0, 1)$ for all k . See (45).

Page 687

Equation (B.19): Insert a comma (,) between μ_j and μ_k ; and also add the condition $j \neq k$ so that (B.19) reads

$$\text{cov} [\mu_j, \mu_k] = -\frac{\alpha_j \alpha_k}{\hat{\alpha}^2 (\hat{\alpha} + 1)}, \quad j \neq k. \quad (192)$$

Page 687

Equation (B.20): The mode (B.20) of the Dirichlet exists “if $\alpha_k > 1$ for all k .”

Page 687

Equation (B.25): The differential operator d should be an upright d.

Page 687

Paragraph –1, Line –1: Since we assume in this report that $\mu_k \in (0, 1)$ for all k , we have no singularity as with the beta distribution. See (45).

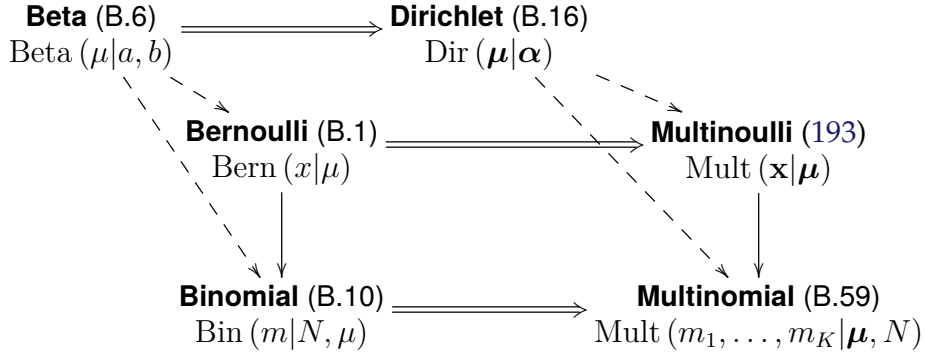


Figure 5 The relationship between discrete distributions and their conjugate priors. Here, “ $A \Rightarrow B$ ” (respectively, “ $A \rightarrow B$ ”) denotes “ A generalizes to B with multiple categories (observations) concerned”; and “ $A \dashrightarrow B$ ” denotes “ A is the conjugate prior for B .” Note also that, as this diagram suggests, there are some inconsistencies in parameterization; it is probably better for consistency to write the binomial as $\text{Bin}(m|\mu, N)$ instead of $\text{Bin}(m|N, \mu)$ and the multinomial as $\text{Mult}(\mathbf{m}|\mu, N)$ instead of $\text{Mult}(m_1, \dots, m_K|\mu, N)$ where $\mathbf{m} = (m_1, \dots, m_K)^T$, for example.

Page 688

Paragraph 1, Line 1: “Gamma” should read “gamma” (without capitalization).

Page 689

Paragraph 1, Line 1: “positive-definite” should read “positive definite” (without hyphenation).

Page 689

Equation (B.49): x in the right hand side should read x_a .

Page 690

Equation (B.52): μ_o in the right hand side should read μ_0 (the subscript should be a zero 0).

Page 690

Equation (B.54): The discrete distribution of the form (B.54), or (2.26), is known as the *categorical* or the *multinoulli* distribution (Murphy, 2012). It is also sometimes called, less precisely, the “multinomial” or the “discrete” distribution. Of these terms, I would prefer the term *multinoulli* because it naturally suggests that it is a generalization of the *Bernoulli* distribution (B.1) to multiple categories $K > 2$ and also a special case of the *multinomial* distribution (B.59) where we have only a single observation $N = 1$. Since we often make use of this discrete distribution, we shall introduce some notation for the right hand side of (B.54). See Figure 5 for the relationship between the discrete distributions found in PRML.

Multinoulli distribution The *multinoulli* distribution is a distribution over the K -dimensional binary variable $\mathbf{x} = (x_1, \dots, x_K)^T$ where $x_k \in \{0, 1\}$ such that $\sum_k x_k = 1$, i.e., we employ the

one-of- K coding scheme for \mathbf{x} . Here, we “overload” the notation (B.59) for the multinomial and write the multinoulli as

$$\text{Mult}(\mathbf{x}|\boldsymbol{\mu}) \equiv \prod_{k=1}^K \mu_k^{x_k} = \exp \left\{ \sum_{k=1}^K x_k \ln \mu_k \right\} \quad (193)$$

where the parameter $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^T$ consists of (normalized) probabilities $\mu_k \in (0, 1)$ such that $\sum_k \mu_k = 1$.

When we identify a multinoulli distribution from its functional form, e.g., in the posterior distribution (158) for the Gaussian mixture model (159) of Section 9.2, one will find it helpful to know that the multinoulli distribution (193) can also be expressed in terms of unnormalized probabilities $\tilde{\mu}_k > 0$, i.e.,

$$\text{Mult}(\mathbf{x}|\boldsymbol{\mu}) \propto \prod_{k=1}^K \tilde{\mu}_k^{x_k} = \exp \left\{ \sum_{k=1}^K x_k \ln \tilde{\mu}_k \right\} \quad (194)$$

where the normalized probabilities μ_k can be found by

$$\mu_k = p(x_k = 1) = \frac{\tilde{\mu}_k}{\sum_j \tilde{\mu}_j}. \quad (195)$$

Page 690

Equation (B.57): Insert a comma (,) between x_j and x_k so that the left hand side of (B.57) reads $\text{cov}[x_j, x_k]$.

Page 691

Paragraph 1, Line 2: We assume in this report that $\mu_k \in (0, 1)$ for all k . See (45).

Page 691

Paragraph 2, Line –3: We assume in this report that $\mu_k \in (0, 1)$ for all k . See (45).

Page 691

Equations (B.59) and (B.63): The *multinomial coefficient* (B.63) should read (46).

Page 691

Equation (B.62): Insert a comma (,) between m_j and m_k so that the left hand side of (B.62) reads $\text{cov}[m_j, m_k]$.

Page 691

The icon for Student’s t-distribution: As we have seen in the erratum for Figure 2.15, the tails of the t-distributions are too high. Figure 4 gives the correct plot.

Equation (B.68): This form of multivariate Student's t-distribution is derived in Section 2.3.7 by marginalizing over the gamma distributed (scalar) variable η in (2.161), but *not* by marginalizing over the $D \times D$ precision matrix Λ that is governed by the Wishart distribution $\mathcal{W}(\Lambda|\mathbf{W}, \nu)$ where $\mathbf{W} \succ 0$ and $\nu > D - 1$, which results in a marginal distribution of the form

$$p(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}, \nu) = \int d\Lambda \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Lambda^{-1}) \mathcal{W}(\Lambda|\mathbf{W}, \nu) d\Lambda. \quad (196)$$

The above marginal (196) is indeed equivalent to (B.68) with some reparameterization. However, this result is not so obvious that I would like to show it here. Note that such marginalization is also used to derive a mixture of Student's t-distributions given by (10.81) in Exercise 10.19.

Multivariate Student's t-distribution as a marginal over Wishart The key idea to evaluating the right hand side of (196) is that the integrand can be identified as an unnormalized Wishart distribution and the marginalization can be done in a symbolic manner. More specifically, we have

$$p(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}, \nu) = \int d\Lambda \frac{|\Lambda|^{1/2}}{(2\pi)^{D/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Lambda (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ \times B(\mathbf{W}, \nu) |\Lambda|^{(\nu-D-1)/2} \exp \left\{ -\frac{1}{2} \text{Tr}(\mathbf{W}^{-1} \Lambda) \right\} \quad (197)$$

$$= \frac{2^{(\nu+1)D/2} \Gamma_D\left(\frac{\nu+1}{2}\right) \left| \mathbf{W}^{-1} + (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \right|^{-(\nu+1)/2}}{(2\pi)^{D/2} 2^{\nu D/2} \Gamma_D\left(\frac{\nu}{2}\right) |\mathbf{W}|^{\nu/2}} \quad (198)$$

where we have used the fact that the Wishart distribution (B.78) is correctly normalized (which will be shown later); and introduced the *multivariate gamma function* $\Gamma_D(\cdot)$ given by (252), by which we can simplify the normalization constant $B(\mathbf{W}, \nu)$ in the form (253). Finally, we obtain

$$p(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}, \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu+1-D}{2}\right)} \frac{|\mathbf{W}|^{1/2}}{\pi^{D/2}} \left[1 + (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{W} (\mathbf{x} - \boldsymbol{\mu}) \right]^{-(\nu+1)/2} \quad (199)$$

where we have used (252) and (C.15). Thus, we see that the marginal distribution of the form (196) is equivalent to the multivariate Student's t-distribution of the form (B.68) or (2.162); they are related by

$$p(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}, \nu) = \text{St}(\mathbf{x}|\boldsymbol{\mu}, (\nu + 1 - D)\mathbf{W}, \nu + 1 - D). \quad (200)$$

If the scale matrix is isotropic, which is common in practice, so that $\mathbf{W} = \widetilde{W}\mathbf{I}$ where $\widetilde{W} > 0$, then the resulting multivariate Student's t-distribution (200) is again isotropic. The same marginal distribution can also be obtained by marginalizing with respect to a univariate Wishart (gamma) prior so that

$$\int_0^\infty \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \widetilde{\lambda}^{-1}\mathbf{I}) \mathcal{W}(\widetilde{\lambda}|\widetilde{W}, \widetilde{\nu}) d\widetilde{\lambda} = \text{St}(\mathbf{x}|\boldsymbol{\mu}, \widetilde{\nu}\widetilde{W}\mathbf{I}, \widetilde{\nu}) \quad (201)$$

where $\tilde{\nu} = \nu + 1 - D > 0$. Note that the “covariance” parameter (79) of the corresponding multivariate Wishart prior $\mathcal{W}(\Lambda|\mathbf{W}, \nu)$ for which we obtain the same marginal (201) is however *not* equal to $\tilde{\sigma}^2 \mathbf{I}$ where $\tilde{\sigma}^2 = (\tilde{\nu} \tilde{W})^{-1}$ is the “covariance” parameter of the univariate Wishart prior $\mathcal{W}(\tilde{\lambda}|\tilde{W}, \tilde{\nu})$, but is given by

$$\Sigma = (\nu \mathbf{W})^{-1} = \frac{\tilde{\nu}}{\tilde{\nu} - 1 + D} \tilde{\sigma}^2 \mathbf{I}. \quad (202)$$

So far, we have observed that a marginal distribution of the form (196) where the marginalization is taken over a matrix-valued random variable Λ is equivalent to a marginal of the form (2.161) or, if the scale matrix is isotropic, of the form (201) where the marginalization is over a scalar random variable η or $\tilde{\lambda}$, respectively. Given that those marginals reduce to an identical multivariate Student’s t-distribution (with some reparameterization), we now have a natural question: *Which form of marginal is better than the other?* I would argue that a marginal with fewer latent variables, i.e., (2.161) or (201), is always better than a marginal with more latent variables, i.e., (196), because fewer latent variables imply less computational space and complexity as well as a tighter bound on the (marginal) likelihood and thus faster convergence when we infer a model involving such marginals with the EM algorithm (see Chapter 9) or variational methods (Chapter 10). Moreover, the marginal of the form (2.161) enjoys even greater modeling flexibility in that it allows us to learn the mean μ and the precision Λ parameters with, e.g., maximum likelihood (see Exercise 12.24) or variational Bayes by introducing a (conditionally) conjugate prior for μ and Λ (Svensén and Bishop, 2005).

Page 692

Paragraph –1, Line 2: Since $\text{Beta}(\mu|1, 1) \equiv \text{U}(\mu|0, 1)$, I would prefer to write the domain of the uniform distribution $\text{U}(x|a, b)$ where $a < b$ as $x \in (a, b)$ for consistency. See also (45).

Page 693

Equations (B.78) through (B.82): Some appropriate citation, e.g., Anderson (2003), is necessary for the Wishart distribution (B.78) or (2.155), which is first introduced in Section 2.3.6 as the conjugate prior for the precision matrix Λ of the multivariate Gaussian $\mathcal{N}(\mathbf{x}|\mu, \Lambda^{-1})$ where the mean μ is assumed to be known (see also Exercise 2.45), because no proof has been given for the normalization constant (B.79) or (2.156). Furthermore, the expectations (B.80) and (B.81) as well as the entropy (B.82), of which we make use in Section 10.2, have not been shown either.

Note however that most multivariate statistics textbooks, including Anderson (2003), motivate the Wishart distribution differently from PRML; they typically introduce the Wishart distribution as the distribution over a symmetric positive-semidefinite matrix (called the *scatter matrix*) of the form

$$\mathbf{S} = \sum_{n=1}^{\nu} \mathbf{x}_n \mathbf{x}_n^T \quad (203)$$

where $\mathbf{x}_1, \dots, \mathbf{x}_\nu$ are samples that have been drawn independently from a zero-mean multivariate Gaussian $\mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma)$. More specifically, it can be shown that the distribution

over the matrix \mathbf{S} is given by

$$p(\mathbf{S}) = \mathcal{W}(\mathbf{S}|\boldsymbol{\Sigma}, \nu) \quad (204)$$

if $\nu \geq D$ where D is the dimensionality of \mathbf{S} .²⁶

The derivation of the Wishart distribution along this line is indirect for our purpose (we are mainly interested in its conjugacy). In the following, I would instead like to show the normalization (B.79) as well as the expectations (B.80) and (B.81) directly just as we have done for the gamma distribution (2.146). To this end, we first introduce some notation for subsets of the space of square matrices such that all the eigenvalues are positive, after which we review an important matrix factorization method called the *Cholesky decomposition* as well as the associated Jacobian. We also introduce the *multivariate gamma function*, which simplifies the form of the normalization constant (B.79). The normalization of the Wishart distribution can be shown through a change of variables similar to the one that we apply for evaluating the multivariate gamma function. The expectations (B.80) and (B.81) are shown by making use of the general identity (58).

Spaces of square matrices with positive eigenvalues To facilitate our discussion, we define two subsets of the space $\mathbb{R}^{D \times D}$ of real square matrices of dimensionality D . Let $\mathcal{S}_+^D \subset \mathbb{R}^{D \times D}$ be the space of symmetric positive-definite matrices of dimensionality D ; and $\mathcal{U}_+^D \subset \mathbb{R}^{D \times D}$ be the space of *upper triangular* matrices of dimensionality D with strictly positive diagonal elements. Here, an upper triangular matrix $\mathbf{U} = (U_{ij})$ is a square matrix such that $U_{ij} \equiv 0$ where $i > j$ so that

$$\mathbf{U} = \begin{pmatrix} U_{11} & U_{12} & \cdots & U_{1D} \\ 0 & U_{22} & \cdots & U_{2D} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & U_{DD} \end{pmatrix} \quad (205)$$

where D is the dimensionality of \mathbf{U} . Similarly, a *lower triangular* matrix $\mathbf{L} = (L_{ij})$ is a square matrix such that $L_{ij} \equiv 0$ where $i < j$ or, equivalently, \mathbf{L}^T is upper triangular.

If $\mathbf{U} \in \mathcal{U}_+^D$, then $\mathbf{U} = (U_{ij})$ is an upper triangular matrix of dimensionality D such that $U_{ii} > 0$ for all i . Note also that $\mathbf{A} \in \mathcal{S}_+^D$ is equivalent to saying “ \mathbf{A} is a symmetric positive-definite matrix of dimensionality D .” Positive definiteness of \mathbf{A} , or $\mathbf{A} \succ 0$, implies by definition that

$$\mathbf{x}^T \mathbf{A} \mathbf{x} > 0 \quad (206)$$

for any $\mathbf{x} \neq \mathbf{0}$.

The plus sign (+) in the subscripts of \mathcal{S}_+^D and \mathcal{U}_+^D indicates that any matrix in either of the sets is such that all the eigenvalues are strictly positive. To see this for $\mathbf{A} \in \mathcal{S}_+^D$, consult the discussion regarding the *eigenvalue decomposition* of real symmetric matrices found in Appendix C of PRML; for $\mathbf{U} \in \mathcal{U}_+^D$, consider its *characteristic equation* (C.30), from which it readily follows that the eigenvalues of $\mathbf{U} = (U_{ij})$ are equal to its diagonal elements $U_{ii} > 0$.

The above observation regarding the eigenvalues of $\mathbf{U} = (U_{ij}) \in \mathcal{U}_+^D$ implies that the (absolute) determinant of \mathbf{U} is given by the product of its diagonal elements $U_{ii} > 0$, i.e.,

$$|\mathbf{U}| = \prod_{i=1}^D U_{ii} \quad (207)$$

²⁶It follows from the definition (203) of \mathbf{S} that the mean is given by $\mathbb{E}[\mathbf{S}] = \nu \boldsymbol{\Sigma}$, showing the identity (B.80) when ν is an integer such that $\nu \geq D$.

which can also be shown directly from the definition (C.10) of the determinant.²⁷

Triangular matrix groups It is worth noting here that \mathcal{U}_+^D forms a *group* with respect to matrix multiplication, i.e., (i) for any pair of the elements $\mathbf{U}, \mathbf{U}' \in \mathcal{U}_+^D$, their product is again in the group \mathcal{U}_+^D so that $\mathbf{U}\mathbf{U}' \in \mathcal{U}_+^D$; and (ii) for any $\mathbf{U} \in \mathcal{U}_+^D$, there exists an inverse $\mathbf{U}^{-1} \in \mathcal{U}_+^D$ such that $\mathbf{U}^{-1}\mathbf{U} = \mathbf{U}\mathbf{U}^{-1} = \mathbf{I}$ where $\mathbf{I} \in \mathcal{U}_+^D$ is the identity matrix.²⁸ Note however that \mathcal{S}_+^D does not form a group.²⁹

To see the first part (i), we use mathematical induction. It is trivial to show the case where $D = 1$. For $D > 1$, let us write $\mathbf{U}, \mathbf{U}' \in \mathcal{U}_+^D$ as partitioned matrices such that

$$\mathbf{U} = \begin{pmatrix} \hat{\mathbf{U}} & \beta \\ \mathbf{0}^T & \alpha \end{pmatrix}, \quad \mathbf{U}' = \begin{pmatrix} \hat{\mathbf{U}}' & \beta' \\ \mathbf{0}^T & \alpha' \end{pmatrix} \quad (210)$$

where $\hat{\mathbf{U}}, \hat{\mathbf{U}}' \in \mathcal{U}_+^{D-1}$ and $\alpha, \alpha' > 0$. Assuming that the product of $\hat{\mathbf{U}}$ and $\hat{\mathbf{U}}'$ is again in the group \mathcal{U}_+^{D-1} so that $\hat{\mathbf{U}}\hat{\mathbf{U}}' \in \mathcal{U}_+^{D-1}$, we have

$$\mathbf{U}\mathbf{U}' = \begin{pmatrix} \hat{\mathbf{U}} & \beta \\ \mathbf{0}^T & \alpha \end{pmatrix} \begin{pmatrix} \hat{\mathbf{U}}' & \beta' \\ \mathbf{0}^T & \alpha' \end{pmatrix} = \begin{pmatrix} \hat{\mathbf{U}}\hat{\mathbf{U}}' & \eta \\ \mathbf{0}^T & \alpha\alpha' \end{pmatrix} \in \mathcal{U}_+^D \quad (211)$$

where

$$\eta = \hat{\mathbf{U}}\beta' + \beta\alpha'. \quad (212)$$

The second part (ii) can be shown similarly by induction. Again, it is trivial to show the case where $D = 1$; and, for $D > 1$, we make use of partitioned matrices. Assuming that the inverse $\hat{\mathbf{U}}^{-1}$ of $\hat{\mathbf{U}} \in \mathcal{U}_+^{D-1}$ is again in the group \mathcal{U}_+^{D-1} so that $\hat{\mathbf{U}}^{-1} \in \mathcal{U}_+^{D-1}$, we easily find \mathbf{U}^{-1} in the form

$$\mathbf{U}^{-1} = \begin{pmatrix} \hat{\mathbf{U}} & \beta \\ \mathbf{0}^T & \alpha \end{pmatrix}^{-1} = \begin{pmatrix} \hat{\mathbf{U}}^{-1} & \xi \\ \mathbf{0}^T & \alpha^{-1} \end{pmatrix} \in \mathcal{U}_+^D \quad (213)$$

where

$$\xi = -\hat{\mathbf{U}}^{-1}\beta\alpha^{-1}. \quad (214)$$

Cholesky decomposition For any symmetric positive-definite matrix $\mathbf{A} \in \mathcal{S}_+^D$, there exists a unique upper triangular matrix $\mathbf{U} \in \mathcal{U}_+^D$ such that

$$\mathbf{A} = \mathbf{U}^T \mathbf{U} \quad (215)$$

²⁷The identity (C.10) is also known as the *Leibniz formula for determinants*.

²⁸The inverse \mathbf{U}^{-1} is unique in general. To see this, suppose that \mathbf{R} is another inverse of \mathbf{U} . Then, we have by associativity that $\mathbf{R} = (\mathbf{U}^{-1}\mathbf{U})\mathbf{R} = \mathbf{U}^{-1}(\mathbf{U}\mathbf{R}) = \mathbf{U}^{-1}$.

²⁹In fact, it is easy to see that the product $\mathbf{P} = \mathbf{A}\mathbf{B}$ of $\mathbf{A}, \mathbf{B} \in \mathcal{S}_+^D$ is not symmetric in general. Furthermore, we do not either have (206) where \mathbf{A} is replaced by \mathbf{P} . To see this, consider, e.g.,

$$\mathbf{A} = \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} \in \mathcal{S}_+^2, \quad \mathbf{B} = \begin{pmatrix} 1 & -2 \\ -2 & 5 \end{pmatrix} \in \mathcal{S}_+^2 \quad (208)$$

the product of which is given by

$$\mathbf{P} = \mathbf{A}\mathbf{B} = \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 1 & -2 \\ -2 & 5 \end{pmatrix} = \begin{pmatrix} -1 & 3 \\ -3 & 8 \end{pmatrix} \notin \mathcal{S}_+^2. \quad (209)$$

If we take, e.g., $\mathbf{x} = (1, 0)^T$, then $\mathbf{x}^T \mathbf{P} \mathbf{x} = -1 < 0$.

or, if the matrices $\mathbf{A} = (A_{ij})$ and $\mathbf{U} = (U_{ij})$ are written element-wise,

$$\begin{pmatrix} A_{11} & A_{12} & \cdots & A_{1D} \\ A_{21} & A_{22} & \cdots & A_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ A_{D1} & A_{D2} & \cdots & A_{DD} \end{pmatrix} = \begin{pmatrix} U_{11} & 0 & \cdots & 0 \\ U_{12} & U_{22} & \ddots & \vdots \\ \vdots & \vdots & \ddots & 0 \\ U_{1D} & U_{2D} & \cdots & U_{DD} \end{pmatrix} \begin{pmatrix} U_{11} & U_{12} & \cdots & U_{1D} \\ 0 & U_{22} & \cdots & U_{2D} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & U_{DD} \end{pmatrix}. \quad (216)$$

This matrix factorization is known as the *Cholesky decomposition* (Anderson, 2003; Press et al., 1992) or *Cholesky factorization* (Golub and Van Loan, 2013; Szeliski, 2010). The upper triangular matrix \mathbf{U} is called the *Cholesky factor* of \mathbf{A} .³⁰

Conversely, for any upper triangular matrix $\mathbf{U} \in \mathcal{U}_+^D$, there exists a unique symmetric positive-definite matrix $\mathbf{A} \in \mathcal{S}_+^D$ such that (215) holds. Therefore, the Cholesky decomposition (215) can be regarded as a bijective function $\mathcal{S}_+^D \rightarrow \mathcal{U}_+^D$. It should also be noted here that, as a necessary condition for a matrix transformation to be bijective, the number of independent parameters in $\mathbf{A} \in \mathcal{S}_+^D$ indeed agrees with that of $\mathbf{U} \in \mathcal{U}_+^D$. More specifically, either of \mathbf{A} and \mathbf{U} has $D(D+1)/2$ independent parameters because \mathbf{A} is a symmetric matrix and \mathbf{U} is an (upper) triangular matrix.

The existence of the Cholesky decomposition (215) can be shown again by induction. First, assume $D = 1$, in which case we have $A > 0$ so that $U = \sqrt{A} > 0$ satisfies (215). Next, we consider the general case where $D > 1$. As we have done in (210), let us write $\mathbf{A} \in \mathcal{S}_+^D$ and $\mathbf{U} \in \mathcal{U}_+^D$ as partitioned matrices such that

$$\mathbf{A} = \begin{pmatrix} \hat{\mathbf{A}} & \mathbf{b} \\ \mathbf{b}^T & a \end{pmatrix}, \quad \mathbf{U} = \begin{pmatrix} \hat{\mathbf{U}} & \boldsymbol{\beta} \\ \mathbf{0}^T & \alpha \end{pmatrix} \quad (218)$$

where $\hat{\mathbf{A}} \in \mathcal{S}_+^{D-1}$ and $\hat{\mathbf{U}} \in \mathcal{U}_+^{D-1}$.³¹ Then, the right hand side of (215) can be written as

$$\mathbf{U}^T \mathbf{U} = \begin{pmatrix} \hat{\mathbf{U}}^T & \mathbf{0} \\ \boldsymbol{\beta}^T & \alpha \end{pmatrix} \begin{pmatrix} \hat{\mathbf{U}} & \boldsymbol{\beta} \\ \mathbf{0}^T & \alpha \end{pmatrix} = \begin{pmatrix} \hat{\mathbf{U}}^T \hat{\mathbf{U}} & \hat{\mathbf{U}}^T \boldsymbol{\beta} \\ \boldsymbol{\beta}^T \hat{\mathbf{U}} & \|\boldsymbol{\beta}\|^2 + \alpha^2 \end{pmatrix}. \quad (219)$$

Equating the right hand side of (219) with the partitioned form (218) of \mathbf{A} gives

$$\hat{\mathbf{A}} = \hat{\mathbf{U}}^T \hat{\mathbf{U}} \quad (220)$$

$$\mathbf{b} = \hat{\mathbf{U}}^T \boldsymbol{\beta} \quad (221)$$

$$a = \|\boldsymbol{\beta}\|^2 + \alpha^2. \quad (222)$$

We see that (220) is the Cholesky decomposition from $\hat{\mathbf{A}} \in \mathcal{S}_+^{D-1}$ to $\hat{\mathbf{U}} \in \mathcal{U}_+^{D-1}$. Assuming that the $(D-1)$ -dimensional Cholesky decomposition (220) exists, we can find $\boldsymbol{\beta}$ and α as

$$\boldsymbol{\beta} = \hat{\mathbf{U}}^{-T} \mathbf{b} \quad (223)$$

$$\alpha = \sqrt{a - \|\boldsymbol{\beta}\|^2} > 0. \quad (224)$$

³⁰The Cholesky decomposition (215) can also be written in the form

$$\mathbf{A} = \mathbf{L} \mathbf{L}^T \quad (217)$$

where \mathbf{L} is a lower triangular matrix with positive diagonals, in which case \mathbf{L} is called the Cholesky factor.

³¹It directly follows from (206) that $\hat{\mathbf{A}} \succ 0$. In fact, substituting $\mathbf{x} = (\hat{\mathbf{x}}^T, 0)^T$ and the partitioned form (218) of \mathbf{A} into (206), we have $\hat{\mathbf{x}}^T \hat{\mathbf{A}} \hat{\mathbf{x}} > 0$ for any $\hat{\mathbf{x}} \neq \mathbf{0}$, showing $\hat{\mathbf{A}} \succ 0$.

Note that (223) is well-defined because $\hat{\mathbf{U}} \in \mathcal{U}_+^{D-1}$ is nonsingular; and so is (224) because we have $a - \|\beta\|^2 > 0$, which can be shown by substituting

$$\mathbf{x} = \begin{pmatrix} \hat{\mathbf{A}}^{-1}\mathbf{b} \\ -1 \end{pmatrix} \quad (225)$$

into (206), giving

$$0 < a - \mathbf{b}^T \hat{\mathbf{A}}^{-1} \mathbf{b} \quad (226)$$

$$= a - \|\beta\|^2 \quad (227)$$

where we have used (220) (with the both sides inverted) and (223).³²

Finally, we observe that the above induction effectively constructs a recursive algorithm to compute $\mathbf{U} \in \mathcal{U}_+^D$ from $\mathbf{A} \in \mathcal{S}_+^D$. The uniqueness of the Cholesky decomposition (215) is implied by the uniqueness of each operation in the algorithm.

Jacobians of matrix transformations Bijectivity of the Cholesky decomposition (215) allows us to *change variables* between a symmetric positive-definite matrix $\mathbf{A} \in \mathcal{S}_+^D$ and its Cholesky factor $\mathbf{U} \in \mathcal{U}_+^D$; this change of variables technique is, as we shall see shortly, useful for evaluating integrals over the space \mathcal{S}_+^D of symmetric positive-definite matrices because the Cholesky factor $\mathbf{U} = (U_{ij}) \in \mathcal{U}_+^D$ has a simpler domain of integration such that

$$U_{ii} \in (0, \infty), \quad U_{ij} \in (-\infty, \infty) \quad (i < j) \quad (228)$$

for diagonal and off-diagonal elements, respectively. In the following, we review some necessary Jacobians for such matrix transformations.

First, let us find the Jacobian for the Cholesky decomposition (215) between $\mathbf{A} \in \mathcal{S}_+^D$ and $\mathbf{U} \in \mathcal{U}_+^D$. To do so, we work with its element-wise representation given by (216). Since $\mathbf{A} = (A_{ij})$ is symmetric so that $A_{ij} \equiv A_{ji}$, we only consider the upper triangular elements A_{ij} of \mathbf{A} such that $i \leq j$. Writing down the upper triangular elements A_{ij} in *lexicographic* (or row-major) order, we have

$$A_{11} = U_{11}^2 \quad (229)$$

$$A_{12} = U_{11}U_{12} \quad (230)$$

$$\vdots$$

$$A_{1D} = U_{11}U_{1D} \quad (231)$$

$$A_{22} = U_{12}^2 + U_{22}^2 \quad (232)$$

$$\vdots$$

$$A_{2D} = U_{12}U_{1D} + U_{22}U_{2D} \quad (233)$$

$$\vdots$$

$$A_{ij} = U_{1i}U_{1j} + \cdots + U_{ii}U_{ij} \quad (i \leq j) \quad (234)$$

$$\vdots$$

$$A_{DD} = U_{1D}^2 + U_{2D}^2 + \cdots + U_{DD}^2. \quad (235)$$

Since A_{ij} depends only on U_{ij} and the preceding elements of \mathbf{U} in lexicographic order,³³ the

³²The inequality (226) can also be shown by noting that the *Schur complement* (282) of $\mathbf{A} \succ 0$ with respect to $\hat{\mathbf{A}}$, given by $\mathbf{A}/\hat{\mathbf{A}} = a - \mathbf{b}^T \hat{\mathbf{A}}^{-1} \mathbf{b}$, is again positive definite.

³³More specifically, A_{ij} depends on U_{1i}, \dots, U_{ii} and U_{1j}, \dots, U_{ij} where $i \leq j$, all of which are contained in the rectangular matrix of dimensionality (i, j) aligned at the upper-left corner of \mathbf{U} .

Jacobian of \mathbf{A} with respect to \mathbf{U} is a lower triangular matrix of the form

$$\frac{\partial \mathbf{A}}{\partial \mathbf{U}} \equiv \frac{\partial (A_{11}, A_{12}, \dots, A_{1D}, A_{22}, \dots, A_{2D}, \dots, A_{DD})}{\partial (U_{11}, U_{12}, \dots, U_{1D}, U_{22}, \dots, U_{2D}, \dots, U_{DD})} \quad (236)$$

$$= \begin{matrix} & U_{11} & U_{12} & \cdots & U_{1D} & U_{22} & \cdots & U_{2D} & \cdots & U_{DD} \\ \begin{matrix} A_{11} \\ A_{12} \\ \vdots \\ A_{1D} \\ A_{22} \\ \vdots \\ A_{2D} \\ \vdots \\ A_{DD} \end{matrix} & \begin{pmatrix} 2U_{11} & 0 & \cdots & 0 & 0 & \cdots & 0 & \cdots & 0 \\ * & U_{11} & \cdots & 0 & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ * & * & \cdots & U_{11} & 0 & \cdots & 0 & \cdots & 0 \\ * & * & \cdots & * & 2U_{22} & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ * & * & \cdots & * & * & \cdots & U_{22} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ * & * & \cdots & * & * & \cdots & * & \cdots & 2U_{DD} \end{pmatrix} \end{matrix} \quad (237)$$

where the diagonal elements are given by

$$\frac{\partial A_{ii}}{\partial U_{ii}} = 2U_{ii} > 0, \quad \frac{\partial A_{ij}}{\partial U_{ij}} = U_{ii} > 0 \quad (i < j) \quad (238)$$

and the elements denoted by $*$ are elements possibly nonzero. Thus, we obtain the (absolute) determinant of the Jacobian in the form

$$\left| \frac{\partial \mathbf{A}}{\partial \mathbf{U}} \right| = 2U_{11}^D \times 2U_{22}^{D-1} \times \cdots \times 2U_{DD} = 2^D \prod_{i=1}^D U_{ii}^{D+1-i}. \quad (239)$$

Another matrix transformation of interest here is a linear transformation between two upper triangular matrices $\mathbf{U}, \mathbf{R} \in \mathcal{U}_+^D$ of the form

$$\mathbf{U} = \mathbf{R}\mathbf{G} \quad (240)$$

where $\mathbf{G} \in \mathcal{U}_+^D$ is a constant (recall that \mathcal{U}_+^D forms a group with respect to multiplication). If these matrices $\mathbf{U} = (U_{ij})$, $\mathbf{R} = (R_{ij})$, and $\mathbf{G} = (G_{ij})$ are written element-wise, then

$$\begin{pmatrix} U_{11} & U_{12} & \cdots & U_{1D} \\ 0 & U_{22} & \cdots & U_{2D} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & U_{DD} \end{pmatrix} = \begin{pmatrix} R_{11} & R_{12} & \cdots & R_{1D} \\ 0 & R_{22} & \cdots & R_{2D} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & R_{DD} \end{pmatrix} \begin{pmatrix} G_{11} & G_{12} & \cdots & G_{1D} \\ 0 & G_{22} & \cdots & G_{2D} \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & G_{DD} \end{pmatrix} \quad (241)$$

so that

$$U_{11} = R_{11}G_{11} \quad (242)$$

$$U_{12} = R_{11}G_{12} + R_{12}G_{22} \quad (243)$$

$$\vdots \\ U_{1D} = R_{11}G_{1D} + R_{12}G_{2D} + \cdots + R_{1D}G_{DD} \quad (244)$$

$$U_{22} = R_{22}G_{22} \quad (245)$$

$$\vdots \\ U_{2D} = R_{22}G_{2D} + \cdots + R_{2D}G_{DD} \quad (246)$$

$$\vdots \\ U_{ij} = R_{ii}G_{ij} + \cdots + R_{ij}G_{jj} \quad (i \leq j) \quad (247)$$

$$\vdots \\ U_{DD} = R_{DD}G_{DD}. \quad (248)$$

Since U_{ij} depends only on R_{ij} and the preceding elements of \mathbf{R} in lexicographic order, the Jacobian matrix $\partial \mathbf{U} / \partial \mathbf{R}$, defined similarly to (236), is again lower triangular. The diagonal elements of $\partial \mathbf{U} / \partial \mathbf{R}$ are given by

$$\frac{\partial U_{ij}}{\partial R_{ij}} = G_{jj} > 0 \quad (249)$$

where $i \leq j$ so that

$$\left| \frac{\partial \mathbf{U}}{\partial \mathbf{R}} \right| = G_{11} \times G_{22}^2 \times \cdots \times G_{DD}^D = \prod_{i=1}^D G_{ii}^i. \quad (250)$$

Multivariate gamma function The *multivariate gamma function* (Anderson, 2003; Olver et al., 2017) is defined by

$$\Gamma_D(a) \equiv \int_{S_+^D} |\mathbf{X}|^{a-(D+1)/2} \exp \{ -\text{Tr}(\mathbf{X}) \} d\mathbf{X} \quad (251)$$

where $a > (D-1)/2$; and the integration is taken over the space \mathcal{S}_+^D of symmetric positive-definite matrices of dimensionality D . It is a generalization of the ordinary (univariate) gamma function $\Gamma(\cdot)$ defined by (1.141). In fact, when $D = 1$, the multivariate gamma function $\Gamma_D(\cdot)$ reduces to the univariate gamma function $\Gamma(\cdot)$ so that $\Gamma_1(a) \equiv \Gamma(a)$.

The multivariate gamma function $\Gamma_D(\cdot)$ can also be written in terms of the univariate gamma function $\Gamma(\cdot)$ as

$$\Gamma_D(a) = \pi^{D(D-1)/4} \prod_{i=1}^D \Gamma \left(a - \frac{i-1}{2} \right) \quad (252)$$

so that we can simplify the normalization constant (B.79) of the Wishart distribution (B.78) in the form

$$B(\mathbf{W}, \nu)^{-1} = 2^{\nu D/2} |\mathbf{W}|^{\nu/2} \Gamma_D \left(\frac{\nu}{2} \right) \quad (253)$$

where $\mathbf{W} \in \mathcal{S}_+^D$ (so that $\mathbf{W} \succ 0$) and $\nu > D-1$. Similarly, if we define the *multivariate digamma function* by

$$\psi_D(a) \equiv \frac{d}{da} \ln \Gamma_D(a) = \sum_{i=1}^D \psi \left(a - \frac{i-1}{2} \right) \quad (254)$$

where $\psi(\cdot)$ is the (univariate) digamma function defined by (101), then the log (absolute) determinant expectation (B.81) can be written as

$$\mathbb{E}[\ln |\mathbf{\Lambda}|] = D \ln 2 + \ln |\mathbf{W}| + \psi_D \left(\frac{\nu}{2} \right). \quad (255)$$

Let us now evaluate the integral (251) to show the identity (252). We first note that, since $\mathbf{X} = (X_{ij}) \in \mathcal{S}_+^D$ is a symmetric matrix, we take the integration only over the $D(D+1)/2$ upper triangular elements X_{ij} where $i \leq j$ so that the differential $d\mathbf{X}$ means

$$d\mathbf{X} \equiv \prod_{1 \leq i \leq j \leq D} dX_{ij} = dX_{11} dX_{12} \cdots dX_{1D} dX_{22} \cdots dX_{2D} \cdots dX_{DD}. \quad (256)$$

The differential $d\mathbf{U}$ of an upper triangular matrix $\mathbf{U} = (U_{ij}) \in \mathcal{U}_+^D$ is defined similarly. Making a change of variables from $\mathbf{X} \in \mathcal{S}_+^D$ to the Cholesky factor $\mathbf{U} \in \mathcal{U}_+^D$ such that $\mathbf{X} = \mathbf{U}^T \mathbf{U}$, we have

$$\Gamma_D(a) = \int_{\mathcal{U}_+^D} |\mathbf{X}|^{a-(D+1)/2} \exp \{ -\text{Tr}(\mathbf{X}) \} \left| \frac{\partial \mathbf{X}}{\partial \mathbf{U}} \right| d\mathbf{U} \quad (257)$$

$$= 2^D \left[\prod_{i=1}^D \int_0^\infty U_{ii}^{2a-i} \exp(-U_{ii}^2) dU_{ii} \right] \left[\prod_{1 \leq i < j \leq D} \int_{-\infty}^\infty \exp(-U_{ij}^2) dU_{ij} \right] \quad (258)$$

$$= \pi^{D(D-1)/4} \prod_{i=1}^D \Gamma \left(a - \frac{i-1}{2} \right) \quad (259)$$

where we have used

$$|\mathbf{X}| = |\mathbf{U}^T \mathbf{U}| = |\mathbf{U}|^2 = \prod_{i=1}^D U_{ii}^2 \quad (260)$$

$$\text{Tr}(\mathbf{X}) = \sum_{i=1}^D X_{ii} = \sum_{i=1}^D (U_{1i}^2 + \cdots + U_{ii}^2) = \sum_{1 \leq i \leq j \leq D} U_{ij}^2 \quad (261)$$

together with the result (239) for the Jacobian of the Cholesky decomposition and the integral identities

$$\int_0^\infty u^{2a-i} \exp(-u^2) du = \frac{1}{2} \Gamma \left(a - \frac{i-1}{2} \right) \quad \text{if } a > \frac{i-1}{2} \quad (262)$$

$$\int_{-\infty}^\infty \exp(-u^2) du = \sqrt{\pi}. \quad (263)$$

Normalization of Wishart Let us next show that the Wishart distribution (B.78) is indeed correctly normalized. Specifically, we show the following integral identity

$$B(\mathbf{W}, \nu)^{-1} = \int_{\mathcal{S}_+^D} |\mathbf{\Lambda}|^{(\nu-D-1)/2} \exp \left\{ -\frac{1}{2} \text{Tr}(\mathbf{W}^{-1} \mathbf{\Lambda}) \right\} d\mathbf{\Lambda} \quad (264)$$

where the normalization constant $B(\mathbf{W}, \nu)$ is given by (253); and we have written the domain of integration explicitly so as to make it clear that $\mathbf{\Lambda} \in \mathcal{S}_+^D$.

To evaluate the right hand side of (264), we successively make two changes of variables of the forms

$$\mathbf{\Lambda} = \mathbf{U}^T \mathbf{U}, \quad \mathbf{U} = \mathbf{R} \mathbf{G} \quad (265)$$

where $\mathbf{U}, \mathbf{R} \in \mathcal{U}_+^D$; and $\mathbf{G} \in \mathcal{U}_+^D$ is the Cholesky factor of $2\mathbf{W} \in \mathcal{S}_+^D$ so that

$$2\mathbf{W} = \mathbf{G}^T \mathbf{G}. \quad (266)$$

The overall Jacobian factor for the two successive changes of variables (265) is given by

$$\left| \frac{\partial \mathbf{\Lambda}}{\partial \mathbf{U}} \right| \left| \frac{\partial \mathbf{U}}{\partial \mathbf{R}} \right| = 2^D |\mathbf{G}|^{D+1} \prod_{i=1}^D R_{ii}^{D+1-i} \quad (267)$$

where we have used the results (239) and (250).

It is worth noting here that the change of variables (265) from $\mathbf{\Lambda} \in \mathcal{S}_+^D$ to $\mathbf{R} \in \mathcal{U}_+^D$ effectively makes the resulting random variables R_{ij} where $i \leq j$ to be independently distributed: R_{ii}^2 is distributed according to a gamma distribution; and R_{ij} where $i < j$ to a Gaussian. More specifically, if $\mathbf{\Lambda} \in \mathcal{S}_+^D$ follows a Wishart distribution so that $p(\mathbf{\Lambda}) = \mathcal{W}(\mathbf{\Lambda} | \mathbf{W}, \nu)$, then it can be shown through the change of variables (265) that

$$p(\{\tau_i\}, \{R_{ij} \mid i < j\}) = \left[\prod_{i=1}^D \text{Gam} \left(\tau_i \mid \frac{\nu+1-i}{2}, 1 \right) \right] \left[\prod_{1 \leq i < j \leq D} \mathcal{N} \left(R_{ij} \mid 0, \frac{1}{2} \right) \right] \quad (268)$$

where we have further made the change variables $R_{ii} = \sqrt{\tau_i}$ for the diagonal elements.³⁴ The above observation is useful for sampling (see Chapter 11) and also shows that the Wishart distribution is indeed correctly normalized. In the following, however, we evaluate the integral (264) directly by identifying the multivariate gamma function in order to show the normalization of the Wishart.

Evaluating the right hand side of (264) by making the change of variables (265), we have

$$\int_{\mathcal{S}_+^D} |\mathbf{\Lambda}|^{(\nu-D-1)/2} \exp \left\{ -\frac{1}{2} \text{Tr}(\mathbf{W}^{-1} \mathbf{\Lambda}) \right\} d\mathbf{\Lambda} \quad (269)$$

$$= |\mathbf{G}|^\nu \int_{\mathcal{U}_+^D} |\mathbf{R}^T \mathbf{R}|^{(\nu-D-1)/2} \exp \{ -\text{Tr}(\mathbf{R}^T \mathbf{R}) \} \left[2^D \prod_{i=1}^D R_{ii}^{D+1-i} \right] d\mathbf{R} \quad (270)$$

$$= 2^{\nu D/2} |\mathbf{W}|^{\nu/2} \Gamma_D \left(\frac{\nu}{2} \right) \quad (271)$$

where we have used (267) and the relation

$$|\mathbf{G}| = 2^{D/2} |\mathbf{W}|^{1/2} \quad (272)$$

between the (absolute) determinants of $\mathbf{W} \in \mathcal{S}_+^D$ and $\mathbf{G} \in \mathcal{U}_+^D$; and identified the integral in the right hand side with an integral form (257) of the multivariate gamma function $\Gamma_D(a)$ where $a = \nu/2$.

³⁴A slightly different form of the change of variables (265) where the scale matrix \mathbf{W} is decomposed into the Cholesky factor of itself as $\mathbf{W} = \mathbf{G}^T \mathbf{G}$, instead of (266), is called the *Bartlett decomposition* (Anderson, 2003), in which case R_{ij} where $i < j$ is distributed according to the zero-mean, unit-variance Gaussian $\mathcal{N}(\cdot | 0, 1)$; and R_{ii}^2 to the *chi-squared* distribution (or χ^2 -distribution) with $\nu + 1 - i$ degrees of freedom or, equivalently, a gamma distribution of the form $\text{Gam}(\cdot | (\nu + 1 - i)/2, 1/2)$.

Mean and log determinant expectation of Wishart Finally, we show the mean (B.80) and the log (absolute) determinant expectation (255) of the Wishart distribution (B.78) by making use of the general identity (58). To this end, we first evaluate the derivatives of the log probability

$$\ln \mathcal{W}(\Lambda|\mathbf{W}, \nu) = \ln B(\mathbf{W}, \nu) + \frac{\nu + D - 1}{2} \ln |\Lambda| - \frac{1}{2} \text{Tr}(\mathbf{W}^{-1}\Lambda) \quad (273)$$

where

$$\ln B(\mathbf{W}, \nu) = -\frac{\nu D}{2} \ln 2 - \frac{\nu}{2} \ln |\mathbf{W}| - \ln \Gamma_D\left(\frac{\nu}{2}\right) \quad (274)$$

with respect to the parameters \mathbf{W} and ν , giving

$$\nabla_{\mathbf{W}} \ln \mathcal{W}(\Lambda|\mathbf{W}, \nu) = -\frac{\nu}{2} \mathbf{W}^{-\text{T}} + \frac{1}{2} \mathbf{W}^{-\text{T}} \Lambda^{\text{T}} \mathbf{W}^{-\text{T}} \quad (275)$$

$$\frac{\partial}{\partial \nu} \ln \mathcal{W}(\Lambda|\mathbf{W}, \nu) = -\frac{D}{2} \ln 2 - \frac{1}{2} \ln |\mathbf{W}| - \frac{1}{2} \psi_D\left(\frac{\nu}{2}\right) + \frac{1}{2} \ln |\Lambda| \quad (276)$$

where we have used (332), (337), and (254). Substituting the above derivatives (275) and (276) into (58), we obtain the expectations (B.80) and (255), respectively.

Note that, since \mathbf{W} is symmetric, we should have, strictly speaking, imposed the symmetry constraint on \mathbf{W} when we evaluate the gradient (275). However, since (the expectation of) the gradient (275) obtained without the symmetry constraint is again symmetric because of the symmetry of \mathbf{W} and Λ , the result (349) allows us to use (275) for evaluating the expectation (B.80) by making use of (58).

Having obtained the expectations (B.80) and (255), we can now evaluate the entropy, which is given by $H[\Lambda] = -\mathbb{E}[\ln \mathcal{W}(\Lambda|\mathbf{W}, \nu)]$ and easily obtained in the form (B.82) by making use of (B.80). The entropy (B.82) can be further simplified to

$$H[\Lambda] = \frac{D(D+1)}{2} \ln 2 + \frac{D+1}{2} \ln |\mathbf{W}| + \ln \Gamma_D\left(\frac{\nu}{2}\right) - \frac{\nu - D - 1}{2} \psi_D\left(\frac{\nu}{2}\right) + \frac{\nu D}{2} \quad (277)$$

where we have used (274) and (255).

Page 693

Paragraph -1, Line -4: “Gamma” should read “gamma” (without capitalization).

Page 693

Paragraph -1, Line -1: $b = 1/2W$ should read $b = 1/(2W)$ for clarity.

Page 696

Equation (C.5): Replacing \mathbf{B}^{T} with \mathbf{A} , we obtain a more general identity

$$(\mathbf{P}^{-1} + \mathbf{A}\mathbf{R}^{-1}\mathbf{B})^{-1} \mathbf{A}\mathbf{R}^{-1} = \mathbf{P}\mathbf{A}(\mathbf{B}\mathbf{P}\mathbf{A} + \mathbf{R})^{-1}. \quad (278)$$

The identity (278) is necessary to show the *push-through identity* (C.6), which in turn can be used to show *Sylvester’s determinant identity* (C.14). As suggested in the text, the above identity (278) can be directly verified by right multiplying both sides by $(\mathbf{B}\mathbf{P}\mathbf{A} + \mathbf{R})$.

However, I would prefer to prove the general push-through identity (278) together with the *Woodbury identity* (C.7) in terms of the inverse of a partitioned matrix, which we have already seen in Section 2.3.1. To this end, we first introduce a square matrix \mathbf{M} that is partitioned into four submatrices so that

$$\mathbf{M} = \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} \quad (279)$$

where \mathbf{A} and \mathbf{D} are square (but not necessarily the same dimension) and then note that \mathbf{M} can be block diagonalized as

$$\begin{pmatrix} \mathbf{I} & \mathbf{O} \\ -\mathbf{CA}^{-1} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} \begin{pmatrix} \mathbf{I} & -\mathbf{A}^{-1}\mathbf{B} \\ \mathbf{O} & \mathbf{I} \end{pmatrix} = \begin{pmatrix} \mathbf{A} & \mathbf{O} \\ \mathbf{O} & \mathbf{M}/\mathbf{A} \end{pmatrix} \quad (280)$$

or

$$\begin{pmatrix} \mathbf{I} & -\mathbf{BD}^{-1} \\ \mathbf{O} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} \begin{pmatrix} \mathbf{I} & \mathbf{O} \\ -\mathbf{D}^{-1}\mathbf{C} & \mathbf{I} \end{pmatrix} = \begin{pmatrix} \mathbf{M}/\mathbf{D} & \mathbf{O} \\ \mathbf{O} & \mathbf{D} \end{pmatrix} \quad (281)$$

if \mathbf{A} or \mathbf{D} is nonsingular, respectively, where we have written the *Schur complement* of \mathbf{M} with respect to \mathbf{A} or \mathbf{D} as

$$\mathbf{M}/\mathbf{A} \equiv \mathbf{D} - \mathbf{CA}^{-1}\mathbf{B} \quad (282)$$

or

$$\mathbf{M}/\mathbf{D} \equiv \mathbf{A} - \mathbf{BD}^{-1}\mathbf{C} \quad (283)$$

respectively.³⁵ The above block diagonalization identities (280) and (281) yield two versions of the inverse partitioned matrix \mathbf{M}^{-1} , i.e.,

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{I} & -\mathbf{A}^{-1}\mathbf{B} \\ \mathbf{O} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{A}^{-1} & \mathbf{O} \\ \mathbf{O} & (\mathbf{M}/\mathbf{A})^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{I} & \mathbf{O} \\ -\mathbf{CA}^{-1} & \mathbf{I} \end{pmatrix} \quad (286)$$

$$= \begin{pmatrix} \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{B}(\mathbf{M}/\mathbf{A})^{-1}\mathbf{CA}^{-1} & -\mathbf{A}^{-1}\mathbf{B}(\mathbf{M}/\mathbf{A})^{-1} \\ -(\mathbf{M}/\mathbf{A})^{-1}\mathbf{CA}^{-1} & (\mathbf{M}/\mathbf{A})^{-1} \end{pmatrix} \quad (287)$$

and

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{I} & \mathbf{O} \\ -\mathbf{D}^{-1}\mathbf{C} & \mathbf{I} \end{pmatrix} \begin{pmatrix} (\mathbf{M}/\mathbf{D})^{-1} & \mathbf{O} \\ \mathbf{O} & \mathbf{D}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{I} & -\mathbf{BD}^{-1} \\ \mathbf{O} & \mathbf{I} \end{pmatrix} \quad (288)$$

$$= \begin{pmatrix} (\mathbf{M}/\mathbf{D})^{-1} & -(\mathbf{M}/\mathbf{D})^{-1}\mathbf{BD}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}(\mathbf{M}/\mathbf{D})^{-1} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}(\mathbf{M}/\mathbf{D})^{-1}\mathbf{BD}^{-1} \end{pmatrix} \quad (289)$$

respectively. Equating the right hand sides, we have, e.g.,

$$(\mathbf{M}/\mathbf{D})^{-1} = \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{B}(\mathbf{M}/\mathbf{A})^{-1}\mathbf{CA}^{-1} \quad (290)$$

and

$$-(\mathbf{M}/\mathbf{A})^{-1}\mathbf{CA}^{-1} = -\mathbf{D}^{-1}\mathbf{C}(\mathbf{M}/\mathbf{D})^{-1}. \quad (291)$$

³⁵Note that the notation for the Schur complement is chosen to suggest that it has a flavor of division (Minka, 2000). In fact, taking the determinant on both sides of (280) and (281), we have from the definition (C.10) of the determinant that

$$\det(\mathbf{M}) = \det(\mathbf{A}) \det(\mathbf{M}/\mathbf{A}) \quad (284)$$

and

$$\det(\mathbf{M}) = \det(\mathbf{D}) \det(\mathbf{M}/\mathbf{D}) \quad (285)$$

respectively.

Substituting (282) and (283) into both sides and replacing \mathbf{D} with $-\mathbf{D}$, we finally have

$$(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1} \quad (292)$$

and

$$(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1} = \mathbf{D}^{-1}\mathbf{C}(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} \quad (293)$$

which are equivalent to (C.7) and (278), respectively.

Page 696

Paragraph 3, Line 2: $\sum_n \alpha_n \mathbf{a}_n = 0$ should read $\sum_n \alpha_n \mathbf{a}_n = \mathbf{0}$ (the right hand side should be a zero vector $\mathbf{0}$).

Pages 696 and 697

Equations (C.8), (C.9), and (C.12): Note that, although determinant $\det(\cdot)$ and trace $\text{Tr}(\cdot)$ only apply to square matrices, the matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} in (C.8) and (C.9) themselves are not necessarily square. On the other hand, in order for the determinant identity (C.12) to hold, both \mathbf{A} and \mathbf{B} must be square.

Page 697

Equation (C.17): It is clear that the definitions for the “vector derivative” (C.17) of a scalar with respect to a vector and that (C.18) of a vector with respect to a vector contradict each other. The vector derivative of the form (C.17) is usually called the *gradient* whereas (C.18) is called the *Jacobian* (Minka, 2000). Note that (C.16) is a special case of (C.18) and thus the Jacobian. In order to avoid ambiguity, we should use a different notation, say, ∇ for the gradient, as defined in the following.

Gradient with respect to a vector Given a vector function $\mathbf{y}(\mathbf{x}) = (y_1(\mathbf{x}), \dots, y_M(\mathbf{x}))^T$ of a vector $\mathbf{x} = (x_1, \dots, x_D)^T$, we write the *gradient* $\nabla_{\mathbf{x}}\mathbf{y}$ of $\mathbf{y}(\mathbf{x})$ with respect to \mathbf{x} as a $D \times M$ matrix of partial derivatives so that

$$\nabla_{\mathbf{x}}\mathbf{y} \equiv \left(\frac{\partial y_j}{\partial x_i} \right) = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_M}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_D} & \cdots & \frac{\partial y_M}{\partial x_D} \end{pmatrix}. \quad (294)$$

As a special case, we see that the gradient $\nabla_{\mathbf{x}}y$ of a scalar function $y(\mathbf{x})$ with respect to a column vector \mathbf{x} is again a column vector of the same dimensionality as \mathbf{x} , corresponding to the right hand side of (C.17), i.e.,

$$\nabla_{\mathbf{x}}y = \left(\frac{\partial y}{\partial x_i} \right) = \begin{pmatrix} \frac{\partial y}{\partial x_1} \\ \vdots \\ \frac{\partial y}{\partial x_D} \end{pmatrix}. \quad (295)$$

Chain rule for gradient Note that the right hand side of the definition of the gradient (294) is identical to the transpose of the *Jacobian* $\partial \mathbf{y} / \partial \mathbf{x} = (\partial y_i / \partial x_j)$ so that

$$\nabla_{\mathbf{x}} \mathbf{y} = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \quad (296)$$

as a consequence of which the chain rule for the gradient is such that the intermediate gradients are built up “towards the left,” i.e.,

$$\nabla_{\mathbf{x}} \mathbf{z}(\mathbf{y}) = \left(\frac{\partial \mathbf{z}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T = \nabla_{\mathbf{x}} \mathbf{y} \nabla_{\mathbf{y}} \mathbf{z}. \quad (297)$$

Since the chain rule (297) is handy when we compute the gradients of composite functions (see below), I would suggest that it should also be pointed out in the “(Vector and) Matrix Derivatives” section of Appendix C.

Taylor series in terms of gradients At this point, one might wonder why we use the two different forms of vector derivative that are identical up to the transposed layout, i.e., the gradient $\nabla_{\mathbf{x}} \mathbf{y}$ and the Jacobian $\partial \mathbf{y} / \partial \mathbf{x}$. As [Minka \(2000\)](#) points out, Jacobians are useful in calculus while gradients are useful in optimization. For instance, we can write down the Taylor series expansion (up to the second order) of a scalar function $f(\mathbf{x})$ succinctly in terms of the gradients as

$$f(\mathbf{x} + \epsilon \boldsymbol{\eta}) = f(\mathbf{x}) + \epsilon \boldsymbol{\eta}^T \mathbf{g}(\mathbf{x}) + \frac{\epsilon^2}{2} \boldsymbol{\eta}^T \mathbf{H}(\mathbf{x}) \boldsymbol{\eta} + O(\epsilon^3) \quad (298)$$

where $\mathbf{g}(\mathbf{x})$ and $\mathbf{H}(\mathbf{x})$ are the gradient vector and the Hessian matrix of $f(\mathbf{x})$, respectively, so that

$$\mathbf{g}(\mathbf{x}) \equiv \nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_D} \end{pmatrix}, \quad \mathbf{H}(\mathbf{x}) \equiv \nabla_{\mathbf{x}} \nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_D \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_D \partial x_D} \end{pmatrix}. \quad (299)$$

Page 697

Equation (C.19): Following the gradient notation (294), we see that (C.19) should read

$$\nabla \{ \mathbf{x}^T \mathbf{a} \} = \nabla \{ \mathbf{a}^T \mathbf{x} \} = \mathbf{a} \quad (300)$$

where we have omitted the subscript \mathbf{x} in what should be $\nabla_{\mathbf{x}}$.

Vector derivative identities Some other useful identities I would suggest to include are

$$\nabla \{ \mathbf{x}^T \mathbf{A} \mathbf{x} \} = \nabla \text{Tr}(\mathbf{x} \mathbf{x}^T \mathbf{A}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \quad (301)$$

$$\nabla \{ \mathbf{B} \mathbf{x} \} = \mathbf{B}^T \quad (302)$$

$$\nabla \{ \phi \mathbf{y} \} = \nabla \phi \mathbf{y}^T + \phi \nabla \mathbf{y} \quad (303)$$

where the matrices \mathbf{A} and \mathbf{B} are constant. Note that the term $\mathbf{x}^T \mathbf{A} \mathbf{x}$ in (301) is a quadratic form and thus the square matrix \mathbf{A} is usually taken to be symmetric so that $\mathbf{A} = \mathbf{A}^T$, in which case we have

$$\nabla \{ \mathbf{x}^T \mathbf{A} \mathbf{x} \} = 2 \mathbf{A} \mathbf{x}. \quad (304)$$

Substituting $\mathbf{A} = \mathbf{I}$ gives

$$\nabla \|\mathbf{x}\|^2 = 2\mathbf{x} \quad (305)$$

where $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}}$ is the norm of \mathbf{x} .

We make use of the above identity (305) when, e.g., we take the gradient (108) of a sum-of-squares error function of the form (107) where we also make use of (297) and (302). The same result (108) can also be obtained by first expanding the square norm in (107) and then differentiating the expanded terms with (304) and (302).

The product rule (303) is used when, e.g., we evaluate the Hessian (5.83) of a nonlinear sum-of-squares error function (5.82), which generally takes the form

$$J = \frac{1}{2} \sum_{n=1}^N \varepsilon_n^2 = \frac{1}{2} \|\boldsymbol{\varepsilon}\|^2 \quad (306)$$

where we have written $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)^T$. The gradient $\nabla_{\boldsymbol{\theta}} J$ and the Hessian $\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} J$ of the error function J with respect to some parameters $\boldsymbol{\theta} = (\theta_1, \dots, \theta_M)^T$ (on which the errors $\boldsymbol{\varepsilon}$ depend nonlinearly) are given by

$$\nabla J = \sum_{n=1}^N \varepsilon_n \nabla \varepsilon_n = (\nabla \boldsymbol{\varepsilon}) \boldsymbol{\varepsilon} \quad (307)$$

$$\nabla \nabla J = \sum_{n=1}^N \left\{ \nabla \varepsilon_n (\nabla \varepsilon_n)^T + \varepsilon_n \nabla \nabla \varepsilon_n \right\} \quad (308)$$

$$= \nabla \boldsymbol{\varepsilon} (\nabla \boldsymbol{\varepsilon})^T + [\nabla \text{vec}(\nabla \boldsymbol{\varepsilon})] (\boldsymbol{\varepsilon} \otimes \mathbf{I}_M) \quad (309)$$

where we have made the subscript $\boldsymbol{\theta}$ of the gradient operators implicit; and written the $M \times M$ identity matrix by \mathbf{I}_M . Here, the *vectorization operator* $\text{vec}(\cdot)$ and the *Kronecker product* \otimes are defined as

$$\text{vec}(\mathbf{V}) \equiv \begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_L \end{pmatrix}, \quad \mathbf{A} \otimes \mathbf{B} \equiv \begin{pmatrix} A_{11}\mathbf{B} & \cdots & A_{1N}\mathbf{B} \\ \vdots & \ddots & \vdots \\ A_{M1}\mathbf{B} & \cdots & A_{MN}\mathbf{B} \end{pmatrix} \quad (310)$$

where $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_L)$; and $\mathbf{A} = (A_{ij})$ is an $M \times N$ matrix (Magnus and Neudecker, 2007).

Note that the second form (309) of the Hessian $\nabla \nabla J$ of the nonlinear error function J , which, however, does not necessarily lead to an efficient implementation (neither does that of the gradient ∇J), can be directly obtained by making use of the general product rule

$$\nabla \{\mathbf{R}\boldsymbol{\phi}\} = \nabla \boldsymbol{\phi} \mathbf{R}^T + \nabla \text{vec}(\mathbf{R}) (\boldsymbol{\phi} \otimes \mathbf{I}_M) \quad (311)$$

where M is the number of rows of \mathbf{R} . One can easily show the product rule (311) through its expanded form

$$\nabla \left\{ \sum_{n=1}^N \phi_n \mathbf{r}_n \right\} = \sum_{n=1}^N \nabla \phi_n \mathbf{r}_n^T + \sum_{n=1}^N \phi_n \nabla \mathbf{r}_n \quad (312)$$

where $\mathbf{R} = (\mathbf{r}_1, \dots, \mathbf{r}_N)$ and $\phi = (\phi_1, \dots, \phi_N)^T$.³⁶ Note also that the identities (302) and (303) are special cases of (311).

Finally, I would like to point out that the gradient of the inner product

$$\nabla \{ \mathbf{u}^T \mathbf{v} \} = \nabla \{ \mathbf{v}^T \mathbf{u} \} = (\nabla \mathbf{u}) \mathbf{v} + (\nabla \mathbf{v}) \mathbf{u} \quad (317)$$

is another special case of the general product rule (311). The identities (300) and (301) are, in turn, special cases of (317).

Page 698

Equation (C.20): Although the Jacobian of a vector with respect to a vector is defined in (C.18), the Jacobian of a matrix with respect to a scalar has not been defined. In the following, we define the Jacobian of a matrix as well as the gradient of a scalar with respect to a matrix. See (Minka, 2000) for more discussions.

Jacobian of a matrix The Jacobian $\partial \mathbf{A} / \partial x$ of a matrix $\mathbf{A} = (A_{ij})$ with respect to a scalar x is defined as a matrix of the same dimensionality as \mathbf{A} so that

$$\frac{\partial \mathbf{A}}{\partial x} \equiv \left(\frac{\partial A_{ij}}{\partial x} \right) \quad (318)$$

which is analogous to (C.18) in that the partial derivatives are laid out according to the numerator, i.e., \mathbf{A} .

Gradient with respect to a matrix On the other hand, the gradient (294) is such that the derivatives are laid out according to the denominator. In a similar analogy, we can define the gradient $\nabla_{\mathbf{A}} y$ of a scalar y with respect to a matrix \mathbf{A} as

$$\nabla_{\mathbf{A}} y \equiv \left(\frac{\partial y}{\partial A_{ij}} \right). \quad (319)$$

³⁶The product rule (311) can also be shown directly. For interested readers, I would like to note that the first term in the right hand side of (311), i.e., the gradient through ϕ , directly follows from (302); and the second term, i.e., the gradient through \mathbf{R} , follows from the identity

$$\mathbf{R} \phi = \text{vec}(\mathbf{R} \phi) = (\phi^T \otimes \mathbf{I}_M) \text{vec}(\mathbf{R}) \quad (313)$$

which itself follows from the property of the vectorization operator and the Kronecker product (Magnus and Neudecker, 2007)

$$\text{vec}(\mathbf{ABC}) = (\mathbf{C}^T \otimes \mathbf{A}) \text{vec}(\mathbf{B}). \quad (314)$$

Taking the gradient of the right hand side of (313), we obtain

$$\nabla \{ (\phi^T \otimes \mathbf{I}_M) \text{vec}(\mathbf{R}) \} = \nabla \text{vec}(\mathbf{R}) (\phi \otimes \mathbf{I}_M) \quad (315)$$

where we have assumed ϕ to be constant and used the transpose identity of the Kronecker product (Magnus and Neudecker, 2007)

$$(\mathbf{A} \otimes \mathbf{B})^T = \mathbf{A}^T \otimes \mathbf{B}^T. \quad (316)$$

Equation (C.22): For this identity to be well-defined, it is necessary that we have $\det(\mathbf{A}) > 0$. We should make this assumption clear. Or, if we adopt the absolute determinant notation (50) for $|\mathbf{A}|$, the identity (C.22) holds, in fact, for any nonsingular \mathbf{A} such that $\det(\mathbf{A}) \neq 0$ as we shall see shortly.

The section named “Eigenvector Equation” of Appendix C gives us a hint for a proof of (C.22) where \mathbf{A} is assumed to be symmetric positive definite so that $\mathbf{A} \succ 0$. Although the restricted proof outlined in PRML is indeed highly instructive, we need a more general proof because we make use of this identity, e.g., in Exercise 2.34 without the assumptions required by the restricted proof.³⁷

Jacobi’s formula To this end, we first show *Jacobi’s formula*, which is an identity that holds for any square matrix \mathbf{A} given by

$$\frac{\partial}{\partial x} \det(\mathbf{A}) = \text{Tr} \left(\mathbf{A}^\dagger \frac{\partial \mathbf{A}}{\partial x} \right) \quad (320)$$

where \mathbf{A}^\dagger is the *adjugate matrix* of \mathbf{A} . The (ij) -th element A_{ij}^\dagger of the adjugate matrix \mathbf{A}^\dagger is given by

$$A_{ij}^\dagger = (-1)^{i+j} \det(\mathbf{A}^{(ji)}) \quad (321)$$

(beware that the superscript (ji) of $\mathbf{A}^{(ji)}$ is *not* (ij) but it is “transposed”) where $\mathbf{A}^{(ij)}$ (the superscript (ij) is *not* “transposed” here) denotes a matrix obtained by removing the i -th row and the j -th column of \mathbf{A} .

From the well-known identity

$$\mathbf{A} \mathbf{A}^\dagger = \mathbf{A}^\dagger \mathbf{A} = \det(\mathbf{A}) \mathbf{I} \quad (322)$$

(consult a linear algebra textbook for a proof), we can write the inverse matrix \mathbf{A}^{-1} in terms of the adjugate matrix \mathbf{A}^\dagger so that

$$\mathbf{A}^{-1} = \frac{\mathbf{A}^\dagger}{\det(\mathbf{A})} \quad (323)$$

if \mathbf{A} is nonsingular so that $\det(\mathbf{A}) \neq 0$. Note also that the above identity (322) implies

$$\det(\mathbf{A}) = \sum_k A_{ik} A_{ki}^\dagger = \sum_k A_{jk}^\dagger A_{kj} \quad (324)$$

for any i and j . Substituting this identity (324) into the left hand side of (320) and noting that, from the definition (321) of the adjugate matrix, A_{ji}^\dagger is independent of A_{ik} nor A_{kj} for

³⁷Note that one can easily extend the restricted proof of (C.22) for symmetric positive-definite matrices \mathbf{A} in terms of the eigenvalue decomposition, outlined in PRML, so as to use instead the *singular value decomposition* or SVD (350) in order to show (C.22) for any nonsingular matrix \mathbf{A} such that $\det(\mathbf{A}) \neq 0$. Here, I would however like to present a proof in terms of *Jacobi’s formula* (320), which is more direct and general. I leave the proof of (C.22) in terms of the SVD as an exercise for the reader. Hint: The right hand side of (C.22) can be written as $\text{Tr}(\mathbf{\Sigma}^{-1} \partial \mathbf{\Sigma} / \partial x)$ where $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ is the SVD of \mathbf{A} because it follows from the *orthonormality* (C.37) of \mathbf{U} that $\text{Tr}(\mathbf{U}^T \partial \mathbf{U} / \partial x) = 0$ and similarly for \mathbf{V} . Finally, note that the absolute determinant of \mathbf{A} is equal to that of $\mathbf{\Sigma}$ and thus to the product of the singular values $\sigma_i > 0$ such that $\mathbf{\Sigma} = \text{diag}(\sigma_i)$, i.e., $|\mathbf{A}| = |\mathbf{\Sigma}| = \prod_i \sigma_i$ where we have used (50).

any k , we have

$$\frac{\partial}{\partial x} \det(\mathbf{A}) = \sum_{ij} \left\{ \frac{\partial}{\partial A_{ij}} \sum_k A_{ik} A_{ki}^\dagger \right\} \frac{\partial A_{ij}}{\partial x} = \sum_{ij} \left\{ \frac{\partial}{\partial A_{ij}} \sum_k A_{jk}^\dagger A_{kj} \right\} \frac{\partial A_{ij}}{\partial x} \quad (325)$$

$$= \sum_{ij} A_{ji}^\dagger \frac{\partial A_{ij}}{\partial x} \quad (326)$$

$$= \text{Tr} \left(\frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^\dagger \right) = \text{Tr} \left(\mathbf{A}^\dagger \frac{\partial \mathbf{A}}{\partial x} \right) \quad (327)$$

which proves the identity (320).

Derivative of log absolute determinant Assuming that \mathbf{A} is nonsingular so that $\det(\mathbf{A}) \neq 0$, we can evaluate the left hand side of (C.22) as

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \frac{1}{\det(\mathbf{A})} \frac{\partial}{\partial x} \det(\mathbf{A}) \quad (328)$$

where we have used the notation (50) for $|\mathbf{A}|$. Substituting (320), we obtain

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left(\frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \right) = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right) \quad (329)$$

where we have used (323).

Page 698

Equations (C.24) through (C.28): Since these derivatives are gradients of a scalar with respect to a matrix \mathbf{A} , the operator $\frac{\partial}{\partial \mathbf{A}}$ should read $\nabla_{\mathbf{A}}$ if we adopt the notation (319).

Matrix derivative identities For example, (C.24) should read

$$\nabla_{\mathbf{A}} \text{Tr}(\mathbf{A}\mathbf{B}) = \nabla_{\mathbf{A}} \text{Tr}(\mathbf{A}^T \mathbf{B}^T) = \mathbf{B}^T \quad (330)$$

where we have used the transpose and the cyclic identities of the trace operator $\text{Tr}(\cdot)$, i.e.,

$$\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^T), \quad \text{Tr}(\mathbf{A}\mathbf{B}) = \text{Tr}(\mathbf{B}\mathbf{A}) \quad (331)$$

respectively. As described in the text, the identity (330) directly follows from (C.23). At this moment, I would like to point out an observation helpful for remembering (330). First, note that the gradient of a scalar with respect to a matrix \mathbf{A} is, by definition (319), a matrix of the same dimensionality as \mathbf{A} . On the other hand, in order for the trace $\text{Tr}(\mathbf{A}\mathbf{B})$ to be meaningful, \mathbf{B} must be of the same dimensionality as \mathbf{A}^T . Thus, (330) passes the “dimensionality test,” meaning that all the matrix operations in (330) are meaningful. Note also that (C.25) and (C.26) are special cases of (330).

Similarly, the gradient of the log (absolute) determinant (C.28) should read

$$\nabla_{\mathbf{A}} \ln |\mathbf{A}| = \mathbf{A}^{-T} \quad (332)$$

where we have used (C.4) and defined

$$\mathbf{A}^{-T} \equiv (\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T. \quad (333)$$

Again, if we adopt the notation (50) for $|\mathbf{A}|$, we see that (332) holds for any nonsingular matrix \mathbf{A} such that $\det(\mathbf{A}) \neq 0$.

The identity (332) can be shown by identifying x with A_{ij} in (329) where A_{ij} is the (ij) -th element of \mathbf{A} ; and then making use of (C.23), which can be stated more suitably for our purpose here as

$$\text{Tr} \left(\frac{\partial \mathbf{A}}{\partial A_{ij}} \mathbf{B} \right) = \text{Tr} \left(\mathbf{B} \frac{\partial \mathbf{A}}{\partial A_{ij}} \right) = B_{ji}. \quad (334)$$

Here, the Jacobian matrix $\partial \mathbf{A} / \partial A_{ij}$ is, by definition (318), a matrix of the same dimensionality as \mathbf{A} such that only the (ij) -th element is one whereas all the other elements are zero, i.e.,

$$\frac{\partial \mathbf{A}}{\partial A_{ij}} = i \begin{pmatrix} & & j \\ & \vdots & \\ \cdots & 1 & \cdots \\ & \vdots & \end{pmatrix} \quad (335)$$

where all the elements omitted or denoted by dots (\cdots) are zero. Note that Svensén and Bishop (2009) effectively make use of (334) in the solution of Exercise 2.34.

In addition to the above mentioned matrix derivative identities, I would suggest to include the following:

$$\nabla_{\mathbf{A}} \text{Tr} (\mathbf{A} \mathbf{B} \mathbf{A}^T \mathbf{C}) = \mathbf{C}^T \mathbf{A} \mathbf{B}^T + \mathbf{C} \mathbf{A} \mathbf{B} \quad (336)$$

$$\nabla_{\mathbf{A}} \text{Tr} (\mathbf{A}^{-1} \mathbf{B}) = -\mathbf{A}^{-T} \mathbf{B}^T \mathbf{A}^{-T}. \quad (337)$$

We use the identities (336) and (337), e.g., when we show (13.113) in Exercise 13.33 and (2.122) in Exercise 2.34, respectively. It should also be noted that (C.27) is a special case of (336).

The identity (336) can be shown as follows. Assuming that \mathbf{B} and \mathbf{C} are constants, we have

$$\frac{\partial}{\partial x} \{ \mathbf{A} \mathbf{B} \mathbf{A}^T \mathbf{C} \} = \frac{\partial \mathbf{A}}{\partial x} \mathbf{B} \mathbf{A}^T \mathbf{C} + \mathbf{A} \mathbf{B} \left(\frac{\partial \mathbf{A}}{\partial x} \right)^T \mathbf{C}. \quad (338)$$

Taking the trace of the both sides gives

$$\frac{\partial}{\partial x} \text{Tr} (\mathbf{A} \mathbf{B} \mathbf{A}^T \mathbf{C}) = \text{Tr} \left(\frac{\partial \mathbf{A}}{\partial x} \mathbf{B} \mathbf{A}^T \mathbf{C} \right) + \text{Tr} \left(\frac{\partial \mathbf{A}}{\partial x} \mathbf{B}^T \mathbf{A}^T \mathbf{C}^T \right) \quad (339)$$

where we have rearranged the factors inside the second trace $\text{Tr}(\cdot)$ in the right hand side by making use of (331). We finally obtain (336) by identifying x with A_{ij} and making use of (334). The identity (337) follows similarly from

$$\frac{\partial}{\partial x} \text{Tr} (\mathbf{A}^{-1} \mathbf{B}) = -\text{Tr} \left(\frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1} \right) \quad (340)$$

which itself follows from (C.21).

Symmetric matrix derivatives So far, we have considered derivatives with respect to a matrix that is not symmetric in general. However, matrices for which we take derivatives in order to, say, perform optimization (e.g., maximum likelihood) or evaluate expectations by making use of (58) are often symmetric. For example, the covariance Σ of the multivariate Gaussian distribution is symmetric positive definite. When we derive the maximum likelihood solution

for Σ in Exercise 2.34, we ignore the symmetry constraint on Σ to calculate the derivatives of the log likelihood with respect to Σ . The maximum likelihood solution Σ_{ML} is obtained by solving necessary conditions that the derivatives should vanish, after which we find Σ_{ML} to be symmetric positive definite. The fact that the solution Σ_{ML} is symmetric is not a fortunate coincidence but a consequence of the symmetry in the necessary conditions solved. In fact, even if we had imposed the symmetry constraint on Σ in the first place, we would have obtained an equivalent set of equations to solve, giving the same solution. We can understand why this is the case by considering derivatives with respect to a symmetric matrix in more general terms as follows.

Let $\phi(\mathbf{A})$ be a scalar function of a square matrix \mathbf{A} where $\mathbf{A} = (A_{ij})$ is not symmetric in general so that $A_{ij} \neq A_{ji}$. As usual, we write the gradient of $\phi(\mathbf{A})$ with respect to \mathbf{A} as

$$\nabla_{\mathbf{A}}\phi(\mathbf{A}) = \left(\frac{\partial}{\partial A_{ij}}\phi(\mathbf{A}) \right). \quad (341)$$

Suppose that we want to evaluate the gradient of $\phi(\mathbf{S})$ where $\mathbf{S} = (S_{ij})$ is symmetric so that $S_{ij} \equiv S_{ji}$. The derivative of $\phi(\mathbf{S})$ with respect to an off-diagonal element S_{ij} where $i \neq j$ consists of two derivatives through A_{ij} and A_{ji} so that

$$\frac{\partial}{\partial S_{ij}}\phi(\mathbf{S}) = \frac{\partial}{\partial A_{ij}}\phi(\mathbf{S}) + \frac{\partial}{\partial A_{ji}}\phi(\mathbf{S}) \quad (342)$$

where we have written

$$\frac{\partial}{\partial A_{ij}}\phi(\mathbf{S}) \equiv \frac{\partial}{\partial A_{ij}}\phi(\mathbf{A}) \Big|_{\mathbf{A}=\mathbf{S}}. \quad (343)$$

The derivative of $\phi(\mathbf{S})$ with respect to a diagonal element S_{ii} is given by

$$\frac{\partial}{\partial S_{ii}}\phi(\mathbf{S}) = \frac{\partial}{\partial A_{ii}}\phi(\mathbf{S}). \quad (344)$$

Thus, we can write

$$\nabla_{\mathbf{S}}\phi(\mathbf{S}) = \nabla_{\mathbf{A}}\phi(\mathbf{S}) + \nabla_{\mathbf{A}}\phi(\mathbf{S})^T - \text{diag}(\nabla_{\mathbf{A}}\phi(\mathbf{S})) \quad (345)$$

where we have written

$$\nabla_{\mathbf{A}}\phi(\mathbf{S}) \equiv \nabla_{\mathbf{A}}\phi(\mathbf{A}) \Big|_{\mathbf{A}=\mathbf{S}}. \quad (346)$$

For example, if \mathbf{A} and \mathbf{B} are both symmetric in (330), we have

$$\nabla_{\mathbf{A}} \text{Tr}(\mathbf{AB}) = 2\mathbf{B} - \text{diag}(\mathbf{B}). \quad (347)$$

The identity (345) is, however, not very useful in practice. A more useful observation can be made by considering equations obtained by setting the derivatives equal to zero. Specifically, it readily follows from (342) and (344) that, if $\nabla_{\mathbf{A}}\phi(\mathbf{S})$ is symmetric (which does hold, say, for the necessary conditions for Σ_{ML} we mentioned above), we have

$$\nabla_{\mathbf{S}}\phi(\mathbf{S}) = \mathbf{O} \iff \nabla_{\mathbf{A}}\phi(\mathbf{S}) = \mathbf{O} \quad (348)$$

which implies that we can solve $\nabla_{\mathbf{S}}\phi(\mathbf{S}) = \mathbf{O}$ without the symmetry constraint on \mathbf{S} , i.e., by simply solving $\nabla_{\mathbf{A}}\phi(\mathbf{S}) = \mathbf{O}$ and then obtain a solution \mathbf{S} that is indeed symmetric.

When we evaluate expectations by making use of (58), we consider equations obtained by setting the expected derivatives equal to zero. With much the same discussion as above, if $\mathbb{E} [\nabla_{\mathbf{A}} \phi(\mathbf{S})]$ is symmetric, we have

$$\mathbb{E} [\nabla_{\mathbf{S}} \phi(\mathbf{S})] = \mathbf{O} \iff \mathbb{E} [\nabla_{\mathbf{A}} \phi(\mathbf{S})] = \mathbf{O}. \quad (349)$$

It should be noted here that the score function (57) occurring in (58) is defined only for independent parameters. Therefore, if the parameters of interest are, say, a symmetric matrix (e.g., the scale matrix \mathbf{W} of the Wishart distribution is symmetric positive definite), we must, strictly speaking, impose the symmetry constraint on the parameters. The equivalence relation (349), however, allows us to safely ignore the symmetry constraint on \mathbf{S} and use $\mathbb{E} [\nabla_{\mathbf{A}} \phi(\mathbf{S})] = \mathbf{O}$ instead.

Page 700

Paragraph 2, Line -1: The determinant of the orthogonal matrix \mathbf{U} can be either positive or negative so that we should write $\det(\mathbf{U}) = \pm 1$ (which is, if the notation (50) is adopted, equivalent to $|\mathbf{U}| = 1$). Although it is possible to take \mathbf{U} such that $\det(\mathbf{U}) = 1$ (one can flip the sign of $\det(\mathbf{U})$ by, say, flipping the sign of any one of the eigenvectors $\{\mathbf{u}_i\}$), there is no point in doing so in practice theoretically nor numerically. In fact, it is easy to see that the following discussion remains valid provided that \mathbf{U} is orthogonal so that we have (C.37) but not necessarily that $\det(\mathbf{U}) = 1$. Moreover, most software implementations of symmetric eigenvalue decomposition only guarantee that \mathbf{U} is orthogonal so that $\det(\mathbf{U}) = \pm 1$.

Singular value decomposition In the special case where the matrix \mathbf{A} is symmetric positive semidefinite or $\mathbf{A} \succeq 0$, we can identify the eigenvalue decomposition (C.43) with the *singular value decomposition* or SVD (Press et al., 1992; Golub and Van Loan, 2013) so that we can use an SVD routine to compute the eigenvalue decomposition of \mathbf{A} . The SVD is generally defined for any real matrix \mathbf{P} not necessarily square, say, of dimensionality $M \times N$, so that the SVD of \mathbf{P} is given by

$$\mathbf{P} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \sum_{i=1}^R \sigma_i \mathbf{u}_i \mathbf{v}_i^T \quad (350)$$

where $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_M)$ and $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_N)$ are orthogonal matrices of dimensionalities $M \times M$ and $N \times N$, respectively; $\mathbf{\Sigma}$ is an $M \times N$ diagonal matrix with nonnegative diagonal elements, called the *singular values*, $\sigma_1 \geq \dots \geq \sigma_R \geq 0$ arranged in descending order; and $R \leq \min(M, N)$ is the rank of \mathbf{P} . Note again that \mathbf{U} and \mathbf{V} are only guaranteed to be orthogonal so that $\det(\mathbf{U}) = \pm 1$ and $\det(\mathbf{V}) = \pm 1$.

Page 700

The text following (C.41): The multiplication by \mathbf{U} can be interpreted as a rotation, a reflection, or a combination of the two.

Equation (D.8): It would be helpful if we make it clear that the left hand side of (D.8) corresponds to the functional derivative so that we should modify (D.8) as

$$\frac{\delta F}{\delta y(x)} \equiv \frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) = 0. \quad (351)$$

Paragraph –1, Line 1: Despite the statement, it is not that straightforward to extend the results obtained here to higher dimensions. Although such an extension is not required in PRML, it is useful when we analyze a particular type of constrained optimization problem commonly found in computer vision applications such as *optical flow* (Horn and Schunck, 1981). Here, I would like to consider an extension of the calculus of variations to a system of D -dimensional Cartesian coordinates $\mathbf{x} = (x_1, \dots, x_D)^T \in \mathbb{R}^D$ and find the form of the functional derivative as well as a more general boundary condition for such a derivative to be well-defined. To this end, we first review some identities concerning the *divergence* (Feynman et al., 1964). The divergence of a vector field

$$\mathbf{p}(\mathbf{x}) = \begin{pmatrix} p_1(\mathbf{x}) \\ \vdots \\ p_D(\mathbf{x}) \end{pmatrix} \in \mathbb{R}^D \quad (352)$$

is a scalar field of the form

$$\text{div } \mathbf{p} = \sum_{i=1}^D \frac{\partial p_i}{\partial x_i} \equiv \nabla \cdot \mathbf{p} \quad (353)$$

where we have omitted the coordinates \mathbf{x} in the function arguments to keep the notation uncluttered. For a differentiable vector field $\mathbf{p}(\mathbf{x})$ defined on some volume $\Omega \subset \mathbb{R}^D$, the *divergence theorem* (Feynman et al., 1964) states that

$$\int_{\Omega} \text{div } \mathbf{p} \, dV = \oint_{\partial\Omega} \mathbf{p} \cdot \mathbf{n} \, dS \quad (354)$$

where the left hand side is the volume integral over the volume Ω ; the right hand side is the surface integral over its boundary $\partial\Omega$; and $\mathbf{n}(\mathbf{x})$ is the outward unit normal vector of $\partial\Omega$. Assuming that the coordinates $\mathbf{x} = (x_1, \dots, x_D)^T$ are Cartesian, we can write the volume element as $dV = dx_1 \cdots dx_D \equiv d\mathbf{x}$ and the inner product as $\mathbf{p} \cdot \mathbf{n} = \mathbf{p}^T \mathbf{n}$. Making use of the divergence theorem (354) together with the following identity

$$\text{div}(\phi \mathbf{p}) = \nabla \phi^T \mathbf{p} + \phi \text{div } \mathbf{p} \quad (355)$$

we obtain a multidimensional version of the “integration by parts” formula

$$\int_{\Omega} \nabla \phi^T \mathbf{p} \, d\mathbf{x} = \oint_{\partial\Omega} \phi \mathbf{p}^T \mathbf{n} \, dS - \int_{\Omega} \phi \text{div } \mathbf{p} \, d\mathbf{x}. \quad (356)$$

Let us now consider a functional of the form

$$E[u(\mathbf{x})] = \int_{\Omega} L(\mathbf{x}, u(\mathbf{x}), \nabla u(\mathbf{x})) \, d\mathbf{x} \quad (357)$$

where $u(\mathbf{x}) \in \mathbb{R}$ is a function (scalar field) defined over some volume $\Omega \subset \mathbb{R}^D$ and $L(\mathbf{x}, f, \mathbf{g}) \in \mathbb{R}$ is a function of $\mathbf{x} \in \Omega$, $f \in \mathbb{R}$, and $\mathbf{g} \in \mathbb{R}^D$. Thus, the functional $E[u(\mathbf{x})] \in \mathbb{R}$ maps $u(\mathbf{x})$ to a real number. As in the ordinary calculus, we can define the derivative of a functional according to the *calculus of variations* (Feynman et al., 1964; Bishop, 2006). In order to find the form of the functional derivative, we consider how $E[u(\mathbf{x})]$ varies upon a small change $\epsilon\eta(\mathbf{x})$ in $u(\mathbf{x})$ where $\eta(\mathbf{x})$ is the “direction” of the change and ϵ is some small constant. The first-order variation of $E[u(\mathbf{x})]$ in the direction of $\eta(\mathbf{x})$ can be evaluated as

$$\delta E[u; \eta] \equiv \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \{E[u + \epsilon\eta] - E[u]\} \quad (358)$$

$$= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \int_{\Omega} \{L(\mathbf{x}, u + \epsilon\eta, \nabla(u + \epsilon\eta)) - L(\mathbf{x}, u, \nabla u)\} d\mathbf{x} \quad (359)$$

$$= \int_{\Omega} \left\{ \eta \frac{\partial L}{\partial f} + \nabla \eta^T \nabla_{\mathbf{g}} L \right\} d\mathbf{x} \quad (360)$$

where we have assumed that $L(\mathbf{x}, f, \mathbf{g})$ is differentiable with respect to both f and \mathbf{g} ; and we have written

$$\frac{\partial L}{\partial f} \equiv \frac{\partial}{\partial f} L(\mathbf{x}, u, \nabla u), \quad \nabla_{\mathbf{g}} L \equiv \nabla_{\mathbf{g}} L(\mathbf{x}, u, \nabla u). \quad (361)$$

By making use of the multidimensional integration by parts (356), we can integrate the second term in the right hand side of (360), giving

$$\delta E[u; \eta] = \int_{\Omega} \eta \left\{ \frac{\partial L}{\partial f} - \text{div}(\nabla_{\mathbf{g}} L) \right\} d\mathbf{x} + \oint_{\partial\Omega} \eta \nabla_{\mathbf{g}} L^T \mathbf{n} dS. \quad (362)$$

In order for the functional derivative to be well-defined, we assume the surface integral term in the variation (362) to vanish so that we have the following boundary condition

$$\oint_{\partial\Omega} \eta \nabla_{\mathbf{g}} L^T \mathbf{n} dS = 0. \quad (363)$$

The boundary condition (363) holds if

$$\eta(\mathbf{x}) = 0 \quad (364)$$

or

$$\nabla_{\mathbf{g}} L^T \mathbf{n}(\mathbf{x}) = 0 \quad (365)$$

for all $\mathbf{x} \in \partial\Omega$. The first condition (364) holds if we assume the *Dirichlet boundary condition* for $u(\mathbf{x})$

$$u(\mathbf{x}) = u_0(\mathbf{x}) \quad (366)$$

where $\mathbf{x} \in \partial\Omega$, i.e., $u(\mathbf{x})$ is assumed to be fixed to some value $u_0(\mathbf{x})$ at the boundary $\partial\Omega$ and so is $u(\mathbf{x}) + \epsilon\eta(\mathbf{x})$ in (358), implying (364). Another common boundary condition for $u(\mathbf{x})$ is the *Neumann boundary condition*

$$\nabla u(\mathbf{x})^T \mathbf{n}(\mathbf{x}) = 0 \quad (367)$$

where $\mathbf{x} \in \partial\Omega$. The Neumann boundary condition (367) is implied by the second condition (365) for the optical-flow energy functional as we shall see shortly. Having assumed that the boundary condition (363) holds, we can write the first order variation (362) in the form

$$\delta E[u; \eta] = \int_{\Omega} \eta \frac{\partial E}{\partial u(\mathbf{x})} d\mathbf{x} \quad (368)$$

where we have written

$$\frac{\partial E}{\partial u(\mathbf{x})} \equiv \frac{\partial L}{\partial f} - \operatorname{div}(\nabla_{\mathbf{g}} L). \quad (369)$$

The volume integral in the right hand side of (368) can be seen as the inner product between $\eta(\mathbf{x})$ and $\partial E/\partial u(\mathbf{x})$, from which we conclude that the quantity $\partial E/\partial u(\mathbf{x})$ is what should be called the functional derivative.³⁸ A stationary point of a functional $E[u(\mathbf{x})]$ is a function $u(\mathbf{x})$ such that the variation $\delta E[u; \eta]$ vanishes in any direction $\eta(\mathbf{x})$ and thus satisfies the *Euler-Lagrange equation* given by

$$\frac{\partial E}{\partial u(\mathbf{x})} = 0. \quad (370)$$

Finally, we present an application of the multidimensional calculus of variations to a dense motion analysis technique called optical flow in the following. Suppose that, given a pair of (grayscale) images $I_0(\mathbf{x})$ and $I_1(\mathbf{x})$ where $\mathbf{x} \in \mathbb{R}^2$ that are taken at some discrete time steps $t = 0$ and $t = 1$, respectively, we wish to find a motion vector field from $I_0(\mathbf{x})$ to $I_1(\mathbf{x})$

$$\mathbf{u}(\mathbf{x}) = \begin{pmatrix} u(\mathbf{x}) \\ v(\mathbf{x}) \end{pmatrix} \quad (371)$$

defined over $\mathbf{x} \in \Omega \subset \mathbb{R}^2$. [Horn and Schunck \(1981\)](#) sought for $\mathbf{u}(\mathbf{x})$ that minimizes an energy functional that takes essentially the same form as

$$J[\mathbf{u}(\mathbf{x})] = J_{\text{data}}[\mathbf{u}(\mathbf{x})] + \alpha J_{\text{smooth}}[\mathbf{u}(\mathbf{x})] \quad (372)$$

where

$$J_{\text{data}}[\mathbf{u}(\mathbf{x})] = \frac{1}{2} \int_{\Omega} (I_1(\mathbf{x} + \mathbf{u}(\mathbf{x})) - I_0(\mathbf{x}))^2 d\mathbf{x} \quad (373)$$

$$J_{\text{smooth}}[\mathbf{u}(\mathbf{x})] = \frac{1}{2} \int_{\Omega} (\|\nabla u(\mathbf{x})\|^2 + \|\nabla v(\mathbf{x})\|^2) d\mathbf{x}. \quad (374)$$

Here, the domain Ω is assumed to be continuous and is typically rectangular. We call the first term $J_{\text{data}}[\mathbf{u}(\mathbf{x})]$ in (372) the data-fidelity term; the second term $J_{\text{smooth}}[\mathbf{u}(\mathbf{x})]$ the smoothness (regularization) term; and the coefficient α the regularization parameter. According to the multidimensional calculus of variations, a stationary point of the optical-flow energy functional (372) satisfies Euler-Lagrange equations of the form

$$\nabla_{\mathbf{u}(\mathbf{x})} J \equiv \begin{pmatrix} \partial J / \partial u(\mathbf{x}) \\ \partial J / \partial v(\mathbf{x}) \end{pmatrix} = \nabla_{\mathbf{u}} \left\{ \frac{\varepsilon(\mathbf{x}, \mathbf{u}(\mathbf{x}))^2}{2} \right\} - \alpha \begin{pmatrix} \operatorname{div}(\nabla u(\mathbf{x})) \\ \operatorname{div}(\nabla v(\mathbf{x})) \end{pmatrix} = \mathbf{0} \quad (375)$$

where we have written

$$\varepsilon(\mathbf{x}, \mathbf{u}) = I_1(\mathbf{x} + \mathbf{u}) - I_0(\mathbf{x}). \quad (376)$$

For the functional derivatives $\partial J/\partial u(\mathbf{x})$ and $\partial J/\partial v(\mathbf{x})$ to be well-defined, let us assume the boundary condition given by (365) for each functional derivative, which implies the Neumann boundary condition for $\mathbf{u}(\mathbf{x})$, i.e.,

$$\nabla u(\mathbf{x})^T \mathbf{n}(\mathbf{x}) = 0, \quad \nabla v(\mathbf{x})^T \mathbf{n}(\mathbf{x}) = 0 \quad (377)$$

³⁸Here we use a notation for the functional derivative that is different from the one used in PRML. The notation $\partial E/\partial u(\mathbf{x})$ employed here is more like an ordinary derivative and can be extended to the case of a vector field $\mathbf{u}(\mathbf{x})$ analogously to the gradient as we shall see in (375).

for all $\mathbf{x} \in \partial\Omega$ where $\partial\Omega$ is the boundary of Ω and $\mathbf{n}(\mathbf{x})$ is the outward unit normal vector of $\partial\Omega$. Thus, solving the above Euler-Lagrange equations (375) with the Neumann boundary condition (377), we obtain the desired motion vector field $\mathbf{u}(\mathbf{x})$. The Euler-Lagrange equations given by (375) are *elliptic partial differential equations* (elliptic PDEs) and can be solved numerically by a type of relaxation method such as the Gauss-Seidel method or the (weighted) Jacobi method or by a more efficient *multigrid* technique (Press et al., 1992; Briggs et al., 2000).

Page 708

Equation (E.3): The right hand side should be a zero vector $\mathbf{0}$ instead of a scalar zero 0.

Page 708

The text after (E.4): $\nabla_{\mathbf{x}}L = 0$ should read $\nabla_{\mathbf{x}}L = \mathbf{0}$ (the right hand side should be a zero vector $\mathbf{0}$).

Page 709

Paragraph -2, Line 5: $\nabla f(\mathbf{x}) = 0$ should read $\nabla f(\mathbf{x}) = \mathbf{0}$ (the right hand side should be a zero vector $\mathbf{0}$).

Page 716

Column 1, Entry -1: “The Feynman Lectures of Physics” should read “The Feynman Lectures on Physics.”

Page 717

Column 2, Entry 7: “John Hopkins University Press” should read “The Johns Hopkins University Press.”

References

- Abramowitz, M. and I. A. Stegun (Eds.) (1964). *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. National Bureau of Standards, U.S. Department of Commerce. 5, 23
- Anderson, T. W. (2003). *An Introduction to Multivariate Statistical Analysis* (Third ed.). Wiley. 53, 56, 59, 61
- Anderson, T. W. and I. Olkin (1985). Maximum-likelihood estimation of the parameters of a multivariate normal distribution. *Linear Algebra and Its Applications* 70, 147–171. 15
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer. 1, 74
- Briggs, W. L., V. E. Henson, and S. F. McCormick (2000). *A Multigrid Tutorial* (Second ed.). SIAM. 76

- Feynman, R. P., R. B. Leighton, and M. Sands (1964). *The Feynman Lectures on Physics*, Volume 2. Addison-Wesley. 73, 74
- Golub, G. H. and C. F. Van Loan (2013). *Matrix Computations* (Fourth ed.). The Johns Hopkins University Press. 13, 56, 72
- Horn, B. K. P. and B. G. Schunck (1981). Determining optical flow. *Artificial Intelligence* 17(1), 185–203. 73, 75
- Kullback, S. and R. A. Leibler (1951). On information and sufficiency. *The Annals of Mathematical Statistics* 22(1), 79–86. 7
- MacKay, D. J. C. (2003). *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press. 6, 7
- Magnus, J. R. and H. Neudecker (2007). *Matrix Differential Calculus with Applications in Statistics and Econometrics* (Third ed.). Wiley. <http://www.janmagnus.nl/misc/mdc2007-3rdedition>. 15, 66, 67
- Mermin, N. D. (1989). What’s wrong with these equations? *Physics Today* 42(10), 9–11. 29
- Minka, T. (2005). Divergence measures and message passing. Technical Report MSR-TR-2005-173, Microsoft Research. <https://www.microsoft.com/en-us/research/publication/divergence-measures-and-message-passing/>. 9
- Minka, T. P. (2000). Old and new matrix algebra useful for statistics. <https://tminka.github.io/papers/matrix/>. 63, 64, 65, 67
- Minka, T. P. (2001). Expectation propagation for approximate bayesian inference. In J. Breese and D. Koller (Eds.), *Proceedings of the Seventeenth International Conference on Artificial Intelligence*, pp. 362–369. Morgan Kaufmann. <https://tminka.github.io/papers/ep/>. 41
- Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press. 50
- Olver, F. W. J., A. B. Olde Daalhuis, D. W. Lozier, B. I. Schneider, R. F. Boisvert, C. W. Clark, B. R. Miller, and B. V. Saunders (Eds.) (2017). *NIST Digital Library of Mathematical Functions*. National Institute of Standards and Technology, U.S. Department of Commerce. <http://dlmf.nist.gov/> (Release 1.0.15 of 2017-06-01). 2, 23, 59
- Press, W. M., S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery (1992). *Numerical Recipes in C: The Art of Scientific Computing* (Second ed.). Cambridge University Press. 22, 28, 56, 72, 76
- Svensén, M. and C. M. Bishop (2005). Robust Bayesian mixture modelling. *Neurocomputing* 64, 235–252. 53
- Svensén, M. and C. M. Bishop (2009). Pattern Recognition and Machine Learning: Solutions to the exercises (web edition). <https://www.microsoft.com/en-us/research/people/cmbishop/#prml-book>. 38, 70

- Svensén, M. and C. M. Bishop (2011). Pattern Recognition and Machine Learning: Errata and additional comments. <https://www.microsoft.com/en-us/research/people/cmbishop/#prml-book>. 1, 36, 37, 41, 44, 46
- Szeliski, R. (2010). *Computer Vision: Algorithms and Applications*. Springer. 22, 56
- Tao, T. (2011). *An Introduction to Measure Theory*. American Mathematical Society. 2, 6
- Tipping, M. E. and A. C. Faul (2003). Fast marginal likelihood maximisation for sparse Bayesian models. In C. M. Bishop and B. J. Frey (Eds.), *Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics*, Key West, Florida. 32
- Watanabe, S. (2010). Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of Machine Learning Research* 11, 3571–3594. 25
- Watanabe, S. (2013). A widely applicable Bayesian information criterion. *Journal of Machine Learning Research* 14, 867–897. 25
- Zhu, H. and R. Rohwer (1995). Information geometric measurements of generalisation. Technical Report NCRG/4350, Aston University. <http://publications.aston.ac.uk/507/>. 9