

1

# PRML Solutions

2

Youming Zhao  
Email: [youming0.zhao@gmail.com](mailto:youming0.zhao@gmail.com)

3

First draft: January 24, 2023    Last update: October 31, 2023

4

## Contents

5

<b>1</b>	<b>Introduction</b>	<b>2</b>
----------	---------------------	----------

6

1.1	Exercises	2
-----	-----------	---

7

<b>2</b>	<b>Chapter 4   Linear Models for Classification</b>	<b>28</b>
----------	---	-----------

8

2.1	Discriminant Functions	28
-----	------------------------	----

9

2.1.1	The derivation of Equation (4.5)	28
-------	----------------------------------	----

10

2.2	Exercises	29
-----	-----------	----

11

## Notations

$$\binom{\alpha}{k} = \frac{\alpha(\alpha-1)(\alpha-2)\cdots(\alpha-k+1)}{k!} \quad (k=1,2,\dots)$$
$$\binom{\alpha}{0} = 1$$

where  $\alpha$  is a nonzero real number. Note that in combinatorics  $\alpha$  is usually a positive integer  $n$ , i.e.,  $\binom{n}{k}$  which is also denoted as  $C_n^k$  with  $C_n^0 = 1$ . In this case, we have

$$C_n^k = \frac{n!}{k!(n-k)!}$$

# 1 Introduction

## 1.1 Exercises

### Exercise 1.1

Consider the sum-of-squares error function given by (1.2),

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

in which the function  $y(x, \mathbf{w})$  is given by the polynomial (1.1),

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \cdots + w_M x^M = \sum_{j=0}^M w_j x^j. \quad (2)$$

Show that the coefficients  $\mathbf{w} = \{w_i\}$  that minimize this error function are given by the solution to the following set of linear equations

$$\sum_{j=0}^M A_{ij} w_j = T_i \quad (3)$$

where

$$A_{ij} = \sum_{n=1}^N (x_n)^{i+j}, \quad T_i = \sum_{n=1}^N (x_n)^i t_n. \quad (4)$$

Here a suffix  $i$  or  $j$  denotes the index of a component, whereas  $(x)^i$  denotes  $x$  raised to the power of  $i$ .

*Proof.* Since  $E(\mathbf{w})$  is a quadratic function, it follows that  $E(\mathbf{w})$  is convex with respect to  $\mathbf{w}$ . Additionally,  $E(\mathbf{w})$  is lower bounded by 0 and its feasible set is the entire space  $\mathbb{R}^M$ , which is convex as well. Hence, the minimum of  $E(\mathbf{w})$  can be achieved at its stationary points  $\mathbf{w}^*$ , i.e.  $\nabla_{\mathbf{w}^*}(E(\mathbf{w}^*)) = \mathbf{0}$ .

We will denote by  $\mathbf{x}$  and  $\mathbf{t}$  the column vectors  $(1, x, x^2, \dots, x^M)^T$  and  $(t_1, t_2, \dots, t_N)^T$ , respectively. Furthermore, we can combine the observations  $\{\mathbf{x}_n\}$  into a data matrix  $\mathbf{X}$  in which the  $n^{\text{th}}$  row of  $\mathbf{X}$  corresponds to the row vector  $\mathbf{x}_n^T$ . Then, we can get the following compact formulations,

$$y(x, \mathbf{w}) = \mathbf{w}^T \mathbf{x}, \quad E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \mathbf{x}_n - t_n\}^2 = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{t})^T (\mathbf{X}\mathbf{w} - \mathbf{t}) \quad (5)$$

Taking gradients of  $E(\mathbf{w})$  w.r.t.  $\mathbf{w}$  gives

$$\nabla_{\mathbf{w}}(E(\mathbf{w})) = \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{t}) = \mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{t} \quad (6)$$

Setting  $\nabla_{\mathbf{w}}(E(\mathbf{w})) = \mathbf{0}$  yields  $\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{t}$ . By expanding this compact result, we get

$$\begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \\ x_1^2 & x_2^2 & \cdots & x_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^M & x_2^M & \cdots & x_n^M \end{pmatrix} \begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^M \\ 1 & x_2 & x_2^2 & \cdots & x_2^M \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^M \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{pmatrix} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \\ x_1^2 & x_2^2 & \cdots & x_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^M & x_2^M & \cdots & x_n^M \end{pmatrix} \begin{pmatrix} t_0 \\ t_1 \\ \vdots \\ t_N \end{pmatrix} \quad (7)$$

$$\Downarrow$$

$$\begin{pmatrix} \sum_{n=1}^N 1^{0+0} & \sum_{n=1}^N x_n^{0+1} & \cdots & \sum_{n=1}^N x_n^{0+M} \\ \sum_{n=1}^N x_n^{1+0} & \sum_{n=1}^N x_n^{1+1} & \cdots & \sum_{n=1}^N x_n^{1+M} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{n=1}^N x_n^{M+0} & \sum_{n=1}^N x_n^{M+1} & \cdots & \sum_{n=1}^N x_n^{M+M} \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{pmatrix} = \begin{pmatrix} \sum_{j=0}^M \sum_{n=0}^N x_n^{0+j} w_j \\ \sum_{j=0}^M \sum_{n=0}^N x_n^{1+j} w_j \\ \vdots \\ \sum_{j=0}^M \sum_{n=0}^N x_n^{M+j} w_j \end{pmatrix} = \begin{pmatrix} \sum_{n=1}^N t_n \\ \sum_{n=1}^N x_n t_n \\ \vdots \\ \sum_{n=1}^N x_n^M t_n \end{pmatrix} \quad (8)$$

$$\Downarrow$$

$$\mathbf{A}\mathbf{w} = \mathbf{T} \quad (9)$$

19 where  $A_{ij} = \sum_{n=1}^N x_n^{i+j}$  and  $\mathbf{T}$  is a column vector with elements  $T_i = \sum_{n=1}^N x_n^i t_n$  for  $i, j = 0, 1, \dots, M$ .  
 20 Note that we omitted the brackets around  $x_n$  for notational brevity. This completes the proof.  $\square$

### Exercise 1.2

Write down the set of coupled linear equations, analogous to (3) ((1.122) in PRML), satisfied by the coefficients  $w_i$  which minimize the regularized sum-of-squares error function given by ((1.4) in PRML)

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2. \quad (10)$$

**Solution:** Taking gradients of  $\tilde{E}(\mathbf{w})$  w.r.t.  $\mathbf{w}$  gives

$$\nabla_{\mathbf{w}}(\tilde{E}(\mathbf{w})) = \mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{t}) + \lambda\mathbf{w} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})\mathbf{w} - \mathbf{X}^T\mathbf{t}. \quad (11)$$

22 Setting  $\tilde{E}(\mathbf{w}) = \mathbf{0}$  yields  $(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})\mathbf{w} = \mathbf{X}^T\mathbf{t}$ . Thus, we get  $\tilde{\mathbf{A}}\mathbf{w} = \mathbf{T}$  where  $T_i = \sum_{n=1}^N x_n^i t_n$   
 23 is identical to the counterpart in Exercise 1.1. Following a similar argument, we obtain  $\tilde{A}_{ij} =$   
 24  $\sum_{n=1}^N (x_n^{i+j} + \lambda\delta_{ij})$  where  $\delta_{ij} = 1$  when  $i = j$  otherwise  $\delta_{ij} = 0$ .  $\square$

### Exercise 1.3

Suppose that you have three coloured boxes  $r$  (red),  $b$  (blue), and  $g$  (green). Box  $r$  contains 3 apples, 4 oranges, and 3 limes, box  $b$  contains 1 apple, 1 orange, and 0 limes, and box  $g$  contains 3 apples, 3 oranges, and 4 limes. If a box is chosen at random with probabilities  $p(r) = 0.2$ ,  $p(b) = 0.2$ ,  $p(g) = 0.6$ , and a piece of fruit is removed from the box with equal probability of selecting any of the items in the box), then what is the probability of selecting an apple? If we observe that the selected fruit is in fact an orange, what is the probability that it came from the green box?

**Solution:** The probabilities of selecting an apple from the red, the blue, or the green box are given by

$$p(F = a|r) = \frac{3}{3+4+3} = 0.3 \quad (12)$$

$$p(F = a|b) = \frac{1}{1+1} = 0.5 \quad (13)$$

$$p(F = a|g) = \frac{3}{3+3+4} = 0.3 \quad (14)$$

respectively. We use the sum and product rules of probability to evaluate the probability of selecting an apple.

$$p(F = a) = p(F = a|r)p(r) + p(F = a|b)p(b) + p(F = a|g)p(g) \quad (15)$$

$$= 0.3 \times 0.2 + 0.5 \times 0.2 + 0.3 \times 0.6 = 0.06 + 0.1 + 0.18 \quad (16)$$

$$= 0.34 \quad (17)$$

By the Bayes' Theorem, the probability of a selected orange that came from the green box is

$$p(g|F = o) = \frac{p(F = o|g)p(g)}{p(F = o)} \quad (18)$$

$$= \frac{p(F = o|g)p(g)}{p(F = o|r)p(r) + p(F = o|b)p(b) + p(F = o|g)p(g)} \quad (19)$$

$$= \frac{0.3 \times 0.6}{0.4 \times 0.2 + 0.5 \times 0.2 + 0.3 \times 0.6} = \frac{0.18}{0.08 + 0.1 + 0.18} \quad (20)$$

$$= 0.5 \quad (21)$$

26

□

#### Exercise 1.4

Consider a probability density  $p_x(x)$  defined over a continuous variable  $x$ , and suppose that we make a nonlinear change of variable using  $x = g(y)$ , so that the density transforms according to ((1.27) in PRML book)

$$p_y(y) = p_x(x)|g'(y)|. \quad (22)$$

By differentiating (22), show that the location  $\hat{y}$  of the maximum of the density in  $y$  is not in general related to the location  $\hat{x}$  of the maximum of the density over  $x$  by the simple functional relation  $\hat{x} = g(\hat{y})$  as a consequence of the Jacobian factor. This shows that the maximum of a probability density (in contrast to a simple function) is dependent on the change of variable. Verify that, in the case of a linear transformation, the location of the maximum transforms in the same way as the variable itself.

27

28

The proof below follows the same logic as the official solution.

*Proof.* Given a function  $f(x)$  and the relation  $x = g(y)$ , we can get a new function

$$\tilde{f}(y) = f(g(y)). \quad (23)$$

Suppose  $f(x)$  achieves its maximum at  $\hat{x}$  so that  $f'(\hat{x}) = 0$ . The corresponding maximum  $\tilde{f}(\hat{y})$  will be obtained by differentiating both sides of (23) w.r.t  $y$

$$\tilde{f}'(\hat{y}) = f'(g(\hat{y}))g'(\hat{y}) = 0. \quad (24)$$

29 Assuming  $g'(\hat{y}) \neq 0$  at the maximum  $\tilde{f}(\hat{y})$ , then  $\tilde{f}'(g(\hat{y})) = 0$ . Since  $f'(\hat{x}) = 0$ , we see that the  
30 locations of the maximum are related by  $\hat{x} = g(\hat{y})$ . Thus, finding a maximum w.r.t  $x$  is equivalent to  
31 first transforming to  $y$ , and then find a maximum w.r.t  $y$ , and then transforming back to  $x$ .

Now consider the behavior of a probability density  $p_x(x)$  under the change of variables  $x = g(y)$ . According to (22), the new density is given by

$$p_y(y) = p_x(x) \left| \frac{dx}{dy} \right| = p_x(g(y))|g'(y)|. \quad (25)$$

Let  $|g'(y)| = sg'(y)$  where  $s \in \{-1, 1\}$ , then

$$p_y(y) = sp_x(g(y))g'(y). \quad (26)$$

Differentiating both sides w.r.t  $y$  yields

$$p'_y(y) = sp'_x(g(y))(g'(y))^2 + sp_x(g(y))g''(y). \quad (27)$$

Due to the presence of the second term on the right hand side of (27), the result  $\hat{x} = g(\hat{y})$  no longer holds. This implies that we can not get the maximum of  $p_x(x)$  by simply transforming to  $p_y(y)$  then maximizing w.r.t  $y$  and then transforming back to  $x$ . In other words, maxima of densities are dependent on the choice of variables. From the above analyses, we see that this is exactly the consequence of the Jacobian factor  $|g'(y)|$ .

In the case of linear transformation,  $g''(y)$  vanishes and  $g'(y)$  is a constant denoted  $c$ , then we have

$$p'_y(y) = sc^2 p'_x(g(y)). \quad (28)$$

which implies  $p'_y(\hat{y}) = p'_x(g(\hat{y})) = p'_x(\hat{x}) = 0$  at the stationarity  $\hat{y}$ . Thus, the location of the maximum transforms according to  $\hat{x} = g(\hat{y})$ . This completes the proof.  $\square$

#### Exercise 1.5

Using the definition ((1.38) in PRML book)

$$\text{var}[f] = \mathbb{E} [(f(x) - \mathbb{E}[f(x)])^2] \quad (29)$$

show that  $\text{var}[f(x)]$  satisfies ((1.39) in PRML book)

$$\text{var}[f] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2. \quad (30)$$

*Proof.* Expanding the right hand side of (29) gives

$$\text{var}[f] = \mathbb{E} [f(x)^2 - 2\mathbb{E}[f(x)]f(x) + \mathbb{E}[f(x)]^2] \quad (31)$$

$$= \mathbb{E}[f(x)^2] - 2\mathbb{E}[f(x)]\mathbb{E}[f(x)] + \mathbb{E}[f(x)]^2 \quad (32)$$

$$= \mathbb{E}[f(x)^2] - 2\mathbb{E}[f(x)]^2 + \mathbb{E}[f(x)]^2 \quad (33)$$

$$= \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2 \quad (34)$$

as desired.  $\square$

#### Exercise 1.6

Show that if two variables  $x$  and  $y$  are independent, then their covariance is zero.

*Proof.* By the definition of covariance, we have

$$\text{cov}[x, y] = \mathbb{E}_{x,y} [\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}] \quad (35)$$

$$= \mathbb{E}_{x,y} [xy - \mathbb{E}[x]y - \mathbb{E}[y]x + \mathbb{E}[x]\mathbb{E}[y]] \quad (36)$$

$$= \mathbb{E}_{x,y} [xy] - \mathbb{E}[x]\mathbb{E}[y] - \mathbb{E}[y]\mathbb{E}[x] + \mathbb{E}[x]\mathbb{E}[y] \quad (37)$$

$$= \mathbb{E}_{x,y} [xy] - \mathbb{E}[x]\mathbb{E}[y] \quad (38)$$

$$= \iint xyp(x, y)dx dy - \mathbb{E}[x]\mathbb{E}[y] \quad (39)$$

$$= \iint xyp(x)p(y)dxdy - \mathbb{E}[x]\mathbb{E}[y] \quad (40)$$

$$= \int xp(x)dx \int yp(y)dy - \mathbb{E}[x]\mathbb{E}[y] \quad (41)$$

$$= \mathbb{E}[x]\mathbb{E}[y] - \mathbb{E}[x]\mathbb{E}[y] = 0 \quad (42)$$

as desired.  $\square$

### Exercise 1.7

In this exercise, we prove the normalization condition ((1.48) in PRML book)

$$\int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) = 1 \quad (43)$$

for the univariate Gaussian. To do this consider, the integral

$$I = \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2}x^2\right) dx \quad (44)$$

which we can evaluate by first writing its square in the form

$$I^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2\sigma^2}x^2 - \frac{1}{2\sigma^2}y^2\right) dxdy. \quad (45)$$

Now make the transformation from Cartesian coordinates  $(x, y)$  to polar coordinates  $(r, \theta)$  and then substitute  $u = r^2$ . Show that, by performing the integrals over  $\theta$  and  $u$ , and then taking the square root of both sides, we obtain

$$I = (2\pi\sigma^2)^{1/2}. \quad (46)$$

Finally, use this result to show that the Gaussian distribution  $\mathcal{N}(x|\mu, \sigma^2)$  is normalized.

*Proof.* By making the transformation from Cartesian coordinates  $(x, y)$  to polar coordinates  $(r, \theta)$  and then substitute  $u = r^2$ , we have

$$I^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\left(-\frac{x^2 + y^2}{2\sigma^2}\right) dxdy \quad (47)$$

$$= \int_0^{2\pi} \int_0^{\infty} \exp\left(-\frac{r^2}{2\sigma^2}\right) r dr d\theta \quad (48)$$

$$= 2\pi\sigma^2 \int_0^{\infty} \exp\left(-\frac{r^2}{2\sigma^2}\right) d\frac{r^2}{2\sigma^2} \quad (49)$$

$$= 2\pi\sigma^2 \int_0^{\infty} \exp\left(-\frac{u}{2\sigma^2}\right) d\frac{u}{2\sigma^2} \quad (50)$$

$$= -2\pi\sigma^2 \exp\left(-\frac{u}{2\sigma^2}\right) \Big|_0^{+\infty} \quad (51)$$

$$= -2\pi\sigma^2(0 - 1) = 2\pi\sigma^2. \quad (52)$$

Thus,

$$I = (2\pi\sigma^2)^{1/2}. \quad (53)$$

Furthermore,

$$\int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) dx = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx = \frac{1}{\sqrt{2\pi}\sigma} I = 1. \quad (54)$$

44 This completes our proof. □

### Exercise 1.8

By using a change of variables, verify that the univariate Gaussian distribution given by ((1.46) in PRML book)

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\} \quad (55)$$

satisfies ((1.49) in PRML book)

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) x dx = \mu. \quad (56)$$

Next, by differentiating both sides of the normalization condition

$$\int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) dx = 1 \quad (57)$$

with respect to  $\sigma^2$ , verify that the Gaussian satisfies ((1.50) in PRML book)

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2. \quad (58)$$

Finally, show that ((1.51) in PRML book)

$$\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2 \quad (59)$$

holds.

45

*Proof.* Let's first verify (56).

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\} x dx \quad (60)$$

$$= \int_{-\infty}^{\infty} \frac{y + \mu}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{y^2}{2\sigma^2}\right\} dy \quad (y = x - \mu) \quad (61)$$

$$= \int_{-\infty}^{\infty} \frac{y}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{y^2}{2\sigma^2}\right\} dy + \underbrace{\mu \int_{-\infty}^{\infty} \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{y^2}{2\sigma^2}\right\} dy}_{=1} \quad (62)$$

$$= 0 + \mu = \mu \quad (63)$$

46 where the first term of the second last line vanishes since the integrand is an odd function with  
47 respect to  $y$  and the region of integration is symmetric about 0.

Next, to derive (58), we first substitute the standard form of Gaussian distribution into (43) and make some rearrangements.

$$\int_{-\infty}^{\infty} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) dx = \sqrt{2\pi\sigma^2}. \quad (64)$$

Before doing differentiation on both sides, we need to explain why we can swap the differentiation and the integration. Define  $f(x, \sigma^2) = \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$  and  $I(\sigma^2) = \int_{-\infty}^{\infty} f(x, \sigma^2)dx$ , then  $f'_{\sigma^2}(x, \sigma^2)$  is given by

$$f'_{\sigma^2}(x, \sigma^2) = \frac{(x-\mu)^2}{2(\sigma^2)^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (65)$$

It is easy to see that  $f(x, \sigma^2)$  and  $f'_{\sigma^2}(x, \sigma^2)$  are continuous on  $(-\infty, +\infty) \times (0, +\infty)$ , and for every  $\sigma^2 \in (0, +\infty)$ ,  $I(\sigma^2)$  converges to  $\sqrt{2\pi\sigma^2}$ <sup>1</sup>. The last thing we need to check is if  $\int_{-\infty}^{\infty} f'_{\sigma^2}(x, \sigma^2)dx$  is uniformly convergent for  $\sigma^2 \in (0, +\infty)$ . Actually, since  $(0, +\infty)$  is open, we only need to see if  $\int_{-\infty}^{\infty} f'_{\sigma^2}(x, \sigma^2)dx$  is uniformly convergent on any closed subset of  $(0, +\infty)$ . To do this, let  $z = (x - \mu)/\sigma^2$ , then

$$\int_{-\infty}^{\infty} f'_{\sigma^2}(x, \sigma^2)dx = \int_{-\infty}^{\infty} \frac{(x-\mu)^2}{2(\sigma^2)^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \quad (66)$$

$$= \frac{4\sigma^2}{2(\sigma^2)^2} \int_{-\infty}^{\infty} \frac{(x-\mu)^2}{4\sigma^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \quad (67)$$

$$< \frac{2}{\sigma^2} \int_{-\infty}^{\infty} \exp\left(-\frac{(x-\mu)^2}{4\sigma^2}\right) \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \quad (68)$$

$$= \frac{2}{\sigma^2} \int_{-\infty}^{\infty} \exp\left(-\frac{(x-\mu)^2}{4\sigma^2}\right) dx \quad (69)$$

$$= \frac{2}{\sigma^2} \sqrt{2\pi(\sqrt{2}\sigma)^2} = \frac{4}{\sigma^2} \sqrt{\pi\sigma^2} \quad (70)$$

where the inequality in the third line follows from  $x < e^x$  for any  $x \in \mathbb{R}$ . Since  $f'_{\sigma^2}(x, \sigma^2) \geq 0$ , according to Weierstrass's test for absolute uniform convergence, the above derivation shows that  $\int_{-\infty}^{\infty} f'_{\sigma^2}(x, \sigma^2)dx$  is uniformly convergent. Thus, we can interchange the differentiation and integral safely.

$$I'(\sigma^2) = \int_{-\infty}^{\infty} f'_{\sigma^2}(x, \sigma^2)dx = \frac{d}{d\sigma^2} \sqrt{2\pi\sigma^2} = \frac{\sqrt{2\pi}}{2\sqrt{\sigma^2}} \quad (71)$$

We can rewrite (67) as

$$\frac{4\sigma^2}{2(\sigma^2)^2} \int_{-\infty}^{\infty} \frac{(x-\mu)^2}{4\sigma^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \quad (72)$$

$$= \frac{4\sigma^2 \sqrt{2\pi\sigma^2}}{2(\sigma^2)^2} \int_{-\infty}^{\infty} \frac{(x-\mu)^2}{\sqrt{2\pi\sigma^2} 4\sigma^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \quad (73)$$

$$= \frac{\sqrt{2\pi\sigma^2}}{2(\sigma^2)^2} \int_{-\infty}^{\infty} \frac{(x-\mu)^2}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \quad (74)$$

$$= \frac{\sqrt{2\pi\sigma^2}}{2(\sigma^2)^2} \text{var}[x]. \quad (75)$$

Combining (71) and (72) yields

$$\frac{\sqrt{2\pi\sigma^2}}{2(\sigma^2)^2} \text{var}[x] = \frac{\sqrt{2\pi}}{2\sqrt{\sigma^2}} \iff \text{var}[x] = \sigma^2 \quad (76)$$

Furthermore, by the definition of variance,

$$\text{var}[x] = \sigma^2 = \int_{-\infty}^{\infty} \frac{(x-\mu)^2}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \quad (77)$$

<sup>1</sup><http://homepages.math.uic.edu/~jyang06/stat411/handouts/InterchangeDiffandIntegral.pdf>



$$= \int_{-\infty}^{\infty} \frac{x^2 - 2\mu x + \mu^2}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) dx \quad (78)$$

$$= \mathbb{E}[x^2] - 2\mu\mathbb{E}[x] + \mu^2 = \mathbb{E}[x^2] - \mu^2 \quad (79)$$

$$\iff \mathbb{E}[x^2] = \mu^2 + \sigma^2. \quad (80)$$

48 The last two claims have been proved together.  $\square$

### Exercise 1.9

Show that the mode (i.e. maximum) of the Gaussian distribution ((1.46) in PRML book)

$$\mathcal{N}(x \mid \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} \quad (81)$$

is given by  $\mu$ . Similarly, show that the mode of the multivariate Gaussian ((1.52) in PRML book)

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\} \quad (82)$$

is given by  $\boldsymbol{\mu}$ . Here,  $\mathbf{x}$  is a  $D$ -dimensional vector of continuous variables.

49

*Proof.* For the univariate case, differentiating the Gaussian density function with respect to  $x$  gives

$$\frac{\partial \mathcal{N}(x \mid \mu, \sigma^2)}{\partial x} = -\frac{x-\mu}{\sigma^2} \cdot \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} \quad (83)$$

50 Setting this to 0 yields  $x = \mu$ . So  $x = \mu$  is the only stationary point. Since  $x \in \mathbb{R}$  and  $\lim_{x \rightarrow \infty} \mathcal{N}(x \mid \mu, \sigma^2) = 0$ , then the mode of  $\mathcal{N}(x \mid \mu, \sigma^2)$  is given by  $\mu$ .

51

Similarly, for the multivariate case, according to the result  $\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x}$  where  $\mathbf{A}$  is a symmetric matrix, differentiating the multivariate Gaussian with respect to  $\mathbf{x}$  gives

$$\frac{\partial \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \mathbf{x}} = -\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) \cdot \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}) \quad (84)$$

52 Setting this to  $\mathbf{0}$  and left-multiplying by  $\boldsymbol{\Sigma}$  yield  $\mathbf{x} = \boldsymbol{\mu}$ . The same argument is applicable here.  $\square$

### Exercise 1.10

Suppose that the two variables  $x$  and  $z$  are statistically independent. Show that the mean and variance of their sum satisfies

$$\mathbb{E}[x+z] = \mathbb{E}[x] + \mathbb{E}[z] \quad (85)$$

$$\text{Var}[x+z] = \text{Var}[x] + \text{Var}[z]. \quad (86)$$

53

*Proof.* We first consider the case when  $x$  and  $z$  are continuous.

$$\mathbb{E}[x+z] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x+z)p(x,z)dx dz \quad (\text{Definition of mean}) \quad (87)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x+z)p(x)p(z)dx dz \quad (x \text{ and } z \text{ are independent}) \quad (88)$$

$$= \int_{-\infty}^{\infty} xp(x)dx + \int_{-\infty}^{\infty} zp(z)dz \quad (89)$$

$$= \mathbb{E}[x] + \mathbb{E}[z] \quad (\text{Definition of mean}) \quad (90)$$

For the variances, since  $x$  and  $z$  are independent,

$$(x + z - \mathbb{E}(x + z))^2 = \quad (91)$$

$$\text{Var}[x + z] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x + z - \mathbb{E}(x + z))^2 p(x, z) dx dz \quad (\text{Definition of variance}) \quad (92)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} ((x - \mathbb{E}[x])^2 + (z - \mathbb{E}[z])^2 \quad (93)$$

$$- 2(x - \mathbb{E}[x])(z - \mathbb{E}[z])) p(x)p(z) dx dz \quad (x \text{ and } z \text{ are independent}) \quad (94)$$

$$= \int_{-\infty}^{\infty} (x - \mathbb{E}[x])^2 p(x) dx + \int_{-\infty}^{\infty} (z - \mathbb{E}[z])^2 p(z) dz \quad (95)$$

$$- 2 \int_{-\infty}^{\infty} (x - \mathbb{E}(x)) p(x) dx \int_{-\infty}^{\infty} (z - \mathbb{E}(z)) p(z) dz \quad (96)$$

$$= \int_{-\infty}^{\infty} (x - \mathbb{E}[x])^2 p(x) dx + \int_{-\infty}^{\infty} (z - \mathbb{E}[z])^2 p(z) dz \quad (97)$$

$$= \text{Var}[x] + \text{Var}[z] \quad (\text{Definition of variance}) \quad (98)$$

54

□

### Exercise 1.11

By setting the derivatives of the log likelihood function ((1.54) in PRML book)

$$\ln p(\mathbf{x} \mid \mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) \quad (99)$$

with respect to  $\mu$  and  $\sigma^2$  equal to zero, verify the results

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n \quad (100)$$

and

$$\sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2. \quad (101)$$

55

*Proof.*

$$\frac{\partial \ln p(\mathbf{x} \mid \mu, \sigma^2)}{\partial \mu} = -\frac{1}{\sigma^2} \sum_{n=1}^N (x_n - \mu) = 0 \quad (102)$$

$$\Downarrow \quad (103)$$

$$\sum_{n=1}^N (x_n - \mu) = 0 \Rightarrow \sum_{n=1}^N x_n - N\mu = 0 \Rightarrow \mu = \frac{1}{N} \sum_{n=1}^N x_n \quad (104)$$

Thus,  $\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$ . Now we plug  $\mu_{\text{ML}}$  into (99) and then take derivatives with respect to  $\sigma^2$ .

$$\frac{\partial \ln p(\mathbf{x} \mid \mu, \sigma^2)}{\partial \sigma^2} = \frac{1}{2(\sigma^2)^2} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2 - \frac{N}{2\sigma^2} = 0 \quad (105)$$

$$\Downarrow \quad (106)$$

$$\frac{1}{\sigma^2} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2 - N = 0 \Rightarrow N\sigma^2 = \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2 \Rightarrow \sigma^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2. \quad (107)$$

56 Hence,  $\sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2$ . This completes the verification.  $\square$

### Exercise 1.12

Using the results in PRML book, i.e. (1.49)

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) x dx = \mu \quad (108)$$

and (1.50)

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2, \quad (109)$$

show that

$$\mathbb{E}[x_n x_m] = \mu^2 + I_{nm} \sigma^2 \quad (110)$$

where  $x_n$  and  $x_m$  denote data points sampled from a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ , and  $I_{nm}$  satisfies  $I_{nm} = 1$  if  $n = m$  and  $I_{nm} = 0$  otherwise. Hence prove that the results (1.57) and (1.58) in PRML book as follows.

$$\mathbb{E}[\mu_{\text{ML}}] = \mu \quad (111)$$

$$\mathbb{E}[\sigma_{\text{ML}}^2] = \left( \frac{N-1}{N} \right) \sigma^2. \quad (112)$$

57 *Proof.* When  $n = m$ ,  $\mathbb{E}[x_n x_m] = \mathbb{E}[x_n^2] = \mu^2 + \sigma^2$ . However, if  $n \neq m$ , since  $x_n$  and  $x_m$  are independent,  $\mathbb{E}[x_n x_m] = \mathbb{E}[x_n] \mathbb{E}[x_m] = \mu^2$ . Thus,  $\mathbb{E}[x_n x_m] = \mu^2 + I_{nm} \sigma^2$  holds.

$$\mathbb{E}[\mu_{\text{ML}}] = \mathbb{E} \left[ \frac{1}{N} \sum_{n=1}^N x_n \right] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[x_n] = \frac{N\mu}{N} = \mu \quad (113)$$

$$\mathbb{E}[\sigma_{\text{ML}}^2] = \mathbb{E} \left[ \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2 \right] \quad (114)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbb{E} [x_n^2 - 2x_n \mu_{\text{ML}} + \mu_{\text{ML}}^2] \quad (115)$$

$$= \frac{1}{N} \sum_{n=1}^N (\mathbb{E}[x_n^2] - 2\mathbb{E}[x_n \mu_{\text{ML}}] + \mathbb{E}[\mu_{\text{ML}}^2]) \quad (116)$$

$$= \frac{1}{N} \sum_{n=1}^N \left( \mu^2 + \sigma^2 - 2\mathbb{E} \left[ x_n \frac{1}{N} \sum_{n=1}^N x_n \right] + \mathbb{E} \left[ \left( \frac{1}{N} \sum_{n=1}^N x_n \right)^2 \right] \right) \quad (117)$$

$$= \frac{1}{N} \sum_{n=1}^N \left( \mu^2 + \sigma^2 - \frac{2}{N} \mathbb{E} \left[ x_n^2 + \sum_{i \neq n} x_n x_i \right] + \frac{1}{N^2} \mathbb{E} \left[ \sum_{i=1}^N x_i^2 + \sum_{i \neq j} x_i x_j \right] \right) \quad (118)$$

$$= \frac{1}{N} \sum_{n=1}^N \left( \mu^2 + \sigma^2 - \frac{2}{N} \left( \mathbb{E}[x_n^2] + \sum_{i \neq n}^N \mathbb{E}[x_n x_i] \right) + \frac{1}{N^2} \left( \sum_{n=1}^N \mathbb{E}[x_n^2] + \sum_{i \neq j}^N \mathbb{E}[x_i x_j] \right) \right) \quad (119)$$

$$= \frac{1}{N} \sum_{n=1}^N \left( \mu^2 + \sigma^2 - \frac{2}{N} (\mu^2 + \sigma^2 + (N-1)\mu^2) + \frac{1}{N^2} (N(\mu^2 + \sigma^2) + N(N-1)\mu^2) \right) \quad (120)$$

$$= \frac{1}{N} \sum_{n=1}^N \left( \mu^2 + \sigma^2 - \frac{1}{N} (\sigma^2 + N\mu^2) \right) = \left( \frac{N-1}{N} \right) \sigma^2 \quad (121)$$

58 which completes the proof.  $\square$

### Exercise 1.13

Suppose that the variance of a Gaussian is estimated using  $\sigma_{\text{ML}}^2$  but with the maximum likelihood estimate  $\mu_{\text{ML}}$  replaced with the true value  $\mu$  of the mean. Show that this estimator has the property that its expectation is given by the true variance  $\sigma^2$ .

59

*Proof.* By replacing  $\mu_{\text{ML}}$  in the proof of Exercise 1.12 with  $\mu$ , we get

$$\mathbb{E}[\sigma_{\text{ML}}^2] = \mathbb{E} \left[ \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2 \right] \quad (122)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbb{E} [x_n^2 - 2x_n\mu + \mu^2] \quad (123)$$

$$= \frac{1}{N} \sum_{n=1}^N (\mathbb{E}[x_n^2] - 2\mu\mathbb{E}[x_n] + \mathbb{E}[\mu^2]) \quad (124)$$

$$= \frac{1}{N} \sum_{n=1}^N (\mu^2 + \sigma^2 - 2\mu^2 + \mu^2) \quad (125)$$

$$= \frac{1}{N} \sum_{n=1}^N \sigma^2 = \frac{N\sigma^2}{N} = \sigma^2 \quad (126)$$

60  $\square$

### Exercise 1.14

Show that an arbitrary square matrix with the elements  $w_{ij}$  can be written in the form  $w_{ij} = w_{ij}^S + w_{ij}^A$  where  $w_{ij}^S$  and  $w_{ij}^A$  are symmetric and anti-symmetric matrices, respectively, satisfying  $w_{ij}^S = w_{ji}^S$  and  $w_{ij}^A = -w_{ji}^A$  for all  $i$  and  $j$ . Now consider the second order term in a higher order polynomial in  $D$  dimensions, given by

$$\sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j. \quad (127)$$

Show that

$$\sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j = \sum_{i=1}^D \sum_{j=1}^D w_{ij}^S x_i x_j \quad (128)$$

so that the contribution from the anti-symmetric matrix vanishes. We therefore see that, without loss of generality, the matrix of coefficients  $w_{ij}$  can be chosen to be symmetric, and so not all of the  $D^2$  elements of this matrix can be chosen independently. Show that the number of independent parameters in the matrix  $w_{ij}^S$  is given by  $D(D+1)/2$ .

61

*Proof.* <sup>2</sup> Given an arbitrary square matrix  $\mathbf{W}$ , let  $\mathbf{S} = (\mathbf{W} + \mathbf{W}^T)/2$  and  $\mathbf{A} = (\mathbf{W} - \mathbf{W}^T)/2$ , then we have  $\mathbf{W} = \mathbf{S} + \mathbf{A}$ , namely,  $w_{ij} = w_{ij}^S + w_{ij}^A$ . Since  $\mathbf{S} = \mathbf{S}^T$  and  $\mathbf{A}^T = -\mathbf{A}$ ,  $\mathbf{S}$  and  $\mathbf{A}$  are symmetric and anti-symmetric matrices, respectively. With this, we have

$$\sum_{i=1}^D \sum_{j=1}^D w_{ij}^A x_i x_j = \frac{1}{2} \sum_{i=1}^D \sum_{j=1}^D (w_{ij} - w_{ji}) x_i x_j \quad (129)$$

$$= \frac{1}{2} \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j - \frac{1}{2} \sum_{i=1}^D \sum_{j=1}^D w_{ji} x_i x_j \quad (130)$$

$$= \frac{1}{2} \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j - \frac{1}{2} \sum_{j=1}^D \sum_{i=1}^D w_{ij} x_j x_i \quad (131)$$

$$= \frac{1}{2} \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j - \frac{1}{2} \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j = 0 \quad (132)$$

It follows that

$$\sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j = \sum_{i=1}^D \sum_{j=1}^D (w_{ij}^S + w_{ij}^A) x_i x_j \quad (133)$$

$$= \sum_{i=1}^D \sum_{j=1}^D w_{ij}^S x_i x_j + \sum_{i=1}^D \sum_{j=1}^D w_{ij}^A x_i x_j \quad (134)$$

$$= \sum_{i=1}^D \sum_{j=1}^D w_{ij}^S x_i x_j \quad (135)$$

For a symmetric matrix, the parameters of only the upper triangle part or the lower triangle part are independent. Therefore, the number of independent parameters can be calculated as follows.

$$1 + 2 + \dots + (D-1) + D = \frac{D(D+1)}{2}. \quad (136)$$

<sup>2</sup>When working on the second part of the proof, I referenced the official solution manual.

62 This completes the proof. □

### Exercise 1.15

In this exercise and the next, we explore how the number of independent parameters in a polynomial grows with the order  $M$  of the polynomial and with the dimension  $D$  of the input space. We start by writing down the  $M^{\text{th}}$  order term for a polynomial in  $D$  dimensions in the form

$$\sum_{i_1=1}^D \sum_{i_2=1}^D \cdots \sum_{i_M=1}^D w_{i_1 i_2 \dots i_M} x_{i_1} x_{i_2} \cdots x_{i_M}. \quad (137)$$

The coefficients  $w_{i_1 i_2 \dots i_M}$  comprise  $D^M$  elements, but the number of independent parameters is significantly fewer due to the many interchange symmetries of the factor  $x_{i_1} x_{i_2} \cdots x_{i_M}$ . Begin by showing that the redundancy in the coefficients can be removed by rewriting this  $M^{\text{th}}$  order term in the form

$$\sum_{i_1=1}^D \sum_{i_2=1}^{i_1} \cdots \sum_{i_M=1}^{i_{M-1}} \tilde{w}_{i_1 i_2 \dots i_M} x_{i_1} x_{i_2} \cdots x_{i_M}. \quad (138)$$

Note that the precise relationship between the  $\tilde{w}$  coefficients and  $w$  coefficients need not be made explicit. Use this result to show that the number of independent parameters  $n(D, M)$ , which appear at order  $M$ , satisfies the following recursion relation

$$n(D, M) = \sum_{i=1}^D n(i, M-1). \quad (139)$$

Next use proof by induction to show that the following result holds

$$\sum_{i=1}^D \frac{(i+M-2)!}{(i-1)!(M-1)!} = \frac{(D+M-1)!}{(D-1)!M!} \quad (140)$$

which can be done by first proving the result for  $D = 1$  and arbitrary  $M$  by making use of the result  $0! = 1$ , then assuming it is correct for dimension  $D$  and verifying that it is correct for dimension  $D + 1$ . Finally, use the two previous results, together with proof by induction, to show

$$n(D, M) = \frac{(D+M-1)!}{(D-1)!M!}. \quad (141)$$

To do this, first show that the result is true for  $M = 2$ , and any value of  $D \geq 1$ , by comparison with the result of Exercise 1.14. Then make use of (139), together with (140), to show that, if the result holds at order  $M - 1$ , then it will also hold at order  $M$ .

*Proof.* <sup>3</sup> The redundancy in (137) arises from the interchange symmetries between the indices  $i_k$ . Enforcing the appearance order of the indices  $i_k$  can remove this redundancy as in (138). To derive (139), the number of independent parameters which appear at order  $M$  can be written as

$$n(D, M) = \sum_{i_1=1}^D \sum_{i_2=1}^{i_1} \cdots \sum_{i_M=1}^{i_{M-1}} 1 \quad (142)$$

<sup>3</sup>I referenced the official manual for the first part of the exercise.

which can be rewritten as

$$n(D, M) = \sum_{i_1=1}^D \left( \sum_{i_2=1}^{i_1} \cdots \sum_{i_M=1}^{i_{M-1}} 1 \right) = \sum_{i_1=1}^D n(i_1, M-1) \quad (143)$$

64 which is equivalent to (139).

Now we use proof by induction to show (140). For  $D = 1$  and arbitrary  $M$ , using the result  $0! = 1$ , we have

$$\sum_{i=1}^1 \frac{(1+M-2)!}{(1-1)!(M-1)!} = \frac{(M-1)!}{0!(M-1)!} = 1 = \frac{1+M-1}{(1-1)!M!}. \quad (144)$$

which shows (140) is correct when  $D = 1$ . Assuming (140) is correct for dimension  $D$ , then for dimension  $D + 1$ , we have

$$\sum_{i=1}^{D+1} \frac{(i+M-2)!}{(i-1)!(M-1)!} = \frac{(D+M-1)!}{(D-1)!M!} + \frac{(D+1+M-2)!}{(D+1-1)!(M-1)!} \quad (145)$$

$$= \frac{(D+M-1)!}{(D-1)!M!} + \frac{(D+M-1)!}{D!(M-1)!} \quad (146)$$

$$= \frac{(D+M-1)!(D+M)}{D!M!} \quad (147)$$

$$= \frac{(D+M)!}{D!M!} \quad (148)$$

as desired. Once again, we show (141) with proof by induction. When  $M = 2$  and  $D \geq 1$ , it is exactly the case of Exercise 1.14. Then we have

$$n(D, 2) = \frac{(D+2-1)!}{(D-1)!2!} = \frac{(D+1)!}{(D-1)!2} = \frac{D(D+1)}{2} \quad (149)$$

which shows that (141) holds when  $M = 2$  and any  $D \geq 1$ . Assuming (141) is correct for  $M - 1$ , then we have

$$n(D, M-1) = \frac{(D+M-2)!}{(D-1)!(M-1)!} \quad (150)$$

Combining this with (139) and (140), we get

$$n(D, M) = \sum_{i=1}^D n(i, M-1) = \sum_{i=1}^D \frac{(i+M-2)!}{(i-1)!(M-1)!} = \frac{(D+M-1)!}{(D-1)!M!}. \quad (151)$$

65 This completes the proof. □

### Exercise 1.16

In Exercise 1.15, we proved the result (139) for the number of independent parameters in the  $M^{\text{th}}$  order term of a  $D$ -dimensional polynomial. We now find an expression for the total number  $N(D, M)$  of independent parameters in all of the terms up to and including the  $M$ th order. First show that  $N(D, M)$  satisfies

$$N(D, M) = \sum_{m=0}^M n(D, m) \quad (152)$$

where  $n(D, m)$  is the number of independent parameters in the term of order  $m$ . Now make use of the result (141), together with proof by induction, to show that

$$N(D, M) = \frac{(D + M)!}{D! M!}. \quad (153)$$

This can be done by first proving that the result holds for  $M = 0$  and arbitrary  $D \geq 1$ , then assuming that it holds at order  $M$ , and hence showing that it holds at order  $M + 1$ . Finally, make use of Stirling's approximation in the form

$$n! \simeq n^n e^{-n} \quad (154)$$

for large  $n$  to show that, for  $D \gg M$ , the quantity  $N(D, M)$  grows like  $D^M$ , and for  $M \gg D$  it grows like  $M^D$ . Consider a cubic ( $M = 3$ ) polynomial in  $D$  dimensions, and evaluate numerically the total number of independent parameters for (i)  $D = 10$  and (ii)  $D = 100$ , which correspond to typical small-scale and medium-scale machine learning applications.

66

*Proof.* We have proved the number of independent parameters in the  $M$ th order term of a  $D$ -dimensional polynomial can be written as (141). Since  $N(D, M)$  represents the total number of independent parameters including the 0th order to the  $M$ th order. Therefore, summing over all the orders gives the total number of independent parameters as follows.

$$N(D, M) = \sum_{m=0}^M n(D, m) \quad (155)$$

We use (141), together with proof by induction to show (153). For  $M = 0$  and arbitrary  $D \geq 0$ ,

$$N(D, 0) = n(D, 0) = \frac{(D + 0 - 1)!}{(D - 1)! 0!} = \frac{(D - 1)!}{(D - 1)! 0!} = 1 = \frac{(D + 0)!}{D! 0!}. \quad (156)$$

The case when  $M = 0$  obviously holds for arbitrary  $D \geq 1$ . Assume (155) holds at order  $M$ , then

$$N(D, M + 1) = \frac{(D + M)!}{D! M!} + n(D, M + 1) \quad (157)$$

$$= \frac{(D + M)!}{D! M!} + \frac{(D + M)!}{(D - 1)! (M + 1)!} \quad (158)$$

$$= \frac{(D + M)! (M + 1)}{D! (M + 1)!} + \frac{(D + M)! D}{D! (M + 1)!} \quad (159)$$

$$= \frac{(D + M)! (M + 1)}{D! (M + 1)!} + \frac{(D + M)! D}{D! (M + 1)!} \quad (160)$$

$$= \frac{(D + M)! (D + M + 1)}{D! (M + 1)!} = \frac{(D + M + 1)!}{D! (M + 1)!} \quad (161)$$



67 as desired.

Finally, when  $D \gg M$ , substituting Stirling's approximation into (153) yields

$$N(D, M) = \frac{(D+M)!}{D! M!} \quad (162)$$

$$\approx \frac{(D+M)^{(D+M)} e^{-(D+M)}}{D^D e^{-D} M!} \quad (163)$$

$$= \frac{(D+M)^{(D+M)} e^{-M}}{D^D M!} \quad (164)$$

$$= \frac{D^M (D+M)^{(D+M)} e^{-M}}{D^{(D+M)} M!} \quad (165)$$

$$= \frac{D^M e^{-M}}{M!} \left(1 + \frac{M}{D}\right)^{(D+M)} \quad (166)$$

$$\approx \frac{D^M e^{-M}}{M!} \left(1 + \frac{M(D+M)}{D}\right) \quad (167)$$

$$\approx \frac{D^M e^{-M}}{M!} (1+M) \quad (168)$$

Thus, for  $D \gg M$ , the quantity  $N(D, M)$  grows like  $D^M$ . Likewise, for  $M \gg D$ , we have

$$N(D, M) \approx \frac{M^D e^{-D}}{D!} (1+D). \quad (169)$$

68 Hence, for  $M \gg D$  it grows like  $M^D$ . When  $M = 3$ , we employ (153) to get  $N(10, 3) = 286$  and  
 69  $N(100, 3) = 176851$ .  $\square$

### Exercise 1.17

The gamma function is defined by

$$\Gamma(x) \equiv \int_0^\infty u^{x-1} e^{-u} du. \quad (170)$$

Using integration by parts, prove the relation  $\Gamma(x+1) = x\Gamma(x)$ . Show also that  $\Gamma(1) = 1$  and hence that  $\Gamma(x+1) = x!$  when  $x$  is an integer.

70

*Proof.*

$$\Gamma(x+1) = \int_0^\infty u^x e^{-u} du \quad (171)$$

$$= - \left[ u^x e^{-u} \right]_0^\infty - x \int_0^\infty u^{x-1} e^{-u} du \quad (172)$$

$$= x \int_0^\infty u^{x-1} e^{-u} du = x\Gamma(x). \quad (173)$$

To calculate  $\Gamma(1)$ , we have

$$\Gamma(1) = \int_0^\infty e^{-u} du = -e^{-u} \Big|_0^\infty = 1. \quad (174)$$

71 Thus,  $\Gamma(x+1) = x!$  when  $x$  is an integer.  $\square$

### Exercise 1.18

We can use the result (46) to derive an expression for the surface area  $S_D$ , and the volume  $V_D$ , of a sphere of a unit radius in  $D$  dimensions. To do this, consider the following result, which is obtained by transforming from Cartesian to polar coordinates

$$\prod_{i=1}^D \int_{-\infty}^{\infty} e^{-x_i^2} dx_i = S_D \int_0^{\infty} e^{-r^2} r^{D-1} dr. \quad (175)$$

Using the definition (170) of the Gamma function, together with (46), evaluate both sides of this equation, and hence show that

$$S_D = \frac{2\pi^{D/2}}{\Gamma(D/2)}. \quad (176)$$

Next, by integrating with respect to radius from 0 to 1, show that the volume of the unit sphere in  $D$  dimensions is given by

$$V_D = \frac{S_D}{D}. \quad (177)$$

Finally, use the results  $\Gamma(1) = 1$  and  $\Gamma(3/2) = \sqrt{\pi}/2$  to show that (176) and (177) reduce to the usual expressions for  $D = 2$  and  $D = 3$ .

Note that it is necessary to clarify that  $S_D$  is not the surface area and  $S_D r^{D-1}$  represents the true surface area. In this exercise, in numerical sense they are equal because they only consider a sphere of a unit radius in  $D$  dimensions. If you are interested in this topic, please check out [https://zhangyk8.github.io/teaching/file/Exercise\\_4\\_insight.pdf](https://zhangyk8.github.io/teaching/file/Exercise_4_insight.pdf).

*Proof.* Given  $\int_{-\infty}^{\infty} \exp(-\frac{1}{2\sigma^2}x^2) dx = \sqrt{2\pi\sigma^2}$ , we have  $\int_{-\infty}^{\infty} \exp(-x_i^2) dx = \sqrt{\pi}$ .

$$\prod_{i=1}^D \int_{-\infty}^{\infty} e^{-x_i^2} dx_i = \pi^{D/2}. \quad (178)$$

Let  $u = r^2$ , then  $dr = \frac{1}{2\sqrt{u}} du$ . Then the right hand side can be evaluated as follows.

$$S_D \int_0^{\infty} e^{-r^2} r^{D-1} dr = S_D \cdot \frac{1}{2} \int_0^{\infty} u^{D/2-1} e^{-u} du = \frac{S_D}{2} \Gamma(D/2) = \pi^{D/2} \quad (179)$$

which shows

$$S_D = \frac{2\pi^{D/2}}{\Gamma(D/2)}. \quad (180)$$

To compute the volume, after observing (175), we drop the terms  $e^{-x_i^2}$  and  $e^{-r^2}$ , and integrate w.r.t radius from 0 to 1 as follows.

$$V_D = \int \cdots \int_{x_1^2 + x_2^2 + \cdots + x_n^2 \leq 1} dx_1 dx_2 \cdots dx_n = S_D \int_0^1 r^{D-1} dr = \frac{S_D}{D} \quad (181)$$

where an excellent derivation for the integrand  $r^{D-1}$  can be found at [https://zhangyk8.github.io/teaching/file/Exercise\\_4\\_insight.pdf](https://zhangyk8.github.io/teaching/file/Exercise_4_insight.pdf). When  $D = 2$ , we have

$$S_2 = \frac{2\pi}{\Gamma(1)} = 2\pi, \quad V_2 = \pi. \quad (182)$$

And when  $D = 3$ , we get

$$S_3 = \frac{2\pi^{3/2}}{\Gamma(3/2)} = \frac{2\pi^{3/2}}{\sqrt{\pi}/2} = 4\pi, \quad V_2 = 4\pi/3. \quad (183)$$

77

□

### Exercise 1.19

Consider a sphere of radius  $a$  in  $D$ -dimensions together with the concentric hypercube of side  $2a$ , so that the sphere touches the hypercube at the centers of each of its sides. By using the results of Exercise 1.18, show that the ratio of the volume of the sphere to the volume of the cube is given by

$$\frac{\text{volume of sphere}}{\text{volume of cube}} = \frac{\pi^{D/2}}{D2^{D-1}\Gamma(D/2)}. \quad (184)$$

Now make use of Stirling's formula in the form

$$\Gamma(x+1) \approx (2\pi)^{1/2} e^{-x} x^{x+1/2} \quad (185)$$

which is valid for  $x \gg 1$ , to show that, as  $D \rightarrow \infty$ , the ratio (184) goes to zero. Show also that the ratio of the distance from the center of the hypercube to one of the corners, divided by the perpendicular distance to one of the sides, is  $\sqrt{D}$ , which therefore goes to  $\infty$  as  $D \rightarrow \infty$ . From these results we see that, in a space of high dimensionality, most of the volume of a cube is concentrated in the large number of corners, which themselves become very long 'spikes'!

78

*Proof.* Without loss of generality, we consider the case when  $a = 1$ . Then we can use the results of Exercise 1.18 directly as follows.

$$\frac{\text{volume of sphere}}{\text{volume of cube}} = \frac{2\pi^{D/2}}{D\Gamma(D/2) \times 2^D} = \frac{\pi^{D/2}}{D2^{D-1}\Gamma(D/2)}. \quad (186)$$

By Stirling's formula, namely, (185), we get

$$\frac{\pi^{D/2}}{D2^{D-1}\Gamma(D/2)} = \frac{\pi^{D/2}}{D2^{D-1}(2\pi)^{1/2}e^{-D/2}(D/2)^{D/2+1/2}} \quad (187)$$

$$= \frac{2\sqrt{2}\pi^{D/2}e^{D/2}}{D^{3/2}4^{D/2}(2\pi)^{1/2}(D/2)^{D/2}} \quad (188)$$

$$= \frac{2\sqrt{2}}{D^{3/2}(2\pi)^{1/2}} \cdot \left(\frac{\pi e}{2D}\right)^{D/2} \quad (189)$$

79 which goes to 0 as  $D \rightarrow \infty$ . The distance from the center of the hypercube to one of the corners is  
80  $\sqrt{1^2 + 1^2 + \dots + 1^2} = \sqrt{D}$ , but the perpendicular distance to one of the sides is 1. Hence, the ratio  
81 between them is  $\sqrt{D}$ . Thus, it goes to  $\infty$  as  $D \rightarrow \infty$ . □

### Exercise 1.20

In this exercise, we explore the behavior of the Gaussian distribution in high-dimensional spaces. Consider a Gaussian distribution in  $D$  dimensions given by

$$p(\mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{D/2}} \exp\left(-\frac{\|\mathbf{x}\|^2}{2\sigma^2}\right). \quad (190)$$

We wish to find the density with respect to radius in polar coordinates in which the direction variables have been integrated out. To do this, show that the integral of the probability density over a thin shell of radius  $r$  and thickness  $\epsilon$ , where  $\epsilon \ll 1$ , is given by  $p(r)\epsilon$  where

$$p(r) = \frac{S_D r^{D-1}}{(2\pi\sigma^2)^{D/2}} \exp\left(-\frac{r^2}{2\sigma^2}\right) \quad (191)$$

where  $S_D$  is the surface area of a unit sphere in  $D$  dimensions. Show that the function  $p(r)$  has a single stationary point located, for large  $D$ , at  $\hat{r} \approx \sqrt{D}\sigma$ . By considering  $p(\hat{r} + \epsilon)$  where  $\epsilon \ll \hat{r}$ , show that for large  $D$ ,

$$p(\hat{r} + \epsilon) = p(\hat{r}) \exp\left(-\frac{\epsilon^2}{\sigma^2}\right) \quad (192)$$

which shows that  $\hat{r}$  is a maximum of the radial probability density and also that  $p(r)$  decays exponentially away from its maximum at  $\hat{r}$  with length scale  $\sigma$ . We have already seen that  $\sigma \ll \hat{r}$  for large  $D$ , and so we see that most of the probability mass is concentrated in a thin shell at large radius. Finally, show that the probability density  $p(\mathbf{x})$  is larger at the origin than at the radius  $\hat{r}$  by a factor of  $\exp(D/2)$ . We therefore see that most of the probability mass in a high-dimensional Gaussian distribution is located at a different radius from the region of high probability density. This property of distributions in spaces of high dimensionality will have important consequences when we consider Bayesian inference of model parameters in later chapters.

82

*Proof.* From Exercise 1.18, the surface area of a sphere of radius  $r$  in  $D$  dimensions can be represented as  $S_D r^{D-1}$ . Given this result and  $\epsilon \ll 1$ , we have

$$\int_{\text{shell}} p(\mathbf{x}) d\mathbf{x} = p(\mathbf{x}) V_{\text{shell}} = \frac{1}{(2\pi\sigma^2)^{D/2}} \exp\left(-\frac{\|r\|^2}{2\sigma^2}\right) S_D r^{D-1} \epsilon = p(r) \epsilon. \quad (193)$$

which implies that  $p(r)$  can be written as

$$p(r) = \frac{S_D r^{D-1}}{(2\pi\sigma^2)^{D/2}} \exp\left(-\frac{r^2}{2\sigma^2}\right). \quad (194)$$

By setting the gradient  $\nabla_r p(r)$  to 0, we get

$$\nabla_r p(r) = \frac{S_D}{(2\pi\sigma^2)^{D/2}} \left[ (D-1) r^{D-2} \exp\left(-\frac{r^2}{2\sigma^2}\right) + \frac{r^D}{\sigma^2} \exp\left(-\frac{r^2}{2\sigma^2}\right) \right] = 0 \iff \hat{r} = \sqrt{D-1}\sigma. \quad (195)$$

Therefore, for large  $D$ ,  $\hat{r} \approx \sqrt{D}\sigma$ . Provided  $\epsilon \ll \hat{r}$ , we have

$$\frac{p(\hat{r} + \epsilon)}{p(\hat{r})} = \frac{(r + \epsilon)^{D-1} \exp\left(-\frac{(r+\epsilon)^2}{2\sigma^2}\right)}{r^{D-1} \exp\left(-\frac{r^2}{2\sigma^2}\right)} \quad (196)$$

$$= \left(1 + \frac{\epsilon}{\hat{r}}\right)^{D-1} \exp\left(-\frac{2\epsilon\hat{r} + \epsilon^2}{2\sigma^2}\right) \quad (197)$$

$$= \exp\left((D-1)\ln\left(1 + \frac{\epsilon}{\hat{r}}\right) - \frac{2\epsilon\hat{r} + \epsilon^2}{2\sigma^2}\right) \quad (198)$$

$$\approx \exp\left((D-1)\left(\frac{\epsilon}{\hat{r}} - \frac{\epsilon^2}{2\hat{r}^2}\right) - \frac{2\epsilon\hat{r} + \epsilon^2}{2\sigma^2}\right) \quad (\ln(1+x) = x - \frac{x^2}{2} + O(x^3)) \quad (199)$$

$$= \exp\left(\frac{2\epsilon\hat{r} - \epsilon^2}{2\sigma^2} - \frac{2\epsilon\hat{r} + \epsilon^2}{2\sigma^2}\right) \quad (D-1 = \frac{\hat{r}^2}{\sigma^2}) \quad (200)$$

$$= \exp\left(-\frac{\epsilon^2}{\sigma^2}\right) \quad (201)$$

which shows that for large  $D$ ,

$$p(\hat{r} + \epsilon) = p(\hat{r}) \exp\left(-\frac{\epsilon^2}{\sigma^2}\right). \quad (202)$$

This indicates that  $\hat{r}$  is a maximum of radial probability density and also that  $p(r)$  decays exponentially away from its maximum at  $\hat{r}$  with length scale  $\sigma$ . At the origin, we have

$$p(\mathbf{0}) = \frac{1}{(2\pi\sigma^2)^{D/2}}. \quad (203)$$

At the radius  $\hat{r}$ , we have

$$p(\|\mathbf{x}\| = \hat{r}) = \frac{1}{(2\pi\sigma^2)^{D/2}} \exp\left(-\frac{\hat{r}^2}{2\sigma^2}\right) \approx \frac{1}{(2\pi\sigma^2)^{D/2}} \exp\left(-\frac{D}{2}\right). \quad (204)$$

Thus,

$$\frac{p(\mathbf{0})}{p(\|\mathbf{x}\| = \hat{r})} = \exp\left(\frac{D}{2}\right). \quad (205)$$

83 This completes the proof. □

### Exercise 1.21

Consider two nonnegative numbers  $a$  and  $b$ , show that, if  $a \leq b$ , then  $a \leq (ab)^{1/2}$ . Use this result to show that, if the decision regions of a two-class classification problem are chosen to minimize the probability of misclassification, this probability will satisfy

$$p(\text{mistake}) \leq \int \{p(\mathbf{x}, \mathcal{C}_1)p(\mathbf{x}, \mathcal{C}_2)\}^{1/2} d\mathbf{x}. \quad (206)$$

84

*Proof.* Given  $a, b \geq 0$ , we have

$$a \leq b \iff a^2 \leq ab \iff a \leq (ab)^{1/2} \quad (207)$$

where the last ‘iff’ follows from the square root function is increasing on  $\mathbb{R}_+$ . Combining this result and (1.78) in PRML book,

$$p(\text{mistake}) = p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \quad (208)$$

$$= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x} \quad (209)$$

$$\leq \int_{\mathcal{R}_1} \{p(\mathbf{x}, \mathcal{C}_1)p(\mathbf{x}, \mathcal{C}_2)\}^{1/2} d\mathbf{x} + \int_{\mathcal{R}_2} \{p(\mathbf{x}, \mathcal{C}_1)p(\mathbf{x}, \mathcal{C}_2)\}^{1/2} d\mathbf{x} \quad (210)$$

$$= \int \{p(\mathbf{x}, \mathcal{C}_1)p(\mathbf{x}, \mathcal{C}_2)\}^{1/2} d\mathbf{x} \quad (211)$$

85 where the inequality follows from  $p(\mathbf{x}, \mathcal{C}_1) > p(\mathbf{x}, \mathcal{C}_2)$  for  $\mathbf{x} \in \mathcal{R}_1$  and  $p(\mathbf{x}, \mathcal{C}_2) > p(\mathbf{x}, \mathcal{C}_1)$  for  $\mathbf{x} \in \mathcal{R}_2$ ,  
 86 since the decision regions are chosen to minimize the probability of misclassification. This completes  
 87 the proof.  $\square$

### Exercise 1.22

Given a loss matrix with elements  $L_{kj}$ , the expected risk is minimized if, for each  $\mathbf{x}$ , we choose the class that minimizes (1.81) in PRML book

$$\sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x}). \quad (212)$$

Verify that, when the loss matrix is given by  $L_{kj} = 1 - I_{kj}$ , where  $I_{kj}$  are the elements of the identity matrix, this reduces to the criterion of choosing the class having the largest posterior probability. What is the interpretation of this form of loss matrix?

88 *Proof.* Plugging  $L_{kj} = 1 - I_{kj}$  into (212) gives

$$\sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x}) = \sum_k (1 - I_{kj}) p(\mathcal{C}_k | \mathbf{x}) \quad (213)$$

$$= \sum_k p(\mathcal{C}_k | \mathbf{x}) - \sum_k I_{kj} p(\mathcal{C}_k | \mathbf{x}) \quad (214)$$

$$= 1 - p(\mathcal{C}_j | \mathbf{x}) \quad (215)$$

89 which implies that minimizing  $\sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$  is equivalent to maximizing  $p(\mathcal{C}_j | \mathbf{x})$ .

90 This kind of loss matrix gives a loss of one if the sample is misclassified and a loss of zero if it is  
 91 classified correctly. The above proof shows that minimizing the expected risk in the sense of this loss  
 92 matrix is equivalent to minimizing the misclassification rate.  $\square$

### Exercise 1.23

Derive the criterion for minimizing the expected loss when there is a general loss matrix and general prior probabilities for the classes.

**Solution:** From (212), minimizing the expected loss is equivalent to minimizing

$$\sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x}) = \frac{1}{p(\mathbf{x})} \sum_k L_{kj} p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k). \quad (216)$$

94  $\square$

### Exercise 1.24

Consider a classification problem in which the loss incurred when an input vector from class  $\mathcal{C}_k$  is classified as belonging to class  $\mathcal{C}_j$  is given by loss matrix  $L_{kj}$ , and for which the loss incurred in selecting the reject option is  $\lambda$ . Find the decision criterion that will give the minimum expected loss. Verify that this reduces to the reject criterion discussed in Section 1.5.3 when the loss matrix is given by  $L_{kj} = 1 - I_{kj}$ . What is the relationship between  $\lambda$  and the rejection threshold  $\theta$ ?

95

*Proof.* From Section 1.5.2, the decision rule that minimizes the expected loss is the one that assigns each new  $\mathbf{x}$  to the class  $j$  for which the quantity

$$\sum_k L_{kj} p(C_k | \mathbf{x}) \quad (217)$$

is a minimum. Equivalently, we should assign a new  $\mathbf{x}$  to class  $j$  for which  $j = \operatorname{argmin}_l L_{kl} p(C_k | \mathbf{x})$  and  $L_{kj} p(C_k | \mathbf{x}) < \lambda$ , otherwise we reject  $\mathbf{x}$ . Given a loss matrix  $L_{kj} = 1 - I_{kj}$ , according to Exercise 1.22,  $\sum_k L_{kj} p(C_k | \mathbf{x}) = 1 - p(C_j)$ . If the smallest  $1 - p(C_j) < \lambda$ , or equivalently the largest  $p(C_j) > 1 - \lambda$ , we assign  $\mathbf{x}$  to  $j$ , otherwise we reject  $\mathbf{x}$ . In other words,  $1 - \lambda = \theta$ .  $\square$

### Exercise 1.25

Consider the generalization of the squared loss function (1.87) in PRML book for a single target variable  $t$  to the case of multiple target variables described by the vector  $\mathbf{t}$  given by

$$\mathbb{E}[L(t, y(\mathbf{x}))] = \iint \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}. \quad (218)$$

Using the calculus of variations, show that the function  $\mathbf{y}(\mathbf{x})$  for which this expected loss is minimized is given by  $\mathbf{y}(\mathbf{x}) = \mathbb{E}_{\mathbf{t}}[\mathbf{t} | \mathbf{x}]$ . Show that this result reduces to  $y(\mathbf{x}) = \mathbb{E}_t[t | \mathbf{x}]$  for the case of a single variable target  $t$ .

*Proof.* Our goal is to choose  $\mathbf{y}(\mathbf{x})$  so as to minimize  $\mathbb{E}[L(t, y(\mathbf{x}))]$ . If we assume a completely flexible function  $y(\mathbf{x})$ , we can do this formally using the calculus of variations to give

$$\frac{\delta \mathbb{E}[L(t, y(\mathbf{x}))]}{\delta \mathbf{y}(\mathbf{x})} = 2 \int (\mathbf{y}(\mathbf{x}) - \mathbf{t}) p(\mathbf{x}, \mathbf{t}) d\mathbf{t} = 0. \quad (219)$$

Solving for  $\mathbf{y}(\mathbf{x})$ , and using the sum and product rules of probability, we obtain

$$\mathbf{y}(\mathbf{x}) = \frac{\int \mathbf{t} p(\mathbf{x}, \mathbf{t}) d\mathbf{t}}{p(\mathbf{x})} = \int \mathbf{t} p(\mathbf{t} | \mathbf{x}) d\mathbf{t} = \mathbb{E}_{\mathbf{t}}[\mathbf{t} | \mathbf{x}] \quad (220)$$

which is the conditional average of  $\mathbf{t}$  conditioned on  $\mathbf{x}$ . When  $t$  is a scalar,

$$y(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int t p(t | \mathbf{x}) dt = \mathbb{E}_t[t | \mathbf{x}] \quad (221)$$

which is equivalent to (1.87) in the PRML book.  $\square$

### Exercise 1.26

By expansion of the square in (218), derive a result analogous to (1.90) in the PRML book, which is

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}_t[t | \mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \operatorname{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}. \quad (222)$$

Using the calculus of variations, show that the function  $\mathbf{y}(\mathbf{x})$  that minimizes the expected squared loss for the case of a vector  $\mathbf{t}$  of target variables is again given by the conditional expectation of  $\mathbf{t}$ .

*Proof.* Following Section 1.5.5, we have

$$\|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 = \|\mathbf{y}(\mathbf{x}) - \mathbb{E}(\mathbf{t} | \mathbf{x}) + \mathbb{E}(\mathbf{t} | \mathbf{x}) - \mathbf{t}\|^2 \quad (223)$$

$$= \|\mathbf{y}(\mathbf{x}) - \mathbb{E}(\mathbf{t}|\mathbf{x})\|^2 + (\mathbf{y}(\mathbf{x}) - \mathbb{E}(\mathbf{t}|\mathbf{x}))^T (\mathbb{E}(\mathbf{t}|\mathbf{x}) - \mathbf{t}) \quad (224)$$

$$+ (\mathbb{E}(\mathbf{t}|\mathbf{x}) - \mathbf{t})^T (\mathbf{y}(\mathbf{x}) - \mathbb{E}(\mathbf{t}|\mathbf{x})) + \|\mathbb{E}(\mathbf{t}|\mathbf{x}) - \mathbf{t}\|^2. \quad (225)$$

Substituting this into (218) and performing the integral over  $\mathbf{t}$ , we see that the cross term vanishes and the last term is exactly the definition of variance  $\mathbb{E}[\mathbf{t}|\mathbf{x}]$ . Thus,

$$\mathbb{E}[L] = \int \{\mathbf{y}(\mathbf{x}) - \mathbb{E}_{\mathbf{t}}[\mathbf{t}|\mathbf{x}]\} p(\mathbf{x}) d\mathbf{x} + \int \text{var}[\mathbf{t}|\mathbf{x}] p(\mathbf{x}) d\mathbf{x}. \quad (226)$$

103

□

### Exercise 1.27

Consider the expected loss for regression problems under the  $L_q$  loss function given by (1.91) in the PRML book, i.e.,

$$\mathbb{E}[L_q] = \iint |y(\mathbf{x}) - t|^q p(\mathbf{x}, t) d\mathbf{x} dt. \quad (227)$$

Write down the condition that  $y(\mathbf{x})$  must satisfy in order to minimize  $\mathbb{E}[L_q]$ . Show that, for  $q = 1$ , this solution represents the conditional median, i.e., the function  $y(\mathbf{x})$  such that the probability mass for  $t < y(\mathbf{x})$  is the same as for  $t \geq y(\mathbf{x})$ . Also show that the minimum expected  $L_q$  loss for  $q \rightarrow 0$  is given by the conditional mode, i.e., by the function  $y(\mathbf{x})$  equal to the value of  $t$  that maximizes  $p(t|\mathbf{x})$  for each  $\mathbf{x}$ .

104

*Proof.* We follow the logic from the official solution manual. (227) can be rewritten as

$$\mathbb{E}[L_q] = \int \left( \int |y(\mathbf{x}) - t|^q p(t|\mathbf{x}) dt \right) p(\mathbf{x}) d\mathbf{x}. \quad (228)$$

Since  $y(\mathbf{x})$  can be chosen independently for each  $\mathbf{x}$ , the minimum of  $\mathbb{E}[L_q]$  can be obtained via minimizing the following integrand

$$\int |y(\mathbf{x}) - t|^q p(t|\mathbf{x}) dt \quad (229)$$

for each  $\mathbf{x}$ . Setting the derivative to 0 yields

$$q \int |y(\mathbf{x}) - t|^{q-1} \text{sign}(y(\mathbf{x}) - t) p(t|\mathbf{x}) dt = 0. \quad (230)$$

Then,

$$\int_{y(\mathbf{x})}^{\infty} |y(\mathbf{x}) - t|^{q-1} p(t|\mathbf{x}) dt = \int_{-\infty}^{y(\mathbf{x})} |y(\mathbf{x}) - t|^{q-1} p(t|\mathbf{x}) dt. \quad (231)$$

When  $q = 1$ , we have

$$\int_{y(\mathbf{x})}^{\infty} p(t|\mathbf{x}) dt = \int_{-\infty}^{y(\mathbf{x})} p(t|\mathbf{x}) dt \quad (232)$$

105 which implies that the conditional median is the solution for the case of  $q = 1$ . As  $q \rightarrow 0$ ,  $|y(\mathbf{x}) - t|^q$   
 106 is close to 1 except in a small neighborhood around  $y(\mathbf{x})$  where it falls to 0. Therefore, the value  
 107 of (229) is close to 1. Normally,  $p(t|\mathbf{x})$  is normalized, but in this case there is a notch at  $t = y(\mathbf{x})$ .  
 108 Therefore, the value of (229) is slightly reduced from 1. To minimize (229), we can choose  $y(\mathbf{x})$  to  
 109 coincide with the maximum  $p(t|\mathbf{x})$ , namely, the conditional mode. □



### Exercise 1.28

In Section 1.6, we introduced the idea of entropy  $h(x)$  as the information gained on observing the value of a random variable  $x$  having distribution  $p(x)$ . We saw that, for independent variables  $x$  and  $y$  for which  $p(x, y) = p(x)p(y)$ , the entropy functions are additive, so that  $h(x, y) = h(x) + h(y)$ . In this exercise, we derive the relation between  $h$  and  $p$  in the form of a function  $h(p)$ . First show that  $h(p^2) = 2h(p)$ , and hence by induction that  $h(p^n) = nh(p)$  where  $n$  is a positive integer. Hence show that  $h(p^{n/m}) = (n/m)h(p)$  where  $m$  is also a positive integer. This implies that  $h(p^x) = xh(p)$  where  $x$  is a positive rational number, and hence by continuity when it is a positive real number. Finally, show that this implies  $h(p)$  must take the form  $h(p) \propto \ln p$ .

*Proof.* Since  $h(x, y) = h(x) + h(y)$ , we have  $h(p^2) = h(p, p) = h(p) + h(p) = 2h(p)$ . Then suppose that  $h(p^{n-1}) = (n-1)h(p)$ , and we get  $h(p^n) = h(p^{n-1}) + h(p) = nh(p)$ . With this, we have

$$h(p^{n/m}) = \frac{mh(p^{n/m})}{m} = \frac{h(p^n)}{m} = \frac{nh(p)}{m}. \quad (233)$$

Hence, by continuity,  $h(p^x) = xh(p)$  holds for any positive real number. Suppose  $p = q^x$ , then we have

$$\frac{h(p)}{\ln p} = \frac{h(q^x)}{\ln q^x} = \frac{xh(q)}{x \ln q} = \frac{h(q)}{\ln q} \quad (234)$$

which implies  $h(p) \propto \ln p$ .  $\square$

### Exercise 1.29

Consider an  $M$ -state discrete random variable  $x$ , and use Jensen's inequality in the form of

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i) \quad (235)$$

where  $\lambda_i \geq 0$  and  $\sum_i \lambda_i = 1$ , for any set of points  $\{x_i\}$ , to show that the entropy of its distribution  $p(x)$  satisfies  $H[x] \leq \ln M$ .

*Proof.* By the definition of entropy, we have

$$H[x] = -\sum_{i=1}^M p(x_i) \ln p(x_i) = \sum_{i=1}^M p(x_i) \ln \frac{1}{p(x_i)}. \quad (236)$$

Since  $\ln \frac{1}{x}$  is concave, by Jensen's inequality, we get

$$H[x] = \sum_{i=1}^M p(x_i) \ln \frac{1}{p(x_i)} \leq \ln \sum_{i=1}^M p(x_i) \frac{1}{p(x_i)} = \ln M. \quad (237)$$

which completes the proof.  $\square$

### Exercise 1.30

Evaluate the Kullback-Leibler divergence

$$\text{KL}(p||q) = \int p(x) \ln \frac{p(x)}{q(x)} dx \quad (238)$$

between two Gaussians,  $p(x) = \mathcal{N}(x|\mu, \sigma^2)$  and  $q(x) = \mathcal{N}(x|m^2, s^2)$ .

114

**Solution:**

$$\text{KL}(p||q) = \int \mathcal{N}(x|\mu, \sigma^2) \ln \frac{\mathcal{N}(x|\mu, \sigma^2)}{\mathcal{N}(x|m^2, s^2)} dx \quad (239)$$

$$= \int \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \ln \frac{\sqrt{s^2/\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{e^{-\frac{(x-m)^2}{2s^2}}} dx \quad (240)$$

$$= \frac{1}{2} \ln \frac{s^2}{\sigma^2} + \int \left( \frac{(x-m)^2}{2s^2} - \frac{(x-\mu)^2}{2\sigma^2} \right) \mathcal{N}(x|\mu, \sigma^2) dx \quad (241)$$

$$= \frac{1}{2} \ln \frac{s^2}{\sigma^2} + \int \frac{(x-m)^2}{2s^2} \mathcal{N}(x|\mu, \sigma^2) dx - \frac{\sigma^2}{2\sigma^2} \quad (242)$$

$$= \frac{1}{2} \ln \frac{s^2}{\sigma^2} + \int \frac{x^2 - 2mx + m^2}{2s^2} \mathcal{N}(x|\mu, \sigma^2) dx - \frac{1}{2} \quad (243)$$

$$= \frac{1}{2} \ln \frac{s^2}{\sigma^2} + \frac{\mu^2 + \sigma^2 - 2m\mu + m^2}{2s^2} - \frac{1}{2}. \quad (244)$$

115

□

### Exercise 1.31

Consider two variables  $\mathbf{x}$  and  $\mathbf{y}$  having joint distribution  $p(\mathbf{x}, \mathbf{y})$ . Show that the differential entropy of this pair of variables satisfies

$$H[\mathbf{x}, \mathbf{y}] \leq H[\mathbf{x}] + H[\mathbf{y}] \quad (245)$$

with equality if, and only if,  $\mathbf{x}$  and  $\mathbf{y}$  are statistically independent.

116

*Proof.* From Exercise 1.41, we have  $I(\mathbf{x}, \mathbf{y}) = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$ . Since the mutual information is a form of KL divergence,  $H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}] \geq 0$  which implies  $H[\mathbf{y}] \geq H[\mathbf{y}|\mathbf{x}]$ . With the relation  $H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}]$  which will be shown in Exercise 1.37, we have

$$H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}] \leq H[\mathbf{y}] + H[\mathbf{x}]. \quad (246)$$

If  $\mathbf{x}$  and  $\mathbf{y}$  are statistically independent, we have  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$ . Thus,

$$H(\mathbf{x}, \mathbf{y}) = - \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}, \mathbf{y}) d\mathbf{x}d\mathbf{y} \quad (247)$$

$$= - \iint p(\mathbf{x})p(\mathbf{y}) \ln p(\mathbf{x})p(\mathbf{y}) d\mathbf{x}d\mathbf{y} \quad (248)$$

$$= - \iint p(\mathbf{x})p(\mathbf{y}) \ln p(\mathbf{x}) d\mathbf{x}d\mathbf{y} - \iint p(\mathbf{x})p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{x}d\mathbf{y} \quad (249)$$

$$= H[\mathbf{x}] + H[\mathbf{y}] \quad (250)$$

117 which shows the sufficiency. For the necessity, if the equality in (246) holds, then  $H[\mathbf{y}|\mathbf{x}] = H[\mathbf{y}]$ .  
 118 According to Exercise 1.41,  $I(\mathbf{x}, \mathbf{y}) = 0$ . By the definition of  $I(\mathbf{x}, \mathbf{y})$ , the KL divergence of two  
 119 distributions is 0 if and only if the two distributions are identical, i.e.  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$ . In other  
 120 words,  $\mathbf{x}$  and  $\mathbf{y}$  are statistically independent. This completes the proof.  $\square$

### Exercise 1.32

Consider a vector  $\mathbf{x}$  of continuous variables with distribution  $p(\mathbf{x})$  and corresponding entropy  $H[\mathbf{x}]$ . Suppose that we make a nonsingular linear transformation of  $\mathbf{x}$  to obtain a new variable  $\mathbf{y} = \mathbf{A}\mathbf{x}$ . Show that the corresponding entropy is given by  $H[\mathbf{y}] = H[\mathbf{x}] + \ln|\mathbf{A}|$  where  $|\mathbf{A}|$  denotes the determinant of  $\mathbf{A}$ .

121

*Proof.* By the change of variables, we have

$$p(\mathbf{x}) = p(\mathbf{y})|\mathbf{A}|. \quad (251)$$

Then,

$$H[\mathbf{y}] = - \int p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{y} \quad (252)$$

$$= - \int p(\mathbf{x}) \ln p(\mathbf{x}) |\mathbf{A}|^{-1} d\mathbf{x} \quad (253)$$

$$= - \int p(\mathbf{x}) (\ln p(\mathbf{x}) - \ln |\mathbf{A}|) d\mathbf{x} \quad (254)$$

$$= - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} + \int p(\mathbf{x}) \ln |\mathbf{A}| d\mathbf{x} \quad (255)$$

$$= H[\mathbf{x}] + \ln |\mathbf{A}| \quad (256)$$

122 where the second equality we used  $p(\mathbf{y})d\mathbf{y} = p(\mathbf{x})d\mathbf{x}$ . This completes the proof.  $\square$

### Exercise 1.37

Using the definition of the conditional entropy of  $\mathbf{y}$  given  $\mathbf{x}$

$$H(\mathbf{y}|\mathbf{x}) = - \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{x} d\mathbf{y} \quad (257)$$

together with the product rule of probability, prove the result

$$H(\mathbf{x}, \mathbf{y}) = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}]. \quad (258)$$

123

*Proof.* By the definition of the conditional entropy, we have

$$H(\mathbf{y}|\mathbf{x}) = - \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{x} d\mathbf{y} \quad (259)$$

$$= - \iint p(\mathbf{x}, \mathbf{y}) \ln \frac{p(\mathbf{y}, \mathbf{x})}{p(\mathbf{x})} d\mathbf{x} d\mathbf{y} \quad (260)$$

$$= - \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{y}, \mathbf{x}) d\mathbf{x} d\mathbf{y} + \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}) d\mathbf{x} d\mathbf{y} \quad (261)$$

$$= H[\mathbf{x}, \mathbf{y}] - H[\mathbf{x}]. \quad (262)$$

After rearrangement, we get

$$H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}]. \quad (263)$$

124

$\square$

### Exercise 1.41

Using the sum and product rules of probability, show that the mutual information  $I(\mathbf{x}, \mathbf{y})$  satisfies the relation

$$I(\mathbf{x}, \mathbf{y}) = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]. \quad (264)$$

*Proof.* By the definition of mutual information,

$$I(\mathbf{x}, \mathbf{y}) \equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) || p(\mathbf{x})p(\mathbf{y})) \quad (265)$$

$$= \iint p(\mathbf{x}, \mathbf{y}) \ln \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} d\mathbf{x}d\mathbf{y} \quad (266)$$

$$= \iint p(\mathbf{x}, \mathbf{y}) \ln \frac{p(\mathbf{x}|\mathbf{y})}{p(\mathbf{x})} d\mathbf{x}d\mathbf{y} \quad (267)$$

$$= \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}|\mathbf{y}) d\mathbf{x}d\mathbf{y} - \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}) d\mathbf{x}d\mathbf{y} \quad (268)$$

$$= - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} - \left( - \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}|\mathbf{y}) d\mathbf{x}d\mathbf{y} \right) \quad (269)$$

$$= H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}]. \quad (270)$$

Similarly, we can get that

$$I(\mathbf{x}, \mathbf{y}) = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]. \quad (271)$$

□

## 2 Chapter 4 Linear Models for Classification

### 2.1 Discriminant Functions

#### 2.1.1 The derivation of Equation (4.5)

Equation (4.5) gives the normal distance from the origin to the decision surface  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0$ . Here is a thorough derivation with an illustration shown in Figure 1.

$$\mathbf{w}^T \mathbf{x} = \|\mathbf{w}\| \|\mathbf{x}\| \cos \theta = \|\mathbf{w}\| \cdot d \implies \frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = d. \quad (272)$$

Since any point  $\mathbf{x}$  on the decision surface satisfies  $\mathbf{w}^T \mathbf{x} + w_0 = 0$ , we have

$$\mathbf{w}^T \mathbf{x} = -w_0 \implies \frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = d = \frac{-w_0}{\|\mathbf{w}\|}. \quad (273)$$

You may have noticed Equation (4.7), i.e.  $r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$ , on Page 182 of the PRML textbook. When the point  $\mathbf{x}$  lies at the origin, namely  $\mathbf{x} = \mathbf{0}$ , it is easy to get

$$r = \frac{y(\mathbf{0})}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|}. \quad (274)$$

which does not contradict the form of  $d$ , though both represent the distance between the origin and the decision surface. Mathematically speaking, they are signed distances. We can follow the wording from the textbook to interpret  $d$  as the normal distance *from the origin to the decision surface* and  $r$  as the perpendicular (orthogonal) distance *from the decision surface to the point  $\mathbf{x}$* .

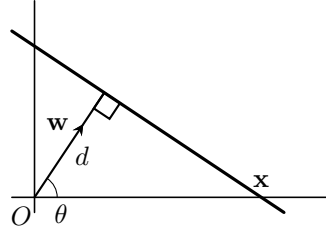


Figure 1: Illustration of the geometry of a linear discriminant function in two dimensions. The direction of  $\mathbf{w}$  depends on the form of the decision surface, shown in thick,  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ . This makes  $\theta$  in range  $[0, \pi]$ .

## 2.2 Exercises

### Exercise 4.1

Given a set of data points  $\mathbf{x}_n$ , we can define the *convex hull* to be the set of all points  $\mathbf{x}$  given by

$$\mathbf{x} = \sum_n a_n \mathbf{x}_n$$

where  $a_n \geq 0$  and  $\sum_n a_n = 1$ . Consider a second set of points  $\mathbf{y}_n$  together with their corresponding convex hull. By definition, the two sets of points will be linearly separable if there exists a vector  $\hat{\mathbf{w}}$  and a scalar  $w_0$  such that  $\hat{\mathbf{w}}^T \mathbf{x}_n + w_0 > 0$  for all  $\mathbf{x}_n$ , and  $\hat{\mathbf{w}}^T \mathbf{y}_n + w_0 < 0$  for all  $\mathbf{y}_n$ . Show that if their convex hulls intersect, the two sets of points cannot be linearly separable, and conversely that if they are linearly separable, their convex hulls do not intersect.

*Proof.* We prove the first part by contradiction. Given the convex hulls of the two sets of points intersect, we need to show that the two sets cannot be linearly separable. Since the two convex hulls intersect, there exists at least one point  $\mathbf{z} = \sum_n a_n \mathbf{x}_n = \sum_n b_n \mathbf{y}_n$  that lies in their intersection, where  $a_n, b_n \geq 0$  and  $\sum_n a_n = \sum_n b_n = 1$ . For the sake of contradiction, suppose that the two sets of points are linearly separable, then there exists a vector  $\hat{\mathbf{w}}$  and a scalar  $w_0$  such that  $\hat{\mathbf{w}}^T \mathbf{x}_n + w_0 > 0$  for all  $\mathbf{x}_n$ , and  $\hat{\mathbf{w}}^T \mathbf{y}_n + w_0 < 0$  for all  $\mathbf{y}_n$ . Furthermore,

$$\begin{aligned} \hat{\mathbf{w}}^T \mathbf{z} &= \hat{\mathbf{w}}^T \left( \sum_n a_n \mathbf{x}_n \right) = \sum_n a_n \hat{\mathbf{w}}^T \mathbf{x}_n = \sum_n a_n (\hat{\mathbf{w}}^T \mathbf{x}_n + w_0 - w_0) \\ &= \sum_n a_n (\hat{\mathbf{w}}^T \mathbf{x}_n + w_0) - w_0 \underbrace{\sum_n a_n}_{=1} \\ &= \sum_n a_n \underbrace{(\hat{\mathbf{w}}^T \mathbf{x}_n + w_0)}_{>0} - w_0 \\ &> -w_0 \end{aligned}$$

Similarly,

$$\begin{aligned} \hat{\mathbf{w}}^T \mathbf{z} &= \hat{\mathbf{w}}^T \left( \sum_n b_n \mathbf{y}_n \right) = \sum_n b_n \hat{\mathbf{w}}^T \mathbf{y}_n = \sum_n b_n (\hat{\mathbf{w}}^T \mathbf{y}_n + w_0 - w_0) \\ &= \sum_n b_n (\hat{\mathbf{w}}^T \mathbf{y}_n + w_0) - w_0 \underbrace{\sum_n b_n}_{=1} \end{aligned}$$

$$\begin{aligned}
&= \sum_n b_n \underbrace{(\hat{\mathbf{w}}^T \mathbf{y}_n + w_0)}_{<0} - w_0. \\
&< -w_0
\end{aligned}$$

136 which gives an obvious contradiction. This implies that the assumption does not hold. In other  
137 words, the two sets of points cannot be linearly separable.

Now we show the second half still by contradiction. Given the two sets are linearly separable, then there exists a vector  $\hat{\mathbf{w}}$  and a scalar  $w_0$  such that  $\hat{\mathbf{w}}^T \mathbf{x}_n + w_0 > 0$  for all  $\mathbf{x}_n$ , and  $\hat{\mathbf{w}}^T \mathbf{y}_n + w_0 < 0$  for all  $\mathbf{y}_n$ . Suppose their convex hulls intersect, then there exists at least one point  $\mathbf{z}$  such that  $\mathbf{z} = \sum_n a_n \mathbf{x}_n = \sum_n b_n \mathbf{y}_n$  with  $\sum_n a_n = \sum_n b_n = 1$  and  $a_n, b_n \geq 0$ . For any  $a_n > 0$ , we have

$$\begin{aligned}
\hat{\mathbf{w}}^T \mathbf{x}_n + w_0 > 0 &\implies a_n \hat{\mathbf{w}}^T \mathbf{x}_n + a_n w_0 > 0 \\
&\implies \hat{\mathbf{w}}^T a_n \mathbf{x}_n + a_n w_0 > 0.
\end{aligned}$$

For  $a_n = 0$ ,  $\hat{\mathbf{w}}^T a_n \mathbf{x}_n + a_n w_0 = 0$ . Summing over  $n$ , we get

$$\sum_n \hat{\mathbf{w}}^T a_n \mathbf{x}_n + a_n w_0 > 0 \implies \hat{\mathbf{w}}^T \underbrace{\sum_n (a_n \mathbf{x}_n)}_{=\mathbf{z}} + w_0 \underbrace{\sum_n a_n}_{=1} > 0 \implies \hat{\mathbf{w}}^T \mathbf{z} + w_0 > 0.$$

Likewise,

$$\sum_n \hat{\mathbf{w}}^T b_n \mathbf{y}_n + b_n w_0 < 0 \implies \hat{\mathbf{w}}^T \underbrace{\sum_n (b_n \mathbf{y}_n)}_{=\mathbf{z}} + w_0 \underbrace{\sum_n b_n}_{=1} < 0 \implies \hat{\mathbf{w}}^T \mathbf{z} + w_0 < 0.$$

138 which leads to a contradiction. This shows their convex hulls do not intersect. This completes the  
139 proof.

140

□