

Brief Article

The Author

September 14, 2013

I I.1

Straightforward. Calculate $\frac{\partial E(\mathbf{w})}{\partial w_j} = 0$ for $j = 0, \dots, M$.

2 I.2

Straightforward. Like I.1.

After reading the solution manual, I think it's much better to present A as a whole matrix.

3 I.3

$P(a) = 0.2 \times 3/10 + 0.2 \times 1/2 + 0.6 \times 3/10 = 0.34$. $P(g|o) = P(g)P(o|g)/P(o) = 0.6 * 0.3/0.36 = 0.5$.

4 I.4

By brute-force differentiation, we have

$$p'_y(y) = sp'_x(g)(g'(y))^2 + sp_x(g)g''(y) \quad (\text{I})$$

Assuming $\hat{x} = g(\hat{y})$, $p'_x(\hat{x}) = p'_y(\hat{y}) = 0$, we have $p'_y(\hat{y}) = sp'_x(\hat{x})(g'(\hat{y}))^2 + sp_x(\hat{x})g''(\hat{y}) = sp_x(\hat{x})g''(\hat{y}) = 0$, and this requires $g''(\hat{y}) = 0$.

In general, this won't be satisfied. However, when $g(\cdot)$ is linear, we have $g''(\hat{y}) = 0$.

In this question, the uniqueness about \hat{x} and \hat{y} is ignored, but this is irrelevant in most cases.

5 1.5

Straightforward.

6 1.6

Use the fact $p(x, y) = p(x)p(y)$ and change the order of integration (or integrate x and y separately)

7 1.7

Here, I just use the result $I = (2\pi\sigma^2)^{1/2}$. Then just use this result in the integration of (1.46), and we are done.

8 1.8

About (1.49): well... I have to admit that this problem is not as complicated as I thought... Let $y = x - \nu$ and notice it's the sum of the integration of a odd function, and the integration of a constant times the integration in 1.7.

About (1.50): I'm so poor at calculus...

9 1.9

Trivial... Perhaps we need that Σ^{-1} is positive-definite...

10 1.10

Trivial...

11 1.11

Trivial...

I2 I.I2

Trivial...

I3 I.I3

Trivial...

I4 I.I4

Trivial... Let $w_{ij}^S = (w_{ij} + w_{ji})/2$, and $w_{ij}^A = (w_{ij} - w_{ji})/2$.

I5 I.I5

(I.I34): given original $x_{i_1}, x_{i_2}, \dots, x_{i_M}$ (D^M in total), we can always arrange i_1, i_2, \dots in decreasing order. So, (I.I34) should suffice.

(I.I35): for each \tilde{w} , consider the corresponding arrangement of i_1, i_2, \dots, i_M . If $i_1 = D$, then other i_m ($M-1$ terms) must be less or equal than D , so for this case, there should be $n(D, M-1)$ possible arrangements of i_1, i_2, \dots, i_M . Similarly, considering the case $i_1 = m$, there're $n(m, M-1)$ terms. So we have (I.I35).

(I.I36): when $D = 1$, it's obviously true. Then it's trivial..

(I.I37): obvious for $M = 2$. For $M > 2$, use (I.I36) and see that its form is just (I.I35), so we're done.

I6 I.I6

(I.I38): trivial.

(I.I39): trivial.

(I.I40): for each case ($D \gg M$ or $M \gg D$), set M or D constant, and expand using Stirling's formula.

Hint: you can multiply the intermediate result by M^D (the thing you desire) and divide it by M^D again. Things can become clear then.

I7 I.I7

trivial.. just remember $\int u dv = uv - \int v du$

18 1.18

(1.142): for left side, use (1.126). for right side, use (1.141) and use $2rdr = dr^2$.

(1.144): here, we conveniently assume for a sphere of radius r , $S_D(r) = S_D(1)r^{D-1}$, and V_D is calculated by summing the volumes of many small sphere shells ($S_D(r)dr$).

19 1.19

(1.145): volume of cube is $(2a)^D$. Letting $a = 1$, we have (1.145).

ratio of (1.146): trivial.

ratio of that in the text: distance to sides is a , and distance from center to corner is $\sqrt{Da^2}$, and the ratio \sqrt{D} follows.

20 1.20

skip for now...

(1.148): I think (45) in reference solution manual is wrong. (1.148) is somewhat obvious. $p(r)$ should be $p(\mathbf{x})$ integrated over a specific sphere of radius r , and this sphere has area $S_D r^{D-1}$, so we get (1.148).

stationary point: differentiate $p(r)$ with respect to r .

stationary point plus ϵ : a lot of approximation... see reference solution.

density at origin and that at \hat{r} : trivial.

21 1.21

$a \leq (ab)^{1/2}$: trivial.

The inequality: $p(\text{mistake})$ is an integration over \mathbf{x} , and so is the right side. At every \mathbf{x} , $p(\text{mistake})$ is the smaller (a) of $p(\mathbf{x}, \mathcal{C}_1)$ and $p(\mathbf{x}, \mathcal{C}_2)$. Then we're done.

22 1.22

Trivial. Interpretation: classification rate.

23 1.23

Trivial.

24 I.24

The phrase “decision criterion” means a rule to follow when making decisions. Given \mathbf{x} , with true label unknown, if we reject it, the loss is λ . If we classify it as class j , the loss is L_{kj} , assuming that its label is k . However, we don’t know its true label, and the distribution of this is given by $p(\mathcal{C}|\mathbf{x})$, so the expected loss if we classify it as class j is $\sum_k L_{kj}p(\mathcal{C}_k|\mathbf{x})$. If the minimum of this value (over j) is greater than λ , we reject. Otherwise, we choose j .

If $L_{kj} = 1 - I_{kj}$, then $\sum_k L_{kj}p(\mathcal{C}_k|\mathbf{x}) = 1 - p(\mathcal{C}_j|\mathbf{x})$. Therefore, the criterion becomes: if the minimum of $1 - p(\mathcal{C}_j|\mathbf{x})$ over j is greater than λ , we reject. Otherwise, we choose j . Here, $p(\mathcal{C}_j|\mathbf{x})$ is largest. So the criterion can be reformulated as: if largest of $p(\mathcal{C}_j|\mathbf{x})$ is less than $1 - \lambda$, we reject. So we have $\theta = 1 - \lambda$.

25 I.25

Skip...

26 I.26

Trivial...

27 I.27

Skip...

28 I.28

I think the wording and notation in page 48 is somewhat confusing. In my understanding, $h(\cdot)$ is a univariate function taking the probability of event x . So, $h(x, y)$ should be written as $h(p(x, y)) = h(p(x)p(y)) = h(p(x)) + h(p(y))$.

Using this notation, we first see $h(p^2) = h(p)h(p)$. Then by induction we trivially have $h(p^n) = h(p^{n-1}) + h(p) = (n-1)h(p) + h(p) = nh(p)$. Then, regarding $p^{1/m}$ as a whole, we have $h(p^{n/m}) = nh(p^{1/m}) = (n/m)mh(p^{1/m}) = n/mh(p)$. By continuity, we have $h(p^x) = xh(p)$. Last, given two positive real numbers $p, q = p^x$, we have

$$\frac{h(q)}{\ln q} = \frac{xh(p)}{x \ln p} = \frac{h(p)}{\ln p}. \quad (2)$$

Thus, we have $h(p) \propto \ln p$.

Hint: sometimes, we should regard different things as the “unit” to be learned, like $p^{1/m}$ and p .

29 I.29

$$H(x) = - \sum_{i=1}^M p(x_i) \ln p(x_i) = \sum_{i=1}^M p(x_i) \ln(1/p(x_i)) \quad (3)$$

Since $\ln(\cdot)$ is concave, the sign in the (1.115) should be reversed.

$$H(x) \leq \ln \sum_{i=1}^M (p(x_i)/p(x_i)) = \ln M. \quad (4)$$

Hint: it's wrong to let $f(x) = -\ln(x)$. This will lead to $H(x) \geq -\ln \sum_{i=1}^M p(x_i)^2 \leq \ln M$, which is true and useless.

30 I.30

Trivial... Just need patience and carefulness.

31 I.31

Trivial... The solution in the manual seems redundant... Use $\text{KL}(p(x, y) \| p(x)p(y)) \geq 0$, and everything follows.

32 I.32

See the manual... I have no idea of Jacobian...

33 I.33

Trivial... But we have to assume that $x_1 \ln x_2 = 0$, whenever $x_1 = 0$ or $x_2 = 0$.

34 I.34

Skip...

35 I.36

Refer to the manual... It's so tricky...

36 I.37

$$H[\mathbf{x}, \mathbf{y}] = - \iint p(\mathbf{x}, \mathbf{y}) \ln p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} \quad (5)$$

$$= - \iint p(\mathbf{x}, \mathbf{y}) \ln(p(\mathbf{y}|\mathbf{x}) + p(\mathbf{x})) d\mathbf{x} d\mathbf{y} \quad (6)$$

By separating the logarithm of joint probability, the two terms become $H(\mathbf{x})$ and $H(\mathbf{y}|\mathbf{x})$ respectively.

37 I.38

Trivial...