

# Bishop Questions

David Ruhe

March 30, 2020

## Chapter 1

### 1.1

Differentiating to  $\{w_i\}$  means differentiating to every single  $w_i$  in  $\mathbf{w}$ . Sum of squares:  $E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$ .

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

$$= x_n^i \sum_{n=1}^N \left\{ \sum_{j=0}^M w_j x_n^j - t_n \right\} = 0 \quad (2)$$

$$= \sum_{n=1}^N \left\{ \sum_{j=0}^M w_j (x_n)^{i+j} - (x_n)^i t_n \right\} = 0 \quad (3)$$

$$= \sum_{n=1}^N \sum_{j=0}^M \{w_j (x_n)^{i+j}\} - \sum_{n=1}^N (x_n)^i t_n = 0 \quad (4)$$

$$= \sum_{n=1}^N \sum_{j=0}^M \{w_j (x_n)^{i+j}\} = \sum_{n=1}^N (x_n)^i t_n \quad (5)$$

$$= \sum_{j=0}^M w_j \sum_{n=1}^N (x_n)^{i+j} = \sum_{n=1}^N (x_n)^i t_n \quad (6)$$

$$= \sum_{j=0}^M w_j A_{ij} = T_i \quad (7)$$

## 1.2

$$\tilde{E} = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t - N\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

$$\frac{\partial \tilde{E}(\mathbf{w})}{\partial w_i} = x_n^i \sum_{n=1}^N \left\{ \sum_{j=1}^M w_j x_n^j - t_n \right\} + \lambda w_i \quad (8)$$

$$= \sum_{n=1}^N \left\{ \sum_{j=0}^M w_j x_n^{j+i} - t_n + \frac{\lambda}{N} w_i \right\} = 0 \quad (9)$$

$$= \sum_{n=1}^N \left\{ \sum_{j=0}^M w_j x_n^{j+i} + \frac{\lambda}{N} w_i \right\} = \sum_{n=1}^N x_n^i t_n \quad (10)$$

$$= \sum_{j=0}^M w_j \left\{ \sum_{n=1}^N x_n^{j+i} + \frac{\lambda}{N} w_i \right\} = \sum_{n=1}^N x_n^i t_n \quad (11)$$

$$= \sum_{j=0}^M w_j A_{ij} = T_i \quad (12)$$

## 1.3

$$p(F = a) = \sum_b p(F = a | B = b) p(B = b) = 0.3 \cdot 0.2 + 0.2 \cdot 0.5 + 0.6 \cdot 0.3 = 0.34 \quad (13)$$

$$p(B = g | F = o) = \frac{p(F = o | B = g) p(B = g)}{p(F = o)} \quad (14)$$

$$= \frac{0.3 \cdot 0.6}{\sum_b p(F = o | B = b) p(B = b)} \quad (15)$$

$$= \frac{0.3 \cdot 0.6}{0.4 \cdot 0.2 + 0.5 \cdot 0.2 + 0.3 \cdot 0.6} = 0.5 \quad (16)$$

## 1.5

$$\text{var}[f] = \mathbb{E} [(f(x) - \mathbb{E}[f(x)])^2] \quad (17)$$

$$= \mathbb{E} [f(x)^2 + \mathbb{E}[f(x)]^2 - 2f(x) \mathbb{E}[f(x)]] \quad (18)$$

$$= \mathbb{E}[f(x)^2] + \mathbb{E}[f(x)]^2 - 2 \mathbb{E}[f(x)] \mathbb{E}[f(x)] \quad (19)$$

$$= \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2 \quad (20)$$

## 1.6

$$\text{cov}[x, y] = \mathbb{E}_{x,y}[\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}] \quad (21)$$

$$= \int_{x,y} p(x, y) [\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}] \quad (22)$$

$$= \int_x \int_y p(x) p(y) [\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}] \quad (23)$$

$$= \int_x p(x) \{x - \mathbb{E}[x]\} \int_y p(y) \{y - \mathbb{E}[y]\} \quad (24)$$

$$= \mathbb{E}[x - \mathbb{E}[x]] \mathbb{E}[y - \mathbb{E}[y]] \quad (25)$$

$$= (\mathbb{E}[x] - \mathbb{E}[x])(\mathbb{E}[y] - \mathbb{E}[y]) \quad (26)$$

## 1.9

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\} \quad (27)$$

$$\frac{\partial \mathcal{N}(x|\mu, \sigma^2)}{\partial \mu} = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\} \sigma^2(x - \mu) = 0 \quad (28)$$

$$x - \mu = 0 \quad (29)$$

$$x = \mu \quad (30)$$

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^2} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right] \quad (31)$$

$$\frac{\partial \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}} = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^2} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right] \left[-\frac{1}{2}(\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}^{-1T})(\mathbf{x} - \boldsymbol{\mu})\right] = 0 \quad (32)$$

$$(33)$$

So only when  $\mathbf{x} = \boldsymbol{\mu}$ . Used Matrix cookbook and diagonality of  $\boldsymbol{\Sigma}^{-1}$ .

## 1.10

$$\mathbb{E}[x + z] = \int_{x+z} p(x + z) [x + z] \quad (34)$$

$$= \int_{x+z} p(x + z) x + \int_{x+z} p(x + z) z \quad (35)$$

$$= \int_{x,z} p(x, z) x + \int_{x,z} p(x, z) z \quad (36)$$

$$= \int_{x,z} p(x) p(z) x + \int_{x,z} p(x) p(z) z \quad (37)$$

$$= \int_x p(x) x + \int_z p(z) z \quad (38)$$

$$= \mathbb{E}[x] + \mathbb{E}[z] \quad (39)$$

$$\text{var}[x + z] = \mathbb{E} \left[ \{(x + z) - \mathbb{E}[x + z]\}^2 \right] \quad (40)$$

$$= \mathbb{E} \left[ (x + z)^2 \right] - \mathbb{E}[x + z]^2 \quad (41)$$

$$= \mathbb{E} \left[ x^2 + z^2 + 2xz \right] - (\mathbb{E}[x] + \mathbb{E}[z])^2 \quad (42)$$

$$= \mathbb{E}[x^2] + \mathbb{E}[z^2] + \mathbb{E}[2xz] - \mathbb{E}[x]^2 - \mathbb{E}[z]^2 - 2\mathbb{E}[x]\mathbb{E}[z] \quad (43)$$

$$= \text{var}[x] + \text{var}[y] + 2 \int_{x,z} p(x, z) xz - 2 \int_x p(x) \int_z p(z) \quad (44)$$

$$= \text{var}[x] + \text{var}[y] + 2 \int_{x,z} p(x) p(z) xz - 2 \int_x p(x) \int_z p(z) \quad (45)$$

$$= \text{var}[x] + \text{var}[y] + 2 \int_x p(x) \int_z p(z) - 2 \int_x p(x) \int_z p(z) \quad (46)$$

## 1.11

Mean:

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln 2\pi \quad (47)$$

$$\frac{\partial \ln p(\mathbf{x}|\mu, \sigma^2)}{\partial \mu} = \sigma^2 \sum_{n=1}^N (x_n - \mu) = 0 \quad (48)$$

$$\sum_{n=1}^N x_n = N\mu \quad (49)$$

$$\mu = \frac{1}{N} \sum_{n=1}^N x_n \quad (50)$$

Variance:

$$-\frac{1}{2\sigma^2} = -(2\sigma^2)^{-1} \quad (51)$$

$$\frac{\partial \ln p(\mathbf{x}|\mu, \sigma^2)}{\partial \sigma^2} = 2 \cdot \frac{1}{4\sigma^4} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2\sigma^2} = 0 \quad (52)$$

$$N \cdot 2\sigma^4 = 2\sigma^2 \sum_{n=1}^N (x_n - \mu)^2 \quad (53)$$

$$\sigma^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2 \quad (54)$$

### 1.13

From 1.56:

$$\mathbb{E} \left[ \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2 \right] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[x^2] - \mathbb{E}[2x_n\mu] + \mathbb{E}[\mu^2] \quad (55)$$

$$= \mathbb{E}[x^2] - 2\mu \mathbb{E}[x] + \mu^2 \quad (56)$$

$$= \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2 \quad (57)$$

### 1.22

$$1 - I_{kj} = \begin{bmatrix} 0 & 1 & \dots & 1 \\ 1 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 0 \end{bmatrix} \quad (58)$$

$$\mathcal{R}_j^* = \arg \min_j L_{kj} p(C_k | \mathbf{x}) \quad (59)$$

$$= \arg \min_j \sum_{j \neq k} L_{kj} p(C_k | \mathbf{x}) \quad (60)$$

$$= \arg \max_j p(C_k | \mathbf{x}) \quad (61)$$

### 1.23

### 1.25

$$\mathbb{E} [L(\mathbf{t}, \mathbf{y}(\mathbf{x}))] = \int \int \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t} \quad (62)$$

We want to find  $\mathbb{E}_t[\mathbf{t} | \mathbf{x}]$  so we switch the order of integrals.

$$= \int \int \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{t} d\mathbf{x} \quad (63)$$

$$F[\mathbf{x}, \mathbf{y}, \mathbf{t}] = \int_{\mathbf{t}} \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) \quad (64)$$

Euler-Lagrange:

$$\frac{\partial F}{\partial \mathbf{y}} - \frac{d}{dx} \left( \frac{\partial F}{\partial \mathbf{y}'} \right) = 0 \quad (65)$$

$$\frac{\partial F}{\partial \mathbf{y}'} = 0 \quad (66)$$

$$\frac{\partial F}{\partial \mathbf{y}} = \int_t (\mathbf{y}(\mathbf{x}) - t) p(\mathbf{x}, t) = 0 \quad (67)$$

$$\int_t \mathbf{y}(\mathbf{x}) p(\mathbf{x}, t) - \int_t t p(\mathbf{x}, t) = 0 \quad (68)$$

$$\mathbf{y}(\mathbf{x}) \int_t p(t, \mathbf{x}) = \int_t t p(\mathbf{x}, t) \quad (69)$$

$$\mathbf{y}(\mathbf{x}) = \frac{1}{p(\mathbf{x})} \int_t t p(\mathbf{x}, t) = \mathbb{E}_t[t|\mathbf{x}] \quad (70)$$

$$(71)$$

## 1.26

First, let's derive the single value case.

$$\{y(\mathbf{x}) - t\}^2 = \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \quad (72)$$

$$= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + \{\mathbb{E}[t|\mathbf{x}] - t\}^2 + 2[y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]]\{\mathbb{E}[t|\mathbf{x}] - t\} \quad (73)$$

$$\mathbb{E}[L] = \underbrace{\mathbb{E}_{x,t} [(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2]}_{(1)} + \underbrace{\mathbb{E}_{x,t} [(\mathbb{E}[t|\mathbf{x}] - t)^2]}_{(2)} + \underbrace{\mathbb{E}_{x,t} [2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) (\mathbb{E}[t|\mathbf{x}] - t)]}_{(3)} \quad (74)$$

(1)

$$\begin{aligned} \mathbb{E}_{x,t} [(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2] &= \mathbb{E}_x \mathbb{E}_t [(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2 | \mathbf{x}] \\ &= \mathbb{E}_x (y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2 \end{aligned} \quad (75)$$

(Since the function is over  $x$ , the expectation over  $t$  vanishes)  
(76)

(2)

$$\mathbb{E}_{x,t} [(\mathbb{E}[t|\mathbf{x}] - t)^2] = \mathbb{E}_x [\mathbb{E}_t [\mathbb{E}[t|\mathbf{x}]^2 + t^2 - 2\mathbb{E}[t|\mathbf{x}]t | \mathbf{x}]] \quad (77)$$

$$\begin{aligned} &= \mathbb{E}_x [\mathbb{E}[t|\mathbf{x}]^2 + \mathbb{E}[t^2|\mathbf{x}] - 2\mathbb{E}_t[t|\mathbf{x}]\mathbb{E}_t[t|\mathbf{x}]] \\ &(\mathbb{E}[t|\mathbf{x}] \text{ can be taken out of the inner expectation since } t \text{ is integrated out.}) \\ &= \mathbb{E}_x [\text{var}[t|\mathbf{x}]] \end{aligned} \quad (78)$$

(3)

$$\mathbb{E}_{x,t} [2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) (\mathbb{E}[t|\mathbf{x}] - t)] = \mathbb{E}_x \mathbb{E}_t [2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) (\mathbb{E}[t|\mathbf{x}] - t) | \mathbf{x}] \quad (79)$$

$$= \mathbb{E}_x [2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) \mathbb{E}_t [(\mathbb{E}[t|\mathbf{x}] - t) | \mathbf{x}]] \quad (80)$$

$$= \mathbb{E}_x [2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) \cdot 0] \quad (81)$$

$$= 0 \quad (82)$$

Now we continue for the vector case.

$$\|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 = \|\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}\|^2 \quad (83)$$

$$= (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^T (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}) \quad (84)$$

$$= (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^2 + (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^2 + 2 [(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})] \quad (85)$$

$$\mathbb{E}[L(\mathbf{t}, \mathbf{y}(\mathbf{x}))] = \mathbb{E}_{x,t} [\underbrace{(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^2}_{(1)} + \underbrace{(\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^2}_{(2)} + \underbrace{2 [(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})]}_{(3)}] \quad (86)$$

(1)

$$\mathbb{E}_{x,t} [(\mathbf{y}(\mathbf{x})) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^2] = \mathbb{E}_x \mathbb{E}_t [(\mathbf{y}(\mathbf{x})) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^2 | \mathbf{x}] \quad (87)$$

$$= \mathbb{E}_x (\mathbf{y}(\mathbf{x})) - \mathbb{E}[\mathbf{t}|\mathbf{x}]^2 \quad (88)$$

(2)

$$\mathbb{E}_{x,t} [(\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^2] = \mathbb{E}_x \mathbb{E}_t [(\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^2 | \mathbf{x}] \quad (89)$$

$$= \mathbb{E}_x \mathbb{E}_t [\mathbb{E}[\mathbf{t}|\mathbf{x}]^T \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbf{t}^T \mathbf{t} - 2 \mathbb{E}[\mathbf{t}|\mathbf{x}]^T \mathbf{t} | \mathbf{x}] \quad (90)$$

$$= \mathbb{E}_x [\mathbb{E}[\mathbf{t}|\mathbf{x}]^T \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbb{E}_t [\mathbf{t}^T \mathbf{t}] - 2 \mathbb{E}[\mathbf{t}|\mathbf{x}]^T \mathbb{E}[\mathbf{t}|\mathbf{x}]] \quad (91)$$

$$= \mathbb{E}_x [\text{var}[\mathbf{t}|\mathbf{x}]] \quad (92)$$

(3)

$$\mathbb{E}_{x,t} \{2 [\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})\} = \mathbb{E}_x \mathbb{E}_t \{2 ([\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}) | \mathbf{x}\} \quad (93)$$

$$= \mathbb{E}_x 2 (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T \mathbb{E}_t [\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}] \quad (94)$$

$$= \mathbf{0} \quad (95)$$

## 1.28

$n = k = 2$

$$h(p^n) = -\log(p(x)^2) = -2 \log p(x) = 2h(p) \quad (96)$$

Induction:

$$h(p^{k+1}) = -\log p(x)^{k+1} = -(k+1) \log(p(x)) = -(k+1) \log p(x) = (k+1)h(p) \quad (97)$$

So it works for  $k+1$ .  $\frac{n}{m}$  is straightforward.

$$h(p) = -\log_z(p) = -\frac{\ln p}{\ln z} \propto \ln p \quad (???)$$

### 1.29

$$H[x] = - \sum_{n=1}^M p(x) \ln p(x) \quad (98)$$

$$= \sum_{n=1}^M p(x) \ln \frac{1}{p(x)} \quad (99)$$

$$\leq \ln \sum_{n=1}^M p(x) \frac{1}{p(x)} \quad (\text{Concavity of } \ln \cdot)$$

$$= \ln M \quad (100)$$

### 1.32

$$H[\mathbf{y}] = - \int p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{y} \quad (101)$$

$$= - \int p(\mathbf{A}\mathbf{x}) \ln p(\mathbf{A}\mathbf{x}) \left| \frac{d\mathbf{y}}{d\mathbf{x}} \right| d\mathbf{x} \quad (\text{Integration by substitution})$$

$$= - \int \frac{p(\mathbf{x})}{|\mathbf{A}|} \ln \frac{p(\mathbf{x})}{|\mathbf{A}|} |\mathbf{A}| d\mathbf{x} \quad (p(\mathbf{x}) = p(\mathbf{y})|\mathbf{A}| \text{ for a function } f(\mathbf{x}) \rightarrow \mathbf{x})$$

$$= - \int p(\mathbf{x}) \ln \frac{p(\mathbf{x})}{|\mathbf{A}|} d\mathbf{x} \quad (102)$$

$$= - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} + \int p(\mathbf{x}) \ln |\mathbf{A}| d\mathbf{x} \quad (103)$$

$$= H[\mathbf{x}] + \ln |\mathbf{A}| \int p(\mathbf{x}) d\mathbf{x} \quad (104)$$

$$= H[\mathbf{x}] + \ln |\mathbf{A}| \quad (105)$$

### 1.34

$$-H[x] = \int p(x) \ln p(x) dx \quad (106)$$

$$= \int p(x) \ln \frac{1}{\sqrt{2\pi\sigma^2}} dx - \int p(x) \frac{(x-\mu)^2}{2\sigma^2} dx \quad (107)$$

$$= \int p(x) \ln 1 dx - \int p(x) \ln \sqrt{2\pi\sigma^2} dx - \int p(x) \frac{(x-\mu)^2}{2\sigma^2} dx \quad (108)$$

$$= -\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \int p(x) (x-\mu)^2 dx \quad (\text{Since } x \text{ is integrated out of } \sigma^2)$$

$$= -\frac{1}{2} \ln 2\pi\sigma^2 - \frac{\sigma^2}{2\sigma^2} \quad (109)$$

$$= -\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2} \quad (110)$$

$$H[x] = \frac{1}{2} [\ln 2\pi\sigma^2 + 1] \quad (111)$$



### 1.36

For  $(a, b) \in \mathbb{R}$  and  $b \geq a$ , convexity means  $f'(b) \geq f'(a)$ . Therefore, We have

$$f''(x) = \frac{f'(b) - f'(a)}{b - a} \geq 0$$

### 1.37

$$-H[x, y] = \int \int p(y, x) \ln p(x, y) dx dy \quad (112)$$

$$= \int \int p(y, x) \ln [p(y|x)p(x)] dx dy \quad (113)$$

$$= \int \int p(y, x) \ln p(y|x) + p(y, x) \ln p(x) dx dy \quad (114)$$

$$= -H[y|x] + \int \int p(y, x) dy \ln p(x) dx \quad (115)$$

$$= -H[y|x] - H[x] \quad (116)$$

### 1.40

$$\ln \prod_{i=1}^n a_i^{1/n} = \frac{1}{n} \ln \prod_{i=1}^n a_i \quad (117)$$

$$= \frac{1}{n} \sum_{i=1}^n \ln a_i \quad (118)$$

$$\geq \ln \frac{1}{n} \sum_{i=1}^n a_i \quad (\text{Jensen})$$

$$(119)$$

By monotonicity of logarithm we have  $\prod_{i=1}^n a_i^{1/n} \geq \frac{1}{n} \sum_{i=1}^n a_i$ .

### 1.41

$$-I[x, y] = \int \int [p(x, y) \ln p(x) + p(x, y) \ln p(y) - p(x, y) \ln p(x|y) - p(x, y) \ln p(y)] dx dy \quad (120)$$

$$= -H[x] + H[x|y] \quad (121)$$

## 2.1

$$p(x|\mu) = \mu^x (1 - \mu)^{1-x} \quad (122)$$

$$\sum_{x=0}^1 p(x|\mu) = 1 - \mu + \mu = 1 \quad (123)$$

$$\mathbb{E}[x] = \sum_{x=0}^1 x \mu^x (1-\mu)^{1-x} \quad (124)$$

$$= \mu(1-\mu)^0 = \mu \quad (125)$$

$$\text{Var}[x] = [(x - \mathbb{E}[x])^2] \quad (126)$$

$$= [(x - \mu)^2] \quad (127)$$

$$= \sum_{i=0}^1 (x - \mu)^2 \mu^x (1-\mu)^{1-x} \quad (128)$$

$$(129)$$

Writing out this sum and the resulting squares results in the desired result.

$$H[x] = \sum i = 0^1 p(x) \ln p(x) = \mu \ln \mu + (1-\mu) \ln(1-\mu) \quad (130)$$

## 2.6

$$p(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \quad (131)$$

$$\int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \quad (132)$$

$$\Gamma(x) = \int_0^\infty \mu^{x-1} e^{-u} du \quad (133)$$

$$\mathbb{E}[\mu] = \int_0^1 p(\mu|a, b) \mu d\mu \quad (134)$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \mu d\mu \quad (135)$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^a (1-\mu)^{b-1} d\mu \quad (136)$$

Since added 1 in the exponent this now looks like a beta distribution with  $p(\mu, a+1, b)$

$$= \frac{a}{a+b} \int_0^1 \frac{\Gamma(a+1+b)}{\Gamma(a+1)\Gamma(b)} \mu^a (1-\mu)^{b-1} d\mu \quad (137)$$

$$= \frac{a}{a+b} \quad (138)$$

Where we used

$$\Gamma(a+1) = a\Gamma(a) \text{ and } \Gamma(a+b+1) = (a+b)\Gamma(a+b) \quad (139)$$

$$\text{Var}[\mu] = \mathbb{E}[\mu^2] - \mathbb{E}[\mu]^2 \quad (140)$$

$$\mathbb{E}[\mu] = \int_0^1 p(\mu|a, b) \mu d\mu^2 \quad (141)$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a+1} (1-\mu)^{b-1} \mu d\mu \quad (142)$$

Which looks like Beta with a+2

$$= \frac{a}{a+b} \int_0^1 \frac{\Gamma(a+1+b)}{\Gamma(a+1)\Gamma(b)} \mu^{a+2} (1-\mu)^{b-1} d\mu \quad (143)$$

$$= \frac{a(a+1)}{(a+b)(a+1+b)} \int_0^1 \frac{\Gamma(a+2+b)}{\Gamma(a+2)\Gamma(b)} \mu^{a+2} (1-\mu)^{b-1} d\mu \quad (144)$$

$$= \frac{a(a+1)}{(a+b)(a+1+b)} \quad (145)$$

$$\text{Var}[\mu] = \frac{a(a+1)}{(a+b)(a+1+b)} - \frac{a^2}{(a+b)^2} \quad (146)$$

$$= \frac{(a+b)a(a+1)}{(a+b)^2(a+1+b)} - \frac{(a+b+1)a^2}{(a+b+1)(a+b)^2} \quad (147)$$

$$= \frac{(a^2+ab)(a+1)}{(a+b)^2(a+1+b)} - \frac{(a+b+1)a^2}{(a+b+1)(a+b)^2} \quad (148)$$

$$= \frac{a^3+a^2b+a^2+ab}{(a+b)^2(a+1+b)} - \frac{a^3+ba^2+a^2}{(a+b+1)(a+b)^2} \quad (149)$$

Which gives you the desired result.

Mode is given where the derivative w.r.t.  $\mu$  is zero:

$$\frac{\partial p(\mu|a, b)}{\partial \mu} = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} [(a-1)\mu^{a-2} - \mu^{a-1}(b-1)(1-\mu)^{b-2}] = 0 \quad (150)$$

$$(a-1)\mu^{a-2}(1-\mu)^{b-1} = \mu^{a-1}(b-1)(1-\mu)^{b-2} \quad (151)$$

$$\frac{a-1}{\mu} = \frac{b-1}{1-\mu} \quad (152)$$

$$\mu = \frac{-a+1}{-a-b+2} = \frac{a-1}{a+b-2} \quad (153)$$

## 2.8

$$\mathbb{E}[x] = \int p(x)xdx \quad (154)$$

$$= \int \int p(x,y)xdxdy \quad (155)$$

$$= \int \int p(x|y)p(y)xdxdy \quad (156)$$

$$= \int p(y) \int p(x|y)xdxdy \quad (157)$$

$$= \int p(y) \mathbb{E}_x[x|y]dy \quad (158)$$

$$= \mathbb{E}_y \mathbb{E}_x[x|y] \quad (159)$$

$$\mathbb{E}_y[\text{Var}_x[x|y]] = \mathbb{E}_y \mathbb{E}_{x|y}[(x - \mathbb{E}[x|y])^2] \quad (160)$$

$$= \mathbb{E}_y \mathbb{E}_{x|y}[(x^2 + \mathbb{E}[x|y]^2 - 2x \mathbb{E}[x|y])] \quad (161)$$

$$= \mathbb{E}_y \mathbb{E}_{x|y} x^2 + \mathbb{E}_y \mathbb{E}_{x|y} \mathbb{E}[x|y]^2 - 2 \mathbb{E}_y \mathbb{E}_{x|y} x \mathbb{E}[x|y]] \quad (162)$$

$$= \mathbb{E}_y \int p(x|y)x^2dx + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \int p(y) \int p(x|y)x \mathbb{E}[x|y]dxdy \quad (163)$$

$$= \int p(y) \int p(x|y)x^2dxdy + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \int \mathbb{E}[x|y] \int p(x|y)xdxdy \quad (164)$$

$$= \int \int p(x,y)x^2dxdy + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \int \mathbb{E}[x|y]^2dy \quad (165)$$

$$= \int \int p(x,y)x^2dxdy + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \mathbb{E}_y \mathbb{E}[x|y]^2] \quad (166)$$

$$= \mathbb{E}[x^2] - \mathbb{E}_y \mathbb{E}[x|y]^2 \quad (167)$$

$$(168)$$

$$\text{Var}_y[\mathbb{E}_x[x|y]] = \mathbb{E}_y [(\mathbb{E}_x[x|y] - \mathbb{E}_y \mathbb{E}_x[x|y])^2] \quad (169)$$

$$= \mathbb{E}_y [(\mathbb{E}_x[x|y] - \mathbb{E}_x[x])^2] \quad (170)$$

$$= \mathbb{E}_y [(\mathbb{E}_x[x|y]^2 + \mathbb{E}_x[x]^2 - 2 \mathbb{E}_x[x] \mathbb{E}_x[x|y])] \quad (171)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] + \mathbb{E}_y[\mathbb{E}_x[x]^2] - 2 \mathbb{E}_y \mathbb{E}_x[x] \mathbb{E}_x[x|y] \quad (172)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] + \mathbb{E}_x[x]^2 - 2 \mathbb{E}_x[x] \mathbb{E}_y \mathbb{E}_x[x|y] \quad (173)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] + \mathbb{E}_x[x]^2 - 2 \mathbb{E}_x[x]^2 \quad (174)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] - \mathbb{E}_x[x]^2 \quad (175)$$

$$(176)$$

Sum of these expressions gives desired result.

## 2.11

## 2.12

$$\int_a^b \frac{1}{b-a} dx = \frac{x}{b-a} - \frac{x}{b-a} \Big|_a^b = \frac{b-a}{b-a} = 1 \quad (177)$$

$$\mathbb{E}[x] = \int_a^b p(x) x dx \quad (178)$$

$$= \int_a^b \frac{x}{b-a} dx \quad (179)$$

$$= \frac{x^2}{2(b-a)} \Big|_a^b \quad (180)$$

$$= \frac{b^2 - a^2}{2(b-a)} = \frac{(b-a)(b+a)}{2(b-a)} = \frac{b+a}{2} \quad (181)$$

$$\text{Var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 \quad (182)$$

$$= \int_a^b \frac{x^2}{b-a} dx - \frac{(b+a)^2}{4} \quad (183)$$

$$= \frac{\frac{1}{3}x^3}{b-a} \Big|_a^b - \frac{(b+a)^2}{4} \quad (184)$$

$$= \frac{\frac{1}{3}b^3}{b-a} - \frac{\frac{1}{3}a^3}{b-a} - \frac{(b+a)^2}{4} \quad (185)$$

$$= \frac{(b-a)(b^2 + a^2 + ab)}{3(b-a)} - \frac{(b+a)^2}{4} \quad (186)$$

$$= \frac{b^2 + a^2 + ab}{3} - \frac{b^2 + a^2 + 2ab}{4} \quad (187)$$

$$= \frac{4b^2 + 4a^2 + 4ab}{12} - \frac{3b^2 + 3a^2 + 6ab}{12} \quad (188)$$

$$= \frac{b^2 + a^2 - 2ab}{12} \quad (189)$$

$$= \frac{(b-a)^2}{12} \quad (190)$$

$$(191)$$

## 2.17

$$A = \frac{1}{2} \underbrace{A + A^T}_{\text{Symmetric}} + \frac{1}{2} \underbrace{A - A^T}_{\text{Non-symmetric}} \quad (192)$$

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp -\frac{1}{2}(x - \mu)^T \frac{1}{2}(\Sigma^{-1} + \Sigma^{-T} + \Sigma^{-1} - \Sigma^{-T})(x - \mu) \quad (193)$$

$$= \exp -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M (x_n - \mu_n)(\Sigma_{nm}^{-1} + \Sigma_{nm}^{-T})(x_m - \mu_m) \quad (194)$$

$$- \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M (x_n - \mu_n) \underbrace{(\Sigma_{nm}^{-1} - \Sigma_{nm}^{-T})}_{=0} (x_m - \mu_m) \quad (195)$$

Which leaves us only with a symmetric covariance matrix.

## 2.21

Symmetry means that if the lower triangular half of the matrix is defined, the other half is defined too. That means that we have to count the amount of entries in a triangular matrix of size  $D$ .

$$\# \text{parameters} = D + (D - 1) + (D - 2) + \dots + D - (D - 2) + D - (D - 1) \quad (196)$$

$$= (D + 1) + (D + 1) + \dots + (D + 1) \quad (197)$$

$$= \frac{D}{2}(D + 1) \quad (198)$$

## 2.22

$$A^{-1}A = A^{-1}A^T \quad (\text{Symmetry})$$

$$A^{-1}AA^{T^{-1}} = A^{-1}A^T A^{T^{-1}} \quad (199)$$

$$A^{-1T} = A^{-1}A^T A^{T^{-1}} \quad (200)$$

$$= A^{-1} \quad (201)$$

## 2.27

$$\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{x} + \mathbf{z}] \quad (202)$$

$$= \int \int p(\mathbf{x}, \mathbf{z})(\mathbf{x} + \mathbf{z}) d\mathbf{x} d\mathbf{z} \quad (203)$$

$$= \int \int p(\mathbf{x}, \mathbf{z})\mathbf{x} + p(\mathbf{x}, \mathbf{z})\mathbf{z} d\mathbf{x} d\mathbf{z} \quad (204)$$

$$= \int \int p(\mathbf{x})p(\mathbf{z})\mathbf{x} + p(\mathbf{z})p(\mathbf{x})\mathbf{z} d\mathbf{x} d\mathbf{z} \quad (205)$$

$$= \int p(\mathbf{z}) \int p(\mathbf{x})\mathbf{x} d\mathbf{x} d\mathbf{z} + \int p(\mathbf{x}) \int p(\mathbf{z})\mathbf{z} d\mathbf{x} d\mathbf{z} \quad (206)$$

$$= \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{z}] \quad (207)$$

## 2.30

$$\mathbb{E}[z] = R^{-1} \begin{pmatrix} \Lambda\mu - A^T Lb \\ Lb \end{pmatrix} \quad (208)$$

$$= \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} A^T \\ A\Lambda^{-1} & L^{-1} + A\Lambda^{-1} A^T \end{pmatrix} \begin{pmatrix} \Lambda\mu - A^T Lb \\ Lb \end{pmatrix} \quad (209)$$

$$= \begin{pmatrix} \Lambda^{-1}(\Lambda\mu - A^T Lb) + \Lambda^{-1} A^T Lb \\ A\Lambda^{-1}(\Lambda\mu - A^T Lb) + \mathcal{L}^{-1} Lb + A\Lambda^{-1} A^T Lb \end{pmatrix} \quad (210)$$

$$= \begin{pmatrix} \mu - \Lambda^{-1} A^T Lb + \Lambda^{-1} A^T Lb \\ A\mu - A\Lambda^{-1} A^T Lb + b + A\Lambda^{-1} A^T Lb \end{pmatrix} \quad (211)$$

Which gives the desired result.

## 2.41

$$\Gamma(x) = \int_0^\infty u^{x-1} \exp(-u) du \quad (212)$$

$$\int_0^\infty \text{Gam}(\lambda|a, b) d\lambda = \int_0^\infty \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) d\lambda \quad (213)$$

$$= \frac{1}{\Gamma(a)} b^a \int_0^\infty \lambda^{a-1} \exp(-b\lambda) d\lambda \quad (214)$$

$$= \frac{1}{\Gamma(a)} b^a \int_0^\infty u^{a-1} b^{-a+1} \exp(-u) b^{-1} du \quad (\lambda = \frac{u}{b})$$

$$= \frac{\Gamma(a)}{\Gamma(a)} = 1 \quad (215)$$

## 2.46

$$a = \nu/2 \quad b = \frac{\nu}{2\lambda} \quad (216)$$

$$p(x|\mu, a, b) = \frac{b^a}{\Gamma(a)} \left( \frac{1}{2a} \right)^{1/2} \left[ b + \frac{1}{2}(x - \mu)^2 \right]^{-a-\frac{1}{2}} \Gamma(a + \frac{1}{2}) \quad (217)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left( \frac{1}{2\pi} \right)^{1/2} \left[ \frac{\nu}{2\lambda} + \frac{1}{2}(x - \mu)^2 \right]^{-\frac{\nu}{2}-\frac{1}{2}} \left( \frac{\nu}{2\lambda} \right)^{\nu/2} \quad (218)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left[ \frac{2\lambda}{\nu} \left( \frac{\nu}{2\lambda} + \frac{1}{2}(x - \mu)^2 \right) \right]^{-\frac{\nu}{2}} \left[ 2\pi \left( \frac{\nu}{2\lambda} + \frac{1}{2}(x - \mu)^2 \right) \right]^{-\frac{1}{2}} \quad (219)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left[ 1 + \frac{\lambda(x - \mu)^2}{\nu} \right]^{-\frac{\nu}{2}} \left[ \frac{\pi\nu}{\lambda} + \pi(x - \mu)^2 \right]^{-\frac{1}{2}} \quad (220)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left[ 1 + \frac{\lambda(x - \mu)^2}{\nu} \right]^{-\frac{\nu}{2}} \left[ 1 + \frac{\lambda}{\nu}(x - \mu)^2 \right]^{-\frac{1}{2}} \left( \frac{\lambda}{\pi\nu} \right)^{1/2} \quad (221)$$

$$(222)$$

## 2.47

$$\text{St}(x|\mu, \lambda, \nu) \propto \left[1 + \frac{\lambda(x - \mu)^2}{\nu}\right]^{-v/2-1/2} \quad (223)$$

$$= \exp \ln \left[1 + \frac{\lambda(x - \mu)^2}{\nu}\right]^{-v/2-1/2} \quad (224)$$

$$= \exp \left[ \frac{-v-1}{2} \ln \left(1 + \frac{\lambda(x - \mu)^2}{\nu}\right) \right] \quad (225)$$

$$= \exp \left[ \frac{-v-1}{2} \left( \frac{\lambda(x - \mu)^2}{\nu} - \mathcal{O} \left[ \frac{\lambda(x - \mu)^2}{2\nu} \right]^2 \right) \right] \quad (226)$$

$$= \exp \left[ \frac{-v-1}{2} \left( \frac{\lambda(x - \mu)^2}{\nu} - \mathcal{O} \left[ \frac{\lambda^2(x - \mu)^4}{4\nu^2} \right] \right) \right] \quad (227)$$

$$\approx \exp \left[ \frac{-1}{2} (\lambda(x - \mu)^2) \right] \quad (v \rightarrow \infty)$$

## 2.48

$$\Gamma(x) = \int \mu^{x-1} e^{-u} du \quad (228)$$

$$|(\eta\Lambda)^{-1}|^{1/2} = [\eta^{-D}|\Lambda^{-1}|]^{1/2} = \eta^{-D/2}|\Lambda|^{-1/2} \quad (229)$$

$$\Gamma(D/2 + \nu/2) = \int \eta^{D/2+\nu/2-1} \exp(-\eta) d\eta \quad (230)$$

$$z = \eta/2(\Delta^2 + \nu) \iff \eta = 2 \frac{z}{\Delta^2 + \nu} \quad (231)$$

$$d\eta = \frac{d\eta}{dz} dz = \frac{z}{\Delta^2 + \nu} dz \quad (232)$$



$$\int \mathcal{N}(x|\mu, (\eta\Lambda)^{-1}) \text{Gam}\left(\eta\left|\frac{\nu}{2}, \frac{\nu}{2}\right.\right) d\eta \quad (233)$$

$$= \int \frac{1}{(2\pi)^{D/2}} \frac{1}{|(\eta\Lambda)^{-1}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T(\eta\Lambda)(\mathbf{x} - \mu)\right\} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (234)$$

$$= \int \frac{1}{(2\pi)^{D/2}} \frac{1}{|(\eta\Lambda)^{-1}|^{1/2}} \exp\left\{-\frac{\eta}{2}\Delta^2\right\} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (235)$$

$$= \int \frac{1}{(2\pi)^{D/2}} \eta^{D/2} |\Lambda|^{1/2} \exp\left\{-\frac{\eta}{2}\Delta^2\right\} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (236)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \int \eta^{D/2} \exp\left\{-\frac{\eta}{2}\Delta^2\right\} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (237)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \int \eta^{D/2+\nu/2-1} \exp\left\{-\frac{\eta}{2}(\Delta^2 + \nu)\right\} d\eta \quad (238)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} 2(\Delta^2 + \nu)^{-D/2-\nu/2} \int z^{1/2D+\nu/2} \exp\{-z\} dz \quad (239)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} 2(\Delta^2 + \nu)^{-D/2-\nu/2} \Gamma(D/2 + \nu/2) \quad (240)$$

$$(241)$$

Which works out algebraically to the desired result.

## 2.50

$$\text{St}(\mathbf{x}|\mu, \Lambda, \nu) \propto \left[1 + \frac{1}{\nu}\Delta^2\right]^{-D/2-\nu/2} \quad (242)$$

$$= \exp \frac{-D-\nu}{2} \ln \left[1 + \frac{1}{\nu}\Delta^2\right] \quad (243)$$

$$= \exp \frac{-D-\nu}{2} \left[\frac{1}{\nu}\Delta^2 - \mathcal{O}(1/\nu^2)\right] \quad (\text{Taylor})$$

$$= \exp \frac{-D-\nu}{2} \frac{1}{\nu}\Delta^2 - \frac{-D-\nu}{2} \mathcal{O}(1/\nu^2) \quad (244)$$

$$= \exp -1/2\Delta^2 \quad (\nu \rightarrow \infty)$$

## 2.51

$$1 = [\cos(A) + i \sin(A)][\cos(A) - i \sin(A)] \quad (245)$$

$$= \cos^2(A) - \sin^2(A) \quad (246)$$

$$\cos(A - B) = \Re \exp[i(A - B)] \quad (247)$$

$$= \Re \exp(iA) \exp(-iB) \quad (248)$$

$$= \Re[\cos(A) + i \sin(A)][\cos(B) - i \sin(B)] \quad (249)$$

$$= \cos(A) \cos(B) - \sin(A) \sin(B) \quad (250)$$

The final question is exactly the same but considering the  $\Im$  part.

## 2.54

$$0 = \sum_{n=1}^N \cos(\theta_0) \sin(\theta_n) - \cos(\theta_n) \sin(\theta_0) \quad (251)$$

$$= \cos(\theta_0) \sum_{n=1}^N \sin(\theta_n) - \sin(\theta_0) \sum_{n=1}^N \cos(\theta_n) \quad (252)$$

$$(253)$$

$$\frac{\sin(\theta_0)}{\cos(\theta_0)} = \frac{\sum_{n=1}^N \sin(\theta_n)}{\sum_{n=1}^N \cos(\theta_n)} \quad (254)$$

$$\tan(\theta_0) = \frac{\sum_{n=1}^N \sin(\theta_n)}{\sum_{n=1}^N \cos(\theta_n)} \quad (255)$$

$$\theta_0 = \arctan \left\{ \frac{\sum_{n=1}^N \sin(\theta_n)}{\sum_{n=1}^N \cos(\theta_n)} \right\} \quad (256)$$

## 2.54

$$\frac{\partial p}{\partial \theta} = -(2\pi I_0(m))^{-1} \exp\{m \cos(\theta - \theta_0)\} (m \sin(\theta - \theta_0)) \quad (257)$$

Setting to 0 and solving gives  $\sin(\theta - \theta_0) = 0$  which resolves to  $\theta^* = \theta_0 + n\pi$   $n \in \mathbb{Z}$  where  $\mathbb{Z}$  are the positive integers.

$$\frac{\partial \partial p}{\partial \theta} = \frac{1}{2\pi I_0(m)} [-\exp(m \cos(\theta - \theta_0))(m \sin(\theta - \theta_0))^2 - m \cos(\theta - \theta_0) \exp(m \cos(\theta - \theta_0))] \quad (258)$$

Left term vanishes since  $\sin(\theta^*) = 0$ . Right term is positive (so maximal) for  $\theta^* = \theta_0 + 0\pi \pmod{2\pi}$ , negative (minimal) for  $\theta^* = \theta_0 + \pi \pmod{2\pi}$ .

## 2.55

$$A(m_{ML}) = \left( \frac{1}{N} \sum_{n=1}^N \cos \theta_n \right) \cos \theta_0^{ML} - \left( \frac{1}{N} \sum_{n=1}^N \sin \theta_n \right) \sin \theta_0^{ML} \quad (259)$$

$$= \bar{r} (\cos \bar{\theta} \cos \theta_0^{ML} - \sin \bar{\theta} \sin \theta_0^{ML}) \quad (260)$$

$$= \bar{r} (\cos(\bar{\theta} - \theta_0^{ML})) \quad (261)$$

$$= \bar{r} \quad (262)$$

## 2.57

For this question, the following knowledge is necessary:

$$a^T B a = \underbrace{B : aa^T}_{\text{Frobenius product (Hadamard \& sum)}} = \text{vec}(B)^T \text{vec}(aa^T) \quad (263)$$

where  $\text{vec}(\cdot)$  is the vectorization operation.

$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (264)$$

$$= (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp \left[ -\frac{1}{2} (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right] \quad (265)$$

$$= (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp(-1/2 \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \exp \left[ \begin{pmatrix} -1/2 \text{vec}(\boldsymbol{\Sigma}^{-1}) \\ -\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \end{pmatrix} \begin{pmatrix} \text{vec}(\mathbf{x} \mathbf{x}^T) & \mathbf{x} \end{pmatrix} \right] \quad (266)$$

$$= (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp(-1/2 \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \exp [\eta^T u(\mathbf{x})] \quad (267)$$

$$(268)$$

Now we only need to rewrite the remaining factors in terms of  $\eta$ . which gives:

$$g(\eta) = |\text{vec}^{-1}(-2\eta_1)|^{-1/2} \eta_2 \text{vec}^{-1}(-2\eta_1) \eta_2^T \quad (269)$$

and

$$h(x) = (2\pi)^{-D/2} \quad (270)$$

## 2.59

$$\int_0^\infty \frac{1}{x} f\left(\frac{1}{x}\right) dx = \int_0^\infty \frac{1}{x} f(y) x dy = \int_0^\infty f(y) dy = 1 \quad (271)$$

## 2.61

### 3.1

$$2\sigma(2a) - 1 = \frac{2}{1 + e^{-2a}} - 1 \quad (272)$$

$$= \frac{e^a}{e^a} \frac{2}{1 + e^{-2a}} - 1 \quad (273)$$

$$= \frac{2e^a}{e^a + e^{-a}} - 1 \quad (274)$$

$$= \frac{2e^a}{e^a + e^{-a}} - \frac{e^a + e^{-a}}{e^a + e^{-a}} \quad (275)$$

$$= \tanh(a) \quad (276)$$

$$a = \frac{x - \mu_j}{s} \quad (277)$$

$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M w_j \sigma(a) \quad (278)$$

$$= w_0 + \sum_{j=1}^M \frac{w_j}{2} \tanh\left(\frac{1}{2}a\right) + w_j \quad (279)$$

$$= w_0 + \sum_{j=1}^M + \sum_{j=1}^M \frac{w_j}{2} \tanh(a') \quad (280)$$

$$= u_0 + \sum_{j=1}^M u_j \tanh(a') \quad (281)$$

### 3.3

$$\frac{dE_D(\mathbf{w})}{d\mathbf{w}} = \sum_{n=1}^N r_n [t_n - \mathbf{w}^T \phi(\mathbf{x}_n)] \phi(\mathbf{x}_n) = 0 \quad (282)$$

$$\sum_{n=1}^N r_n t_n \phi(\mathbf{x}_n) = \sum_{n=1}^N \mathbf{w}^T \phi(\mathbf{x}_n) \phi(\mathbf{x}_n) \quad (283)$$

$$= \sum_{n=1}^N \phi(\mathbf{x}_n)^T \mathbf{w} \phi(\mathbf{x}_n) \quad (284)$$

$$= \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \mathbf{w} \quad (285)$$

$$= \Phi^T \Phi \mathbf{w} \quad (286)$$

$$([\mathbf{r} \cdot \mathbf{t}]^T \Phi)^T = \Phi^T \Phi \mathbf{w} \quad (287)$$

$$\Phi^T [\mathbf{r} \cdot \mathbf{t}] = \Phi^T \Phi \mathbf{w} \quad (288)$$

$$(\Phi^T \Phi)^{-1} \Phi^T [\mathbf{r} \cdot \mathbf{t}] = \mathbf{w} \quad (289)$$

(i) ?

(ii)  $r_n > 0$  essentially replicates data-points that otherwise would have been summed.

### 3.4

$$\mathbb{E}[E_d(\mathbf{w})] = \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2\right] \quad (290)$$

$$= \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (w_0 + \sum_{i=1}^D w_i (x_i + \epsilon_i) - t_n)^2\right] \quad (291)$$

$$= \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) + \sum_{i=1}^D w_i \epsilon_i - t_n)^2\right] \quad (292)$$

$$= \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2 + \underbrace{\left(\sum_{i=1}^D w_i \epsilon_i\right)^2}_{=0(\mathbb{E}[\epsilon_i]=0)} + 2(y(\mathbf{x}_n, \mathbf{w}) - t_n) \sum_{i=1}^D w_i \epsilon_i\right] \quad (293)$$

Since  $\mathbb{E}[w_i w_j] = \mathbb{E}[w_i] \mathbb{E}[w_j] = 0$  for all  $i \neq j$ :

$$\mathbb{E}\left(\sum_{i=1}^D w_i \epsilon_i\right)^2 = \mathbb{E}\left[\sum_{i=1}^D \sum_{j=1}^D w_j \epsilon_j w_i \epsilon_i\right] = \sum_{i=1}^D w_i^2 \mathbb{E} \epsilon_i^2 = \sum_{i=1}^D w_i^2 \mathbb{E} \epsilon_i^2 = \sum_{i=1}^D w_i^2 (\sigma^2 + \underbrace{E[\epsilon_i^2]}_{=0}) \quad (294)$$

Which gives us our desired result.

### 3.5

$$\mathcal{L}(\mathbf{w}, \lambda) = \frac{1}{2} \sum_{n=1}^N [t_n - \mathbf{w}^T \phi(\mathbf{x}_n)]^2 + \lambda \left[ \sum_{j=1}^M |w_j|^q - \eta \right] \quad (295)$$

Which has the same dependence on  $\mathbf{w}$  up to a scaling factor.

### 3.6

$$p(T|\mathbf{W}, \Sigma) = \prod_{n=1}^N p(t_n|\mathbf{W}, \Sigma) \quad (296)$$

$$\ln p(\mathbf{T}|\mathbf{W}, \mathbf{\Sigma}) = \ln \prod_{n=1}^N p(\mathbf{t}_n|\mathbf{W}, \mathbf{\Sigma}) = \sum_{n=1}^N \ln p(\mathbf{t}_n|\mathbf{W}, \mathbf{\Sigma}) \quad (297)$$

$$= \sum_{n=1}^N \ln \left( (2\pi)^{-D/2} |\mathbf{\Sigma}|^{-1/2} \right) + \sum_{n=1}^N \frac{1}{2} (\mathbf{y}(\mathbf{x}, \mathbf{W}) - \mathbf{t})^T \mathbf{\Sigma}^{-1} (\mathbf{y}(\mathbf{x}, \mathbf{W}) - \mathbf{t}) \quad (298)$$

$$\frac{\partial p}{\partial \mathbf{W}} = \frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial \mathbf{W}} [(\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t})^T \mathbf{\Sigma}^{-1} (\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t})] \quad (299)$$

$$= \frac{1}{2} \sum_{n=1}^N (\mathbf{\Sigma}^{-1} + \mathbf{\Sigma}^{-1T}) (\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t}) \phi(\mathbf{x})^T = 0 \quad (300)$$

$$= \sum_{n=1}^N \mathbf{\Sigma}^{-1} (\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t}) \phi(\mathbf{x})^T \quad (301)$$

$$\mathbf{\Sigma}^{-1} \sum_{n=1}^N \mathbf{W}^T \phi(\mathbf{x}) \phi(\mathbf{x})^T = \mathbf{\Sigma}^{-1} \sum_{n=1}^N \mathbf{t} \phi(\mathbf{x})^T \quad (302)$$

$$\mathbf{W}^T \mathbf{\Phi}^T \mathbf{\Phi} = \mathbf{T} \mathbf{\Phi}^T \quad (303)$$

$$\mathbf{W}^T = \mathbf{T} \mathbf{\Phi}^T (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \quad (304)$$

$$\mathbf{W} = (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \mathbf{\Phi}^T \mathbf{T} \quad (305)$$

$$\frac{d \ln p(\mathbf{T}|\mathbf{x}, \mathbf{W}, \mathbf{\Sigma})}{d \mathbf{\Sigma}} = -\frac{N}{2} \frac{d}{d \mathbf{\Sigma}} \ln |\mathbf{\Sigma}| - \frac{1}{2} \sum_{n=1}^N \frac{d}{d \mathbf{\Sigma}} (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \mathbf{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) \quad (306)$$

$$= \frac{N}{2} \frac{d}{d \mathbf{\Sigma}} \ln |\mathbf{\Sigma}|^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (307)$$

$$= \frac{N}{2} \frac{d}{d \mathbf{\Sigma}} \ln |\mathbf{\Sigma}^{-1}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (308)$$

$$= \frac{N}{2} \mathbf{\Sigma} - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (309)$$

$$(310)$$

$$\mathbf{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (311)$$

### 3.7

Hint 1: use  $x^T A x + x^T b + c = (x - h)^T A (x - h) + k$  where  $h = -(1/2)A^{-1}b$  and  $k = c - \frac{1}{4}b^T A^{-1}b$  if  $A$  is symmetric.

$$p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta^{-1}) \quad (312)$$

$$= \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{S}_0) \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (313)$$

$$= (2\pi)^{-D/2} |\mathbf{S}_0|^{-1/2} \exp[-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0)] \prod_{n=1}^N \left(\frac{\beta}{2\pi}\right)^{\frac{1}{2}} \exp[-\frac{\beta}{2}(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2] \quad (314)$$

$$= (2\pi)^{-D/2} |\mathbf{S}_0|^{-1/2} \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \exp[-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N -\frac{\beta}{2}(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2] \quad (315)$$

$$(316)$$

Now we have to get the exponent right.

$$-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N -\frac{\beta}{2}(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 \quad (317)$$

$$= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) - \frac{\beta}{2} \sum_{n=1}^N (t_n^2 + \mathbf{w}^T \phi(\mathbf{x}_n) \mathbf{w}^T \phi(\mathbf{x}_n) - 2\mathbf{w}^T \phi(\mathbf{x}_n) t_n) \quad (318)$$

$$= -\frac{1}{2}(\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{w} - \mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{m}_0 - \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{w} + \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0) - \frac{\beta}{2}(\mathbf{t}^T \mathbf{t} + \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{t}) \quad (319)$$

$$= -\frac{1}{2}(\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{w} - 2\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{m}_0 + \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0 - \beta \mathbf{t}^T \mathbf{t} + \beta \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\beta \mathbf{w}^T \Phi^T \mathbf{t}) \quad (320)$$

$$= \mathbf{w}^T (-\frac{1}{2}(\mathbf{S}_0^{-1} + \beta \Phi^T \Phi)) \mathbf{w} + \mathbf{w}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}) + c \quad (321)$$

$$= \mathbf{w}^T (-\frac{1}{2} \mathbf{S}_N^{-1}) \mathbf{w} + \mathbf{w}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}) + c \quad (322)$$

Now completing the square.

$$\mathbf{w}^T (-\frac{1}{2} \mathbf{S}_N^{-1}) \mathbf{w} + \mathbf{w}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}) + c \quad (323)$$

$$= (\mathbf{w} + \frac{1}{2}(-\frac{1}{2} \mathbf{S}_N^{-1})^{-1}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}))^T (-\frac{1}{2} \mathbf{S}_N^{-1})(\mathbf{w} + \frac{1}{2}(-\frac{1}{2} \mathbf{S}_N^{-1})^{-1}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t})) + k \quad (324)$$

$$= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^T \mathbf{S}_N^{-1}(\mathbf{w} - \mathbf{m}_N) + k \quad (325)$$

### 3.15

Hint 1: Use 3.95 and 3.92

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N \quad (326)$$

$$= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\gamma(\mathbf{m}_N^T \mathbf{m}_N)^{-1}}{2} \mathbf{m}_N^T \mathbf{m}_N \quad (327)$$

$$= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\gamma}{2} \quad (328)$$

$$(329)$$

$$\beta = (N - \gamma) \|\mathbf{t} - \Phi^T \mathbf{m}_n\|^{-2} \quad (330)$$

$$E(\mathbf{m}_N) = \frac{1}{2}(\mathcal{N} - \gamma + \gamma) = \frac{N}{2} \quad (331)$$

### 3.17

Evidence function:

$$p(D|M_i) = \int p(D|\mathbf{w}, M_i) p(\mathbf{w}, M_i) d\mathbf{w} \quad (332)$$

For Linear regression we have  $M_i = (\alpha, \beta)$

$$p(D|\alpha, \beta) = \int p(\mathbf{w}|\alpha^{-1}I) \prod_{n=1}^N p(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (333)$$

$$= \int (2\pi)^{-M/2} |\alpha^{-1}I|^{-1/2} \exp[-1/2(\mathbf{w}^T(\alpha^{-1}I)^{-1}\mathbf{w})] \prod_{n=1}^N \frac{1}{\sqrt{2\pi\beta^{-1}}} \exp[\beta/2(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2] \quad (334)$$

$$= \int (2\pi)^{-M/2} (\alpha^{-1})^{-M/2} \frac{\beta^{N/2}}{2\pi} \exp[-\alpha/2\mathbf{w}^T \mathbf{w} - \beta/2\|\mathbf{t} - \Phi \mathbf{w}\|^2] \quad (335)$$

### 4.4

$$\mathcal{L} = \mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1) + \lambda(\mathbf{w}^T \mathbf{w} - 1) \quad (336)$$

$$\nabla_{\mathbf{w}} \mathcal{L} = (\mathbf{m}_2 - \mathbf{m}_1) + 2\lambda \mathbf{w} = 0 \quad (337)$$

$$\mathbf{w} = -\frac{1}{2\lambda}(\mathbf{m}_2 - \mathbf{m}_1) \propto (\mathbf{m}_2 - \mathbf{m}_1) \quad (338)$$



## 4.5

It's quite straightforward if we fill in the given equations. Numerator:

$$(m_2 - m_1)^2 = (\mathbf{w}^T \mathbf{m}_2 - \mathbf{w}^T \mathbf{m}_1)^2 \quad (339)$$

$$= \mathbf{w}^T \mathbf{m}_2 \mathbf{w}^T \mathbf{m}_2 - \mathbf{w}^T \mathbf{m}_2 \mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{m}_2 + \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{m}_1 \quad (340)$$

$$= \mathbf{w}^T (\mathbf{m}_2 \mathbf{m}_2^T - \mathbf{m}_2 \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{m}_2^T + \mathbf{m}_1 \mathbf{m}_1^T) \mathbf{w} \quad (341)$$

$$= \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w} \quad (342)$$

The denominator uses exactly the same approach. I'll show for  $s_1^2$ .

$$s_1^2 = \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m}_1)^2 \quad (343)$$

$$= \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n \mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x}_n \mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{x}_n + \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{m}_1) \quad (344)$$

$$= \sum_{n \in C_1} \mathbf{w}^T (\mathbf{x}_n \mathbf{x}_n^T - \mathbf{x}_n \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{x}_n^T + \mathbf{m}_1 \mathbf{m}_1^T) \mathbf{w} \quad (345)$$

$$= \sum_{n \in C_1} \mathbf{w}^T (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T \mathbf{w} \quad (346)$$

## 4.6

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0 \quad (347)$$

$$= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m} - t_n) \mathbf{x}_n \quad (348)$$

$$\mathbf{w}^T \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m}) \mathbf{x}_n = \sum_{n=1}^N t_n \mathbf{x}_n \quad (349)$$

This gives us the three terms that we have to solve for.

$$\sum_{n=1}^N t_n \mathbf{x}_n = \sum_{n \in C_1}^{N_1} t_n \mathbf{x}_n - \sum_{n \in C_2}^{N_2} t_n \mathbf{x}_n \quad (350)$$

$$= \sum_{n \in C_1}^{N_1} \frac{N}{N_1} \mathbf{x}_n - \sum_{n \in C_2}^{N_2} \frac{N}{N_2} \mathbf{x}_n \quad (351)$$

$$= N \left( \sum_{n \in C_1}^{N_1} \frac{1}{N_1} \mathbf{x}_n - \sum_{n \in C_2}^{N_2} \frac{1}{N_2} \mathbf{x}_n \right) \quad (352)$$

$$= N(\mathbf{m}_1 - \mathbf{m}_2) \quad (353)$$

$$-\mathbf{w}^T \mathbf{m} \sum_{n=1}^N \mathbf{x}_n = -\frac{1}{N} \mathbf{w}^T (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) \quad (354)$$

$$= -\frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)^T \mathbf{w} \quad (355)$$

$$= -\frac{1}{N} (N_1^2 \mathbf{m}_1 \mathbf{m}_1^T + N_1 N_2 \mathbf{m}_1 \mathbf{m}_2^T + N_2 N_1 \mathbf{m}_2 \mathbf{m}_1^T + N_2^2 \mathbf{m}_2 \mathbf{m}_2^T) \mathbf{w} \quad (356)$$

$$= -\frac{1}{N} ((N - N_2) N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_1 N_2 \mathbf{m}_1 \mathbf{m}_2^T + N_2 N_1 \mathbf{m}_2 \mathbf{m}_1^T + (N - N_1) N_2 \mathbf{m}_2 \mathbf{m}_2^T) \mathbf{w} \quad (357)$$

$$= \left( -N_1 \mathbf{m}_1 \mathbf{m}_1^T + \frac{N_1 N_2}{N} \mathbf{m}_1 \mathbf{m}_1^T - \frac{N_1 N_2}{N} \mathbf{m}_1 \mathbf{m}_2^T - \frac{N_2 N_1}{N} \mathbf{m}_2 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \frac{N_1 N_2}{N} \mathbf{m}_2 \mathbf{m}_2^T \right) \mathbf{w} \quad (358)$$

$$= \left( -N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \frac{N_1 N_2}{N} (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \right) \mathbf{w} \quad (359)$$

$$= \left( -N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} \quad (360)$$

$$(361)$$

We add the remaining terms to the final term:

$$\left( -N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} \quad (362)$$

$$= \left( N_1 \mathbf{m}_1 \mathbf{m}_1^T - 2N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T - 2N_2 \mathbf{m}_2 \mathbf{m}_2^T + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} \quad (363)$$

$$= \left( N_1 \mathbf{m}_1 \mathbf{m}_1^T - \sum_{n \in C_1} \mathbf{x}_n \mathbf{m}_1^T - \mathbf{m}_1 \sum_{n \in C_1} \mathbf{x}_n^T + \dots + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} \quad (364)$$

$$= \left( \sum_{n \in C_1} \mathbf{m}_1 \mathbf{m}_1^T - \mathbf{x}_n \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{x}_n^T + \mathbf{x}_n \mathbf{x}_n^T + \dots \right) \mathbf{w} \quad (365)$$

Which gives the product we need. "..." denotes symmetric steps but for the class 2.

## 4.7

$$\sigma(-a) = \frac{1}{e^x + 1} = \frac{e^{-x}}{1 + e^{-x}} = \frac{e^{-x} + 1}{e^{-x} + 1} - \frac{1}{e^{-x} + 1} \quad (366)$$

$$y = \frac{1}{1 + e^{-x}} \quad (367)$$

$$e^{-x} = \frac{1 - y}{y} \quad (368)$$

$$y = e^x (1 - y) \quad (369)$$

$$e^x = y/(1-y) \quad (370)$$

$$x = \ln[y/(1-y)] \quad (371)$$

## 4.8

$$p(C_1|\mathbf{x}) = \sigma(a) \quad (372)$$

So we have to show:  $a = \mathbf{w}^T \mathbf{x} + w_0$

$$a = \ln \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} + \ln \frac{p(C_1)}{p(C_2)} \quad (373)$$

$$\ln \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} = \ln \left[ \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left( -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) \right) \right] \quad (374)$$

$$- \ln \left[ \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left( -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \right) \right] \quad (375)$$

$$= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \quad (376)$$

$$= \frac{1}{2} [2\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 - 2\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_2 + \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2] \quad (377)$$

$$(378)$$

From which the result easily follows.

## 4.9

Hint 1: Using a Lagrange multiplier, make sure  $\sum_{j=1}^k \pi_j = 1$  before optimizing.

$$\ln p(\mathbf{X}|\mathbf{T}) = \sum_{n=1}^N \ln \prod_{j=1}^K (\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma))^{t_j} \quad (379)$$

$$= \sum_{n=1}^N \sum_{j=1}^K t_j \ln(\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma)) \quad (380)$$

Now we optimize with constraint  $\sum_{j=1}^K \pi_j = 1$ .

$$\mathcal{L}(\boldsymbol{\pi}, \lambda) = \ln(p(\mathbf{X}, \mathbf{T})) - \lambda \left( \sum_{j=1}^K \pi_j - 1 \right) \quad (381)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_j} = \sum_{n=1}^N \frac{t_j}{\pi_j} - \lambda = 0 \quad (382)$$

$$\lambda = \sum_{n=1}^N \frac{t_j}{\pi_j} = N \frac{t_j}{\pi_j} = \frac{N_j}{\pi_j} \quad (383)$$

$$\pi_j = \frac{N_j}{\lambda} \quad (384)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sum_{j=1}^K \pi_j = 1 \quad (385)$$

Plugging this into 380.

$$\sum_{j=1}^K \frac{N_j}{\lambda} = 1 \iff \lambda = N \quad (386)$$

Gives us the desired result.

## 4.12

$$\frac{d\sigma}{da} = (1 + e^{-a})^{-2} e^{-a} \quad (387)$$

$$= \frac{1}{1 + e^{-a}} \frac{1}{1 + e^{-a}} e^{-a} \quad (388)$$

$$= \sigma(a) \frac{e^{-a}}{1 + e^{-a}} \quad (389)$$

$$= \sigma(a) \left[ \frac{1 + e^{-a}}{1 + e^{-a}} - \frac{1}{1 + e^{-a}} \right] \quad (390)$$

## 4.13

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = - \sum_{n=1}^N \frac{t_n}{y_n} \nabla_{\mathbf{w}} \sigma(\mathbf{w}^T \phi_n) - \frac{1 - t_n}{1 - y_n} \nabla_{\mathbf{w}} \sigma(\mathbf{w}^T \phi_n) \quad (391)$$

$$= - \sum_{n=1}^N \frac{t_n}{y_n} \sigma(\mathbf{w}^T \phi_n) (1 - \sigma(\mathbf{w}^T \phi_n)) \phi_n - \frac{1 - t_n}{1 - y_n} \sigma(\mathbf{w}^T \phi_n) (1 - \sigma(\mathbf{w}^T \phi_n)) \phi_n \quad (392)$$

$$= - \sum_{n=1}^N t_n (1 - y_n) \phi_n - (1 - t_n) y_n \phi_n \quad (393)$$

$$= - \sum_{n=1}^N (t_n - y_n) \phi_n \quad (394)$$

## 4.14

Hint 1: approach it with an argument, using that we have a perfect decision boundary at  $\mathbf{w}^T \phi = 0$ .

We know that if  $C_1$  is labelled with  $t_{C_1} = 1$  and  $C_2$  is labelled with  $t_{C_2} = 0$  then we want  $p(C_1|\phi) = \sigma(\mathbf{w}^T \phi) > 0.5$  and  $p(C_2|\phi) = \sigma(\mathbf{w}^T \phi) < 0.5$  which happens if the decision boundary perfectly separates them at  $\mathbf{w}^T \phi = 0$ . Now the binary cross entropy will be minimal as  $p(C_1|\phi) \rightarrow 1$  which happens when  $\mathbf{w} \rightarrow \infty$ . And vice versa.

## 4.16

$$p(\mathbf{t}, \mathbf{w}) = \prod_{n=1}^N y_n^{\pi_n} [1 - y_n]^{1-t_n} \quad (395)$$

$$\ln p = \sum_{n=1}^N \pi_n \ln y_n + (1 - \pi_n) \ln(1 - y_n) \quad (396)$$

## 4.17

$$p(C_k | \phi) = y_k = \frac{\exp a_k}{\sum_{j=1} \exp a_j} \quad (397)$$

$$\frac{\partial y_k}{\partial a_j} = -\exp a_k \left( \sum_j \exp(a_j) \right)^{-2} \exp(a_j) \quad (398)$$

$$= \begin{cases} y_k(0 - y_j) & j \neq k \\ y_k(1 - y_j) & j = k \end{cases} \quad (399)$$

$$= y_k(I_{kj} - y_j) \quad (400)$$

## 4.18

$$\nabla_{\mathbf{w}_j} E(\mathbf{W}) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \nabla_{\mathbf{w}_j} \ln y_{nk} \quad (401)$$

$$\nabla_{\mathbf{w}_j} \ln(y_{nk}) = -(I_{kj} - y_j) \phi_n \quad (402)$$

$$\nabla_{\mathbf{w}_j} E(\mathbf{W}) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} (I_{kj} - y_j) \phi_n \quad (403)$$

$$= \phi \sum_{n=1}^N \sum_{k=1}^K t_{nk} y_{jn} \phi - t_{nk} I_{kj} \quad (404)$$

$$= \phi \sum_{n=1}^N t_{nj} - y_{jn} \phi \underbrace{\sum_{k=1}^K t_{nk}}_{=1} \quad (405)$$

$$= \sum_{n=1}^N \phi (y_{jn} - t_{nj}) \quad (406)$$

## 4.19

Hint 1: Use binary cross netropy and 4.114 as the activation function. Use the fundamental theorem of calculus.

$$p(\mathbf{t}, \mathbf{w}) = \sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)] \quad (407)$$

$$\nabla_{\mathbf{w}} \sum_{n=1}^N [t_n \ln \Phi(a) + (1 - t_n) \ln(1 - \Phi(a))] = \sum_{n=1}^N \left( \frac{t_n}{\Phi(a)} - \frac{1 - t_n}{1 - \Phi(a)} \right) \Phi(a) \phi_n \quad (408)$$

## 4.21

$$\Phi(a) = \int_0^a \mathcal{N}(0, 1) d\theta \quad (409)$$

$$= \int_0^a 1/\sqrt{2\pi} \exp(-\frac{1}{2}\theta^2) d\theta \quad (410)$$

$$= \frac{1}{2} + \int_{-\infty}^a \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}\theta^2) d\theta \quad (411)$$

$$= \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \frac{\sqrt{\pi}}{2} \int_{-\infty}^a \frac{2}{\sqrt{\pi}} \exp(-\frac{1}{2}\theta^2) d\theta \quad (412)$$

$$= \frac{1}{2} \left( 1 + \frac{1}{\sqrt{2}} \text{erf}(a) \right) \quad (413)$$

## 4.22

$$\ln p(D) = \ln \left[ f(z_0) \frac{(2\pi)^{M/2}}{|A|^{1/2}} \right] \quad (414)$$

$$= \ln f(z_0) = \frac{M}{2} \ln 2\pi - \frac{1}{2} \ln |A| \quad (415)$$

$z_0$  is the location of the  $\boldsymbol{\theta}_{MAP}$  estimate so

$$\ln f(\boldsymbol{\theta}) \Big|_{z_0} = \ln f(\boldsymbol{\theta}_{MAP}) = \ln p(D|\boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) \quad (416)$$

## 5.2

$$p(\mathbf{T}, \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \mathcal{N}(\mathbf{t}_n | y(\mathbf{x}_n, \mathbf{w}), \beta^{-1} I) \quad (417)$$

$$= \prod_{n=1}^N \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\beta^{-1} I|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))^T \beta^{-1} I (\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))\right) \quad (418)$$

Now it's obvious that if we take the the log likelihood it cancels the exp and we end up with

$$\left(-\frac{1}{2}(\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))^T \beta^{-1} I(\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))\right) = -\frac{1}{2\beta} \|\mathbf{y}(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2. \quad (419)$$

## 5.5

$$p(\mathbf{T}, \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \prod_{k=1}^K y_k(\mathbf{x}_n, \mathbf{w})^{t_n^k} (1 - y_k)(\mathbf{x}_n, \mathbf{w})^{1-t_n^k} \quad (420)$$

Taking the log becomes the cross-entropy function.

## 5.6

$$\frac{\partial E(\mathbf{w})}{\partial a_k} = \sum_{n=1}^N \left( \frac{\partial}{\partial a_k} t_n \ln y_n + \frac{\partial}{\partial a_k} (1 - t_n) \ln(1 - y_n) \right) \quad (421)$$

$$= -\frac{t_k}{y_k} \frac{\partial}{\partial a_k} y_k - \frac{1 - t_k}{1 - y_k} \frac{\partial}{\partial a_k} (1 - y_k) \quad (422)$$

$$= -t_k(1 - y_k) + (1 - t_k)y_k \quad (423)$$

$$= y_k - t_k \quad (424)$$

## 5.7

$$E(\mathbf{w}) = \sum_{k=1}^K t_k \ln y_k(\mathbf{x}, \mathbf{w}) \quad (425)$$

$$\frac{\partial E}{\partial a_j} = -\sum_{k=1}^K \frac{t_k}{y_k} y_k (I_{kj} - y_j) \quad (426)$$

$$= -\sum_{k=1}^K t_k (I_{kj} - y_j) \quad (427)$$

$$= -t_j + y_j \quad (428)$$

$$\frac{d \tanh}{da} = (e^a - e^{-a}) \frac{d}{da} (e^a + e^{-a})^{-1} + (e^a + e^{-a})^{-1} \frac{d}{da} (e^a - e^{-a}) \quad (429)$$

$$= -(e^a - e^{-a})(e^a + e^{-a})^{-2} \frac{d}{da} (e^a + e^{-a}) + (e^a + e^{-a})^{-1} (e^a - e^{-a}) \quad (430)$$

$$= -(e^a - e^{-a})(e^a + e^{-a})^{-2} (e^a - e^{-a}) + 1 \quad (431)$$

$$= -h^2(a) + 1 \quad (432)$$

## 5.9

Still Bernoulli, so

$$p(t|\mathbf{x}, \mathbf{w}) = \left(\frac{1+y}{2}\right)^{\frac{1+t}{2}} \left(\frac{1-y}{2}\right)^{\frac{1-t}{2}} \quad (433)$$

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}) \quad (434)$$

$$-\ln(p(\mathbf{t}|\mathbf{X}, \mathbf{w})) = -\sum_{n=1}^N \ln p(t_n|\mathbf{x}_n, \mathbf{w}) \quad (435)$$

$$= -\sum_{n=1}^N \left( \left(\frac{1+t}{2}\right) \ln\left(\frac{1+y}{2}\right) + \left(\frac{1-t}{2}\right) \ln\left(\frac{1-y}{2}\right) \right) \quad (436)$$

We can use tanh as activation function.

## 5.10

$$u_i^T H u_i = u_i^T \lambda u_i = \delta_i i \lambda_i = \lambda_i \quad (437)$$

We started with 5.37, so this will always be positive.

The converse direction is the one in the book.

## 5.13

Hint 1:  $\mathbf{b}$  has  $W$  parameters and  $H$  is  $W \times W$ . We already know  $H$  has  $\frac{N(N+1)}{2}$  parameters and  $b$  has  $N$  parameters.

$$\frac{N(N+1)}{2} + N = \frac{N(N+3)}{2} \quad (438)$$

## 5.14

Hint 1: Taylor expansion on both terms in numerator.

Taylor:

$$E_n(w_j + \epsilon) = E_n(w_{ji}) + \epsilon \frac{\partial E_n}{\partial w_{ji}} + \frac{\epsilon^2}{2} \frac{\partial^2 E_n}{\partial w_{ji}^2} + O(\epsilon^3) \quad (439)$$

$$E_n(w_j - \epsilon) = E_n(w_{ji}) - \epsilon \frac{\partial E_n}{\partial w_{ji}} + \frac{\epsilon^2}{2} \frac{\partial^2 E_n}{\partial w_{ji}^2} - O(\epsilon^3) \quad (440)$$

Subtracting these and solving for the partial derivative shows that the second order terms cancel.



## 5.16

$$E = \frac{1}{2} \sum_{n=1}^N \|\mathbf{y}_n - \mathbf{t}_n\|^2 = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M (\mathbf{y}_n - \mathbf{t}_n)^2 \quad (441)$$

$$\nabla E = \sum_{n=1}^N \sum_{m=1}^M (\mathbf{y}_n - \mathbf{t}_n) \nabla \mathbf{y}_n \quad (442)$$

$$\nabla \nabla E = \sum_{n=1}^N \sum_{m=1}^M \nabla \mathbf{y}_n \nabla (\mathbf{y}_n - \mathbf{t}_n) + (\mathbf{y}_n - \mathbf{t}_n) \nabla \nabla \mathbf{y}_n \approx \sum_{n=1}^N \sum_{m=1}^M \nabla \mathbf{y}_n \nabla \mathbf{y}_n^T \quad (443)$$

## 5.17

Hint 1: Use  $y(\mathbf{x}, \mathbf{w}) = \int t p(t|x) dt$

$$\frac{\partial E}{\partial w_r} = \int \int (y(x, w) - t) \frac{\partial y}{\partial w_r} p(x, t) dx dt \quad (444)$$

$$\frac{\partial^2 E}{\partial w_r \partial w_s} = \int \int \left[ (y(x, w) - t) \frac{\partial^2 y(x, w)}{\partial w_r \partial w_s} + \frac{\partial y(x, w)}{\partial w_r} \frac{\partial y}{\partial w_s} \right] p(x, t) dx dt \quad (445)$$

$$\int \int (y(x, w) - t) p(x, t) dx dt = \int \int (y(x, w) - t) p(t|x) p(x) dx dt \quad (446)$$

$$= \int p(x) \left( y(x, w) - \underbrace{\int t p(t|x)}_{=y(x, w)} \right) dx dt = 0 \quad (447)$$

The remaining integral is the answer.

## 5.18

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left( \sum_{j=1}^M w_{kj} h \left( \sum_{i=1}^D w_{ji} x_i + w_{j0} \right) + \sum_{l=1}^D w_l x_l \right) \quad (448)$$

Finding the derivatives to these skip weights is straightforward.

## 5.19

$$E = - \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \quad (449)$$

$$\nabla E = \sum_{n=1}^N (y_n - t_n) \nabla a_n \quad (\text{Taken from earlier solutions})$$

$$\nabla \nabla E = \sum_{n=1}^N y_n (1 - y_n) \nabla a_n \nabla a_n^T + (y_n - t_n) \nabla \nabla a_n \quad (450)$$

## 5.20

$$\nabla_{w_j} E(\mathbf{W}) = \sum_{n=1}^N (y_{n_j} - bt_{n_j}) \nabla a_j \quad (451)$$

$$\nabla \nabla \mathbf{w}_j \approx \sum_{n=1} y_k (I - y_j) \nabla a_j \nabla a_j^T \quad (452)$$

## 5.24

$$\sum_i \frac{1}{a} w_{ji} (ax_i + b) + w_{j0} - \frac{b}{a} \sum_i w_{ji} = \sum_i w_{ji} x + \frac{b}{a} w_{ji} + w_{j0} - \frac{b}{a} \sum_i w_{ji} \quad (453)$$

$$= \sum_i w_{ji} x_i + w_{j0} \quad (454)$$

$y_k$  scaling is similar.

## 5.28

If we normally have  $y = \sum_{j=0}^M w_{kj} z_j$  we now have  $y_k = \sum_{j=0}^M w_{kj} z_j$ . Therefore the backprop becomes  $\frac{\partial}{\partial w_k} = \sum_{j=0}^M z_j$ . I.e., the weights are updated according to the outputs that the generated for all receptive fields and summed.

## 5.29

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \quad (455)$$

$$\frac{\partial \tilde{E}}{\partial w_i} = \frac{\partial E}{\partial w_i} + \frac{\partial}{\partial w_i} \lambda \Omega(\mathbf{w}) \quad (456)$$

$$\frac{\partial}{\partial w_i} \lambda \Omega(\mathbf{w}) = -\lambda \frac{\partial}{\partial w_i} \ln \left( \sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right) \quad (457)$$

$$\frac{\partial}{\partial w_i} = \frac{1}{\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2)} \frac{\partial}{\partial w_i} \ln \left( \sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right) \quad (458)$$

$$\frac{\partial}{\partial w_i} \ln \left( \sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right) = \sum_{j=1}^M \pi_j \frac{\partial}{\partial w_i} \mathcal{N}(w_i | \mu_j, \sigma_j^2) \quad (459)$$

$$\frac{\partial}{\partial w_i} \mathcal{N}(w_i | \mu_j, \sigma_j^2) = \frac{1}{\sqrt{2\pi\sigma_j^2} \exp(-\frac{1}{2} \frac{(w_i - \mu_j)^2}{\sigma_j^2})} \frac{\partial}{\partial w_i} \left( -\frac{1}{2} \frac{(w_i - \mu_j)^2}{\sigma_j^2} \right) \quad (460)$$

$$\frac{\partial}{\partial w_i} = -\sigma_j^{-2} (w_i - \mu_j) \quad (461)$$

Plugging everything in gives

$$\frac{\partial}{\partial w_i} = \frac{\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j)}{\sum_{k=1}^K \mathcal{N}(w_i | \mu_j, \sigma_j^2)} \quad (462)$$

### 5.30

$$\frac{\partial \tilde{E}}{\partial \mu_j} = -\lambda \frac{\partial}{\partial \mu_j} \sum_i \ln \sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j) \quad (463)$$

$$= -\lambda \sum_i \frac{1}{\sum_j \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j)} \quad (464)$$

$$= -\lambda \sum_i \frac{1}{\sum_j \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j)} \sum_j \pi_j \frac{d}{d \mu_j} \mathcal{N}(w_i | \mu_j, \sigma_j) \quad (465)$$

$$\frac{d}{d \mu_j} \mathcal{N}(w_i | \mu_j, \sigma_j) = -\frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{1}{2\sigma_j^2}(w_i - \mu_j)^2\right) \left[\frac{1}{\sigma_j^2}(w_i - \mu_j)\right] \quad (466)$$

Plugging in gives the result.

### 5.31

$$\frac{\partial}{\partial \sigma_j} \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j) \quad (467)$$

$$= \pi \left[ \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \frac{\partial}{\partial \sigma_j} \frac{1}{\sqrt{2\pi\sigma_j^2}} + \frac{1}{\sqrt{2\pi\sigma_j^2}} \frac{\partial}{\partial \sigma_j} \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \right] \quad (468)$$

$$\frac{\partial}{\partial \sigma_j} (2\pi\sigma_j^2)^{-\frac{1}{2}} = -\frac{1}{2} (2\pi\sigma_j^2)^{-\frac{3}{2}} \frac{\partial}{\partial \sigma_j} (2\pi\sigma_j^2) \quad (469)$$

$$= -\frac{1}{2} (2\pi\sigma_j^2)^{-\frac{3}{2}} 4\pi\sigma_j \quad (470)$$

$$= -\frac{1}{\sigma_j} \frac{1}{\sqrt{2\pi\sigma_j^2}} \quad (471)$$

$$\frac{\partial}{\partial \sigma_j} \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) = \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \frac{\partial}{\partial \sigma_j} \left(-\frac{1}{2\sigma_j^2}(\mu_j - w_i)^2\right) \quad (472)$$

$$= \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \frac{1}{\sigma_j^3} (\mu_j - w_i)^2 = \quad (473)$$

Plugging these values in gives the result.

### 5.33

If we start with  $\mathbf{v} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  then using trigonometry:

$$\mathbf{y} = \begin{pmatrix} \cos(\pi - \theta)L_1 + v_1 \\ \sin(\pi - \theta)L_1 + v_2 \end{pmatrix} \quad (474)$$

$$\mathbf{x} = \begin{pmatrix} y_1 + \cos(\theta_1 + \theta_2 - \pi)L_2 \\ y_2 + \sin(\theta_1 + \theta_2 - \pi)L_2 \end{pmatrix} \quad (475)$$

### 5.34

$$\frac{\partial E_n}{\partial a_{kl}^\pi} = - \frac{1}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (476)$$

$$\frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) = \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_k \quad (477)$$

$$\frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_k = \sum_{l=1} \pi_k (I_{kl} - \pi_l) \quad (478)$$

$$= \pi_k - \pi_k \sum_{l=1} \pi_l \quad (479)$$

$$\frac{\partial E_n}{\partial a_{kl}^\pi} = - \frac{\pi_k \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) - \pi_k \sum_{l=1} \pi_l \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \quad (480)$$

$$= -\gamma_k(w) + \pi_k \quad (481)$$

### 5.35

$$\frac{\partial E_n}{\partial a_{kl}^\mu} = - \frac{1}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \frac{\partial}{\partial a_k^\mu} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (482)$$

$$\frac{\partial}{\partial a_k^\mu} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (483)$$

$$= \pi_k(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \frac{\partial}{\partial a_k^\mu} - \frac{1}{2} (\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1} (\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \quad (484)$$

$$(485)$$

$$\frac{\partial}{\partial a_k^\mu} - \frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \quad (486)$$

$$= -\sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \underbrace{\frac{\partial \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w})}{\partial a_k^\mu}}_{=1} \quad (487)$$

$$\frac{\partial E_n}{\partial a_{kl}^\mu} = \gamma_{nk} \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \quad (488)$$

### 5.36

$$\frac{\partial E_n}{\partial a_{kl}^\sigma} = - \frac{1}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \frac{\partial}{\partial a_k^\sigma} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (489)$$

$$\frac{\partial}{\partial a_k^\sigma} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) = \pi_k(\mathbf{x}_n, \mathbf{w}) \frac{\partial}{\partial a_k^\sigma} \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (490)$$

$$\frac{\partial}{\partial a_k^\sigma} \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) = \frac{\partial}{\partial a_k^\sigma} \frac{1}{(2\pi)^{\frac{L}{2}}} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) \quad (491)$$

$$= \frac{1}{(2\pi)^{\frac{L}{2}}} \left( \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) \frac{\partial}{\partial a_k^\sigma} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \right. \quad (492)$$

$$\left. + \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \frac{\partial}{\partial a_k^\sigma} \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) \right) \quad (493)$$

$$\frac{\partial}{\partial a_k^\sigma} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} = \frac{\partial}{\partial a_k^\sigma} \frac{1}{\sigma_k(\mathbf{x}, \mathbf{w})^L} = -\frac{L}{\sigma^2(\mathbf{x}, \mathbf{w})^{L-1}} = \frac{L}{\sigma_k} \frac{1}{\sigma_k(\mathbf{x}, \mathbf{w})^L} = \frac{L}{\sigma_k} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \quad (494)$$

$$\frac{\partial}{\partial a_k^\sigma} \left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) = \frac{\|\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}\|^2}{\sigma_k^3} \quad (495)$$

Plugging in these values gives the correct result. The book's answer has a typo.

### 5.37

### 5.38

### 5.39

Using the approximation:

$$p(D|\alpha, \beta) = \int p(D|w, \beta)p(w|\alpha)dw \approx p(D, w_{map})p(w_{map})\frac{(2\pi)^{W/2}}{|A|^{\frac{1}{2}}} \quad (496)$$

Taking the log gives the result.

### 5.40

Use a softmax activation.

### 6.3

$$\|x - x_n\| = x^T x + x_n^T x - 2x_n^T x \quad (497)$$

$$= k(x, x) + k(x_n, x) - 2k(x_n, x) \quad (498)$$

### 6.4

$$\begin{pmatrix} 2 & 0 \\ -1 & 3 \end{pmatrix} \quad (499)$$

### 6.5

$$ck(\mathbf{x}, \mathbf{x}') = c\psi^T \psi = (c^{\frac{1}{2}}\psi)^T (c^{\frac{1}{2}}\psi) = \phi^T \phi \quad (500)$$

$$f(\mathbf{x})k(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') = f(\mathbf{x})\psi^T \psi f(\mathbf{x}') = (f(\mathbf{x})\psi^T)^T (\psi f(\mathbf{x}')) = \phi^T \phi \quad (501)$$

### 6.6

$$q(k(x, x')) = \sum_{n=1}^N \beta_n k(x, x')^n \quad (502)$$

$$= \sum_{n=1}^N \beta_0 k(x, x') \quad 6.18 \quad (503)$$

$$= \sum_{n=1}^N k(x, x') \quad 6.13 \quad (504)$$

$$= k(x, x') \quad 6.17 \quad (505)$$

6.16 follows from the previous proof and the definition of the exponential function.

## 6.7

$$\psi(x)^T \psi(x') + \theta(x)^T \theta(x') = \sum_{n=1}^N (\psi(x_n) + \theta(x_n))(\psi(x'_n) + \theta(x'_n)) - \psi(x_n)\theta(x'_n) - \theta(x_n)\psi(x'_n) \quad (506)$$

$$= \begin{bmatrix} \dots & \psi(x_n) + \theta(x_n) & -\psi(x_n) & -\theta(x_n) & \dots \end{bmatrix} \begin{bmatrix} \vdots \\ \psi(x'_n) + \theta(x'_n) \\ \theta(x'_n) \\ \psi(x'_n) \\ \vdots \end{bmatrix} \quad (507)$$

$$= \phi(x)^T \phi(x') \quad (508)$$

$$k_1(x, x') k_2(x, x') = \theta(x)^T \theta(x') (\psi(x)^T \psi(x')) \quad (509)$$

$$= \sum_{n=1}^N \theta(x_n) \theta(x'_n) \psi(x_n) \psi(x'_n) \quad (510)$$

$$= \sum_{n=1}^N \theta(x_n) \psi(x_n) \theta(x'_n) \psi(x'_n) \quad (511)$$

$$= \phi(x)^T \phi(x') \quad (512)$$

$$= k(x, x') \quad (513)$$

## 6.8

$$k_3(\phi(x), \phi(x')) = \psi(\phi(x))^T \psi(\phi(x')) = k(x, x') \quad (514)$$

$$x^T A x' = x^T U^T U x' \quad (\text{Symmetry})$$

$$= (Ux)^T Ux' \quad (515)$$

$$= \phi(x)^T \phi(x') \quad (516)$$

$$= k(x, x') \quad (517)$$

## 6.9

$$\sum_{n=1}^N \psi(x_{na})\psi(x'_{na}) + \sum_{m=1}^M \theta(x_{mb})\theta(x'_{mb}) \quad (518)$$

$$= \begin{bmatrix} \psi(x_{a1}) & \dots & \psi(x_{aN}) & \theta(x_{b1}) & \dots & \theta(x_{bM}) \end{bmatrix} \begin{bmatrix} \psi(x'_{a1}) \\ \vdots \\ \psi(x'_{aN}) \\ \theta(x'_{b1}) \\ \vdots \\ \theta(x'_{bM}) \end{bmatrix} = k(x, x') \quad (519)$$

$$x \in \mathbb{R}^n \quad (520)$$

$g$  bijective,  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

$$y = f(g(x)) \quad (521)$$

$$y'(x) = \nabla_{g(x)} f(g(x)) J_x(g(x)) \quad (522)$$

$$k_a(x_a, x'_a)k(x_b, x'_b) = \sum_{n=1}^N \theta(x_{an})\theta(x'_{an}) \sum_{m=1}^M \psi(x_{bm})\psi(x'_{bm}) \quad (523)$$

$$= \sum_{n=1}^N \sum_{m=1}^M \theta(x_{an})\theta(x'_{an})\psi(x_{bm})\psi(x'_{bm}) \quad (524)$$

$$= \sum_{n=1}^N \sum_{m=1}^M \theta(x_{an})\psi(x_{bm})\theta(x'_{an})\psi(x'_{bm}) \quad (525)$$

$$= \phi(x)^T \phi(x') = k(x, x') \quad (526)$$

## 6.10

$$y(x) = \sum_{n=1}^N f(x_n)f(x_n)(K + \lambda I_N)^{-1}t \quad (527)$$

which is proportional to  $f$ ?

## 6.11

Just by plugging in we see that we observe an infinite amount of terms in the kernel dot-product, and therefore the vectors are of infinite dimensionality.



## 6.13

If  $\phi$  is an invertible differentiable transformation of  $\theta$ .

$$g(\theta, \mathbf{x}) = J_\theta(\phi) \nabla_\phi \ln p(\mathbf{x}, \phi) = J_\theta(\phi) g(\phi, \mathbf{x}) \quad (528)$$

$$g(\phi, \mathbf{x}) = J_\theta^{-1}(\phi) g(\theta, \mathbf{x}) \quad (529)$$

$$\mathbf{F}' = J_\theta^{-1}(\phi) \mathbb{E}_x[g(\theta, \mathbf{x})g(\theta, \mathbf{x})^T]J_\theta^{-T}(\phi) = J_\theta^{-1}(\phi) \mathbf{F} J_\theta^{-T}(\phi) \quad (530)$$

$$k'(\mathbf{x}, \mathbf{x}') = g(\phi, \mathbf{x})^T \mathbf{F}'^{-1} g(\phi, \mathbf{x}') \quad (531)$$

$$= (J_\theta^{-1}(\phi) g(\theta, \mathbf{x}))^T (J_\theta^{-1}(\phi) \mathbf{F} J_\theta^{-T}(\phi))^{-1} J_\theta^{-1}(\phi) g(\theta, \mathbf{x}') \quad (532)$$

$$= g(\theta, \mathbf{x})^T J_\theta^{-T}(\phi) J_\theta^T(\phi) \mathbf{F} J_\theta(\phi) J_\theta^{-1}(\phi) g(\theta, \mathbf{x}') \quad (533)$$

$$= g(\theta, \mathbf{x})^T \mathbf{F} g(\theta, \mathbf{x}') \quad (534)$$

$$= k(\mathbf{x}, \mathbf{x}') \quad (535)$$

Therefore, the Fisher kernel is invariant.

## 6.14

$$\nabla_\theta (\ln p(\mathbf{x}, \theta)) = \nabla_\mu \ln p(\mathbf{x}, \mu) \quad (536)$$

$$= S^{-1}(\mathbf{x} - \mu) \quad (537)$$

$$\mathbf{F} = \mathbb{E}[S^{-1}(\mathbf{x} - \mu)(S^{-1}(\mathbf{x} - \mu))^T] \quad (538)$$

$$= S^{-1} \mathbb{E}_x[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] S^{-1} \quad (539)$$

$$= S^{-1} S S^{-1} = S^{-1} \quad (540)$$

$$k(\mathbf{x}, \mathbf{x}') = (S^{-1}(\mathbf{x} - \mu))^T \mathbf{F}^{-1} (S^{-1}(\mathbf{x}' - \mu)) \quad (541)$$

$$= (\mathbf{x} - \mu)^T S^{-1} S S^{-1} (\mathbf{x}' - \mu) \quad (542)$$

$$= (\mathbf{x} - \mu)^T S^{-1} (\mathbf{x}' - \mu) \quad (543)$$

## 6.15

Since the gram matrix is positive semi-definite we have:

$$|K| = k(x_1, x_1)k(x_2, x_2) - k(x_2, x_1)k(x_1, x_2) \geq 0 \quad (544)$$

from which the result follows.

## 6.18

$$\mathbf{z} = (x, t) \quad (545)$$

$$p(\mathbf{z}) = \frac{\sum_n \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)}{\sum_m \int \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I) dt} = \frac{\sum_n \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)}{\sum_m \mathcal{N}(x - x_n, \sigma^2 I)} \quad (546)$$

$$\frac{\sum_n \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)}{\sum_m \mathcal{N}(x - x_n, \sigma^2 I)} = \frac{\sum_n \frac{1}{2\pi\sigma^2} \exp[-\frac{1}{2\sigma^2}(\mathbf{z} - \mathbf{z}_n)^T(\mathbf{z} - \mathbf{z}_n)]}{\sum_m \frac{1}{\sqrt{2\pi}\sigma} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]} \quad (547)$$

$$= \frac{\sum_n \frac{1}{2\pi\sigma^2} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2 + -\frac{1}{2\sigma^2}(t - t_n)^2]}{\sum_m \frac{1}{\sqrt{2\pi}\sigma} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]} \quad (548)$$

$$= \sum_n \frac{\frac{1}{\sqrt{2\pi}\sigma^2} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]}{\sum_m \frac{1}{\sqrt{2\pi}\sigma} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]} \frac{1}{\sqrt{2\pi}\sigma^2} \exp[-\frac{1}{2\sigma^2}(t - t_n)^2] \quad (549)$$

$$= \sum_n \pi_n \mathcal{N}(t|t_n, \sigma^2) \quad (550)$$

## 6.24

$$u^T W u = u^T \sqrt{W} \sqrt{W} u \quad (551)$$

$$= (\sqrt{W} u)^T \sqrt{W} u \quad (552)$$

$$> 0 \quad \forall u \in \mathbb{R} \setminus \mathbf{0} \quad (553)$$

$$u^T (W + V) u = u^T W u + u^T V u > 0 \quad (554)$$

## 6.25

$$a_N = a_N - \nabla \nabla \Psi(a_N) \nabla \Psi(a_n) \quad (555)$$

$$= a_N + (W_N + C^{-1})^{-1} [t_N - \sigma_N - C^{-1} a_N] \quad (556)$$

$$= (W_N + C_N^{-1}) [W_N a_N - \sigma_N + t_N] \quad (557)$$

$$= C_N ((W_N + C_N^{-1}) C_N)^{-1} [W_N a_N - \sigma_N + t_N] \quad (558)$$

$$= C_N (C_N W_N + I)^{-1} [W_N a_N - \sigma_N + t_N] \quad (559)$$

## 6.26

## 7.2

$$t_n(\mathbf{w}^T \phi(\mathbf{x}) + b) = \gamma \quad (560)$$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \gamma) \quad (561)$$

$$\nabla_{\mathbf{w}} = \mathbf{w} - \sum_{n=1}^N \nabla a_n t_n \mathbf{w}^T \phi(\mathbf{x}_n) \quad (562)$$

$$= \mathbf{w} - \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) = 0 \quad (563)$$

$$\iff \mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \quad (564)$$

$$\frac{\partial}{\partial b} = - \sum_{n=1}^N a_n t_n = 0 \quad (565)$$

$$= \sum_{n=1}^N a_n t_n \quad (566)$$

$$\tilde{L} = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n t_n a_m t_m \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) - \sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \gamma) \quad (567)$$

$$\sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \gamma) = \sum_{n=1}^N 0 - a_n \gamma \quad (568)$$

## 7.6

Since we have  $p(t = 1) = \sigma(y)$  and  $p(t = -1) = \sigma(-y)$  and  $t \in \{-1, 1\}$  we have  $p(t|y) = \sigma(ty)$ . Therefore

$$-\ln p(t_n) = - \sum_{n=1} \ln \sigma(t_n y_n) \quad (569)$$

which is the cross-entropy error function.

## 7.7

It is quite straightforward if you use the following steps.

$$\sum_{n=1}^N \xi_n (C - \mu_n - a_n) + \sum_{n=1}^N \hat{\xi}_n (C - \hat{\mu}_n - \hat{a}_n) = 0 \quad (570)$$

and

$$- \sum_{n=1}^N a_n (y_n) - \sum_{n=1}^N \hat{a}_n (-y_n) = - \sum_{n=1}^N (a_n - \hat{a}_n) y_n = 0 \quad (571)$$

## 7.8

This follows from 7.67 and 7.68.

## 7.9

$$p(\mathbf{w}, \mathbf{t}) = \mathcal{N}(\mathbf{w}|_N, \mathbf{S}_N) \quad (572)$$

$$\mathbf{S}_N = \mathbf{S}_0(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\mathbf{\Phi}^T\mathbf{t}) \quad (573)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta\mathbf{\Phi}^T\mathbf{\Phi} \quad (574)$$

We have

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta\mathbf{\Phi}^T\mathbf{\Phi} = \mathbf{A} + \beta\mathbf{\Phi}^T\mathbf{\Phi} = \Sigma^{-1} \quad (575)$$

Considering  $\mathbf{m}_0 = \mathbf{0}$ ,  $\mathbf{m}_N$  follows directly.

## 7.16

## 7.18

$$\nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) = \nabla \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) - \mathbf{A}\mathbf{w} \quad (576)$$

$$\nabla \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) = \sum_{n=1}^N t_n \nabla \ln y_n + (1 - t_n) \nabla \ln(1 - y_n) \quad (577)$$

$$\nabla \ln y_n = (1 - y_n)\phi(\mathbf{x}_n) \quad (578)$$

$$\nabla \ln(1 - y_n) = -y_n\phi(\mathbf{x}_n) \quad (579)$$

$$\nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) = \sum_{n=1}^N \phi(\mathbf{x}_n)[t_n - y_n t_n - y_n + y_n t_n] - \mathbf{A}\mathbf{w} \quad (580)$$

$$= \mathbf{\Phi}^T[\mathbf{t} - \mathbf{y}] - \mathbf{A}\mathbf{w} \quad (581)$$

$$\nabla \nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) = \nabla \mathbf{\Phi}^T[\mathbf{t} - \mathbf{y}] - \mathbf{A} \quad (582)$$

$$= \nabla \sum_{n=1}^N [t_n - y_n]\phi_n - \mathbf{A} \quad (583)$$

$$= - \sum_{n=1}^N \phi_n \nabla^T y_n - \mathbf{A} \quad (584)$$

$$= - \sum_{n=1}^N \phi_n y_n (1 - y_n) \phi_n^T - \mathbf{A} \quad (585)$$

$$= -\mathbf{\Phi}^T \mathbf{B} \mathbf{\Phi} - \mathbf{A} \quad (586)$$

Which gives the result.

## 8.1

$$\int p(\mathbf{x})d\mathbf{x} = \int \int p(x_K|pa_K) \prod_{k=1}^{K-1} p(x_k|pa_k)dx_K dx_1 \dots dx_{K-1} \quad (587)$$

$$= \int \int p(x_K|pa_K)dx_K \prod_{k=1}^{K-1} p(x_k|pa_k)dx_1 \dots dx_{K-1} \quad (588)$$

$$= \int \prod_{k=1}^{K-1} p(x_k|pa_k)dx_1 \dots dx_{K-1} \quad (589)$$

$$\vdots \quad (590)$$

$$= \int p(x_1)dx_1 = 1 \quad (591)$$

## 8.2

For an acyclic graph, if you number any graph (backwards or forwards) and (by accident) encounter a connection that connects a node to a lower-numbered node you can always swap the numbers (and adjust the rest of the graph accordingly) and continue.

## 8.5

## 8.6

The constraint  $\sum_i \mu_i = 1$  ensures that any  $\mu_i : i > 0$  increases the probability. This means that setting the cutoff point at which we predict  $y = 1$  at  $\mu_0$  ensures the OR-function. The  $\mu_i :> 0$  control the increase in probability for that  $x_i$ .

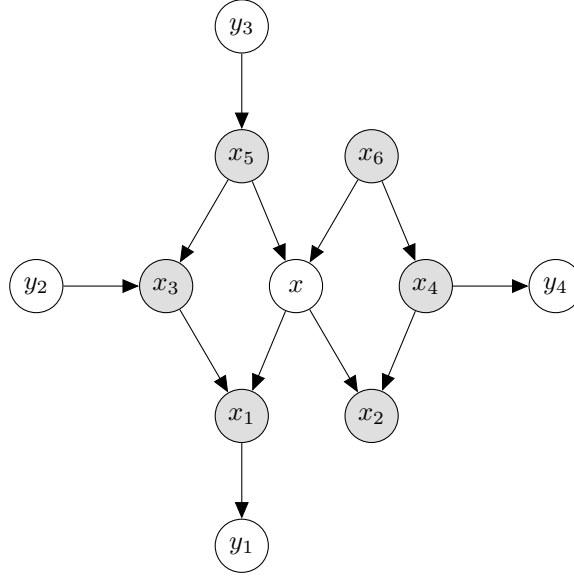
## 8.8

$$p(a, b, c|d) = \frac{p(a, b, c)}{p(d)} = \frac{p(a|b, c)p(b)p(c)}{p(d)} = \frac{p(a)p(b)p(c)}{p(d)} \quad (592)$$

$$\int \frac{p(a)p(b)p(c)}{p(d)}dc = p(a|d)p(b|d) \quad (593)$$

Which is what we needed to prove.

## 8.9



$\{y_1, y_2, y_3, y_4\}$  are all the possible connections for the Markov blanket. The path from  $x$  to  $y_1$  via  $x_1$  is blocked since  $x_1$  is in  $C$  and the path meets head to tail. Same path through  $x_5$  is blocked since they meet tail to tail and  $x_5$  is in  $C$ . The path from  $x$  to  $y_2$  is blocked by  $x_3$  since the arrows meet head to tail. Same for  $y_3$  to  $x$ . All paths have either a head to tail or tail to tail node with an observed variable in it.

## 8.10

$$p(a, b, c, d) = p(a)p(b)p(c|a, b)p(d|c) \quad (594)$$

$$p(a, b) = \int \int p(a, b, c, d) dc dd = p(a)p(b) \quad (595)$$

Second part:

$$p(a, b, c|d) = \frac{p(a, b, c, d)}{p(d)} = \frac{p(a)p(b)p(c|a, b)p(d|c)}{p(d)} \quad (596)$$

$$p(a, b|d) = \int p(a, b, c|d) dc = \frac{p(a)p(b)p(d|c)}{p(d)} \quad (597)$$

This does not factor into  $p(a|d)p(b|d)$ .

## 8.12

In an undirected graph we can remove or add a link between each node and every other node and that will create a new graph. So  $2^{\# \text{Links}}$  graphs. Any node can connect with any other node, giving  $N(N-1)$  pairs, but we cannot count the reverse paths so we divide by two.

### 8.13

$$E(x, y)_{x_k=1} - E(x, y)_{x_k=-1} = h \sum_i x_i - \beta \sum_{i,j} x_i x_j - \eta \sum_i x_i y_i + h - \beta x_k \sum_j x_j - \eta x_k \sum_i y_i \quad (598)$$

$$- h \sum_i x_i - \beta \sum_{i,j} x_i x_j - \eta \sum_i x_i y_i + h - \beta x_k \sum_j x_j - \eta x_k \sum_i y_i \quad (599)$$

$$= 2h - 2\beta x_k \sum_j x_j - \eta x_k \sum_i y_i \quad (600)$$

### 8.14

$$p(x, y) = \frac{1}{Z} \exp[-E(x, y)] \quad (601)$$

$$\ln p(x, y) = -\ln(Z) - E(x, y) = -\ln(Z) + \eta \sum_i x_i y_i \quad (602)$$

This is maximized when  $x_i = y_i$ , i.e.  $-1 \cdot -1 = 1$  or  $1 \cdot 1 = 1$

### 8.20

### 8.22

### 8.26

### 9.1

The loss function clearly is convex. Moreover, 9.2 and 9.4 both are guaranteed to lower the function (arg min and an analytical solution for  $\mu_k$ . Therefore, it will always converge.

### 9.2

### 9.3

$$p(\mathbf{x}, \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} \pi_k^{z_k} \quad (603)$$

$$p(\mathbf{z}) = \sum_{\mathbf{z}} \prod_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} \quad (604)$$

Since the product is 1 for every  $z_k \neq k$  this becomes  $\sum_{\mathbf{z}} \pi_{z_k} \mathcal{N}(\mathbf{x} | \mu_{z_k}, \Sigma_{z_k})$ , which is equal to the required equation 9.7.

## 9.4

Log posterior:  $\ln p(\boldsymbol{\theta}|\mathbf{x}) \propto \ln p(\mathbf{x}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})$ . For the e-step: Evaluate  $p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old})$ , this function only depends on the likelihood part of the objective, so by definition will be the same. For the m-step:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{old})_{MAP} = \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\boldsymbol{\theta}, \mathbf{z}|\mathbf{x}) \quad (605)$$

$$\propto \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\boldsymbol{\theta}, \mathbf{z}|\mathbf{x}) \quad (606)$$

$$= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \quad (607)$$

$$= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) \quad (608)$$

## 9.5

This is obvious, as there simply is no connection between  $z_m$  and  $x_n$

## 9.7

$$\nabla_{\mu} \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{kn} \nabla_{\mu} \ln \mathcal{N}(\mathbf{x}_n, \mu_k, \Sigma_k) \quad (609)$$

$$\nabla \ln \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) = -\frac{1}{2} \nabla (\mathbf{x}_n - \mu_k)^T \sigma_k^{-1} (\mathbf{x}_n - \mu_k) = \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \quad (610)$$

Since  $z$  is 1-of- $k$  this gradient only concerns distribution  $k$  for every  $n$ .

$$\mathcal{L}(\pi_k, \lambda) \ln p(\mathbf{X}, \mathbf{Z}|\mu_k \Sigma_k \pi_k) + \lambda[-1 + \sum_k \pi_k] \quad (611)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \sum_{n=1}^N z_{nk} \frac{1}{\pi_k} - \lambda \quad (612)$$

$$= \sum_{n=1}^N z_{nk} - \pi_k \lambda \quad (613)$$

$$= \sum_{n=1}^N \sum_{k=1}^K z_{nk} - \pi_k \lambda \quad (614)$$

$$= N - \lambda = 0 \iff \lambda = N \quad (615)$$

Substituting back:

$$\sum_{n=1}^N z_{nk} \frac{1}{\pi_k} = N \iff \pi_k = \frac{1}{N} \sum_{n=1}^N z_{nk} = \frac{N_k}{N} \quad (616)$$



$$\nabla_{\mu_k} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) [\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \quad (617)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \nabla \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (618)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (619)$$

$$\Rightarrow \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (620)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n - \underbrace{\sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\mu}_k}_{N_k} = 0 \quad (621)$$

$$(622)$$

Which leads to the correct answer.

## 9.9

$$\frac{\partial}{\partial \boldsymbol{\Sigma}_k} = \sum_{n=1}^N \gamma(z_{nk}) \left[ \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) - \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \right] \quad (623)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[ \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T - \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln |\boldsymbol{\Sigma}^{-1}|^{\frac{1}{2}} \right] \quad (624)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[ \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T - \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{2} \ln |\boldsymbol{\Sigma}^{-1}| \right] \quad (625)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[ \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T - \frac{1}{2} \boldsymbol{\Sigma} \right] = 0 \quad (626)$$

$$\Leftrightarrow \sum_{n=1}^N \gamma(z_{nk}) \left[ \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right] = \sum_{n=1}^N \gamma(z_{nk}) \frac{1}{2} \boldsymbol{\Sigma} = N_k \frac{1}{2} \boldsymbol{\Sigma}_k \quad (627)$$

From which the answer is easily seen.

## 9.11

$$\mathbb{E}_{\mathbf{z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(\mathbf{z}_{nk}) [\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \quad (628)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \gamma(\mathbf{z}_{nk}) [\ln \pi_k + \ln \frac{1}{(2\pi\epsilon)^{M/2}} - \frac{1}{2\epsilon} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2] \quad (629)$$

$$\propto \sum_{n=1}^N \sum_{k=1}^K \gamma(\mathbf{z}_{nk}) [\epsilon \ln \pi_k + \epsilon \ln \frac{1}{(2\pi\epsilon)^{M/2}} - \frac{1}{2} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2] \quad (630)$$

$$= \lim_{\epsilon \rightarrow 0} -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 + C \quad (631)$$

## 9.12

$$\mathbb{E}[\mathbf{x}] = \int \mathbf{x} p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\pi}) d\mathbf{x} \quad (632)$$

$$= \int \mathbf{x} \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} \quad (633)$$

$$= \sum_{k=1}^K \pi_k \int \mathbf{x} p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} \quad (634)$$

$$= \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k \quad (635)$$

$$\text{cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \quad (636)$$

$$= \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T \quad (637)$$

$$(638)$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \int \mathbf{x}\mathbf{x}^T \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} \quad (639)$$

$$= \sum_{k=1}^K \pi_k \int \mathbf{x}\mathbf{x}^T p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} = \sum_{k=1}^K \pi_k \mathbb{E}_k[\mathbf{x}\mathbf{x}^T] \quad (640)$$

$$(641)$$

$$\mathbb{E}_k[\mathbf{x}\mathbf{x}^T] = \mathbb{E}_k[(\mathbf{x} - \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{x}])^T] \quad (642)$$

$$= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T + (\mathbf{x} - \mathbb{E}[\mathbf{x}]) \mathbb{E}[\mathbf{x}]^T + \mathbb{E}[\mathbf{x}] (\mathbf{x} - \mathbb{E}[\mathbf{x}])^T + \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T] \quad (643)$$

$$= \underbrace{\text{cov}_k(\mathbf{x})}_{=\boldsymbol{\Sigma}_k} + \boldsymbol{\mu}\boldsymbol{\mu}^T \quad (644)$$

## 9.14

$$p(\mathbf{x}, \mathbf{z}, \boldsymbol{\mu})p(\mathbf{z}|\boldsymbol{\pi}) = \prod_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \pi_k^{z_k} \quad (645)$$

$$= \sum_{\mathbf{z}} \prod_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \pi_k^{z_k} \quad (646)$$

$$(647)$$

Since  $\mathbf{z}$  is 1-of-K and the inner product only returns when  $k = z_k$  this becomes  $\sum_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \pi_k^{z_k}$

## 9.15

$$\frac{\partial}{\partial \mu_{ki}} = \sum_{n=1}^N \gamma(z_{nk}) \left[ \frac{\partial}{\partial \mu_{ki}} x_{ni} \ln \mu_{ki} + \frac{\partial}{\partial \mu_{ki}} (1 - x_{ni}) \ln(1 - \mu_{ki}) \right] \quad (648)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[ \frac{x_{ni}}{\mu_{ki}} - \frac{1 - x_{ni}}{1 - \mu_{ki}} \right] \quad (649)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[ \frac{x_{ni}(1 - \mu_{ki}) - \mu_{ki}(1 - x_{ni})}{\mu_{ki}(1 - \mu_{ki})} \right] = 0 \quad (650)$$

$$\Rightarrow \sum_{n=1}^N \gamma(z_{nk}) \mu_{ki} = \sum_{n=1}^N \gamma(z_{nk}) x_n \quad (651)$$

$$= N_k \mu_{ki} = \sum_{n=1}^N \gamma(z_{nk}) x_n \quad (652)$$

## 9.16

This solution is exactly the same as exercise 9.7's.

## 9.17

By observing equation 9.51 we see that since  $p(\mathbf{x}_n|\boldsymbol{\mu}_k) \leq 1$  and  $\sum_k \pi_k = 1$  the maximum of the ln is 0.

## 9.20

$$\frac{\partial}{\partial \alpha} = \frac{M}{2} \frac{1}{\alpha} - \frac{1}{2} \mathbb{E}[\mathbf{w}^T \mathbf{w}] = 0 \quad (653)$$

$$\Rightarrow \frac{M}{\alpha} = \mathbb{E}[\mathbf{w}^T \mathbf{w}] \quad (654)$$

$$\alpha = \frac{M}{\mathbb{E}[\mathbf{w}^T \mathbf{w}]} \quad (655)$$

## 9.24

$$\ln p(\mathbf{x}|\boldsymbol{\theta}) = \sum_z q(\mathbf{z}) \ln p(\mathbf{x}|\boldsymbol{\theta}) \quad (656)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \quad (657)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \frac{q(\mathbf{z})}{q(\mathbf{z})} \quad (658)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} + \sum_z q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \quad (659)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} - \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} \quad (660)$$

## 9.25

This directly follows from the fact that the KL-divergence reduces to 0 if the distributions are equal.

## 9.26

### 10.1

We've shown this in the previous chapter.

$$\ln p(\mathbf{x}) = \int q(\mathbf{z}) \ln p(\mathbf{x}) d\mathbf{z} \quad (661)$$

$$= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \quad (662)$$

$$= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z})q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})q(\mathbf{z})} \quad (663)$$

$$= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} - \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{x}|\mathbf{z})} \quad (664)$$

$$(665)$$

### 10.2

This is easily seen by simply filling in the values. The result is a solution to the equations.

## 10.8

Filling in:

$$\mathbb{E}[\tau] = \frac{a}{b} = (a_0 + \frac{N}{2})(b_0 + \frac{1}{2} \mathbb{E}_\mu[\sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^2 + \lambda_0(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^2])^{-1} \underset{N \rightarrow \infty}{\approx} b_N^{-1} \quad (666)$$

Which follows from the fact that the denominator grows proportionally to the numerator.

Similarly, the variance  $\text{Var}[\tau] = \frac{a}{b^2}$  is derived but this time the denominator grows quadratically and therefore goes to zero in the limit  $N \rightarrow \infty$ .

## 10.10

$$\ln p(\mathbf{x}) = \sum_m \sum_{\mathbf{z}} q(\mathbf{z}, m) \ln p(\mathbf{x}) \quad (667)$$

$$= \sum_m \sum_{\mathbf{z}} q(\mathbf{z}, m) \ln \frac{p(\mathbf{x}, m, \mathbf{z})}{p(m, \mathbf{z}|\mathbf{x})} \frac{q(\mathbf{z}, m)}{q(\mathbf{z}, m)} \quad (668)$$

Which results in 10.35 after rearrangement.

## 10.15

We have:

$$\mathbb{E}[\pi_k] = \frac{\alpha_k}{\sum_{k=1}^K \alpha_k} \quad (669)$$

and  $\alpha_k = \alpha_0 + N_k$ .

$$\mathbb{E}[\pi_k] = \frac{\alpha_0 + N_k}{\sum_{k=1}^K \alpha_0 + N_k} \quad (670)$$

$$= K\alpha_0 + \underbrace{\sum_{k=1}^K N_k}_{=N} \quad (671)$$

## 10.21

This is easily observed by the fact that if we fix 1 distribution we still have  $K - 1$  combinations. If we continue this cascade you end up with  $K \cdot K - 1 \dots 2 \cdot 1$  distributions, which is  $K!$

## 10.29

$$\frac{dd \ln(x)}{dx} = -x^{-2} \quad (672)$$

which is negative everywhere so  $\ln(x)$  is concave.

Taylor approximation:

$$y(x) = \ln(\xi) + \frac{1}{\xi}(x - \xi) = \lambda x - \ln(\lambda) - 1 \quad \lambda = \frac{1}{\xi} \quad (673)$$

$$g(\lambda) = \min_x [\lambda x - f(x)] \quad (674)$$

$$\frac{d}{dx} = \lambda - \frac{1}{x} = 0 \iff x = \frac{1}{\lambda} \quad (675)$$

Therefore  $g(\lambda) = 1 - \ln \frac{1}{\lambda} = 1 + \ln \lambda$ .

$$\frac{d}{d\lambda} \lambda x - g(\lambda) = x - \frac{1}{\lambda} \iff x = \frac{1}{\lambda} \quad (676)$$

Plugging this into  $y(x)$  gives the result  $\ln x$

## 10.30

$$\frac{d}{dx} = \frac{1}{1 + e^{-x}} e^{-x} = \sigma(x) e^{-x} \quad (677)$$

$$\frac{dd}{dx} = \sigma(x)(1 - \sigma(x))e^{-x} - e^{-x}\sigma(x) = -\sigma^2(x)e^{-x} \quad (678)$$

Both functions are positive everywhere, showing that the function itself is negative, therefore concave.

$$\frac{d}{dx} = \frac{1}{1 + e^{-x}} e^{-x} \quad (679)$$

$$\frac{dd}{dx} = \sigma(x)(1 - \sigma(x))e^{-x} - e^{-x}\sigma(x) = -\sigma^2(x)e^{-x} \quad (680)$$

Taylor:

$$f(x) = f(\xi) + f'(\xi)(x - \xi) + \mathcal{O}(\xi^2) \quad (681)$$

Since the approximation is linear and  $f(x)$  is concave it must be that the LHS is smaller-equal to the RHS.

$$f(x) \leq -\ln(1 + e^{-xi}) + \sigma(\xi)e^{-\xi}(x - \xi) + \mathcal{O}(\xi^2) \quad (682)$$

$$= -\ln(1 + e^{-xi}) + \sigma(\xi)e^{-\xi}x - \sigma(\xi)e^{-\xi}\xi + \mathcal{O}(\xi^2) \quad (683)$$

$$= \lambda x - g(\lambda) \quad \lambda = \sigma(\xi)e^{-\xi} \quad (684)$$

## 10.33

$$\frac{d}{d\xi_n} = (1 - \sigma(\xi_n)) - \frac{1}{2} - \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n \lambda'(\xi_n) + 2\xi_n \lambda(\xi_n) + \lambda'(\xi_n) \xi_n^2 \quad (685)$$

$$= -2\xi_n \lambda(\xi_n) - \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n \lambda'(\xi_n) + 2\xi_n \lambda(\xi_n) + \lambda'(\xi_n) \xi_n^2 \quad (686)$$

$$= -\lambda'(\xi_n)(\xi_n^2 - \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n) = 0 \quad (687)$$

$$\iff \xi_n^2 = \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n \quad (688)$$

## 10.37

$$q^{\setminus j}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_j(\boldsymbol{\theta})} \quad (689)$$

$$Z_j = \int q^{\setminus j}(\boldsymbol{\theta}) f_j(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int q^{\setminus j}(\boldsymbol{\theta}) \tilde{f}_j(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int (q(\boldsymbol{\theta}) d\boldsymbol{\theta}) = 1 \quad (690)$$

$$\tilde{f}_j(\boldsymbol{\theta}) = \frac{q^{new}}{q^{\setminus j}(\boldsymbol{\theta})} = \tilde{f}_j = f_j \quad (691)$$

## 11.1

$$\mathbb{E}[\hat{f}] = \mathbb{E}\left[\frac{1}{L} \sum_{l=1}^L f(z^l)\right] \quad (692)$$

$$= \frac{1}{L} \sum_p (z) f(z^l) dz \quad (693)$$

$$= \frac{1}{L} \sum_{l=1}^L \mathbb{E}[f] = \mathbb{E}[f] \quad (694)$$

For the variance, we have to note that

$$\mathbb{E}[f(z^l), f(z^k)] = \text{Var}(f(z)) + \mathbb{E}[f(z)]^2 \quad (695)$$

And  $\text{Var}(f(z)) = 0$  for  $k \neq l$ .

$$\text{Var}[\hat{f}] = \mathbb{E}[\hat{f}^2] - \mathbb{E}[\hat{f}]^2 \quad (696)$$

$$= \frac{1}{L^2} \left( \sum_{l=1}^L \sum_{k=1}^L \mathbb{E}[f(z^l) f(z^k)] \right) - \mathbb{E}[f]^2 \quad (697)$$

$$= \frac{1}{L^2} \left( \sum_{l=1}^L \sum_{k=1}^L \mathbb{E}[f(z^l), f(z^k)] \right) - \mathbb{E}[f]^2 \quad (698)$$

$$= \frac{1}{L^2} L \text{Var}(f(z)) + L^2 \mathbb{E}[f]^2 - \mathbb{E}[f]^2 \quad (699)$$

$$= \frac{1}{L} \text{Var}(f(z)) \quad (700)$$

## 11.2

What we need to show is that we can transform  $z$  into any distribution if we use 11.6.

$$p(y) = 1 \cdot \left| \frac{dz}{dy} \right| = 1 \cdot p(y) \quad (701)$$

### 11.3

$$z = \left( \frac{1}{\pi} \int_{-\infty}^y \frac{1}{1 + \hat{y}^2} \right) \quad (702)$$

$$= \left( \frac{1}{\pi} [\arctan(y) - \arctan(-\infty)] \right) \quad (703)$$

$$= \left( \frac{1}{\pi} \arctan(y) + \frac{1}{2} \right) \quad (704)$$

$$y = h^{-1}(z) \quad (705)$$

$$\iff \pi z - \frac{\pi}{2} = \arctan(y) \quad (706)$$

$$\iff h^{-1}(z) = \tan(\pi z - \frac{\pi}{2}) \quad (707)$$

### 11.5

Hint 1: Show that the expectation and covariance are equal to  $\mu$  and  $\Sigma$ .

$$\mathbb{E}[y] = \mathbb{E}_z[\mu] + \mathbb{E}_z[\mathbf{L}z] \quad (708)$$

$$= \mathbb{E}_z[\mu] + \mathbf{L} \mathbb{E}_z[z] \quad (709)$$

$$= \mathbb{E}_z[\mu] + \mathbf{L}\mathbf{0} \quad (710)$$

$$= \mu \quad (711)$$

$$\text{cov}[y] = \mathbb{E}[(y - \mu)(y - \mu)^T] \quad (712)$$

$$= \mathbb{E}[yy^T - \mu y^T - y \mu^T + \mu \mu^T] \quad (713)$$

$$= \mathbb{E}[(\mu + \mathbf{L}z)(\mu + \mathbf{L}z)^T - \mu(\mu + \mathbf{L}z)^T - (\mu + \mathbf{L}z)\mu^T + \mu \mu^T] \quad (714)$$

$$= \mathbb{E}[\mu \mu^T + \mu(\mathbf{L}z)^T + \mathbf{L}z \mu^T + \mathbf{L}z(\mathbf{L}z)^T - \mu \mu^T - \mu(\mathbf{L}z)^T - \mu \mu^T - \mathbf{L}z \mu^T + \mu \mu^T] \quad (715)$$

$$= \mathbb{E}[\mathbf{L}z(\mathbf{L}z)^T] \quad (716)$$

$$= \mathbf{L} \mathbb{E}[zz^T] \mathbf{L}^T \quad (717)$$

Now, we know that  $\text{Var}[z] = \mathbb{E}[zz^T] + \mu \mu^T$ , therefore  $\mathbb{E}[zz^T] = \mathbf{I}$ . Then the result follows.

### 11.7

$$p(y) = p(z) \left| \frac{dz}{dy} \right| \quad (718)$$

We have  $z \sim \mathcal{U}(0, 1)$  and therefore  $p(z) = \frac{1}{1-0} = 1$ . Inverting the given equation:

$$y = b \tan z + c \iff \frac{y - c}{b} = \tan z \iff \arctan \frac{y - c}{b} = z \quad (719)$$



$$\frac{dz}{dy} = \frac{1}{1 + \left(\frac{y-c}{b}\right)^2} \frac{1}{b} \quad (720)$$

$$(721)$$

Multiplying these two gives the desired result without the scaling constant  $k$ .

## 11.10

## 11.12

Since there are regions with zero conditional probability, Gibbs sampling will not be ergodic.

## 11.14

Hint: calculate  $\mathbb{E}[z'_i]$  and  $E[(z'_i - \mu_i)^2]$ .

$$\mathbb{E}_{z,\nu}[z'_i] = \mathbb{E}[\mu_i] + \mathbb{E}[\alpha(z_i - \mu_i)] + \mathbb{E}[\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu] \quad (722)$$

$$= \mu_i \quad (723)$$

$$\text{Var}_{z,\nu} = \mathbb{E}[z'^2_i] - \mathbb{E}[z'_i]^2 \quad (724)$$

$$= \mathbb{E}[(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)] - \mu_i^2 \quad (725)$$

$$= \mathbb{E}[(\mu_i + \alpha(z_i - \mu_i))(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu) + (\mu_i + \alpha(z_i - \mu_i))\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu] \quad (726)$$

$$+ (\mu_i + \alpha(z_i - \mu_i))(\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu) + (\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)(\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)] - \mu_i^2 \quad (727)$$

$$\vdots \quad (728)$$

$$(729)$$

## 11.15

$$\frac{dH}{dr_i} = \frac{1}{Z} \frac{d}{dr_i} r_i^2 = r_i \quad (730)$$

$$\frac{dr_i}{d\tau} = -\frac{dE(z)}{dz_i} = -\frac{dH}{dz_i} \quad (731)$$

## 11.16

$$p(\mathbf{r}|\mathbf{z}) = \frac{p(\mathbf{r}, \mathbf{z})}{p(\mathbf{z})} = \frac{Z_p}{Z_H} \exp(-K(\mathbf{r})) \quad (732)$$

## 11.17

## 12.3

$$\|u_i\|^2 = [\frac{1}{(N\lambda_i)^{1/2}} X^T v_i]^T [\frac{1}{(N\lambda_i)^{1/2}} X^T v_i] \quad (733)$$

$$= (N\lambda_i)^{-1} [X^T v_i]^T [X^T v_i] \quad (734)$$

$$= (N\lambda_i)^{-1} v^T X X^T v_i \quad (735)$$

$$= \lambda_i^{-1} v^T \lambda_i v_i \quad (736)$$

$$= 1 \quad (737)$$

And therefore  $\|u_i\| = 1$

## 12.4

This problem can be solved by using 2.113-2.115 and just filling in the variables.

## 12.9

$$\frac{\partial \ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2)}{\partial \boldsymbol{\mu}} = -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \quad (738)$$

$$= \sum_{n=1}^N \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = 0 \quad (739)$$

$$\Rightarrow 0 = -N\boldsymbol{\mu} + \sum_{n=1}^N \mathbf{x}_n \quad (740)$$

$$\Rightarrow \boldsymbol{\mu} = \bar{\mathbf{x}} \quad (741)$$

## 12.14

$$D(D-1) + 1 - (D-1)(D-2)/2 = D^2 - D + 1 - \frac{1}{2}D^2 + \frac{3}{2}D - 1 = \frac{1}{2}D^2 + \frac{1}{2}D = D(D+1)/2 \quad (742)$$

## 12.18

$$M \times D + 1 \quad (743)$$

where  $W \in \mathbb{R}^{M \times D}$  and sigma is 1-dimensional.

## 13.1

$x_{n+2}$  is d-separated from  $x_n$  Since  $x_{n+1}$  is observed and the nodes meet head-to-tail.

### 13.3

Using the d-separation criterion, we see that there is always a path connecting any two observed variables  $x_n$  and  $x_m$  via the latent variables, and that this path is never blocked. Thus the predictive distribution  $p(x_{n+1}|x_1, \dots, x_n)$  for observation  $x_{n+1}$  given all previous observations does not exhibit any conditional independence properties, and so our predictions for  $x_{n+1}$  depends on all previous observations.

### 13.6

If  $\pi_k$  is 0, there is no probability density for latent variable  $z_k$ . This means that  $z_k$  will always be 0 and therefore the update steps  $\gamma(z_{nk})$  and  $\xi(z_{n-1}, z_{nk})$  will also be 0.

### 13.7

$$\frac{\partial Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})}{\boldsymbol{\mu}_k} = \sum_{n=1}^N \gamma(z_{nk}) \frac{1}{p(\mathbf{x}_n|\boldsymbol{\phi}_k)} \frac{\partial}{\partial \boldsymbol{\mu}_k} p(\mathbf{x}_n|\boldsymbol{\phi}_k) \quad (744)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \frac{1}{p(\mathbf{x}_n|\boldsymbol{\phi}_k)} p(\mathbf{x}_n|\boldsymbol{\phi}_k) \frac{\partial}{\partial \boldsymbol{\mu}_k} - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad (745)$$

$$= 0 \iff \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (746)$$

$$\iff \boldsymbol{\mu} = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \quad (747)$$

$$(748)$$

$$\frac{\partial Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})}{\boldsymbol{\Sigma}_k} = \sum_{n=1}^N \gamma(z_{nk}) \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln p(\mathbf{x}|\boldsymbol{\phi}_k) \quad (749)$$

$$(750)$$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln p(\mathbf{x}|\boldsymbol{\phi}_k) = \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \left[ \ln \frac{1}{(2\pi)^{D/2}} + \ln \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} - \frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (751)$$

$$= \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{2} \ln |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \frac{\partial}{\partial \boldsymbol{\Sigma}_k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) \quad (752)$$

$$= \frac{1}{2} \boldsymbol{\Sigma} - \frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \quad (753)$$

$$\frac{\partial}{\partial \Sigma_k} = \sum_{n=1}^N \gamma(z_{nk}) \left( \frac{1}{2} \Sigma - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right) = 0 \quad (754)$$

$$\iff \Sigma_k = \frac{\sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma(z_{nk})} \quad (755)$$

## 13.17

$$h(z_1) = p(z_1|u_1)p(x_1|z_1, u_1) \quad (756)$$

$$f_n(z_{n-1}, z_n) = p(z_n|z_{n-1}, u_n)p(x_n|z_n, u_n) \quad (757)$$

## 13.19

This follows from:

$$\frac{d}{dz} \ln p(\mathbf{Z}) = \frac{d}{dz} \ln p(z_0)p(z_1|z_0) \times \cdots \times p(z_n|z_{n-1}) = \begin{bmatrix} \frac{\partial}{\partial z_0} \ln p(z_0) \\ \frac{\partial}{\partial z_1} \ln p(z_1|z_0) \\ \vdots \\ \frac{\partial}{\partial z_n} p(z_n|z_{n-1}) \end{bmatrix} \quad (758)$$

So individually optimizing the conditional distributions will optimize the total distribution.

## 13.27

If the noise goes to zero,  $\Sigma = \mathbf{0}$ . Considering  $C = I$ : Note that  $K_1 = V_0(V_0 + \Sigma) = I$  Therefore  $\mu_1 = \mu_0 + x_1 - \mu_0 = x_1$

$V_n = (I - K_n C)P_{n-1} = 0$ . Therefore  $P_{n-1} = \Gamma$ . Now every term in 13.86 and 13.87 is directly dependent on  $x$ , except for  $\Gamma$ . I don't see how this term doesn't influence the distribution.

## 14.2

$$\mathbb{E} \left[ \left( \frac{1}{M} \sum_{m=1}^M \epsilon_m(x) \right)^2 \right] = \mathbb{E} \left[ \left( \frac{1}{M^2} \sum_{m=1}^M \sum_{l=1}^M \epsilon_m(x) \epsilon_l(x) \right) \right] \quad (759)$$

$$= \mathbb{E} \left[ \frac{1}{M^2} \left( \sum_{m=1}^M \sum_{\substack{l=1 \\ l \neq m}}^M \epsilon_m(x) \epsilon_l(x) + \sum_{l=1}^M \epsilon_l(x)^2 \right) \right] \quad (760)$$

$$= \frac{1}{M^2} \sum_{l=1}^M \mathbb{E}[\epsilon_l(x)^2] \quad (761)$$

$$= \frac{1}{M} E_{AV} \quad (762)$$

### 14.3

Through Jensen's inequality:

$$\sum_{m=1}^M \frac{1}{M} \epsilon_m(x)^2 \geq \left( \sum_{m=1}^M \frac{1}{M} \epsilon_m(x) \right)^2 \quad (763)$$

Therefore:

$$\mathbb{E}_x \left[ \sum_{m=1}^M \frac{1}{M} \epsilon_m(x)^2 \right] \geq \mathbb{E} \left[ \left( \sum_{m=1}^M \frac{1}{M} \epsilon_m(x) \right)^2 \right] = E_{com} \quad (764)$$

### 14.6

$$\frac{d}{d\alpha_m} = (e^{\alpha_m/2} + e^{-\alpha/2}) \sum_{n=1}^N w_n^{(m)} I(y_m(x_n) = t_n) = e^{\alpha_m/2} \sum_{n=1}^N w_n^{(m)} = 0 \quad (765)$$

$$\frac{e^{-\alpha_m/2}}{e^{\alpha_m/2} + e^{-\alpha_m/2}} = \epsilon_m \quad (766)$$

Rearranging this results in the desired formula.

[14.7](#) [14.8](#)

### 14.9

$$\mathcal{L}_m(\mathbf{x}_n) = (y_n - \sum_{l=1}^M \alpha_l \hat{y}_l)^2 \quad (767)$$

$$= (y - \underbrace{\sum_{l=1}^{M-1} \alpha_l \hat{y}_l}_{\text{residual}} + \alpha_M \hat{y}_M)^2 \quad (768)$$

$$(769)$$

### 14.10

$$L = \frac{1}{2} \sum_{n=1}^N (t_n - t)^2 \quad (770)$$

$$\frac{dL}{dt} = - \sum_{n=1}^N (t_n - t) = 0 \iff Nt = \sum_{n=1}^N t_n \iff t = \frac{1}{N} \sum_{n=1}^N t_n \quad (771)$$

### 14.13

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) = \prod_{n=1}^N \prod_{k=1}^K (\pi_k \mathcal{N}(t_n | \mathbf{w}_k^T \boldsymbol{\phi}_n, \beta^{-1})^{z_{nk}} \quad (772)$$

Taking the log of this will yield the result.

### 14.14

$$\frac{dL}{d\pi_k} = \sum_{n=1}^N \gamma_{nk} / \pi_k - \lambda = 0 \quad (773)$$

$$\iff \sum_k \sum_{n=1}^N \gamma_{nk} = \sum_k \pi_k \lambda \quad (774)$$

$$(775)$$

Since summing over  $k$  is equal to marginalizing  $p(z | \theta_k)$  we get  $\sum_{n=1}^N 1 = \lambda \iff \lambda = N$  Substituting back results in  $\pi_k = \frac{1}{N} \sum_{n=1}^N \gamma_{nk}$ .

### 14.15

$$\mathbb{E}[t | \hat{\boldsymbol{\phi}}, \boldsymbol{\theta}] = \sum_{k=1}^K \pi_k \mathbb{E}[t | \hat{\boldsymbol{\phi}}, \mathbf{w}_k, \beta] \quad (776)$$