

Bishop Questions

David Ruhe

July 20, 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

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

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

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

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

1.4

$$p'(y) = p'(g(y)) \cdot g'(y) \cdot |g'(y)| + p(g(y)) \frac{d|g'(y)|}{dy} \quad (17)$$

If we maximize $p(x)$ we get that the first term is 0. In the nonlinear case, the second term is not equal to 0 however, so maximizing one distribution does not mean maximizing the other distribution.

Derivation of the change of variables formula. We have a random variable X with pdf $f(x)$. We define $Y=g(X)$ where $g(\cdot)$ is a *monotone* function, the pdf of Y is obtained as follows.

We start with the CDF of y .

$$P(Y \leq y) = P(g(X) \leq y) \quad (18)$$

$$= P(X \leq g^{-1}(y)) \quad \text{Because } g \text{ is monotonically increasing.} \quad (19)$$

$$(20)$$

This is equivalent to $F(y) = F(g^{-1}(y))$. Therefore, if we differentiate both sides and use the chain rule we get

$$f(y) = f(g^{-1}(y)) \frac{d}{dy} g^{-1}(y) \quad (21)$$

If the function is monotonically decreasing we get

$$P(Y \leq y) = P(g(X) \leq y) \quad (22)$$

$$= P(X \geq g^{-1}(y)) \quad (23)$$

$$= 1 - F(g^{-1}(y)) \quad (24)$$

Differentiating both sides leads to

$$f(y) = f(g^{-1}(y)) \frac{d}{dy} g^{-1}(y) \quad (25)$$

Combining these formulas gives the formula with the absolute value. This works for both monotonically decreasing and increasing. If g is not monotonic, we need to split the function into monotonic parts and do this transformation for them individually, then sum.

1.5

$$\text{var}[f] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] \quad (26)$$

$$= \mathbb{E}[f(x)^2 + \mathbb{E}[f(x)]^2 - 2f(x)\mathbb{E}[f(x)]] \quad (27)$$

$$= \mathbb{E}[f(x)^2] + \mathbb{E}[f(x)]^2 - 2\mathbb{E}[f(x)]\mathbb{E}[f(x)] \quad (28)$$

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

1.6

$$\text{cov}[x, y] = \mathbb{E}_{x,y}[\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}] \quad (30)$$

$$= \int_{x,y} p(x, y) [\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}] \quad (31)$$

$$= \int_x \int_y p(x)p(y) [\{x - \mathbb{E}[x]\}\{y - \mathbb{E}[y]\}] \quad (32)$$

$$= \int_x p(x) \{x - \mathbb{E}[x]\} \int_y p(y) \{y - \mathbb{E}[y]\} \quad (33)$$

$$= \mathbb{E}[x - \mathbb{E}[x]] \mathbb{E}[y - \mathbb{E}[y]] \quad (34)$$

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

1.7

We have:

$$I = \int \exp(-\frac{1}{2\sigma^2}x^2)dx \quad (36)$$

Making use of $\int f(x)dx \int g(y)dy = \int \int f(x)g(y)dxdy$:

$$I^2 = \int \int \exp(-\frac{1}{2\sigma^2}(x^2 + y^2))dxdy \quad (37)$$

Using Pythagoras we know $x^2 + y^2 = r^2$.

$$I^2 = \int \int \exp(-\frac{1}{2\sigma^2}r^2)dxdy \quad (38)$$

Changing the integration to polar (i.e., over θ and r) works as follows. We now integrate over a tiny part of a circle ($d\theta$ and dr). The length of the corresponding arc is equal to $S = d\theta r$, where θ is expressed in radians. This follows directly from the fact that full circle circumference is $2\pi r$, where 2π is the ‘angle’ of a full circle expressed in radians. We approximate the area between the arc and dr as a rectangle, giving that the tiny parts we integrate with (corresponding to $dxdy$ in Cartesian plane) are of size $rd\theta dr$. Thus, the integral becomes:

$$I^2 = \int_0^\infty \int_0^{2\pi} \exp(-\frac{1}{2\sigma^2}r^2)rdrd\theta \quad (39)$$

Since one of the integrands is expressed in radians we integrate from 0 to 2π , and the radius is nonnegative.

$$I^2 = \int_0^{2\pi} d\theta \int_0^\infty \exp(-\frac{1}{2\sigma^2}r^2)rdr \quad (40)$$

$$= 2\pi(-\sigma^2 \exp(-\frac{1}{2\sigma^2}r^2)) \Big|_0^\infty \quad (41)$$

$$= 2\pi(-\sigma^2 \exp(-\frac{1}{2\sigma^2}\infty^2)) - 2\pi(-\sigma^2 \exp(-\frac{1}{2\sigma^2}0^2)) \quad (42)$$

$$= 0 + 2\pi\sigma^2 \quad (43)$$

So

$$I = \sqrt{2\pi\sigma^2} \quad (44)$$

which exactly cancels the first factor in the Gaussian equation, so the result is 1.

1.8

$$\mathbb{E}[x] = \int \frac{1}{\sqrt{2\pi\sigma^2}} \exp[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2]xdx \quad (45)$$

Using $y = \frac{x-\mu}{\sigma}$ we get $\frac{dy}{dx} = \frac{1}{\sigma} \Rightarrow dx = dy\sigma$ and $x = y\sigma + \mu$.

$$\mathbb{E}[x] = \int \frac{1}{\sqrt{2\pi\sigma}}(y\sigma + \mu) \exp[-\frac{1}{2}y^2]dy \quad (46)$$

$$= \int \frac{1}{\sqrt{2\pi}}y\sigma \exp[-\frac{1}{2}y^2] + \frac{1}{\sqrt{2\pi}}\mu \exp[-\frac{1}{2}y^2]dy \quad (47)$$

$$= -\frac{\sigma}{\sqrt{2\pi}} \exp[-\frac{1}{2}y^2] \Big|_{-\infty}^\infty + \mu \quad (48)$$

$$= 0 - 0 + \mu \quad (49)$$

Next,

$$\frac{1}{\sqrt{2\pi}\sigma} \int \exp\left[-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right] dx = 1 \quad (50)$$

$$\Rightarrow \int \exp\left[-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right] dx = \sqrt{2\pi}\sigma^2 \quad (51)$$

$$(52)$$

Now differentiating both sides.

$$\int \frac{1}{2} (x - \mu)^2 \sigma^{-4} \exp\left[-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right] dx = \frac{1}{2} (2\pi\sigma^2)^{-1/2} 2\pi \quad (53)$$

$$= \frac{1}{2} (2\pi)^{1/2} \sigma^{-1} \quad (54)$$

$$\sigma^{-4} \int (x - \mu)^2 \exp\left[-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right] dx = (2\pi)^{1/2} \sigma^{-1} \quad (55)$$

$$\Rightarrow \sigma^{-1} (2\pi)^{-1/2} \int (x - \mu)^2 \exp\left[-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right] dx = \sigma^2 \quad (56)$$

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

$$\mathbb{E}[x^2] - 2\mu \mathbb{E}[x] + \mu^2 = \sigma^2 \quad (58)$$

$$\Rightarrow \mathbb{E}[x^2] - 2\mu^2 + \mu^2 = \sigma^2 \quad (59)$$

$$\Rightarrow \mathbb{E}[x^2] = \sigma^2 + \mu^2 \quad (60)$$

The last equality is easily found using the previous results.

1.9

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2} (x - \mu)^2\right\} \quad (61)$$

$$\frac{\partial \mathcal{N}(x|\mu, \sigma^2)}{\partial \mu} = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2} (x - \mu)^2\right\} \sigma^2 (x - \mu) = 0 \quad (62)$$

$$x - \mu = 0 \quad (63)$$

$$x = \mu \quad (64)$$

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^2} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right] \quad (65)$$

$$\frac{\partial \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}} = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^2} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right] \left[-\frac{1}{2} (\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}^{-1T}) (\mathbf{x} - \boldsymbol{\mu})\right] = 0 \quad (66)$$

$$(67)$$

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

1.10

$$\mathbb{E}[x + z] = \int_{x+z} p(x + z) [x + z] \quad (68)$$

$$= \int_{x+z} p(x + z)x + \int_{x+z} p(x + z)z \quad (69)$$

$$= \int_{x,z} p(x, z)x + \int_{x,z} p(x, z)z \quad (70)$$

$$= \int_{x,z} p(x)p(z)x + \int_{x,z} p(x)p(z)z \quad (71)$$

$$= \int_x p(x)x + \int_z p(z)z \quad (72)$$

$$= \mathbb{E}[x] + \mathbb{E}[z] \quad (73)$$

$$\text{var}[x + z] = \mathbb{E} [\{(x + z) - \mathbb{E}[x + z]\}^2] \quad (74)$$

$$= \mathbb{E} [(x + z)^2] - \mathbb{E} [x + z]^2 \quad (75)$$

$$= \mathbb{E} [x^2 + z^2 + 2xz] - (\mathbb{E}[x] + \mathbb{E}[z])^2 \quad (76)$$

$$= \mathbb{E}[x^2] + \mathbb{E}[z^2] + \mathbb{E}[2xz] - \mathbb{E}[x]^2 - \mathbb{E}[z]^2 - 2\mathbb{E}[x]\mathbb{E}[z] \quad (77)$$

$$= \text{var}[x] + \text{var}[y] + 2 \int_{x,z} p(x, z)xz - 2 \int_x p(x) \int_z p(z) \quad (78)$$

$$= \text{var}[x] + \text{var}[y] + 2 \int_{x,z} p(x)p(z)xz - 2 \int_x p(x) \int_z p(z) \quad (79)$$

$$= \text{var}[x] + \text{var}[y] + 2 \int_x p(x) \int_z p(z) - 2 \int_x p(x) \int_z p(z) \quad (80)$$

1.11

Mean:

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln 2\pi \quad (81)$$

$$\frac{\partial \ln p(\mathbf{x}|\mu, \sigma^2)}{\partial \mu} = \sigma^2 \sum_{n=1}^N (x_n - \mu) = 0 \quad (82)$$

$$\sum_{n=1}^N x_n = N\mu \quad (83)$$

$$\mu = \frac{1}{N} \sum_{n=1}^N x_n \quad (84)$$

Variance:

$$-\frac{1}{2\sigma^2} = -(2\sigma^2)^{-1} \quad (85)$$

$$\frac{\partial \ln p(\mathbf{x}|\mu, \sigma^2)}{\partial \sigma^2} = 2 \cdot \frac{1}{4\sigma^4} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2\sigma^2} = 0 \quad (86)$$

$$N \cdot 2\sigma^4 = 2\sigma^2 \sum_{n=1}^N (x_n - \mu)^2 \quad (87)$$

$$\sigma^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2 \quad (88)$$

1.12

If $x_n = x_m$ we have $\mathbb{E}[x_n^2] = \mu^2 + \sigma^2$, otherwise we have $\mathbb{E}[x_n x_m] = \mathbb{E}[x_n] \mathbb{E}[x_m] = \mu^2$. So we only get σ^2 if $n = m$, which leads to the given equation.

$$\mathbb{E}[\mu_{ML}] = \frac{1}{N} \mathbb{E}\left[\sum_{n=1}^N x_n\right] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[x_n] = \frac{1}{N} N\mu \quad (89)$$

$$\mathbb{E}[\sigma_{ML}^2] = \frac{1}{N} \mathbb{E}\left[\sum_{n=1}^N (x_n - \mu_{ML})^2\right] \quad (90)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbb{E}\left[x_n^2 - 2x_n \frac{1}{N} \sum_{m=1}^N x_m + \frac{1}{N^2} \sum_{l=1}^N \sum_{k=1}^N x_l x_k\right] \quad (91)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbb{E}[\mu^2 + \sigma^2 - 2\frac{1}{N} \sum_{m=1}^N x_m x_n + \frac{1}{N^2} \sum_{l=1}^N \sum_{k=1}^N x_l x_k] \quad (92)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbb{E}[\mu^2 + \sigma^2] - 2\frac{1}{N} \sum_{m=1}^N \mathbb{E}[x_m x_n] + \frac{1}{N^2} \sum_{l=1}^N \sum_{k=1}^N \mathbb{E}[x_l x_k] \quad (93)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbb{E}[\mu^2 + \sigma^2] - 2\frac{1}{N} \sum_{m=1}^N (\mu^2 + I_{nm}\sigma^2) + \frac{1}{N^2} \sum_{l=1}^N \sum_{k=1}^N (\mu^2 + I_{kl}\sigma^2) \quad (94)$$

$$= \mu^2 + \sigma^2 - 2\mu^2 - 2\frac{1}{N}\sigma^2 + \mu^2 + \frac{1}{N}\sigma^2 \quad (95)$$

$$= \frac{N}{N}\sigma^2 - \frac{1}{N}\sigma^2 \quad (96)$$

1.13

From 1.56:

$$\mathbb{E}\left[\frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2\right] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[x^2] - \mathbb{E}[2x_n \mu] + \mathbb{E}[\mu^2] \quad (97)$$

$$= \mathbb{E}[x^2] - 2\mu \mathbb{E}[x] + \mu^2 \quad (98)$$

$$= \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2 \quad (99)$$

1.14

We can write any matrix W as a sum of a symmetric matrix S where the off-diagonals $S_{ij} = \frac{W_{ij}+W_{ji}}{2}$ and the diagonals $S_{ii} = W_{ii}$. The antisymmetric matrix diagonal will be $\text{Diag}(0)$. The off-diagonal entries will be $W_{ij} - S_{ij}$. Thus we can have

$$x^T W x = x^T (S + A) x = x^T S x + x^T A x = x^T S x + \sum_{i=1}^D \sum_{j=1}^D x_i x_j w_{ij} \quad (100)$$

The second term will be, because of the antisymmetry, 0. This means that only the symmetric part contributes.

1.17

$$\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du \quad (101)$$

$$\Rightarrow \Gamma(x+1) = \int_0^\infty u^x e^{-u} du \quad (102)$$

Using $\int f(x)g'(x)dx = f(x)g(x) - \int f'(x)g(x)dx$ and using $f(x) = u^x, f'(x) = xu^{x-1}, g'(x) = e^{-u}, g(x) = -e^{-u}$ we get

$$\int_0^\infty u^x e^{-u} = \left(u^x e^{-u} + \int_0^\infty xu^{x-1} e^{-u} \right) \Big|_0^\infty \quad (103)$$

$$= \left(u^x e^{-u} + x \int_0^\infty u^{x-1} e^{-u} du \right) \Big|_0^\infty \quad (104)$$

$$= \left(u^x e^{-u} + x \int_0^\infty u^{x-1} e^{-u} du \right) \Big|_0^\infty \quad (105)$$

$$= \left(u^x e^{-u} + x\Gamma(x) \right) \Big|_0^\infty \quad (106)$$

$$= u^x e^{-u} \Big|_0^\infty + x\Gamma(x) \quad (107)$$

$$(108)$$

Where the second term is no function of u anymore, so we don't need to evaluate it for u . Since $u^x e^{-u} = \frac{u^x}{e^u} = \frac{f(x)}{g(x)}$ with the indefinite form as we take the limit $\frac{\infty}{\infty}$ we can use l'Hopital.

Using l'Hopital we see that $\lim_{u \rightarrow \infty} \frac{u^x}{e^u} = \lim_{u \rightarrow \infty} \frac{(x-1)u^{x-1}}{e^u}$

If we keep doing this recursively for any integer x we will obtain $\lim_{u \rightarrow \infty} \frac{x!}{e^u}$ which clearly converges to zero. However, even for a (negative) non integer y there will be a larger x , meaning a larger numerator, for which the limit still converges to zero, showing that the result still holds for y .

1.18

The right integral simplifies:

$$\int_0^\infty e^{-r^2} r^{D-1} dr = \frac{1}{2} \int_0^\infty e^{-u} \sqrt{D-1} \sqrt{u}^{-1} du \quad u = r^2 \Rightarrow \frac{du}{dr} = 2r = 2\sqrt{u} \Rightarrow dr = \frac{du}{2\sqrt{u}} \quad (109)$$

$$= \frac{1}{2} \int_0^\infty e^{-u} \sqrt{u}^{D-1} (\sqrt{u})^{-1} du \quad (110)$$

$$= \frac{1}{2} \int_0^\infty e^{-u} (\sqrt{u})^{D-2} du \quad (111)$$

$$= \frac{1}{2} \int_0^\infty e^{-u} u^{1/2 D - 1} du \quad (112)$$

$$= \frac{1}{2} \Gamma\left(\frac{1}{2} D\right) \quad (113)$$

Using $\sigma = 1$ we see from 1.126 that $I^D = \prod_{i=1}^D I = (2\pi)^{D/2}$. Filling in these two results gives the answer.

The next part is a bit trickier. We need to integrate over the radius. The total volume will be a sum of 'slices' (spherical shells). E.g., the volume of an onion is the sum of the volume of its rings. For example, the volume of a sphere is $4/3\pi r^3$. What happens to the area if we increase r a bit? We get that $dV = \frac{4}{3}(r+dr)^3 - \frac{4}{3}\pi r^3$. Working this out gives $dV = 4r^2 dr + r dr^2 + dr^3$. We see that the second and third term become arbitrarily small as $dr \rightarrow 0$. Therefore, dV can be well approximated by $4r^2 dr$, using the formula for its surface area. Integrating dV will give V . This generalizes: $\int_0^1 S_D r^{D-1} dr$. Evaluating this integral gives the desired result.

1.19

$$\frac{S_D}{2^D D} = \frac{2\pi^{D/2}}{\Gamma(D/2)} \frac{1}{2^D D} = \frac{2\pi^{D/2}}{\Gamma(D/2) 2^D D} = \frac{2^{1-D} \pi^{D/2}}{\Gamma(D/2) D} = \frac{\pi^{D/2}}{\Gamma(D/2) 2^D D} \quad (114)$$

If we take $\Gamma(D/2) = \Gamma(D/2 - 1 + 1)$ and apply Stirling we get $\Gamma(D/2) = 2\pi^{1/2} e^{-D/2-1} (D/2 - 1)^{D/2+1/2}$. Plugging this into the fraction and plotting shows that this function goes to infinity. An ambitious person could apply l'Hopital.

For the final problem: start with a square ($D=2$). The distance to the corner using Pythagoras will be $\sqrt{a^2 + a^2} = \sqrt{2a^2} = \sqrt{Da^2}$. If we add a dimension the new distance to the corner will be applying pythagoras again, using the previous result. $\sqrt{\sqrt{a^2 + a^2}^2 + a^2} = \sqrt{3a^2}$. Using induction you could show that this generalizes to $\sqrt{Da^2}$. Taking the square root and dividing over a gives the result.

1.20

1.21

$$a \leq b \quad (115)$$

$$\Rightarrow \sqrt{a} \leq \sqrt{b} \quad (\text{Nonnegativity}) \quad (116)$$

$$\Rightarrow a \leq \sqrt{a}\sqrt{b} \quad (117)$$

We know that

$$\int_{R_1} p(x, C_2) \leq \left(\int_{R_1} p(x, C_1) dx \int_{R_1} p(x, C_2) dx \right)^{1/2} \quad (118)$$

Thus, by symmetry we get

$$P(\text{mistake}) = \left(\int_{R_1} p(x, C_1) dx \int_{R_1} p(x, C_2) dx \right)^{1/2} + \left(\int_{R_2} p(x, C_1) dx \int_{R_2} p(x, C_2) dx \right)^{1/2} \quad (119)$$

I do not see how this reduces to the required solution.

1.22

$$1 - I_{kj} = \begin{bmatrix} 0 & 1 & \dots & 1 \\ 1 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 0 \end{bmatrix} \quad (120)$$

$$\mathcal{R}_j^* = \arg \min_j L_{kj} p(C_k | \mathbf{x}) \quad (121)$$

$$= \arg \min_j \sum_{j \neq k} L_{kj} p(C_k | \mathbf{x}) \quad (122)$$

$$= \arg \max_j p(C_k | \mathbf{x}) \quad (123)$$

1.23

1.24

If the probability (output of our model) of vector x belonging to C_k times the loss we will occur if we select that class is smaller than λ we take the bet.

1.25

$$\mathbb{E}[L(\mathbf{t}, \mathbf{y}(\mathbf{x}))] = \int \int \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t} \quad (124)$$

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

$$= \int \int \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{t} d\mathbf{x} \quad (125)$$

$$F[\mathbf{x}, \mathbf{y}, \mathbf{t}] = \int_{\mathbf{t}} \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) \quad (126)$$

Euler-Lagrange:

$$\frac{\partial F}{\partial \mathbf{y}} - \frac{d}{d\mathbf{x}} \left(\frac{\partial F}{\partial \mathbf{y}'} \right) = 0 \quad (127)$$

$$\frac{\partial F}{\partial \mathbf{y}'} = 0 \quad (128)$$

$$\frac{\partial F}{\partial \mathbf{y}} = \int_t (\mathbf{y}(\mathbf{x}) - t)p(\mathbf{x}, t) = 0 \quad (129)$$

$$\int_t \mathbf{y}(\mathbf{x})p(\mathbf{x}, t) - \int_t tp(\mathbf{x}, t) = 0 \quad (130)$$

$$\mathbf{y}(\mathbf{x}) \int_t p(t, \mathbf{x}) = \int_t tp(\mathbf{x}, t) \quad (131)$$

$$\mathbf{y}(\mathbf{x}) = \frac{1}{p(\mathbf{x})} \int_t tp(\mathbf{x}, t) = \mathbb{E}_t[t|\mathbf{x}] \quad (132)$$

$$(133)$$

1.26

First, let's derive the single value case.

$$\{y(\mathbf{x}) - t\}^2 = \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \quad (134)$$

$$= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + \{\mathbb{E}[t|\mathbf{x}] - t\}^2 + 2[y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]]\{\mathbb{E}[t|\mathbf{x}] - t\} \quad (135)$$

$$\mathbb{E}[L] = \underbrace{\mathbb{E}_{x,t}[(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2]}_{(1)} + \underbrace{\mathbb{E}_{x,t}[(\mathbb{E}[t|\mathbf{x}] - t)^2]}_{(2)} + \underbrace{\mathbb{E}_{x,t}[2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) (\mathbb{E}[t|\mathbf{x}] - t)]}_{(3)} \quad (136)$$

(1)

$$\begin{aligned} \mathbb{E}_{x,t}[(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2] &= \mathbb{E}_x \mathbb{E}_t[(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2 | \mathbf{x}] \\ &= \mathbb{E}_x (y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2 \end{aligned} \quad (137)$$

(Since the function is over x , the expectation over t vanishes)
(138)

(2)

$$\begin{aligned} \mathbb{E}_{x,t}[(\mathbb{E}[t|\mathbf{x}] - t)^2] &= \mathbb{E}_x [\mathbb{E}_t [\mathbb{E}[t|\mathbf{x}]^2 + t^2 - 2\mathbb{E}[t|\mathbf{x}]t | \mathbf{x}]] \\ &= \mathbb{E}_x [\mathbb{E}[t|\mathbf{x}]^2 + \mathbb{E}[t^2|\mathbf{x}] - 2\mathbb{E}_t[t|\mathbf{x}]\mathbb{E}_t[t|\mathbf{x}]] \end{aligned} \quad (139)$$

($\mathbb{E}[t|\mathbf{x}]$ can be taken out of the inner expectation since t is integrated out.)

$$= \mathbb{E}_x [\text{var}[t|\mathbf{x}]] \quad (140)$$

(3)

$$\mathbb{E}_{x,t}[2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) (\mathbb{E}[t|\mathbf{x}] - t)] = \mathbb{E}_x \mathbb{E}_t[2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) (\mathbb{E}[t|\mathbf{x}] - t) | \mathbf{x}] \quad (141)$$

$$= \mathbb{E}_x [2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) \mathbb{E}_t[(\mathbb{E}[t|\mathbf{x}] - t) | \mathbf{x}]] \quad (142)$$

$$= \mathbb{E}_x [2(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]) \cdot 0] \quad (143)$$

$$= 0 \quad (144)$$

Now we continue for the vector case.

$$\|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 = \|\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}\|^2 \quad (145)$$

$$= (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^T (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}) \quad (146)$$

$$= (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^2 + (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^2 + 2 [(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})] \quad (147)$$

$$\mathbb{E}[L(\mathbf{t}, \mathbf{y}(\mathbf{x}))] = \mathbb{E}_{x,t} [\underbrace{(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^2}_{(1)} + \underbrace{(\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^2}_{(2)} + 2 \underbrace{[(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})]}_{(3)}] \quad (148)$$

(1)

$$\mathbb{E}_{x,t} [(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^2] = \mathbb{E}_x \mathbb{E}_t [(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^2 | \mathbf{x}] \quad (149)$$

$$= \mathbb{E}_x (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^2 \quad (150)$$

(2)

$$\mathbb{E}_{x,t} [(\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^2] = \mathbb{E}_x \mathbb{E}_t [(\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})^2 | \mathbf{x}] \quad (151)$$

$$= \mathbb{E}_x \mathbb{E}_t [\mathbb{E}[\mathbf{t}|\mathbf{x}]^T \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbf{t}^T \mathbf{t} - 2 \mathbb{E}[\mathbf{t}|\mathbf{x}]^T \mathbf{t} | \mathbf{x}] \quad (152)$$

$$= \mathbb{E}_x [\mathbb{E}[\mathbf{t}|\mathbf{x}]^T \mathbb{E}[\mathbf{t}|\mathbf{x}] + \mathbb{E}_t [\mathbf{t}^T \mathbf{t}] - 2 \mathbb{E}[\mathbf{t}|\mathbf{x}]^T \mathbb{E}[\mathbf{t}|\mathbf{x}]] \quad (153)$$

$$= \mathbb{E}_x [\text{var}[\mathbf{t}|\mathbf{x}]] \quad (154)$$

(3)

$$\mathbb{E}_{x,t} \{2 [(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})]\} = \mathbb{E}_x \mathbb{E}_t \{2 [(\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T (\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t})] | \mathbf{x}\} \quad (155)$$

$$= \mathbb{E}_x 2 (\mathbf{y}(\mathbf{x}) - \mathbb{E}[\mathbf{t}|\mathbf{x}])^T \mathbb{E}_t [\mathbb{E}[\mathbf{t}|\mathbf{x}] - \mathbf{t}] \quad (156)$$

$$= \mathbf{0} \quad (157)$$

You observe from the first term that the function that minimizes this expression is $\mathbf{y}(\mathbf{x}) = \mathbb{E}[\mathbf{t}|\mathbf{x}]$

1.27

Using a graphing calculator one can show that derivatives of absolute value polynomials require an additional **sign** component to account for the absolute parts.

Therefore

$$\frac{d}{dy(\mathbf{x})} = q \int [y(\mathbf{x} - t)]^{q-1} \text{sign}(y(\mathbf{x}) - t) p(\mathbf{x}, t) dt = 0 \quad (158)$$

We see that the **sign** expression is negative for $t > y(\mathbf{x})$ and vice versa. Therefore we can split the integral like so:

$$q \int [y(\mathbf{x} - t)]^{q-1} \text{sign}(y(\mathbf{x}) - t) p(\mathbf{x}, t) dt = q \int_{-\infty}^{y(\mathbf{x})} [y(\mathbf{x} - t)]^{q-1} \text{sign}(y(\mathbf{x}) - t) p(\mathbf{x}, t) dt \quad (159)$$

$$+ q \int_{y(\mathbf{x})}^{\infty} [y(\mathbf{x} - t)]^{q-1} \text{sign}(y(\mathbf{x}) - t) p(\mathbf{x}, t) dt \quad (160)$$

$$(161)$$

The sign of the second term will always be negative and the first term positive:

$$q \int [y(\mathbf{x} - t)]^{q-1} \text{sign}(y(\mathbf{x}) - t) p(\mathbf{x}, t) = q \int_{-\text{inf}ty}^{y(\mathbf{x})} [y(\mathbf{x} - t)]^{q-1} (y(\mathbf{x}) - t) p(\mathbf{x}, t) \quad (162)$$

$$- q \int_{y(\mathbf{x})}^{\infty} [y(\mathbf{x} - t)]^{q-1} (y(\mathbf{x}) - t) p(\mathbf{x}, t) \quad (163)$$

$$= 0 \Rightarrow \int_{y(\mathbf{x})}^{\infty} |y(\mathbf{x}) - t|^{q-1} p(t|x) dt = \int_{y(\mathbf{x})}^{\infty} |y(\mathbf{x}) - t|^{q-1} p(t|x) dt \quad (164)$$

In the case of $q = 1$ this reduces to an equality of two integrals over the pdf. Since the total area is 1 it must be that the both terms are 0.5, which shows that $y(\mathbf{x})$ is the median of the distribution.

1.28

$$n = k = 2$$

$$h(p^n) = -\log(p(x)^2) = -2\log p(x) = 2h(p) \quad (165)$$

Induction:

$$h(p^{k+1}) = -\log p(x)^{k+1} = -(k+1)\log(p(x)) = -(k+1)\log p(x) = (k+1)h(p) \quad (166)$$

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

$$H[x] = -\sum_{n=1}^M p(x) \ln p(x) \quad (167)$$

$$= \sum_{n=1}^M p(x) \ln \frac{1}{p(x)} \quad (168)$$

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

$$= \ln M \quad (169)$$

1.30

$$KL[p||z] = \int p(x) \ln \frac{q(x)}{p(x)} dx \quad (170)$$

$$\ln \frac{q(x)}{p(x)} = \ln q(x) - \ln p(x) \quad (171)$$

$$= \ln \left(\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(x-\mu)^2\right] \right) - \ln \left(\frac{1}{\sqrt{2\pi}s^2} \exp\left[-\frac{1}{2s^2}(x-m)^2\right] \right) \quad (172)$$

$$= -\ln(\sqrt{2\pi}) - \ln \sigma + \left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) + \ln(\sqrt{2\pi}) + \ln s - \left(-\frac{1}{2s^2}(x-m)^2\right) \quad (173)$$

$$= \ln \frac{s}{\sigma} - \frac{1}{2\sigma^2}(x^2 - 2x\mu + \mu^2) + \frac{1}{2s^2}(x^2 - 2xm + m^2) \quad (174)$$

$$= x^2 \left(\frac{1}{2s^2} - \frac{1}{2\sigma^2} \right) + x \left(\frac{\mu}{\sigma^2} - \frac{m}{s^2} \right) + \left(\ln \frac{s}{\sigma} - \frac{\mu^2}{2\sigma^2} + \frac{m^2}{2s^2} \right) \quad (175)$$

The integral now becomes

$$\left(\frac{1}{2s^2} - \frac{1}{2\sigma^2} \right) \int p(x)x^2 dx + \left(\frac{\mu}{\sigma^2} - \frac{m}{s^2} \right) \int xp(x)dx + \left(\ln \frac{s}{\sigma} - \frac{\mu^2}{2\sigma^2} + \frac{m^2}{2s^2} \right) \int p(x)dx \quad (176)$$

$$= \left(\frac{1}{2s^2} - \frac{1}{2\sigma^2} \right) (\mu^2 + \sigma^2) + \left(\frac{\mu}{\sigma^2} - \frac{m}{s^2} \right) \mu + \ln \frac{s}{\sigma} - \frac{\mu^2}{2\sigma^2} + \frac{m^2}{2s^2} \quad (177)$$

Which reduces to

$$\ln \frac{s}{\sigma} + \frac{\sigma^2 + (\mu - m)^2}{2s^2} - \frac{1}{2} \quad (178)$$

1.31

$$H[x, y] = - \int \int p(x, y) \log p(x, y) dx dy \quad (179)$$

$$= - \int \int p(x, y) \log p(x) p(y|x) dx dy \quad (180)$$

$$= - \int \int p(x, y) \log p(x) dx dy - \int \int p(x, y) \log p(y|x) dx dy \quad (181)$$

$$= - \int \log p(x) \int p(x, y) dy dx + H[y|x] \quad (182)$$

$$= - \int \log p(x) p(x) + H[y|x] = H[x] + H[y|x] \quad (183)$$

$$H[x] + H[y] - H[x, y] = H[x] + H[y] - H[x] - H[y|x] \quad (184)$$

$$= H[y] - H[y|x] = I(x; y) \geq 0 \quad (185)$$

1.32

$$H[\mathbf{y}] = - \int p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{y} \quad (186)$$

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

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

$$= - \int p(\mathbf{x}) \ln \frac{p(\mathbf{x})}{|\mathbf{A}|} d\mathbf{x} \quad (187)$$

$$= - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} + \int p(\mathbf{x}) \ln |\mathbf{A}| d\mathbf{x} \quad (188)$$

$$= H[\mathbf{x}] + \ln |\mathbf{A}| \int p(\mathbf{x}) d\mathbf{x} \quad (189)$$

$$= H[\mathbf{x}] + \ln |\mathbf{A}| \quad (190)$$

1.33

$$H[y|x] = 0 \quad (191)$$

$$= - \int \int p(x, y) \log p(y|x) dy dx \quad (192)$$

$$= - \int \int p(x, y) \log p(x, y) dy dx + \int \int p(x, y) \log p(x) dx dy \quad (193)$$

$$\Rightarrow H(x) = H(y, x) \quad (194)$$

From this we see that the information gained when we obtain x is equal to the information we gain when we see both x and y . Therefore y must be a function of x . We can also directly see it from the definition. If we use the product rule $p(x, y) = p(y|x)p(x)$ and $\lim_{p(y|x) \rightarrow 0} \ln p(y|x) = 0$ we see that the equality will only hold if $p(y|x) = 0$ or $p(y|x) = 1$ if $p(x) > 1$ which is given in the exercise. This, together with the requirement that $p(y|x)$ should integrate to 1 means that the only solution is there's 1 $p(y|x) = 1$ and all the others are $p(y|x) = 0$.

[1.34](#)

1.35

$$-H[x] = \int p(x) \ln p(x) dx \quad (195)$$

$$= \int p(x) \ln \frac{1}{\sqrt{2\pi\sigma^2}} dx - \int p(x) \frac{(x-\mu)^2}{2\sigma^2} \quad (196)$$

$$= \int p(x) \ln 1 dx - \int p(x) \ln \sqrt{2\pi\sigma^2} dx - \int p(x) \frac{(x-\mu)^2}{2\sigma^2} \quad (197)$$

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

$$= -\frac{1}{2} \ln 2\pi\sigma^2 - \frac{\sigma^2}{2\sigma^2} \quad (198)$$

$$= -\frac{1}{2} \ln 2\pi\sigma^2 - \frac{1}{2} \quad (199)$$

$$H[x] = \frac{1}{2} [\ln 2\pi\sigma^2 + 1] \quad (200)$$

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

$$-H[x, y] = \int \int p(y, x) \ln p(x, y) dx dy \quad (201)$$

$$= \int \int p(y, x) \ln [p(y|x)p(x)] dx dy \quad (202)$$

$$= \int \int p(y, x) \ln p(y|x) + p(y, x) \ln p(x) dx dy \quad (203)$$

$$= -H[y|x] + \int \int p(y, x) dy \ln p(x) dx \quad (204)$$

$$= -H[y|x] - H[x] \quad (205)$$

1.38

We know from 1.114 that $k = 2$ is true:

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i) \quad (206)$$

With $\{x_i\} = \{a, b\}$ and λ_1 must be $1 - \lambda$ due to the constraint. So $M = k = 2$, the bsse case, checks out. Now $M = k + 1$:

$$f\left(\sum_{i=1}^{k+1} \lambda_i x_i\right) = f\left(\sum_{i=1}^k \lambda_i x_i + \lambda_{k+1} x_{k+1}\right) \quad (207)$$

$$= f\left((1 - \lambda_{k+1}) \frac{1}{1 - \lambda_{k+1}} \sum_{i=1}^K \lambda_i x_i + \lambda_{k+1} x_{k+1}\right) \quad (208)$$

This is of the form $M = 2$ so we can apply the rule:

$$\leq (1 - \lambda_{k+1}) f\left(\frac{1}{1 - \lambda_{k+1}} \sum_{i=1}^K \lambda_i x_i\right) + \lambda_{k+1} x_{k+1} \quad (209)$$

$$= (1 - \lambda_{k+1}) f\left(\sum_{i=1}^K \frac{\lambda_i}{1 - \lambda_{k+1}} x_i\right) + \lambda_{k+1} x_{k+1} \quad (210)$$

Now we have to observe that when $\sum_{i=1}^{K+1} \lambda_i = \sum_{i=1}^K \lambda_i + \lambda_{k+1} = 1 \iff \sum_{i=1}^K \lambda_i = 1 - \lambda_{k+1}$ and then divide both sides by $1 - \lambda_{k+1}$ to see that the factor in the sum adds up to one and therefore we can apply the inequality again.

$$\leq (1 - \lambda_{k+1}) \sum_{i=1}^K \frac{\lambda_i}{1 - \lambda_{k+1}} f(x_i) + \lambda_{k+1} x_{k+1} \quad (211)$$

$$= \sum_{i=1}^{k+1} \lambda_i f(x_i) \quad (212)$$

1.40

$$\ln \prod_{i=1}^n a_i^{1/n} = \frac{1}{n} \ln \prod_{i=1}^n a_i \quad (213)$$

$$= \frac{1}{n} \sum_{i=1}^n \ln a_i \quad (214)$$

$$\geq \ln \frac{1}{n} \sum_{i=1}^n a_i \quad (\text{Jensen}) \quad (215)$$

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

$$-I[x, y] = \int \int \left[p(x, y) \ln p(x) + p(x, y) \ln p(y) - p(x, y) \ln p(x|y) - p(x, y) \ln p(y) \right] dx dy \quad (216)$$

$$= -H[x] + H[x|y] \quad (217)$$

2.1

$$p(x|\mu) = \mu^x(1-\mu)^{1-x} \quad (218)$$

$$\sum_{x=0}^1 p(x|\mu) = 1 - \mu + \mu = 1 \quad (219)$$

$$\mathbb{E}[x] = \sum_{x=0}^1 x\mu^x(1-\mu)^{1-x} \quad (220)$$

$$= \mu(1-\mu)^0 = \mu \quad (221)$$

$$\text{Var}[x] = [(x - \mathbb{E}[x])^2] \quad (222)$$

$$= [(x - \mu)^2] \quad (223)$$

$$= \sum_{i=0}^1 (x - \mu)^2 \mu^x(1-\mu)^{1-x} \quad (224)$$

$$(225)$$

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

$$H[x] = \sum i = 0^1 p(x) \ln p(x) = \mu \ln \mu + (1-\mu) \ln(1-\mu) \quad (226)$$

2.2

$$\sum_x p(x|\mu) = \sum_x \left(\frac{1-\mu}{2}\right)^{(1-x)/2} \left(\frac{1+\mu}{2}\right)^{(1+x)/2} \quad (227)$$

$$= \frac{1-\mu}{2} + \frac{1+\mu}{2} = 1 \quad (228)$$

where $\mathcal{X} = \{0, 1\}$

$$\mathbb{E}[x] = \sum_x xp(x) = \frac{-1+\mu}{2} + \frac{1+\mu}{2} = \mu \quad (229)$$

$$\mathbb{E}[(x - \mathbb{E}[x])^2] = \mathbb{E}[x^2] - 2\mathbb{E}[x\mathbb{E}[x]] + \mathbb{E}[\mathbb{E}[x]^2] \quad (230)$$

$$= 1 - 2\mu\mathbb{E}[x] + \mu^2 \quad (231)$$

$$= 1 - 2\mu^2 + \mu^2 = 1 - \mu^2 \quad (232)$$

$$H[x] = -\sum_x p(x) \ln p(x) = \frac{1+\mu}{2} \ln \frac{1+\mu}{2} + \frac{1-\mu}{2} \ln \frac{1-\mu}{2} \quad (233)$$

2.3

$$\frac{N!}{(N-m)!m!} + \frac{N!}{(N-(m-1))!(m-1)!} \quad (234)$$

$$= \frac{N!}{(N-m+1)!m!(N-m+1)^{-1}} + \frac{N!}{(N-m+1)!m!m^{-1}} \quad (235)$$

$$= \frac{N!(N-m+1)}{(N-m+1)!m!} + \frac{N!m}{(N-m+1)!m!} \quad (236)$$

$$= \frac{N!(N+1)}{(N-m+1)!m!} = \frac{(N+1)!}{(N-m+1)!m!} \quad (237)$$

Now the induction part. The general form is:

$$(x+y)^N = \sum_{m=0}^N \binom{N}{m} x^{N-m} y^m \quad (238)$$

Base case: $N = 0 \Rightarrow 1 = 1$ is trivially true.

$$(x+y)^{N+1} = \sum_{m=0}^{N+1} \binom{N+1}{m} x^{N+1-m} y^m \quad (239)$$

Continuing with the LHS.

$$(x+y)^{N+1} = (x+y) \sum_{m=0}^N \binom{N}{m} x^{N-m} y^m \quad \text{Using the base case} \quad (240)$$

$$= \sum_{m=0}^N \binom{N}{m} x^{N-m+1} y^m + \sum_{m=0}^N \binom{N}{m} x^{N-m} y^{m+1} \quad (241)$$

We can equalize the exponents by changing the sum range.

$$= \sum_{m=0}^N \binom{N}{m} x^{N-m+1} y^m + \sum_{m=1}^{N+1} \binom{N}{m-1} x^{N-(m-1)} y^{(m-1)+1} \quad (242)$$

Now let's take out the terms that they not have in common. I.e., $m = 0$ from the first term and $m = N+1$ from the second.

$$= \binom{N}{0} x^{N-0+1} y^0 + \binom{N}{N} x^{N-N} y^{N+1} + \sum_{m=1}^N \binom{N}{m} x^{N-m+1} y^m + \sum_{m=1}^N \binom{N}{m-1} x^{N-m+1} y^m \quad (243)$$

$$= x^{N+1} + y^{N+1} + \sum_{m=1}^N (x^{N-m+1} y^m) \left(\binom{N}{m} + \binom{N}{m-1} \right) \quad (244)$$

Now we can use our derived formula.

$$= x^{N+1} + y^{N+1} + \sum_{m=1}^N x^{N-m+1} y^m \binom{N+1}{m} \quad (245)$$

Now, see that $x^{N+1} = \binom{N+1}{m} x^{N-m+1} y^m \iff m = 0$ which can be plugged into the sum.

$$= y^{N+1} + \sum_{m=0}^N x^{N-m+1} y^m \binom{N+1}{m} \quad (246)$$

Similary, $y^{N+1} = x^{N-m+1} y^m \binom{N+1}{N+1} \iff m = N+1$

$$= \sum_{m=0}^{N+1} x^{N-m+1} y^m \binom{N+1}{m} \quad (247)$$

Which is equal to the RHS.

Now for the normalization. Using the binomial theorem the LHS is $(\mu + 1 - \mu)^N$

$$(\mu + 1 - \mu)^N = 1^N = 1 \quad (248)$$

2.4

$$\sum_{m=0}^N \binom{N}{m} \frac{m}{\mu} \mu^m (1-\mu)^{N-m} - \frac{N-m}{1-\mu} (1-\mu)^{N-m} \mu^m = 0 \quad (249)$$

$$\iff \sum_{m=0}^N \frac{m}{\mu} \binom{N}{m} \mu^m (1-\mu)^{N-m} = \sum_{m=0}^N \frac{N-m}{1-\mu} \binom{N}{m} (1-\mu)^{N-m} \mu^m \quad (250)$$

$$\iff \frac{1}{\mu} \mathbb{E}[m] = \frac{N}{1-\mu} \sum_{m=0}^N \binom{N}{m} (1-\mu)^{N-m} \mu^m - \frac{1}{1-\mu} \sum_{m=0}^N m \binom{N}{m} (1-\mu)^{N-m} \mu^m \quad (251)$$

$$= \frac{N}{1-\mu} \sum_{m=0}^N \text{Bin}(m|N, \mu) - \frac{1}{1-\mu} \mathbb{E}[m] \quad (252)$$

$$\iff \mathbb{E}[m] = \frac{\mu N}{1-\mu} \sum_{m=0}^N \text{Bin}(m|N, \mu) - \frac{\mu}{1-\mu} \mathbb{E}[m] \quad (253)$$

$$\iff 1 = \frac{\mu N}{\mathbb{E}[m](1-\mu)} \sum_{m=0}^N \text{Bin}(m|N, \mu) - \frac{\mu}{1-\mu} \quad (254)$$

$$\iff \frac{1}{1-\mu} = \frac{\mu N}{\mathbb{E}[m](1-\mu)} \sum_{m=0}^N \text{Bin}(m|N, \mu) \quad (255)$$

2.4

$$\sum_{m=0}^N \binom{N}{m} \frac{m}{\mu} \mu^m (1-\mu)^{N-m} - \frac{N-m}{1-\mu} (1-\mu)^{N-m} \mu^m = 0 \quad (256)$$

$$= \sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} \left(\frac{m}{\mu} - \frac{N-m}{1-\mu} \right) = 0 \quad (257)$$

$$= \sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} \left(\frac{(1-\mu)m}{(1-\mu)\mu} - \frac{\mu(N-m)}{\mu(1-\mu)} \right) = 0 \quad (258)$$

$$= \sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} (m - \mu N) = 0 \quad (259)$$

$$= \mathbb{E}[m] - \mu N \sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} = 0 \quad (260)$$

The sum in the second term is the integral of the binomial distribution (2.264 in the book), and therefore is equal to 1.

2.5

Moving the integral inside

$$\Gamma(a)\Gamma(b) = \int_0^\infty \int_0^\infty \exp(-x) x^{a-1} \exp(-y) y^{b-1} dx dy \quad (261)$$

And using $t = y + x \iff dy = dt$

$$= \int_0^\infty \int_{t=x=0}^\infty \exp(-x) x^{a-1} \exp(x-t) (t-x)^{b-1} dt dx \quad (262)$$

$$= \int_0^\infty \int_{t=x=0}^{x=t=\infty} x^{a-1} \exp(-t) (t-x)^{b-1} dt dx \quad (263)$$

Using the second CoV $x = t\mu \iff dx = d\mu t$

$$= \int_{t=0}^\infty \int_0^1 \exp(-t) t^{a-1} \mu^{a-1} (t-t\mu)^{b-1} t d\mu dt \quad (264)$$

$$= \int_{t=0}^\infty \int_0^1 \exp(-t) t^{a-1} \mu^{a-1} (1-\mu)^{b-1} t^{b-1} t d\mu dt \quad (265)$$

$$= \int_{t=0}^\infty \exp(-t) t^{a-1} t^{b-1} t dt \int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu \quad (266)$$

$$= \Gamma(a+b) \int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu \quad (267)$$

$$(268)$$

2.6

$$p(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \quad (269)$$

$$\int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \quad (270)$$

$$\Gamma(x) = \int_0^\infty \mu^{x-1} e^{-\mu} d\mu \quad (271)$$

$$\mathbb{E}[\mu] = \int_0^1 p(\mu|a, b) \mu d\mu \quad (272)$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \mu d\mu \quad (273)$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^a (1-\mu)^{b-1} d\mu \quad (274)$$

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

$$= \frac{a}{a+b} \int_0^1 \frac{\Gamma(a+1+b)}{\Gamma(a+1)\Gamma(b)} \mu^a (1-\mu)^{b-1} d\mu \quad (275)$$

$$= \frac{a}{a+b} \quad (276)$$

Where we used

$$\Gamma(a+1) = a\Gamma(a) \text{ and } \Gamma(a+b+1) = (a+b)\Gamma(a+b) \quad (277)$$

$$\text{Var}[\mu] = \mathbb{E}[\mu^2] - \mathbb{E}[\mu]^2 \quad (278)$$

$$\mathbb{E}[\mu] = \int_0^1 p(\mu|a, b) \mu d\mu^2 \quad (279)$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a+1} (1-\mu)^{b-1} \mu d\mu \quad (280)$$

Which looks like Beta with $a+2$

$$= \frac{a}{a+b} \int_0^1 \frac{\Gamma(a+1+b)}{\Gamma(a+1)\Gamma(b)} \mu^{a+2} (1-\mu)^{b-1} d\mu \quad (281)$$

$$= \frac{a(a+1)}{(a+b)(a+1+b)} \int_0^1 \frac{\Gamma(a+2+b)}{\Gamma(a+2)\Gamma(b)} \mu^{a+2} (1-\mu)^{b-1} d\mu \quad (282)$$

$$= \frac{a(a+1)}{(a+b)(a+1+b)} \quad (283)$$

$$\text{Var}[\mu] = \frac{a(a+1)}{(a+b)(a+1+b)} - \frac{a^2}{(a+b)^2} \quad (284)$$

$$= \frac{(a+b)a(a+1)}{(a+b)^2(a+1+b)} - \frac{(a+b+1)a^2}{(a+b+1)(a+b)^2} \quad (285)$$

$$= \frac{(a^2+ab)(a+1)}{(a+b)^2(a+1+b)} - \frac{(a+b+1)a^2}{(a+b+1)(a+b)^2} \quad (286)$$

$$= \frac{a^3+a^2b+a^2+ab}{(a+b)^2(a+1+b)} - \frac{a^3+ba^2+a^2}{(a+b+1)(a+b)^2} \quad (287)$$

Which gives you the desired result.

Mode is given where the derivative w.r.t. μ is zero:

$$\frac{\partial p(\mu|a,b)}{\partial \mu} = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} [(a-1)\mu^{a-2} - \mu^{a-1}(b-1)(1-\mu)^{b-2}] = 0 \quad (288)$$

$$(a-1)\mu^{a-2}(1-\mu)^{b-1} = \mu^{a-1}(b-1)(1-\mu)^{b-2} \quad (289)$$

$$\frac{a-1}{\mu} = \frac{b-1}{1-\mu} \quad (290)$$

$$\mu = \frac{-a+1}{-a-b+2} = \frac{a-1}{a+b-2} \quad (291)$$

2.7

The ML estimate is $\mu_{ML} = \frac{m}{N} = \frac{m}{N}$ and the prior mean is $\frac{a}{a+b}$. The posterior is $p(x=1|D) = \frac{m+a}{m+a+n+b}$.

Following the tip in the exercise.

$$\lambda\left(\frac{a}{a+b}\right) + (1-\lambda)\frac{m}{m+n} = \frac{m+a}{m+a+n+b} \quad (292)$$

Solving for lambda gives

$$\lambda = \frac{1}{1+(n+m)/(a+b)} \quad (293)$$

Which must lie in (0,1).

2.8

$$\mathbb{E}[x] = \int p(x)xdx \quad (294)$$

$$= \int \int p(x,y)xdxdy \quad (295)$$

$$= \int \int p(x|y)p(y)xdxdy \quad (296)$$

$$= \int p(y) \int p(x|y)xdxdy \quad (297)$$

$$= \int p(y) \mathbb{E}_x[x|y]dy \quad (298)$$

$$= \mathbb{E}_y \mathbb{E}_x[x|y] \quad (299)$$

$$\mathbb{E}_y[\text{Var}_x[x|y]] = \mathbb{E}_y \mathbb{E}_{x|y}[(x - \mathbb{E}[x|y])^2] \quad (300)$$

$$= \mathbb{E}_y \mathbb{E}_{x|y}[(x^2 + \mathbb{E}[x|y]^2 - 2x \mathbb{E}[x|y])] \quad (301)$$

$$= \mathbb{E}_y \mathbb{E}_{x|y} x^2 + \mathbb{E}_y \mathbb{E}_{x|y} \mathbb{E}[x|y]^2 - 2 \mathbb{E}_y \mathbb{E}_{x|y} x \mathbb{E}[x|y]] \quad (302)$$

$$= \mathbb{E}_y \int p(x|y)x^2dx + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \int p(y) \int p(x|y)x \mathbb{E}[x|y]dxdy \quad (303)$$

$$= \int p(y) \int p(x|y)x^2dxdy + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \int \mathbb{E}[x|y] \int p(x|y)xdxdy \quad (304)$$

$$= \int \int p(x,y)x^2dxdy + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \int \mathbb{E}[x|y]^2dy] \quad (305)$$

$$= \int \int p(x,y)x^2dxdy + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \mathbb{E}_y \mathbb{E}[x|y]^2] \quad (306)$$

$$= \mathbb{E}[x^2] - \mathbb{E}_y \mathbb{E}[x|y]^2 \quad (307)$$

$$(308)$$

$$\text{Var}_y[\mathbb{E}_x[x|y]] = \mathbb{E}_y [(\mathbb{E}_x[x|y] - \mathbb{E}_y \mathbb{E}_x[x|y])^2] \quad (309)$$

$$= \mathbb{E}_y [(\mathbb{E}_x[x|y] - \mathbb{E}_x[x])^2] \quad (310)$$

$$= \mathbb{E}_y [(\mathbb{E}_x[x|y]^2 + \mathbb{E}_x[x]^2 - 2 \mathbb{E}_x[x] \mathbb{E}_x[x|y])] \quad (311)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] + \mathbb{E}_y[\mathbb{E}_x[x]^2] - 2 \mathbb{E}_y \mathbb{E}_x[x] \mathbb{E}_x[x|y] \quad (312)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] + \mathbb{E}_x[x]^2 - 2 \mathbb{E}_x[x] \mathbb{E}_y \mathbb{E}_x[x|y] \quad (313)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] + \mathbb{E}_x[x]^2 - 2 \mathbb{E}_x[x]^2 \quad (314)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] - \mathbb{E}_x[x]^2 \quad (315)$$

$$(316)$$

Sum of these expressions gives desired result.

2.10

$$\mathbb{E}[\mu_j] = \int \mu_j \frac{\Gamma(a_0)}{\Gamma(a_1) \dots \Gamma(a_K)} \prod_{k=1}^K \mu_k^{a_k-1} d\boldsymbol{\mu} \quad (317)$$

$$= \frac{\Gamma(a_0)}{\Gamma(a_1) \dots \Gamma(a_K)} \int \prod_{k=1}^{j-1} \mu_k^{a_k-1} \mu_j^{a_j} \prod_{k=j+1}^K \mu_k^{a_k-1} d\boldsymbol{\mu} \quad (318)$$

Now we use $a_j := b_j - 1$ and $a_k := b_k$, then we get $b_0 = \sum b_k = \sum a_k + 1 = a_0 + 1$

$$= \frac{1/a_0 \Gamma(a_0 + 1)}{1/a_j \Gamma(a_1) \dots \Gamma(a_j + 1) \dots \Gamma(a_K)} \quad (319)$$

$$= \frac{a_j}{a_0} \frac{\Gamma(b_0)}{\Gamma(b_0) \dots \Gamma(b_K)} \int \prod_{k=1}^K \mu_k^{b_k-1} d\boldsymbol{\mu} \quad (320)$$

$$= \frac{a_j}{a_0} \int \text{Dir}(\boldsymbol{\mu}, \mathbf{b}) d\boldsymbol{\mu} \quad (321)$$

$$= \frac{a_j}{a_0} \quad (322)$$

$$\text{var}[\mu_j] = \mathbb{E}[(\mu_j - \mathbb{E}[\mu_j])^2] \quad (323)$$

$$= \mathbb{E}[\mu^2] - \mathbb{E}[\mu]^2 \quad (324)$$

$$= \mathbb{E}[\mu^2] - \left(\frac{a_j}{a_0}\right)^2 \quad (325)$$

$$\mathbb{E}[\mu^2] = \int \mu_j^2 \frac{\Gamma(a_0)}{\Gamma(a_1) \dots \Gamma(a_K)} \prod_{k=1}^K \mu_k^{a_k-1} d\boldsymbol{\mu} \quad (326)$$

$$= \frac{\Gamma(a_0)}{\Gamma(a_1) \dots \Gamma(a_K)} \int \prod_{k=1}^{j-1} \mu_k^{a_k-1} \mu_j^{a_j+1} \prod_{k=j+1}^K \mu_k^{a_k-1} d\boldsymbol{\mu} \quad (327)$$

We apply the same trick: $a_j + 1 = b_j - 1$ and $a_k = b_k$, thus $b_0 = \sum b_k = \sum a_k + 2 = a_0 + 2$

$$= \frac{1/a_0 \Gamma(a_0 + 1)}{1/a_j \Gamma(a_1) \dots \Gamma(a_j + 1) \dots \Gamma(a_K)} \int \prod_{k=1}^{j-1} \mu_k^{a_k-1} \mu_j^{a_j+1} \prod_{k=j+1}^K \mu_k^{a_k-1} d\boldsymbol{\mu} \quad (328)$$

$$= \frac{\frac{1}{a_0(a_0+1)} \Gamma(a_0 + 2)}{\frac{1}{a_j(a_j+1)} \Gamma(a_1) \dots \Gamma(a_j + 2) \dots \Gamma(a_K)} \int \prod_{k=1}^{j-1} \mu_k^{a_k-1} \mu_j^{a_j+1} \prod_{k=j+1}^K \mu_k^{a_k-1} d\boldsymbol{\mu} \quad (329)$$

$$= \frac{a_j(a_j + 1)}{a_0(a_0 + 1)} \int \text{Dir}(\boldsymbol{\mu}, \mathbf{b}) d\boldsymbol{\mu} \quad (330)$$

$$= \frac{a_j(a_j + 1)}{a_0(a_0 + 1)} \quad (331)$$

$$\text{Var}[\mu_j] = \frac{a_j(a_j + 1)}{a_0(a_0 + 1)} - \frac{a_j^2}{a_0^2} \quad (332)$$

$$= \frac{a_0^2 a_j(a_j + 1)}{a_0^2(a_0^2 + a_0)} - \frac{a_j^2(a_0^2 + a_0)}{a_0^2(a_0^2 + a_0)} \quad (333)$$

$$= \frac{a_0^2 a_j(a_j + 1) - a_j^2(a_0^2 + a_0)}{a_0^2(a_0^2 + a_0)} \quad (334)$$

$$= \frac{a_0^2 a_j^2 + a_0^2 a_j - a_j^2 a_0^2 - a_j^2 a_0}{a_0^2(a_0^2 + a_0)} \quad (335)$$

$$= \frac{a_0^2 a_j - a_j^2 a_0}{a_0^2(a_0^2 + a_0)} \quad (336)$$

$$= \frac{a_0 a_j(a_0 - a_j)}{a_0^2(a_0^2 + a_0)} \quad (337)$$

$$= \frac{a_j(a_0 - a_j)}{a_0(a_0^2 + a_0)} \quad (338)$$

$$(339)$$

2.11

2.12

$$\int_a^b \frac{1}{b-a} dx = \frac{x}{b-a} - \frac{x}{b-a} \Big|_a^b = \frac{b-a}{b-a} = 1 \quad (340)$$

$$\mathbb{E}[x] = \int_a^b p(x) x dx \quad (341)$$

$$= \int_a^b \frac{x}{b-a} dx \quad (342)$$

$$= \frac{x^2}{2(b-a)} \Big|_a^b \quad (343)$$

$$= \frac{b^2 - a^2}{2(b-a)} = \frac{(b-a)(b+a)}{2(b-a)} = \frac{b+a}{2} \quad (344)$$

$$\text{Var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 \quad (345)$$

$$= \int_a^b \frac{x^2}{b-a} dx - \frac{(b+a)^2}{4} \quad (346)$$

$$= \left. \frac{\frac{1}{3}x^3}{b-a} \right|_a^b - \frac{(b+a)^2}{4} \quad (347)$$

$$= \frac{\frac{1}{3}b^3}{b-a} - \frac{\frac{1}{3}a^3}{b-a} - \frac{(b+a)^2}{4} \quad (348)$$

$$= \frac{(b-a)(b^2+a^2+ab)}{3(b-a)} - \frac{(b+a)^2}{4} \quad (349)$$

$$= \frac{b^2+a^2+ab}{3} - \frac{b^2+a^2+2ab}{4} \quad (350)$$

$$= \frac{4b^2+4a^2+4ab}{12} - \frac{3b^2+3a^2+6ab}{12} \quad (351)$$

$$= \frac{b^2+a^2-2ab}{12} \quad (352)$$

$$= \frac{(b-a)^2}{12} \quad (353)$$

$$(354)$$

2.13

$$KL[p||q] = \int p(\mathbf{x}) \ln \frac{p(\mathbf{x})}{q(\mathbf{x})} \quad (355)$$

$$= \int p(\mathbf{x}) \ln \frac{|\Sigma|^{-1/2} \exp[-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)]}{|\mathbf{L}|^{-1/2} \exp[-\frac{1}{2}(\mathbf{x}-\mathbf{m})^T \mathbf{L}^{-1}(\mathbf{x}-\mathbf{m})]} \quad (356)$$

$$= \int p(\mathbf{x}) \ln \frac{|\mathbf{L}|^{1/2}}{|\Sigma|^{1/2}} + \frac{1}{2}(\mathbf{x}-\mathbf{m})^T \mathbf{L}^{-1}(\mathbf{x}-\mathbf{m}) - \frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu) \quad (357)$$

$$= \ln \frac{|\mathbf{L}|^{1/2}}{|\Sigma|^{1/2}} + \frac{1}{2} \int p(\mathbf{x})(\mathbf{x}-\mathbf{m})^T \mathbf{L}^{-1}(\mathbf{x}-\mathbf{m}) - \frac{1}{2} \int p(\mathbf{x})(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu) \quad (358)$$

$$= \ln \frac{|\mathbf{L}|^{1/2}}{|\Sigma|^{1/2}} + \frac{1}{2} \mathbb{E}[(\mathbf{x}-\mathbf{m})^T \mathbf{L}^{-1}(\mathbf{x}-\mathbf{m})] - \frac{1}{2} \mathbb{E}[(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)] \quad (359)$$

Now we use Matrix Cookook (15-11-2012) equation 318

$$= \ln \frac{|\mathbf{L}|^{1/2}}{|\Sigma|^{1/2}} + \frac{1}{2}(\text{Tr}(\mathbf{L}^{-1}\mathbf{L}) + \mathbb{E}[\mathbf{x}]^T \mathbf{L}^{-1} \mathbb{E}[\mathbf{x}]) - \frac{1}{2}(\text{Tr}(\Sigma^{-1}\Sigma) + \mathbb{E}[\mathbf{x}]^T \Sigma^{-1} \mathbb{E}[\mathbf{x}]) \quad (360)$$

$$= \ln \frac{|\mathbf{L}|^{1/2}}{|\Sigma|^{1/2}} + \mathbf{m}^T \mathbf{L}^{-1} \mathbf{m} + \mu^T \Sigma^{-1} \mu \quad (361)$$

2.14

I don't see where the Tr comes from in the solution.

2.15

$$H[\mathbf{x}] = - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \quad (362)$$

$$= - \int p(\mathbf{x}) \left[-\ln \left[(2\pi)^{D/2} \right] - \ln \left[|\Sigma|^{1/2} \right] - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (363)$$

$$(364)$$

$$\int p(\mathbf{x}) \ln \left[(2\pi)^{D/2} \right] = \frac{D}{2} \ln 2\pi \quad (365)$$

$$\int p(\mathbf{x}) \ln \left[|\Sigma|^{1/2} \right] = \frac{1}{2} \int p(\mathbf{x}) \ln |\Sigma| = \frac{1}{2} \ln |\Sigma| \quad (366)$$

$$\frac{1}{2} \int p(\mathbf{x}) (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \frac{1}{2} \mathbb{E} \left[(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (367)$$

$$= \frac{1}{2} \text{Tr}(\Sigma^{-1} \Sigma) + \frac{1}{2} \mathbb{E}[\mathbf{x} - \boldsymbol{\mu}]^T \Sigma^{-1} \mathbb{E}[\mathbf{x} - \boldsymbol{\mu}] \quad (368)$$

$$= \frac{1}{2} D + \frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu})^T \Sigma^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}) \quad (369)$$

$$= \frac{D}{2} \quad (370)$$

Adding the terms gives the desired result.

2.17

$$A = \frac{1}{2} \underbrace{A + A^T}_{\text{Symmetric}} + \frac{1}{2} \underbrace{A - A^T}_{\text{Non-symmetric}} \quad (371)$$

$$\mathcal{N}(\boldsymbol{\mu}, \Sigma) \propto \exp - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \frac{1}{2} (\Sigma^{-1} + \Sigma^{-T} + \Sigma^{-1} - \Sigma^{-T}) (\mathbf{x} - \boldsymbol{\mu}) \quad (372)$$

$$= \exp - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M (x_n - \mu_n) (\Sigma_{nm}^{-1} + \Sigma_{nm}^{-T}) (x_m - \mu_m) \quad (373)$$

$$- \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M (x_n - \mu_n) \underbrace{(\Sigma_{nm}^{-1} - \Sigma_{nm}^{-T})}_{=0} (x_m - \mu_m) \quad (374)$$

Which leaves us only with a symmetric covariance matrix.

2.19

From (2.45) we have

$$\mathbf{\Sigma}\mathbf{u}_i = \lambda_i\mathbf{u}_i \quad (375)$$

$$\Rightarrow \mathbf{u}_i^T \mathbf{\Sigma}\mathbf{u}_i = \mathbf{u}_i^T \lambda_i \mathbf{u}_i \quad (376)$$

$$\Rightarrow \mathbf{U}^T \mathbf{\Sigma}\mathbf{U} = \mathbf{U}^T \mathbf{\Lambda}\mathbf{U} \quad (377)$$

$$\Rightarrow \mathbf{\Sigma} = \mathbf{\Lambda} \quad (378)$$

Therefore we have

$$\mathbf{\Sigma} = \mathbf{\Lambda} = \mathbf{\Lambda}\mathbf{U}\mathbf{U}^T = \sum_i \lambda_i \mathbf{u}_i \mathbf{u}_i^T \quad (379)$$

$$\mathbf{\Sigma}^{-1} = (\mathbf{\Lambda}\mathbf{U}\mathbf{U}^T)^{-1} = \mathbf{U}\mathbf{U}^T \mathbf{\Lambda}^{-1} = \sum_i \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad (380)$$

2.20

Since the eigenvectors of $\mathbf{\Sigma}$ form an orthonormal set, they form a basis of \mathbb{R}^D , and therefore $\mathbf{a} = \sum_i a_i \mathbf{u}_i$

$$\mathbf{a}^T \mathbf{\Sigma}\mathbf{a} = \left(\sum_i a_i \mathbf{u}_i\right)^T \mathbf{\Sigma} \left(\sum_i a_i \mathbf{u}_i\right) \quad (381)$$

$$= \left(\sum_i a_i \mathbf{u}_i\right)^T \sum_i a_i \mathbf{\Sigma}\mathbf{u}_i \quad (382)$$

$$= \left(\sum_i a_i \mathbf{u}_i\right)^T \sum_i a_i \lambda_i \mathbf{u}_i \quad (383)$$

$$(384)$$

Example:

$$(a_1 \mathbf{u}_1 + a_2 \mathbf{u}_2)^T (a_1 \lambda_1 \mathbf{u}_1 + a_2 \lambda_2 \mathbf{u}_2) = a_1^2 \lambda_1 \mathbf{u}_1^T \mathbf{u}_1 + a_1 a_2 \lambda_2 \mathbf{u}_1^T \mathbf{u}_2 + a_2 a_1 \lambda_1 \mathbf{u}_2^T \mathbf{u}_1 + a_2^2 \lambda_2 \mathbf{u}_2^T \mathbf{u}_2 \quad (385)$$

$$= a_1^2 \lambda_1 + a_2^2 \lambda_2 \quad (386)$$

Since $\mathbf{u}_i^T \mathbf{u}_j = 0$ for $i \neq j$

So the cross terms vanish, and therefore we are left with an expression that will be positive for all \mathbf{a} only if $\lambda_i > 0$.

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 .

$$\#parameters = D + (D-1) + (D-2) + \dots + D - (D-2) + D - (D-1) \quad (387)$$

$$= (D+1) + (D+1) + \dots + (D+1) \quad (388)$$

$$= \frac{D}{2}(D+1) \quad (389)$$

2.22

$$\begin{aligned} A^{-1}A &= A^{-1}A^T && \text{(Symmetry)} \\ A^{-1}AA^{T^{-1}} &= A^{-1}A^TA^{T^{-1}} && (390) \end{aligned}$$

$$A^{-1T} = A^{-1}A^TA^{T^{-1}} \quad (391)$$

$$= A^{-1} \quad (392)$$

2.23

2.24

The left side will return the identity matrix, but it's not trivial to show as block matrix inversion cannot be done block by block.

The RHS can be multiplied block by block.

Upper left:

$$AM - BD^{-1}CM = (A - BD^{-1}C)M = I \quad (393)$$

Upper right:

$$-AMBD^{-1} + BD^{-1} + BD^{-1}CMBD^{-1} \quad (394)$$

$$= (-AM + I + BD^{-1}CM)BD^{-1} \quad (395)$$

$$= (-AM + AM - BD^{-1}CM + BD^{-1}CM)BD^{-1} = 0 \quad (396)$$

Bottom left:

$$CM - DD^{-1}CM = 0 \quad (397)$$

Bottom right:

$$-CMBD^{-1} + D(D^{-1} + D^{-1}CMBD^{-1}) \quad (398)$$

$$= -CMBD^{-1} + I + CMBD^{-1} = I \quad (399)$$

2.25

We have $p(x_a|x_b) = \frac{1}{p(x_b)} \int p(x_a, x_b, x_c) dx_c$ From section 2.3.2 we know that after marginalizing out

we have $\mathbb{E} \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}$

And then by section 2.3.1 we have $\boldsymbol{\mu}_{a|b}$ equal to (2.81) and $\boldsymbol{\Sigma}_{a|b}$ is equal to (2.82).

2.26

The LHS trivially becomes the identity matrix.

$$(A + BCD)(A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}) \quad (400)$$

$$= I - B(C^{-1} + DA^{-1}B)^{-1}DA^{-1} + BCDA^{-1} - BCDA^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1} \quad (401)$$

$$(402)$$

$$BCDA^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1} \quad (403)$$

$$= BC(C^{-1} - C^{-1} + DA^{-1}B)(C^{-1} + DA^{-1}B)^{-1}DA^{-1} \quad (404)$$

$$= BC(C^{-1} + DA^{-1}B)^{-1}DA^{-1} - BC(C^{-1} + DA^{-1}B)(C^{-1} + DA^{-1}B)^{-1}DA^{-1} \quad (405)$$

$$= BC(C^{-1} + DA^{-1}B)^{-1}DA^{-1} - BCDA^{-1} \quad (406)$$

Plugging this into the previous result gives I .

2.27

$$\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{x} + \mathbf{z}] \quad (407)$$

$$= \int \int p(\mathbf{x}, \mathbf{z})(\mathbf{x} + \mathbf{z})d\mathbf{x}d\mathbf{z} \quad (408)$$

$$= \int \int p(\mathbf{x}, \mathbf{z})\mathbf{x} + p(\mathbf{x}, \mathbf{z})\mathbf{z}d\mathbf{x}d\mathbf{z} \quad (409)$$

$$= \int \int p(\mathbf{x})p(\mathbf{z})\mathbf{x} + p(\mathbf{z})p(\mathbf{x})\mathbf{z}d\mathbf{x}d\mathbf{z} \quad (410)$$

$$= \int p(\mathbf{z}) \int p(\mathbf{x})\mathbf{x}d\mathbf{x}d\mathbf{z} + \int p(\mathbf{x}) \int p(\mathbf{z})\mathbf{z}d\mathbf{x}d\mathbf{z} \quad (411)$$

$$= \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{z}] \quad (412)$$

2.29

It's quite straightforward.

$$M = (A - BD^{-1}C)^{-1} = (\Lambda + A^T LA - ALL^{-1}LA)^{-1} \quad (413)$$

$$= (\Lambda + A^T LA - A^T LA)^{-1} \quad (414)$$

$$= \Lambda^{-1} \quad (415)$$

$$-MBD^{-1} = \Lambda^{-1}A^T LL^{-1} = \Lambda^{-1}A^T \quad (416)$$

$$D^{-1}CM = -L^{-1}LA^{\Lambda^{-1}} = A\lambda^{-1} \quad (417)$$

$$D^{-1} + D^{-1}CMBD^{-1} = L^{-1} + L^{-1}LA\Lambda^{-1}A^TLL^{-1} = L^{-1} + A\Lambda^{-1}A^T \quad (418)$$

2.30

$$\mathbb{E}[z] = R^{-1} \begin{pmatrix} \Lambda\mu - A^T Lb \\ Lb \end{pmatrix} \quad (419)$$

$$= \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1}A^T \\ A\Lambda^{-1} & L^{-1} + A\Lambda^{-1}A^T \end{pmatrix} \begin{pmatrix} \Lambda\mu - A^T Lb \\ Lb \end{pmatrix} \quad (420)$$

$$= \begin{pmatrix} \Lambda^{-1}(\Lambda\mu - A^T Lb) + \Lambda^{-1}A^T Lb \\ A\Lambda^{-1}(\Lambda\mu - A^T Lb) + \mathcal{L}^{-1}Lb + A\Lambda^{-1}A^T Lb \end{pmatrix} \quad (421)$$

$$= \begin{pmatrix} \mu - \Lambda^{-1}A^T Lb + \Lambda^{-1}A^T Lb \\ A\mu - A\Lambda^{-1}A^T Lb + b + A\Lambda^{-1}A^T Lb \end{pmatrix} \quad (422)$$

Which gives the desired result.

2.31

You plug in the values according to equation 2.115.

2.34

The derivative splits into two terms. The first one:

$$-\frac{N}{2} \frac{d}{d\Sigma} \ln |\Sigma| = -\frac{N}{2} \Sigma^{-T} = \frac{N}{2} \Sigma^{-1} \quad (423)$$

The second part:

$$-\frac{1}{2} \sum_{n=1}^N \frac{d}{d\Sigma} (\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \quad (424)$$

By Matrix Cookbook November 15, 2012 eq. 61 we get

$$= \frac{1}{2} \sum_{n=1}^N \Sigma^{-T} (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-T} \quad (425)$$

$$= \frac{1}{2} \sum_{n=1}^N \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-1} \quad (426)$$

$$= \frac{N}{2} \Sigma^{-1} \left(\frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^T \right) \Sigma^{-1} \quad (427)$$

$$(428)$$

Setting the sum of the terms to zero gives

$$\frac{N}{2} \Sigma^{-1} \left(\frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T \right) \Sigma^{-1} = \frac{N}{2} \Sigma^{-1} \quad (429)$$

$$\iff \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T = \Sigma \quad (430)$$

2.35

First we use

$$xx^T = (x - \mathbb{E}[x])(x - \mathbb{E}[x])^T - \mathbb{E}[x] \mathbb{E}[x]^T + x \mathbb{E}[x]^T + \mathbb{E}[x] x^T \quad (431)$$

Then we have

$$\mathbb{E}[xx^T] = \Sigma - \mu\mu^T + \mathbb{E}[x]\mu^T + \mu\mathbb{E}[x]^T = \Sigma + \mu\mu^T \quad (432)$$

The result 2.291 follows from the proof in 2.61 where we see that Σ only appear when $i = j$.

$$\mathbb{E}[\Sigma_{ML}] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[(\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T] \quad (433)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbb{E}[(\mathbf{x}_n - \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n)(\mathbf{x}_n - \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n)^T] \quad (434)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbb{E}[\mathbf{x}_n \mathbf{x}_n^T - \mathbf{x}_n \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n^T - \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \right) \mathbf{x}_n^T + \frac{1}{N^2} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T] \quad (435)$$

$$= \mu\mu^T + \Sigma - \frac{2}{N^2} \sum_{n=1}^N \sum_{n=1}^N \mathbb{E}[\mathbf{x}_n \mathbf{x}_n^T] + \frac{1}{N} (\mu\mu^T + \Sigma) \quad (436)$$

$$= \mu\mu^T + \Sigma - 2\mu\mu^T - \frac{2}{N} \Sigma + \mu\mu^T + \frac{1}{N} \Sigma \quad (437)$$

From which the result follows. We get $\frac{1}{N}$ before the variances since it only appears when $i = j$ in the sum.

2.36

$$\sigma_{ML}^{2(N)} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu^{(N)})^2 \quad (438)$$

$$= \frac{1}{N} (x_n - \mu^{(N)})^2 + \frac{1}{N} \sum_{n=1}^N (x_n - \mu^{(N)})^2 \quad (439)$$

$$= \frac{1}{N} (x_n - \mu^{(N)})^2 + \frac{N-1}{N} \sigma_{ML}^{2(N-1)} \quad (440)$$

$$= \sigma_{ML}^{2(ML)} + \frac{1}{N} \left((x_n - \mu^{(N)})^2 - \sigma_{ML}^{2(N-1)} \right) \quad (441)$$

2.37

$$\Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \quad (442)$$

$$= \frac{1}{N} (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T + \frac{1}{N} \sum_{n=1}^{N-1} (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \quad (443)$$

$$= \frac{1}{N} (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T + \frac{N-1}{N} \Sigma_{ML}^{(N-1)} \quad (444)$$

$$= \Sigma_{ML}^{(N-1)} + \frac{1}{N} \left((\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T - \Sigma_{ML}^{(N-1)} \right) \quad (445)$$

2.39

2.40

The easiest way is to deduce the parameters by examining the proportionality of the log-posterior. The exponent becomes

$$\sum_{n=1}^N -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) - \frac{1}{2} (\mu - \mu_0)^T \Sigma_0^{-1} (\mu - \mu_0) \quad (446)$$

We multiply everything through to get

$$= -\frac{N}{2} \mu^T \Sigma^{-1} \mu - \frac{1}{2} \sum_i^N x^T \Sigma^{-1} x + \sum_i^N x^T \Sigma^{-1} \mu - \frac{1}{2} \mu^T \Sigma_0^{-1} \mu + \mu^T \Sigma_0^{-1} \mu_0 - \frac{1}{2} \mu_0^T \Sigma_0^{-1} \mu_0 \quad (447)$$

Analogously to 2.71 and the text thereafter we exclude the terms not depending on μ and group them together

$$= -\frac{1}{2} \mu^T (N \Sigma^{-1} + \Sigma_0^{-1}) \mu + \mu^T (\Sigma_0^{-1} \mu_0 + \Sigma^{-1} \sum_i^N x_i) + C \quad (448)$$

Now we use the completing-the-square formula $x^T A x + x^T b + c = (x - h)^T A (x - h) + k$ with $h = \frac{1}{2} A^{-1} b$ to see that the mean of the posterior becomes

$$(449)$$

$$h = -\frac{1}{2} A^{-1} b = (N \Sigma^{-1} + \Sigma_0^{-1})^{-1} (\Sigma_0^{-1} \mu_0 + \Sigma^{-1} \sum_i^N x_i) \quad (450)$$

and the variance

$$A = N \Sigma^{-1} + \Sigma_0^{-1} \quad (451)$$

2.41

$$\Gamma(x) = \int_0^\infty u^{x-1} \exp(-u) du \quad (452)$$

$$\int_0^\infty \text{Gam}(\lambda|a, b) d\lambda = \int_0^\infty \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) d\lambda \quad (453)$$

$$= \frac{1}{\Gamma(a)} b^a \int_0^\infty \lambda^{a-1} \exp(-b\lambda) d\lambda \quad (454)$$

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

$$= \frac{\Gamma(a)}{\Gamma(a)} = 1 \quad (455)$$

2.42

$$\mathbb{E}[\lambda] = \int \lambda \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) d\lambda \quad (456)$$

$$= \int \frac{1}{\Gamma(a)} b^a \lambda^a \exp(-b\lambda) d\lambda \quad (457)$$

$$= \frac{1}{\Gamma(a)} b^a \int \lambda^a \exp(-b\lambda) d\lambda \quad (458)$$

Which looks a lot like the gamma function $\Gamma(x) := \int u^{x-1} e^{-u} du$, so we do the change of variables $u = b\lambda \iff d\lambda = \frac{du}{b}$

$$(459)$$

$$= \frac{1}{\Gamma(a)} b^a \int b^{-a} u^a \exp(-u) \frac{1}{b} du \quad (460)$$

$$= \frac{1}{b\Gamma(a)} \Gamma(a+1) \quad (461)$$

$$= \frac{a}{b} \quad (462)$$

$$\mathbb{E}[\lambda^2] = \frac{b^a}{\Gamma(a)} \int \lambda^{a+1} \exp(-b\lambda) d\lambda \quad (463)$$

$$= \frac{b^a}{\Gamma(a)} \int u^{a+1} b^{-a-1} \exp(-u) \frac{1}{b} du \quad (464)$$

$$= \frac{1}{\Gamma(a)b^2} \int u^{a+1} \exp(-u) du \quad (465)$$

$$= \frac{1}{\Gamma(a)b^2} \Gamma(a+2) \quad (466)$$

$$= \frac{a+1}{\Gamma(a)b^2} \Gamma(a+1) \quad (467)$$

$$= \frac{a(a+1)}{b^2} \quad (468)$$

And therefore

$$\mathbb{E}[\lambda^2] - \mathbb{E}[\lambda]^2 = \frac{a(a+1) - a^2}{b} = \frac{a}{b^2} \quad (469)$$

$$\frac{d}{d\lambda} \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) \quad (470)$$

$$= \frac{b^a}{\Gamma(a)} (a-1) \lambda^{a-2} \exp(-b\lambda) - b \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda) = 0 \quad (471)$$

$$\iff (a-1) \lambda^{a-2} = b \lambda^{a-1} \quad (472)$$

$$\iff \lambda = \frac{a-1}{b} \quad (473)$$

$$(474)$$

2.44

I think τ^{-1} should be λ ?

2.46

$$a = \nu/2 \quad b = \frac{\nu}{2\lambda} \quad (475)$$

$$p(x|\mu, a, b) = \frac{b^a}{\Gamma(a)} \left(\frac{1}{2a} \right)^{1/2} \left[b + \frac{1}{2}(x - \mu)^2 \right]^{-a - \frac{1}{2}} \Gamma(a + \frac{1}{2}) \quad (476)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left(\frac{1}{2\pi} \right)^{1/2} \left[\frac{\nu}{2\lambda} + \frac{1}{2}(x - \mu)^2 \right]^{-\frac{\nu}{2} - \frac{1}{2}} \left(\frac{\nu}{2\lambda} \right)^{\nu/2} \quad (477)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left[\frac{2\lambda}{\nu} \left(\frac{\nu}{2\lambda} + \frac{1}{2}(x - \mu)^2 \right) \right]^{-\frac{\nu}{2}} \left[2\pi \left(\frac{\nu}{2\lambda} + \frac{1}{2}(x - \mu)^2 \right) \right]^{-\frac{1}{2}} \quad (478)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left[1 + \frac{\lambda(x - \mu)^2}{\nu} \right]^{-\frac{\nu}{2}} \left[\frac{\pi\nu}{\lambda} + \pi(x - \mu)^2 \right]^{-\frac{1}{2}} \quad (479)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left[1 + \frac{\lambda(x - \mu)^2}{\nu} \right]^{-\frac{\nu}{2}} \left[1 + \frac{\lambda}{\nu}(x - \mu)^2 \right]^{-\frac{1}{2}} \left(\frac{\lambda}{\pi\nu} \right)^{1/2} \quad (480)$$

$$(481)$$

2.47

$$\text{St}(x|\mu, \lambda, \nu) \propto \left[1 + \frac{\lambda(x - \mu)^2}{\nu}\right]^{-v/2-1/2} \quad (482)$$

$$= \exp \ln \left[1 + \frac{\lambda(x - \mu)^2}{\nu}\right]^{-v/2-1/2} \quad (483)$$

$$= \exp \left[\frac{-v-1}{2} \ln \left(1 + \frac{\lambda(x - \mu)^2}{\nu}\right) \right] \quad (484)$$

$$= \exp \left[\frac{-v-1}{2} \left(\frac{\lambda(x - \mu)^2}{\nu} - \mathcal{O} \left[\frac{\lambda(x - \mu)^2}{2\nu} \right]^2 \right) \right] \quad (485)$$

$$= \exp \left[\frac{-v-1}{2} \left(\frac{\lambda(x - \mu)^2}{\nu} - \mathcal{O} \left[\frac{\lambda^2(x - \mu)^4}{4\nu^2} \right] \right) \right] \quad (486)$$

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

2.48

$$\Gamma(x) = \int \mu^{x-1} e^{-u} du \quad (487)$$

$$|(\eta\Lambda)^{-1}|^{1/2} = [\eta^{-D}|\Lambda^{-1}|]^{1/2} = \eta^{-D/2}|\Lambda|^{-1/2} \quad (488)$$

$$\Gamma(D/2 + \nu/2) = \int \eta^{D/2+\nu/2-1} \exp(-\eta) d\eta \quad (489)$$

$$z = \eta/2(\Delta^2 + \nu) \iff \eta = 2 \frac{z}{\Delta^2 + \nu} \quad (490)$$

$$d\eta = \frac{d\eta}{dz} dz = \frac{z}{\Delta^2 + \nu} dz \quad (491)$$

$$\int \mathcal{N}(x|\mu, (\eta\Lambda)^{-1}) \text{Gam}\left(\eta\left|\frac{\nu}{2}, \frac{\nu}{2}\right.\right) d\eta \quad (492)$$

$$= \int \frac{1}{(2\pi)^{D/2}} \frac{1}{|(\eta\Lambda)^{-1}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T(\eta\Lambda)(\mathbf{x} - \mu)\right\} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (493)$$

$$= \int \frac{1}{(2\pi)^{D/2}} \frac{1}{|(\eta\Lambda)^{-1}|^{1/2}} \exp\left\{-\frac{\eta}{2}\Delta^2\right\} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (494)$$

$$= \int \frac{1}{(2\pi)^{D/2}} \eta^{D/2} |\Lambda|^{1/2} \exp\left\{-\frac{\eta}{2}\Delta^2\right\} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (495)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \int \eta^{D/2} \exp\left\{-\frac{\eta}{2}\Delta^2\right\} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (496)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \int \eta^{D/2+\nu/2-1} \exp\left\{-\frac{\eta}{2}(\Delta^2 + \nu)\right\} d\eta \quad (497)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} 2(\Delta^2 + \nu)^{-D/2-\nu/2} \int z^{1/2D+\nu/2} \exp\{-z\} dz \quad (498)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} 2(\Delta^2 + \nu)^{-D/2-\nu/2} \Gamma(D/2 + \nu/2) \quad (499)$$

$$(500)$$

Which works out algebraically to the desired result.

2.49

$$\mathbb{E}_x\left[\int_0^\infty \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, (\eta\boldsymbol{\Lambda})^{-1}) \Gamma(\eta|\nu/2, \nu/2) d\eta\right] = \int_0^\infty \Gamma(\eta|\nu/2, \nu/2) d\eta \mathbb{E}_x[\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, (\eta\boldsymbol{\Lambda})^{-1})] \quad (501)$$

$$= \mathbf{1} \cdot \boldsymbol{\mu} \quad (502)$$

$$\text{cov}[\mathbf{x}] = \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T = \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T \quad (503)$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \int_0^\infty \text{Gam}(\eta|\nu/2, \nu/2) \int \mathbf{x}\mathbf{x}^T \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, (\eta\boldsymbol{\Lambda})^{-1}) d\mathbf{x} d\eta \quad (504)$$

We use 2.62 for the inner expectation

$$= \int_0^\infty \text{Gam}(\eta|\nu/2, \nu/2)(\mu\mu^T + (\eta\Lambda)^{-1})d\eta \quad (505)$$

$$= \mu\mu^T + \Lambda^{-1} \int_0^\infty \text{Gam}(\eta|\nu/2, \nu/2)\eta^{-1}d\eta \quad (506)$$

$$= \mu\mu^T + \Lambda^{-1} \int_0^\infty \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\nu/2-1} \eta^{-1} \exp(-\nu/2\eta)d\eta \quad (507)$$

$$= \mu\mu^T + \Lambda^{-1} \int_0^\infty \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\nu/2-2} \exp(-\nu/2\eta)d\eta \quad (508)$$

$$= \mu\mu^T + \Lambda^{-1} \int_0^\infty \frac{1}{(\nu/2-1)\Gamma(\nu/2-1)} (\nu/2)^{\nu/2} \eta^{\nu/2-2} \exp(-\nu/2\eta)d\eta \quad (509)$$

$$= \mu\mu^T + \Lambda^{-1} \frac{\nu/2}{\nu/2-1} \int_0^\infty \frac{1}{\Gamma(\nu/2-1)} (\nu/2)^{\nu/2-1} \eta^{\nu/2-2} \exp(-\nu/2\eta)d\eta \quad (510)$$

$$= \mu\mu^T + \Lambda^{-1} \frac{\nu/2}{\nu/2-1} \int_0^\infty \text{Gam}(\nu/2-1, \nu/2)d\nu \quad (511)$$

$$= \mu\mu^T + \Lambda^{-1} \frac{\nu}{\nu-2} \quad (512)$$

Plugging this into the covariance formula gives the desired result.

Finding the mode is done by setting the derivative of 2.162 and setting it to zero.

2.50

$$\text{St}(\mathbf{x}|\mu, \Lambda, \nu) \propto \left[1 + \frac{1}{\nu}\Delta^2\right]^{-D/2-\nu/2} \quad (513)$$

$$= \exp \frac{-D-\nu}{2} \ln \left[1 + \frac{1}{\nu}\Delta^2\right] \quad (514)$$

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

$$= \exp \frac{-D-\nu}{2} \frac{1}{\nu}\Delta^2 - \frac{-D-\nu}{2} \mathcal{O}(1/\nu^2) \quad (515)$$

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

2.51

$$1 = [\cos(A) + i \sin(A)][\cos(A) - i \sin(A)] \quad (516)$$

$$= \cos^2(A) - \sin^2(A) \quad (517)$$

$$\cos(A - B) = \Re \exp[i(A - B)] \quad (518)$$

$$= \Re \exp(iA) \exp(-iB) \quad (519)$$

$$= \Re[\cos(A) + i \sin(A)][\cos(B) - i \sin(B)] \quad (520)$$

$$= \cos(A) \cos(B) - \sin(A) \sin(B) \quad (521)$$

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

2.52

Plugging in the Taylor expansion in the exponent, we see that the exponent becomes proportional to $-\frac{m}{2}(\theta - \theta_0)^2 = -\frac{1}{2\sigma^2}(\theta - \theta_0)$ if we define $m = 1/\sigma^2$ the precision.

2.54

$$0 = \sum_{n=1}^N \cos(\theta_0) \sin(\theta_n) - \cos(\theta_n) \sin(\theta_0) \quad (522)$$

$$= \cos(\theta_0) \sum_{n=1}^N \sin(\theta_n) - \sin(\theta_0) \sum_{n=1}^N \cos(\theta_n) \quad (523)$$

$$(524)$$

$$\frac{\sin(\theta_0)}{\cos(\theta_0)} = \frac{\sum_{n=1}^N \sin(\theta_n)}{\sum_{n=1}^N \cos(\theta_n)} \quad (525)$$

$$\tan(\theta_0) = \frac{\sum_{n=1}^N \sin(\theta_n)}{\sum_{n=1}^N \cos(\theta_n)} \quad (526)$$

$$\theta_0 = \arctan \left\{ \frac{\sum_{n=1}^N \sin(\theta_n)}{\sum_{n=1}^N \cos(\theta_n)} \right\} \quad (527)$$

2.54

$$\frac{\partial p}{\partial \theta} = -(2\pi I_0(m))^{-1} \exp\{m \cos(\theta - \theta_0)\} (m \sin(\theta - \theta_0)) \quad (528)$$

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

$$\frac{\partial \partial p}{\partial \theta} = \frac{1}{2\pi I_0(m)} [-\exp(m \cos(\theta - \theta_0))(m \sin(\theta - \theta_0))^2 - m \cos(\theta - \theta_0) \exp(m \cos(\theta - \theta_0))] \quad (529)$$

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

$$A(m_{ML}) = \left(\frac{1}{N} \sum_{n=1}^N \cos \theta_n \right) \cos \theta_0^{ML} - \left(\frac{1}{N} \sum_{n=1}^N \sin \theta_n \right) \sin \theta_0^{ML} \quad (530)$$

$$= \bar{r} (\cos \bar{\theta} \cos \theta_0^{ML} - \sin \bar{\theta} \sin \theta_0^{ML}) \quad (531)$$

$$= \bar{r} (\cos(\bar{\theta} - \theta_0^{ML})) \quad (532)$$

$$= \bar{r} \quad (533)$$

2.56

$$\beta(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \quad (534)$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \exp \log[\mu^{a-1} (1-\mu)^{b-1}] \quad (535)$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} (\exp[(a-1) \log \mu + (b-1) \log[1-\mu]]) \quad (536)$$

$$(537)$$

And so we have $h(\mu) = 1, g(\eta) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}, \eta^T = [a-1, b-1]$ and $u(\mu) = [\ln(\mu), \ln(1-\mu)]^T$

$$\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) \quad (538)$$

$$= \Gamma(a)^{-1} b^a \exp \log[\lambda^{a-1}] \exp(-b\lambda) \quad (539)$$

$$= \Gamma(a)^{-1} b^a \exp(a-1) \log[\lambda] \exp(-b\lambda) \quad (540)$$

$$= \Gamma(a)^{-1} b^a \exp[(a-1) \log[\lambda] - b\lambda] \quad (541)$$

$$(542)$$

And so we have $h(\lambda) = 1, g(\eta) = \Gamma(a)^{-1} b^a, \eta^T = [a-1, -b]$ and $u(\lambda) = [\log \lambda, \lambda]^T$

The von Mises distribution is obvious. We only have to expand the cos in the exponent: $\cos(\theta - \theta_0) = \cos(\theta) \cos(\theta_0) + \sin(\theta) \sin(\theta_0)$ We then have $u(\theta) = [\cos(\theta), \sin(\theta)]^T$ and $\eta^T = [m \cos(\theta_0), m \sin(\theta_0)]$

2.57

For this question, the following knowledge is necessary:

$$a^T B a = \underbrace{B : aa^T}_{\text{Frobenius product (Hadamard \& sum)}} = \text{vec}(B)^T \text{vec}(aa^T) \quad (543)$$

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

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (544)$$

$$= (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp \left[-\frac{1}{2} (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right] \quad (545)$$

$$= (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp(-1/2 \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \exp \left[\begin{pmatrix} -1/2 \text{vec}(\boldsymbol{\Sigma}^{-1}) \\ -\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \end{pmatrix} \begin{pmatrix} \text{vec}(\mathbf{x} \mathbf{x}^T) & \mathbf{x} \end{pmatrix} \right] \quad (546)$$

$$= (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp(-1/2 \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \exp [\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})] \quad (547)$$

$$(548)$$

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

$$g(\boldsymbol{\eta}) = |\text{vec}^{-1}(-2\boldsymbol{\eta}_1)|^{-1/2} \boldsymbol{\eta}_2 \text{vec}^{-1}(-2\boldsymbol{\eta}_1) \boldsymbol{\eta}_2^T \quad (549)$$

and

$$h(\mathbf{x}) = (2\pi)^{-D/2} \quad (550)$$

2.59

$$\int_0^\infty \frac{1}{x} f\left(\frac{1}{x}\right) dx = \int_0^\infty \frac{1}{x} f(y) x dy = \int_0^\infty f(y) dy = 1 \quad (551)$$

2.60

Maybe I'm mistaken, but shouldn't 'volume' be domain here? I think the volume of a region i is equal to the domain Δ_i times it's density h_i . We then have the constraint $\sum_i \Delta_i h_i = 1$

Solving $\sum_{n=1}^N \log h(\mathbf{x}) + \lambda(\sum_i \Delta_i h_i - 1)$ will give the solution.

2.61

3.1

$$2\sigma(2a) - 1 = \frac{2}{1 + e^{-2a}} - 1 \quad (552)$$

$$= \frac{e^a}{e^a} \frac{2}{1 + e^{-2a}} - 1 \quad (553)$$

$$= \frac{2e^a}{e^a + e^{-a}} - 1 \quad (554)$$

$$= \frac{2e^a}{e^a + e^{-a}} - \frac{e^a + e^{-a}}{e^a + e^{-a}} \quad (555)$$

$$= \tanh(a) \quad (556)$$

$$a = \frac{x - \mu_j}{s} \quad (557)$$

$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M w_j \sigma(a) \quad (558)$$

$$= w_0 + \sum_{j=1}^M \frac{w_j}{2} \tanh\left(\frac{1}{2}a\right) + w_j \quad (559)$$

$$= w_0 + \sum_{j=1}^M \frac{w_j}{2} \tanh(a') \quad (560)$$

$$= u_0 + \sum_{j=1}^M u_j \tanh(a') \quad (561)$$

3.2

Proving that some vector is in the column space of a matrix comes down to proving that that vector is the result of a multiplication with that matrix. We have $\Phi(\Phi^T\Phi)^{-1}\Phi^T v$ which can be written as Φw and so the result of the projection operator $(\Phi(\Phi^T\Phi)^{-1}\Phi^T)$ will lie in \mathcal{S} . Now you can see that $y = \Phi w_{ML} = \Phi(\Phi^T\Phi)^{-1}\Phi^T t$ is a projection of t onto \mathcal{S} . Also note that a projection of a vector from that same matrix is the same vector (projection onto its own column space doesn't change anything): $\Phi(\Phi^T\Phi)^{-1}\Phi^T\Phi = \Phi I = \Phi$. Now note by looking at the figure that connects the heads of y and t is the vector $y - t$. If the projection is orthogonal, this vector must be orthogonal to any vector in \mathcal{S} .

$$(y - t)^T \Phi = y^T \Phi - t^T \Phi \quad (562)$$

$$= (\Phi w_{ML})^T \Phi - t^T \Phi \quad (563)$$

$$= w_{ML}^T \Phi^T \Phi - t^T \Phi \quad (564)$$

$$= ((\Phi^T \Phi)^{-1} \Phi^T t)^T \Phi^T \Phi - t^T \Phi \quad (565)$$

$$= t^T \Phi (\Phi^T \Phi)^{-1} \Phi^T \Phi - t^T \Phi \quad (566)$$

$$= t^T \Phi - t^T \Phi \quad (567)$$

which shows that the projection is orthogonal.

3.3

$$\frac{dE_D(\mathbf{w})}{d\mathbf{w}} = \sum_{n=1}^N r_n [t_n - \mathbf{w}^T \phi(\mathbf{x}_n)] \phi(\mathbf{x}_n) = 0 \quad (568)$$

$$\sum_{n=1}^N r_n t_n \phi(\mathbf{x}_n) = \sum_{n=1}^N \mathbf{w}^T \phi(\mathbf{x}_n) \phi(\mathbf{x}_n) \quad (569)$$

$$= \sum_{n=1}^N \phi(\mathbf{x}_n)^T \mathbf{w} \phi(\mathbf{x}_n) \quad (570)$$

$$= \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \mathbf{w} \quad (571)$$

$$= \Phi^T \Phi \mathbf{w} \quad (572)$$

$$([\mathbf{r} \cdot \mathbf{t}]^T \Phi)^T = \Phi^T \Phi \mathbf{w} \quad (573)$$

$$\Phi^T [\mathbf{r} \cdot \mathbf{t}] = \Phi^T \Phi \mathbf{w} \quad (574)$$

$$(\Phi^T \Phi)^{-1} \Phi^T [\mathbf{r} \cdot \mathbf{t}] = \mathbf{w} \quad (575)$$

(i) ?

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

3.4

$$\mathbb{E}[E_d(\mathbf{w})] = \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2\right] \quad (576)$$

$$= \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (w_0 + \sum_{i=1}^D w_i (x_i + \epsilon_i) - t_n)^2\right] \quad (577)$$

$$= \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) + \sum_{i=1}^D w_i \epsilon_i - t_n)^2\right] \quad (578)$$

$$= \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2 + \underbrace{\left(\sum_{i=1}^D w_i \epsilon_i\right)^2}_{=0(\mathbb{E}[\epsilon_i]=0)} + 2(y(\mathbf{x}_n, \mathbf{w}) - t_n) \sum_{i=1}^D w_i \epsilon_i\right] \quad (579)$$

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

$$\mathbb{E}\left(\sum_{i=1}^D w_i \epsilon_i\right)^2 = \mathbb{E}\left[\sum_{i=1}^D \sum_{j=1}^D w_j \epsilon_j w_i \epsilon_i\right] = \sum_{i=1}^D w_i^2 \mathbb{E} \epsilon_i^2 = \sum_{i=1}^D w_i^2 \mathbb{E} \epsilon_i^2 = \sum_{i=1}^D w_i^2 (\sigma^2 + \underbrace{E[\epsilon_i]^2}_{=0}) \quad (580)$$

Which gives us our desired result.

3.5

$$\mathcal{L}(\mathbf{w}, \lambda) = \frac{1}{2} \sum_{n=1}^N [t_n - \mathbf{w}^T \phi(\mathbf{x}_n)]^2 + \lambda \left[\sum_{j=1}^M |w_j|^q - \eta \right] \quad (581)$$

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

3.6

$$p(\mathbf{T}|\mathbf{W}, \Sigma) = \prod_{n=1}^N p(\mathbf{t}_n|\mathbf{W}, \Sigma) \quad (582)$$

$$\ln p(\mathbf{T}|\mathbf{W}, \Sigma) = \ln \prod_{n=1}^N p(\mathbf{t}_n|\mathbf{W}, \Sigma) = \sum_{n=1}^N \ln p(\mathbf{t}_n|\mathbf{W}, \Sigma) \quad (583)$$

$$= \sum_{n=1}^N \ln \left((2\pi)^{-D/2} |\Sigma|^{-1/2} \right) + \sum_{n=1}^N \frac{1}{2} (\mathbf{y}(\mathbf{x}, \mathbf{W}) - \mathbf{t})^T \Sigma^{-1} (\mathbf{y}(\mathbf{x}, \mathbf{W}) - \mathbf{t}) \quad (584)$$

$$\frac{\partial p}{\partial \mathbf{W}} = \frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial \mathbf{W}} [(\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t})^T \Sigma^{-1} (\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t})] \quad (585)$$

$$= \frac{1}{2} \sum_{n=1}^N (\Sigma^{-1} + \Sigma^{-1T}) (\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t}) \phi(\mathbf{x})^T = 0 \quad (586)$$

$$= \sum_{n=1}^N \Sigma^{-1} (\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t}) \phi(\mathbf{x})^T \quad (587)$$

$$\Sigma^{-1} \sum_{n=1}^N \mathbf{W}^T \phi(\mathbf{x}) \phi(\mathbf{x})^T = \Sigma^{-1} \sum_{n=1}^N \mathbf{t} \phi(\mathbf{x})^T \quad (588)$$

$$\mathbf{W}^T \Phi^T \Phi = \mathbf{T} \Phi^T \quad (589)$$

$$\mathbf{W}^T = \mathbf{T} \Phi^T (\Phi^T \Phi)^{-1} \quad (590)$$

$$\mathbf{W} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T} \quad (591)$$

$$\frac{d \ln p(\mathbf{T}|\mathbf{x}, \mathbf{W}, \boldsymbol{\Sigma})}{d\boldsymbol{\Sigma}} = -\frac{N}{2} \frac{d}{d\boldsymbol{\Sigma}} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N \frac{d}{d\boldsymbol{\Sigma}} (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \boldsymbol{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) \quad (592)$$

$$= \frac{N}{2} \frac{d}{d\boldsymbol{\Sigma}} \ln |\boldsymbol{\Sigma}|^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (593)$$

$$= \frac{N}{2} \frac{d}{d\boldsymbol{\Sigma}} \ln |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (594)$$

$$= \frac{N}{2} \boldsymbol{\Sigma} - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (595)$$

$$(596)$$

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (597)$$

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.

$$p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta^{-1}) \quad (598)$$

$$= \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{S}_0) \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (599)$$

$$= (2\pi)^{-D/2} |\mathbf{S}_0|^{-1/2} \exp[-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0)] \prod_{n=1}^N \left(\frac{\beta}{2\pi} \right)^{\frac{1}{2}} \exp[-\frac{\beta}{2}(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2] \quad (600)$$

$$= (2\pi)^{-D/2} |\mathbf{S}_0|^{-1/2} \left(\frac{\beta}{2\pi} \right)^{\frac{N}{2}} \exp[-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N -\frac{\beta}{2}(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2] \quad (601)$$

$$(602)$$

Now we have to get the exponent right.

$$-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N -\frac{\beta}{2} (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 \quad (603)$$

$$= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) - \frac{\beta}{2} \sum_{n=1}^N (t_n^2 + \mathbf{w}^T \phi(\mathbf{x}_n) \mathbf{w}^T \phi(\mathbf{x}_n) - 2\mathbf{w}^T \phi(\mathbf{x}_n) t_n) \quad (604)$$

$$= -\frac{1}{2}(\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{w} - \mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{m}_0 - \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{w} + \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0) - \frac{\beta}{2} (\mathbf{t}^T \mathbf{t} + \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{t}) \quad (605)$$

$$= -\frac{1}{2} (\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{w} - 2\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{m}_0 + \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0 - \beta \mathbf{t}^T \mathbf{t} + \beta \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\beta \mathbf{w}^T \Phi^T \mathbf{t}) \quad (606)$$

$$= \mathbf{w}^T (-\frac{1}{2}(\mathbf{S}_0^{-1} + \beta \Phi^T \Phi)) \mathbf{w} + \mathbf{w}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}) + c \quad (607)$$

$$= \mathbf{w}^T (-\frac{1}{2} \mathbf{S}_N^{-1}) \mathbf{w} + \mathbf{w}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}) + c \quad (608)$$

Now completing the square.

$$\mathbf{w}^T (-\frac{1}{2} \mathbf{S}_N^{-1}) \mathbf{w} + \mathbf{w}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}) + c \quad (609)$$

$$= (\mathbf{w} + \frac{1}{2}(-\frac{1}{2} \mathbf{S}_N^{-1})^{-1}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}))^T (-\frac{1}{2} \mathbf{S}_N^{-1}) (\mathbf{w} + \frac{1}{2}(-\frac{1}{2} \mathbf{S}_N^{-1})^{-1}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t})) + k \quad (610)$$

$$= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^T \mathbf{S}_N^{-1} (\mathbf{w} - \mathbf{m}_N) + k \quad (611)$$

3.8

$$\ln p(\mathbf{w}|\mathbf{t})^{(N)} \propto \ln p(\mathbf{t}|\mathbf{X}, \mathbf{w})^{(N)} + \ln p(\mathbf{w}|\mathbf{t})^{(N-1)} \quad (612)$$

$$= \ln \mathcal{N}(\mathbf{t}|\mathbf{w}^T \phi(\mathbf{x}), \beta^{-1} I) + \ln \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \quad (613)$$

$$= C - \frac{\beta}{2} (\mathbf{t} - \mathbf{w}^T \phi(\mathbf{x}))^T (\mathbf{t} - \mathbf{w}^T \phi(\mathbf{x})) - \frac{1}{2} (\mathbf{w} - \mathbf{m}_N)^T \mathbf{S}_N^{-1} (\mathbf{w} - \mathbf{m}_N) \quad (614)$$

$$= C - \frac{1}{2} \mathbf{w}^T (\beta \phi(\mathbf{x}) \phi(\mathbf{x})^T + \mathbf{S}_N^{-1}) \mathbf{w} - \mathbf{w}^T (\beta \phi(\mathbf{x}) \mathbf{t} + \mathbf{S}_N^{-1} \mathbf{m}_N) \quad (615)$$

We added terms not dependent on \mathbf{w} to C

$$(616)$$

From the completing-the-square formula we can read off the result:

$$\mathbf{S}_{N+1}^{-1} = \mathbf{S}_N^{-1} + \beta \phi(\mathbf{x}) \phi(\mathbf{x})^T \quad (617)$$

$$\mathbf{m}_{N+1} = \mathbf{S}_{N+1} (\mathbf{S}_N^{-1} \mathbf{m}_N + \beta \phi(\mathbf{x}) \mathbf{t}) \quad (618)$$

3.9

We have

$$p(x|y) = \mathcal{N}(x|\Sigma[A^T L(y-b) + \Lambda\mu], \Sigma) \quad (619)$$

with $x = w, \mu = m_N, \Lambda^{-1} = S_N y = t, (Ax + b) = \phi(x)^T w, L^{-1} = \beta^{-1}I$, therefore $\Sigma = (S_N^{-1} + \beta\phi(x)\phi(x)^T)$

$$p(x|y) = \mathcal{N}(w|(S_N^{-1} + \beta\phi(x)\phi(x)^T)[\phi(x)^T \beta I t] + S_N^{-1} m_N, (S_N^{-1} + \beta\phi\phi^T)) \quad (620)$$

$$= \mathcal{N}(w|S_{N+1}^{-1}(\beta\Phi I t + S_N^{-1} m_N, S_{N+1}^{-1})) \quad (621)$$

$$= \mathcal{N}(w|m_{N+1}, S_{N+1}^{-1}) \quad (622)$$

3.10

We have

1. $Ax + b \Rightarrow y(x, w) = \phi(x)^T w$
2. $L^{-1} \Rightarrow \beta^{-1}$
3. $y \Rightarrow t$
4. $\mu \Rightarrow m_N$
5. $\Lambda^{-1} \Rightarrow S_N$

Directly filling in these values gives

$$p(t|x, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\phi(x)^T m_N, \beta^{-1} + \phi(x)^T S_N \phi(x)) \quad (623)$$

3.11

$$S_{N+1} = (S_N^{-1} + \beta(\phi(x_{N+1})\phi(x_{N+1})^T))^{-1} \quad (624)$$

$$= S_N - \left(\frac{S_N \phi(x_{N+1})(\phi(x_{N+1})^T S_N)}{1 + \phi(x)^T S_N \phi(x)} \right) \quad (625)$$

thus

$$\phi(x)^T S_{N+1} \phi(x) = \phi(x)^T S_N \phi(x) - \phi(x)^T A \phi(x) \quad (626)$$

with A as the fraction.

Then we have

$$\sigma_N^2 - \sigma_{N+1}^2 = \phi(x)^T A \phi(x) \geq 0 \quad (627)$$

since (denominator) S_N is positive-semidefinite (i.e., all its quadratic forms are nonnegative) and we get $(\phi(x)^T S_N \phi(x))^2$ in the numerator.

3.12

3.13

The predictive distribution is the integrated posterior multiplied with the likelihood for a new datapoint. In our case the posterior is $p(w, \beta|t)$ and the likelihood $p(t|w, x, \beta)$ (3.113) is our posterior and the likelihood used in 3.12 was $\mathcal{N}(w^T \phi(x), \beta^{-1})$. The predictive distribution thus becomes

$$\int \int \mathcal{N}(w|m_N, \beta^{-1} S_N) \text{Gam}(\beta|a_N, b_N) \mathcal{N}(w^T \phi(x), \beta^{-1}) dw d\beta \quad (628)$$

Noticing that this is the convolution of two Gaussians we can use (2.115). Then we're left with something that looks much like 2.161, and thus a student t.

3.14

3.15

Hint 1: Use 3.95 and 3.92

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N \quad (629)$$

$$= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\gamma(\mathbf{m}_N^T \mathbf{m}_N)^{-1}}{2} \mathbf{m}_N^T \mathbf{m}_N \quad (630)$$

$$= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\gamma}{2} \quad (631)$$

$$(632)$$

$$\beta = (N - \gamma) \|\mathbf{t} - \Phi^T \mathbf{m}_n\|^{-2} \quad (633)$$

$$E(\mathbf{m}_N) = \frac{1}{2} (\mathcal{N} - \gamma + \gamma) = \frac{N}{2} \quad (634)$$

3.17

Evidence function:

$$p(D|M_i) = \int p(D|\mathbf{w}, M_i) p(\mathbf{w}, M_i) d\mathbf{w} \quad (635)$$

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

$$p(D|\alpha, \beta) = \int p(\mathbf{w}|\alpha^{-1} I) \prod_{n=1}^N p(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (636)$$

$$= \int (2\pi)^{-M/2} |\alpha^{-1} I|^{-1/2} \exp[-1/2(\mathbf{w}^T (\alpha^{-1} I)^{-1} \mathbf{w})] \prod_{n=1}^N \frac{1}{\sqrt{2\pi\beta^{-1}}} \exp[\beta/2(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2] \quad (637)$$

$$= \int (2\pi)^{-M/2} (\alpha^{-1})^{-M/2} \frac{\beta^{N/2}}{2\pi} \exp[-\alpha/2 \mathbf{w}^T \mathbf{w} - \beta/2 \|\mathbf{t} - \Phi \mathbf{w}\|^2] \quad (638)$$

3.18

$$\frac{\beta}{2}(\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \quad (639)$$

$$= \frac{\beta}{2} \mathbf{t}^T \mathbf{t} - \beta \mathbf{w}^T \Phi^T \mathbf{t} + \frac{\beta}{2} \mathbf{w}^T \Phi^T \Phi \mathbf{w} + \mathbf{w}^T \frac{\alpha}{2} \mathbf{I} \mathbf{w} \quad (640)$$

$$= \frac{1}{2} \mathbf{w}^T (\beta \Phi^T \Phi + \alpha \mathbf{I}) \mathbf{w} + \beta \mathbf{w}^T \Phi^T \mathbf{t} + \frac{\beta}{2} \mathbf{t}^T \mathbf{t} \quad (641)$$

$$(642)$$

Using the symmetry of $\mathbf{A} = \Phi^T \Phi + \alpha \mathbf{I}$ we can use the matrix completing the square formula with $\mathbf{h} = -\frac{1}{2}(\mathbf{A}^{-1} \Phi^T \mathbf{t})$ and $k = \frac{\beta}{2} \mathbf{t}^T \mathbf{t} - \frac{1}{4} \mathbf{t}^T \Phi \mathbf{A}^{-1} \Phi^T \mathbf{t}$ to get

$$(\mathbf{w} - \mathbf{h})^T \mathbf{A} (\mathbf{w} - \mathbf{h}) + k \quad (643)$$

and so we only have to show that $k = E(m_N)$

3.19

The integrand in the RHS of the first equation, together with $\mathbf{A} = \mathbf{S}_N^{-1}$, is part of a MV Gaussian.

$$\int \exp \left\{ -\frac{1}{2} (\mathbf{w} - m_N)^T \mathbf{A} (\mathbf{w} - m_N) \right\} d\mathbf{w} \quad (644)$$

$$= (2\pi)^{M/2} |\mathbf{A}|^{-1/2} \int (2\pi)^{M/2} |\mathbf{A}|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{w} - m_N)^T \mathbf{A} (\mathbf{w} - m_N) \right\} d\mathbf{w} \quad (645)$$

$$= (2\pi)^{M/2} |\mathbf{A}|^{-1/2} \quad (646)$$

3.20

I think the most interesting part here is to realise that the determinant of a matrix is the product of its eigenvalues and that we have

$$\beta \Phi^T \Phi = \lambda_i u_i \quad (647)$$

Since $\Phi^T \Phi$ is square

$$\iff \beta \Phi^T \Phi u_i + \alpha I u_i = \lambda_i u_i + \alpha I u_i \quad (648)$$

$$\iff (\beta \Phi^T \Phi + \alpha I) u_i = \lambda_i u_i + \alpha u_i = (\lambda_i + \alpha) u_i \quad (649)$$

$$(650)$$

And so the eigenvalues can quite easily be computed.

3.21

Considering the result in the previous exercise (eigenvalue is $\lambda_i + \alpha$, we get

$$\frac{d}{d\alpha} \ln |A| = \frac{d}{d\alpha} \ln \prod_{i=1}^M (\alpha + \lambda_i) \quad (651)$$

$$= \frac{d}{d\alpha} \sum_{i=1}^M \ln(\alpha + \lambda_i) \quad (652)$$

$$= \sum_{i=1}^M \frac{1}{\alpha + \lambda_i} \quad (653)$$

$$= \text{Tr} A^{-1} \quad \text{By spectral theorem, C.46 and C.48} \quad (654)$$

$$= \text{Tr}(A^{-1} \frac{d}{d\alpha} A) \quad (655)$$

Where the last equation is since $\frac{d}{d\alpha} A = \frac{d}{d\alpha} \alpha I + \beta \Phi^T \Phi = I$

3.22

The solutions is pretty detailed in the book. However, I had to realise that if ϵ_i is an eigenvalue of $\Phi^T \Phi$, then $\lambda_i = \beta \epsilon_i$, and therefore $\frac{d\lambda_i}{d\beta} = \epsilon_i = \lambda_i / \beta$

3.23

Like 3.12, this involves multiplying the distributions and then integrating, making use of the fact that the resulting distribution is Gaussian and thus normalized.

3.24

4.4

$$\mathcal{L} = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) + \lambda (\mathbf{w}^T \mathbf{w} - 1) \quad (656)$$

$$\nabla_{\mathbf{w}} \mathcal{L} = (\mathbf{m}_2 - \mathbf{m}_1) + 2\lambda \mathbf{w} = 0 \quad (657)$$

$$\mathbf{w} = -\frac{1}{2\lambda} (\mathbf{m}_2 - \mathbf{m}_1) \propto (\mathbf{m}_2 - \mathbf{m}_1) \quad (658)$$

4.5

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

$$(m_2 - m_1)^2 = (\mathbf{w}^T \mathbf{m}_2 - \mathbf{w}^T \mathbf{m}_1)^2 \quad (659)$$

$$= \mathbf{w}^T \mathbf{m}_2 \mathbf{w}^T \mathbf{m}_2 - \mathbf{w}^T \mathbf{m}_2 \mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{m}_2 + \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{m}_1 \quad (660)$$

$$= \mathbf{w}^T (\mathbf{m}_2 \mathbf{m}_2^T - \mathbf{m}_2 \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{m}_2^T + \mathbf{m}_1 \mathbf{m}_1^T) \mathbf{w} \quad (661)$$

$$= \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w} \quad (662)$$

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

$$s_1^2 = \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m}_1)^2 \quad (663)$$

$$= \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n \mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x}_n \mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{x}_n + \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{m}_1) \quad (664)$$

$$= \sum_{n \in C_1} \mathbf{w}^T (\mathbf{x}_n \mathbf{x}_n^T - \mathbf{x}_n \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{x}_n^T + \mathbf{m}_1 \mathbf{m}_1^T) \mathbf{w} \quad (665)$$

$$= \sum_{n \in C_1} \mathbf{w}^T (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T \mathbf{w} \quad (666)$$

4.6

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0 \quad (667)$$

$$= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m} - t_n) \mathbf{x}_n \quad (668)$$

$$\mathbf{w}^T \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m}) \mathbf{x}_n = \sum_{n=1}^N t_n \mathbf{x}_n \quad (669)$$

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

$$\sum_{n=1}^N t_n \mathbf{x}_n = \sum_{n \in C_1}^{N_1} t_n \mathbf{x}_n - \sum_{n \in C_2}^{N_2} t_n \mathbf{x}_n \quad (670)$$

$$= \sum_{n \in C_1}^{N_1} \frac{N}{N_1} \mathbf{x}_n - \sum_{n \in C_2}^{N_2} \frac{N}{N_2} \mathbf{x}_n \quad (671)$$

$$= N \left(\sum_{n \in C_1}^{N_1} \frac{1}{N_1} \mathbf{x}_n - \sum_{n \in C_2}^{N_2} \frac{1}{N_2} \mathbf{x}_n \right) \quad (672)$$

$$= N(\mathbf{m}_1 - \mathbf{m}_2) \quad (673)$$

$$-\mathbf{w}^T \mathbf{m} \sum_{n=1}^N \mathbf{x}_n = -\frac{1}{N} \mathbf{w}^T (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) \quad (674)$$

$$= -\frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)^T \mathbf{w} \quad (675)$$

$$= -\frac{1}{N} (N_1^2 \mathbf{m}_1 \mathbf{m}_1^T + N_1 N_2 \mathbf{m}_1 \mathbf{m}_2^T + N_2 N_1 \mathbf{m}_2 \mathbf{m}_1^T + N_2^2 \mathbf{m}_2 \mathbf{m}_2^T) \mathbf{w} \quad (676)$$

$$= -\frac{1}{N} ((N - N_2) N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_1 N_2 \mathbf{m}_1 \mathbf{m}_2^T + N_2 N_1 \mathbf{m}_2 \mathbf{m}_1^T + (N - N_1) N_2 \mathbf{m}_2 \mathbf{m}_2^T) \mathbf{w} \quad (677)$$

$$= \left(-N_1 \mathbf{m}_1 \mathbf{m}_1^T + \frac{N_1 N_2}{N} \mathbf{m}_1 \mathbf{m}_1^T - \frac{N_1 N_2}{N} \mathbf{m}_1 \mathbf{m}_2^T - \frac{N_2 N_1}{N} \mathbf{m}_2 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \frac{N_1 N_2}{N} \mathbf{m}_2 \mathbf{m}_2^T \right) \mathbf{w} \quad (678)$$

$$= \left(-N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \frac{N_1 N_2}{N} (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \right) \mathbf{w} \quad (679)$$

$$= \left(-N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} \quad (680)$$

$$(681)$$

We add the remaining terms to the final term:

$$\left(-N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} \quad (682)$$

$$= \left(N_1 \mathbf{m}_1 \mathbf{m}_1^T - 2N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T - 2N_2 \mathbf{m}_2 \mathbf{m}_2^T + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} \quad (683)$$

$$= \left(N_1 \mathbf{m}_1 \mathbf{m}_1^T - \sum_{n \in C_1} \mathbf{x}_n \mathbf{m}_1^T - \mathbf{m}_1 \sum_{n \in C_1} \mathbf{x}_n^T + \dots + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} \quad (684)$$

$$= \left(\sum_{n \in C_1} \mathbf{m}_1 \mathbf{m}_1^T - \mathbf{x}_n \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{x}_n^T + \mathbf{x}_n \mathbf{x}_n^T + \dots \right) \mathbf{w} \quad (685)$$

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

4.7

$$\sigma(-a) = \frac{1}{e^x + 1} = \frac{e^{-x}}{1 + e^{-x}} = \frac{e^{-x} + 1}{e^{-x} + 1} - \frac{1}{e^{-x} + 1} \quad (686)$$

$$y = \frac{1}{1 + e^{-x}} \quad (687)$$

$$e^{-x} = \frac{1 - y}{y} \quad (688)$$

$$y = e^x (1 - y) \quad (689)$$

$$e^x = y/(1-y) \quad (690)$$

$$x = \ln[y/(1-y)] \quad (691)$$

4.8

$$p(C_1|\mathbf{x}) = \sigma(a) \quad (692)$$

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

$$a = \ln \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} + \ln \frac{p(C_1)}{p(C_2)} \quad (693)$$

$$\ln \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} = \ln \left[\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) \right) \right] \quad (694)$$

$$- \ln \left[\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \right) \right] \quad (695)$$

$$= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \quad (696)$$

$$= \frac{1}{2} [2\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 - 2\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_2 + \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2] \quad (697)$$

$$(698)$$

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.

$$\ln p(\mathbf{X}|\mathbf{T}) = \sum_{n=1}^N \ln \prod_{j=1}^K (\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma))^{t_j} \quad (699)$$

$$= \sum_{n=1}^N \sum_{j=1}^K t_j \ln(\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma)) \quad (700)$$

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

$$\mathcal{L}(\boldsymbol{\pi}, \lambda) = \ln(p(\mathbf{X}, \mathbf{T})) - \lambda \left(\sum_{j=1}^K \pi_j - 1 \right) \quad (701)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_j} = \sum_{n=1}^N \frac{t_j}{\pi_j} - \lambda = 0 \quad (702)$$

$$\lambda = \sum_{n=1}^N \frac{t_j}{\pi_j} = N \frac{t_j}{\pi_j} = \frac{N_j}{\pi_j} \quad (703)$$

$$\pi_j = \frac{N_j}{\lambda} \quad (704)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sum_{j=1}^K \pi_j = 1 \quad (705)$$

Plugging this into 380.

$$\sum_{j=1}^K \frac{N_j}{\lambda} = 1 \iff \lambda = N \quad (706)$$

Gives us the desired result.

4.12

$$\frac{d\sigma}{da} = (1 + e^{-a})^{-2} e^{-a} \quad (707)$$

$$= \frac{1}{1 + e^{-a}} \frac{1}{1 + e^{-a}} e^{-a} \quad (708)$$

$$= \sigma(a) \frac{e^{-a}}{1 + e^{-a}} \quad (709)$$

$$= \sigma(a) \left[\frac{1 + e^{-a}}{1 + e^{-a}} - \frac{1}{1 + e^{-a}} \right] \quad (710)$$

4.13

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = - \sum_{n=1}^N \frac{t_n}{y_n} \nabla_{\mathbf{w}} \sigma(\mathbf{w}^T \phi_n) - \frac{1 - t_n}{1 - y_n} \nabla_{\mathbf{w}} \sigma(\mathbf{w}^T \phi_n) \quad (711)$$

$$= - \sum_{n=1}^N \frac{t_n}{y_n} \sigma(\mathbf{w}^T \phi_n) (1 - \sigma(\mathbf{w}^T \phi_n)) \phi_n - \frac{1 - t_n}{1 - y_n} \sigma(\mathbf{w}^T \phi_n) (1 - \sigma(\mathbf{w}^T \phi_n)) \phi_n \quad (712)$$

$$= - \sum_{n=1}^N t_n (1 - y_n) \phi_n - (1 - t_n) y_n \phi_n \quad (713)$$

$$= - \sum_{n=1}^N (t_n - y_n) \phi_n \quad (714)$$

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

$$p(\mathbf{t}, \mathbf{w}) = \prod_{n=1}^N y_n^{\pi_n} [1 - y_n]^{1-t_n} \quad (715)$$

$$\ln p = \sum_{n=1}^N \pi_n \ln y_n + (1 - \pi_n) \ln(1 - y_n) \quad (716)$$

4.17

$$p(C_k | \phi) = y_k = \frac{\exp a_k}{\sum_{j=1} \exp a_j} \quad (717)$$

$$\frac{\partial y_k}{\partial a_j} = -\exp a_k \left(\sum_j \exp(a_j) \right)^{-2} \exp(a_j) \quad (718)$$

$$= \begin{cases} y_k(0 - y_j) & j \neq k \\ y_k(1 - y_j) & j = k \end{cases} \quad (719)$$

$$= y_k(I_{kj} - y_j) \quad (720)$$

4.18

$$\nabla_{\mathbf{w}_j} E(\mathbf{W}) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \nabla_{\mathbf{w}_j} \ln y_{nk} \quad (721)$$

$$\nabla_{\mathbf{w}_j} \ln(y_{nk}) = -(I_{kj} - y_j) \phi_n \quad (722)$$

$$\nabla_{\mathbf{w}_j} E(\mathbf{W}) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} (I_{kj} - y_j) \phi_n \quad (723)$$

$$= \phi \sum_{n=1}^N \sum_{k=1}^K t_{nk} y_{jn} \phi - t_{nk} I_{kj} \quad (724)$$

$$= \phi \sum_{n=1}^N t_{nj} - y_{jn} \phi \underbrace{\sum_{k=1}^K t_{nk}}_{=1} \quad (725)$$

$$= \sum_{n=1}^N \phi (y_{jn} - t_{nj}) \quad (726)$$

4.19

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

$$p(\mathbf{t}, \mathbf{w}) = \sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)] \quad (727)$$

$$\nabla_{\mathbf{w}} \sum_{n=1}^N [t_n \ln \Phi(a) + (1 - t_n) \ln(1 - \Phi(a))] = \sum_{n=1}^N \left(\frac{t_n}{\Phi(a)} - \frac{1 - t_n}{1 - \Phi(a)} \right) \Phi(a) \phi_n \quad (728)$$

4.21

$$\Phi(a) = \int_0^a \mathcal{N}(0, 1) d\theta \quad (729)$$

$$= \int_0^a \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\theta^2\right) d\theta \quad (730)$$

$$= \frac{1}{2} + \int_{-\infty}^a \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\theta^2\right) d\theta \quad (731)$$

$$= \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \frac{\sqrt{\pi}}{2} \int_{-\infty}^a \frac{2}{\sqrt{\pi}} \exp\left(-\frac{1}{2}\theta^2\right) d\theta \quad (732)$$

$$= \frac{1}{2} \left(1 + \frac{1}{\sqrt{2}} \text{erf}(a) \right) \quad (733)$$

4.22

$$\ln p(D) = \ln \left[f(z_0) \frac{(2\pi)^{M/2}}{|A|^{1/2}} \right] \quad (734)$$

$$= \ln f(z_0) = \frac{M}{2} \ln 2\pi - \frac{1}{2} \ln |A| \quad (735)$$

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

$$\ln f(\boldsymbol{\theta}) \Big|_{z_0} = \ln f(\boldsymbol{\theta}_{MAP}) = \ln p(D|\boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) \quad (736)$$

5.2

$$p(\mathbf{T}, \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \mathcal{N}(\mathbf{t}_n | y(\mathbf{x}_n, \mathbf{w}), \beta^{-1} I) \quad (737)$$

$$= \prod_{n=1}^N \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\beta^{-1} I|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))^T \beta^{-1} I (\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))\right) \quad (738)$$

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

$$\left(-\frac{1}{2}(\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))^T \beta^{-1} I(\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))\right) = -\frac{1}{2\beta} \|y(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2. \quad (739)$$

5.5

$$p(\mathbf{T}, \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \prod_{k=1}^K y_k(\mathbf{x}_n, \mathbf{w})^{t_n^k} (1 - y_k)(\mathbf{x}_n, \mathbf{w})^{1-t_n^k} \quad (740)$$

Taking the log becomes the cross-entropy function.

5.6

$$\frac{\partial E(\mathbf{w})}{\partial a_k} = \sum_{n=1}^N \left(\frac{\partial}{\partial a_k} t_n \ln y_n + \frac{\partial}{\partial a_k} (1 - t_n) \ln(1 - y_n) \right) \quad (741)$$

$$= -\frac{t_k}{y_k} \frac{\partial}{\partial a_k} y_k - \frac{1 - t_k}{1 - y_k} \frac{\partial}{\partial a_k} (1 - y_k) \quad (742)$$

$$= -t_k(1 - y_k) + (1 - t_k)y_k \quad (743)$$

$$= y_k - t_k \quad (744)$$

5.7

$$E(\mathbf{w}) = \sum_{k=1}^K t_k \ln y_k(\mathbf{x}, \mathbf{w}) \quad (745)$$

$$\frac{\partial E}{\partial a_j} = -\sum_{k=1}^K \frac{t_k}{y_k} y_k (I_{kj} - y_j) \quad (746)$$

$$= -\sum_{k=1}^K t_k (I_{kj} - y_j) \quad (747)$$

$$= -t_j + y_j \quad (748)$$

$$\frac{d \tanh}{da} = (e^a - e^{-a}) \frac{d}{da} (e^a + e^{-a})^{-1} + (e^a + e^{-a})^{-1} \frac{d}{da} (e^a - e^{-a}) \quad (749)$$

$$= -(e^a - e^{-a})(e^a + e^{-a})^{-2} \frac{d}{da} (e^a + e^{-a}) + (e^a + e^{-a})^{-1} (e^a - e^{-a}) \quad (750)$$

$$= -(e^a - e^{-a})(e^a + e^{-a})^{-2} (e^a - e^{-a}) + 1 \quad (751)$$

$$= -h^2(a) + 1 \quad (752)$$

5.9

Still Bernoulli, so

$$p(t|\mathbf{x}, \mathbf{w}) = \left(\frac{1+y}{2}\right)^{\frac{1+t}{2}} \left(\frac{1-y}{2}\right)^{\frac{1-t}{2}} \quad (753)$$

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}) \quad (754)$$

$$-\ln(p(\mathbf{t}|\mathbf{X}, \mathbf{w})) = -\sum_{n=1}^N \ln p(t_n|\mathbf{x}_n, \mathbf{w}) \quad (755)$$

$$= -\sum_{n=1}^N \left(\left(\frac{1+t}{2}\right) \ln\left(\frac{1+y}{2}\right) + \left(\frac{1-t}{2}\right) \ln\left(\frac{1-y}{2}\right) \right) \quad (756)$$

We can use tanh as activation function.

5.10

$$u_i^T H u_i = u_i^T \lambda u_i = \delta_i i \lambda_i = \lambda_i \quad (757)$$

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.

$$\frac{N(N+1)}{2} + N = \frac{N(N+3)}{2} \quad (758)$$

5.14

Hint 1: Taylor expansion on both terms in numerator.

Taylor:

$$E_n(w_j + \epsilon) = E_n(w_{ji}) + \epsilon \frac{\partial E_n}{\partial w_{ji}} + \frac{\epsilon^2}{2} \frac{\partial^2 E_n}{\partial w_{ji}^2} + O(\epsilon^3) \quad (759)$$

$$E_n(w_j - \epsilon) = E_n(w_{ji}) - \epsilon \frac{\partial E_n}{\partial w_{ji}} + \frac{\epsilon^2}{2} \frac{\partial^2 E_n}{\partial w_{ji}^2} - O(\epsilon^3) \quad (760)$$

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

5.16

$$E = \frac{1}{2} \sum_{n=1}^N \|\mathbf{y}_n - \mathbf{t}_n\|^2 = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M (\mathbf{y}_n - \mathbf{t}_n)^2 \quad (761)$$

$$\nabla E = \sum_{n=1}^N \sum_{m=1}^M (\mathbf{y}_n - \mathbf{t}_n) \nabla \mathbf{y}_n \quad (762)$$

$$\nabla \nabla E = \sum_{n=1}^N \sum_{m=1}^M \nabla \mathbf{y}_n \nabla (\mathbf{y}_n - \mathbf{t}_n) + (\mathbf{y}_n - \mathbf{t}_n) \nabla \nabla \mathbf{y}_n \approx \sum_{n=1}^N \sum_{m=1}^M \nabla \mathbf{y}_n \nabla \mathbf{y}_n^T \quad (763)$$

5.17

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

$$\frac{\partial E}{\partial w_r} = \int \int (y(x, w) - t) \frac{\partial y}{\partial w_r} p(x, t) dx dt \quad (764)$$

$$\frac{\partial^2 E}{\partial w_r \partial w_s} = \int \int \left[(y(x, w) - t) \frac{\partial^2 y(x, w)}{\partial w_r \partial w_s} + \frac{\partial y(x, w)}{\partial w_r} \frac{\partial y}{\partial w_s} \right] p(x, t) dx dt \quad (765)$$

$$\int \int (y(x, w) - t) p(x, t) dx dt = \int \int (y(x, w) - t) p(t|x) p(x) dx dt \quad (766)$$

$$= \int p(x) \left(y(x, w) - \underbrace{\int t p(t|x)}_{=y(x, w)} \right) dx dt = 0 \quad (767)$$

The remaining integral is the answer.

5.18

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left(\sum_{j=1}^M w_{kj} h \left(\sum_{i=1}^D w_{ji} x_i + w_{j0} \right) + \sum_{l=1}^D w_l x_l \right) \quad (768)$$

Finding the derivatives to these skip weights is straightforward.

5.19

$$E = - \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \quad (769)$$

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

$$\nabla \nabla E = \sum_{n=1}^N y_n (1 - y_n) \nabla a_n \nabla a_n^T + (y_n - t_n) \nabla \nabla a_n \quad (770)$$

5.20

$$\nabla_{w_j} E(\mathbf{W}) = \sum_{n=1}^N (y_{n_j} - bt_{n_j}) \nabla a_j \quad (771)$$

$$\nabla \nabla \mathbf{w}_j \approx \sum_{n=1} y_k (I - y_j) \nabla a_j \nabla a_j^T \quad (772)$$

5.24

$$\sum_i \frac{1}{a} w_{ji} (ax_i + b) + w_{j0} - \frac{b}{a} \sum_i w_{ji} = \sum_i w_{ji} x + \frac{b}{a} w_{ji} + w_{j0} - \frac{b}{a} \sum_i w_{ji} \quad (773)$$

$$= \sum_i w_{ji} x_i + w_{j0} \quad (774)$$

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

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \quad (775)$$

$$\frac{\partial \tilde{E}}{\partial w_i} = \frac{\partial E}{\partial w_i} + \frac{\partial}{\partial w_i} \lambda \Omega(\mathbf{w}) \quad (776)$$

$$\frac{\partial}{\partial w_i} \lambda \Omega(\mathbf{w}) = -\lambda \frac{\partial}{\partial w_i} \ln \left(\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right) \quad (777)$$

$$\frac{\partial}{\partial w_i} = \frac{1}{\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2)} \frac{\partial}{\partial w_i} \ln \left(\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right) \quad (778)$$

$$\frac{\partial}{\partial w_i} \ln \left(\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right) = \sum_{j=1}^M \pi_j \frac{\partial}{\partial w_i} \mathcal{N}(w_i | \mu_j, \sigma_j^2) \quad (779)$$

$$\frac{\partial}{\partial w_i} \mathcal{N}(w_i | \mu_j, \sigma_j^2) = \frac{1}{\sqrt{2\pi\sigma_j^2} \exp(-\frac{1}{2} \frac{(w_i - \mu_j)^2}{\sigma_j^2})} \frac{\partial}{\partial w_i} \left(-\frac{1}{2} \frac{(w_i - \mu_j)^2}{\sigma_j^2} \right) \quad (780)$$

$$\frac{\partial}{\partial w_i} = -\sigma_j^{-2} (w_i - \mu_j) \quad (781)$$

Plugging everything in gives

$$\frac{\partial}{\partial w_i} = \frac{\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j)}{\sum_{k=1}^K \mathcal{N}(w_i | \mu_j, \sigma_j^2)} \quad (782)$$

5.30

$$\frac{\partial \tilde{E}}{\partial \mu_j} = -\lambda \frac{\partial}{\partial \mu_j} \sum_i \ln \sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j) \quad (783)$$

$$= -\lambda \sum_i \frac{1}{\sum_j \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j)} \quad (784)$$

$$= -\lambda \sum_i \frac{1}{\sum_j \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j)} \sum_j \pi_j \frac{d}{d \mu_j} \mathcal{N}(w_i | \mu_j, \sigma_j) \quad (785)$$

$$\frac{d}{d \mu_j} \mathcal{N}(w_i | \mu_j, \sigma_j) = -\frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{1}{2\sigma_j^2}(w_i - \mu_j)^2\right) \left[\frac{1}{\sigma_j^2}(w_i - \mu_j)\right] \quad (786)$$

Plugging in gives the result.

5.31

$$\frac{\partial}{\partial \sigma_j} \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j) \quad (787)$$

$$= \pi \left[\exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \frac{\partial}{\partial \sigma_j} \frac{1}{\sqrt{2\pi\sigma_j^2}} + \frac{1}{\sqrt{2\pi\sigma_j^2}} \frac{\partial}{\partial \sigma_j} \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \right] \quad (788)$$

$$\frac{\partial}{\partial \sigma_j} (2\pi\sigma_j^2)^{-\frac{1}{2}} = -\frac{1}{2} (2\pi\sigma_j^2)^{-\frac{3}{2}} \frac{\partial}{\partial \sigma_j} (2\pi\sigma_j^2) \quad (789)$$

$$= -\frac{1}{2} (2\pi\sigma_j^2)^{-\frac{3}{2}} 4\pi\sigma_j \quad (790)$$

$$= -\frac{1}{\sigma_j} \frac{1}{\sqrt{2\pi\sigma_j^2}} \quad (791)$$

$$\frac{\partial}{\partial \sigma_j} \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) = \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \frac{\partial}{\partial \sigma_j} \left(-\frac{1}{2\sigma_j^2}(\mu_j - w_i)^2\right) \quad (792)$$

$$= \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \frac{1}{\sigma_j^3} (\mu_j - w_i)^2 = \quad (793)$$

Plugging these values in gives the result.

5.33

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

$$\mathbf{y} = \begin{pmatrix} \cos(\pi - \theta)L_1 + v_1 \\ \sin(\pi - \theta)L_1 + v_2 \end{pmatrix} \quad (794)$$

$$\mathbf{x} = \begin{pmatrix} y_1 + \cos(\theta_1 + \theta_2 - \pi)L_2 \\ y_2 + \sin(\theta_1 + \theta_2 - \pi)L_2 \end{pmatrix} \quad (795)$$

5.34

$$\frac{\partial E_n}{\partial a_{kl}^\pi} = - \frac{1}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (796)$$

$$\frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) = \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_k \quad (797)$$

$$\frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_k = \sum_{l=1} \pi_k (I_{kl} - \pi_l) \quad (798)$$

$$= \pi_k - \pi_k \sum_{l=1} \pi_l \quad (799)$$

$$\frac{\partial E_n}{\partial a_{kl}^\pi} = - \frac{\pi_k \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) - \pi_k \sum_{l=1} \pi_l \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \quad (800)$$

$$= -\gamma_k(w) + \pi_k \quad (801)$$

5.35

$$\frac{\partial E_n}{\partial a_{kl}^\mu} = - \frac{1}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \frac{\partial}{\partial a_k^\mu} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (802)$$

$$\frac{\partial}{\partial a_k^\mu} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (803)$$

$$= \pi_k(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \frac{\partial}{\partial a_k^\mu} - \frac{1}{2} (\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1} (\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \quad (804)$$

$$(805)$$

$$\frac{\partial}{\partial a_k^\mu} - \frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \quad (806)$$

$$= -\sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \underbrace{\frac{\partial \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w})}{\partial a_k^\mu}}_{=1} \quad (807)$$

$$\frac{\partial E_n}{\partial a_{kl}^\mu} = \gamma_{nk} \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \quad (808)$$

5.36

$$\frac{\partial E_n}{\partial a_{kl}^\sigma} = - \frac{1}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \frac{\partial}{\partial a_k^\sigma} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (809)$$

$$\frac{\partial}{\partial a_k^\sigma} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) = \pi_k(\mathbf{x}_n, \mathbf{w}) \frac{\partial}{\partial a_k^\sigma} \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (810)$$

$$\frac{\partial}{\partial a_k^\sigma} \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) = \frac{\partial}{\partial a_k^\sigma} \frac{1}{(2\pi)^{\frac{L}{2}}} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) \quad (811)$$

$$= \frac{1}{(2\pi)^{\frac{L}{2}}} \left(\exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) \frac{\partial}{\partial a_k^\sigma} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \right. \quad (812)$$

$$\left. + \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \frac{\partial}{\partial a_k^\sigma} \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) \right) \quad (813)$$

$$\frac{\partial}{\partial a_k^\sigma} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} = \frac{\partial}{\partial a_k^\sigma} \frac{1}{\sigma_k(\mathbf{x}, \mathbf{w})^L} = -\frac{L}{\sigma^2(\mathbf{x}, \mathbf{w})^{L-1}} = \frac{L}{\sigma_k} \frac{1}{\sigma_k(\mathbf{x}, \mathbf{w})^L} = \frac{L}{\sigma_k} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \quad (814)$$

$$\frac{\partial}{\partial a_k^\sigma} \left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) = \frac{\|\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}\|^2}{\sigma_k^3} \quad (815)$$

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

5.37

5.38

5.39

Using the approximation:

$$p(D|\alpha, \beta) = \int p(D|w, \beta)p(w|\alpha)dw \approx p(D, w_{map})p(w_{map})\frac{(2\pi)^{W/2}}{|A|^{\frac{1}{2}}} \quad (816)$$

Taking the log gives the result.

5.40

Use a softmax activation.

6.3

$$\|x - x_n\| = x^T x + x_n^T x - 2x_n^T x \quad (817)$$

$$= k(x, x) + k(x_n, x) - 2k(x_n, x) \quad (818)$$

6.4

$$\begin{pmatrix} 2 & 0 \\ -1 & 3 \end{pmatrix} \quad (819)$$

6.5

$$ck(\mathbf{x}, \mathbf{x}') = c\psi^T \psi = (c^{\frac{1}{2}}\psi)^T (c^{\frac{1}{2}}\psi) = \phi^T \phi \quad (820)$$

$$f(\mathbf{x})k(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') = f(\mathbf{x})\psi^T \psi f(\mathbf{x}') = (f(\mathbf{x})\psi^T)^T (\psi f(\mathbf{x}')) = \phi^T \phi \quad (821)$$

6.6

$$q(k(x, x')) = \sum_{n=1}^N \beta_n k(x, x')^n \quad (822)$$

$$= \sum_{n=1}^N \beta_0 k(x, x') \quad 6.18 \quad (823)$$

$$= \sum_{n=1}^N k(x, x') \quad 6.13 \quad (824)$$

$$= k(x, x') \quad 6.17 \quad (825)$$

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

6.7

$$\psi(x)^T \psi(x') + \theta(x)^T \theta(x') = \sum_{n=1}^N (\psi(x_n) + \theta(x_n))(\psi(x'_n) + \theta(x'_n)) - \psi(x_n)\theta(x'_n) - \theta(x_n)\psi(x'_n) \quad (826)$$

$$= \begin{bmatrix} \dots & \psi(x_n) + \theta(x_n) & -\psi(x_n) & -\theta(x_n) & \dots \end{bmatrix} \begin{bmatrix} \vdots \\ \psi(x'_n) + \theta(x'_n) \\ \theta(x'_n) \\ \psi(x'_n) \\ \vdots \end{bmatrix} \quad (827)$$

$$= \phi(x)^T \phi(x') \quad (828)$$

$$k_1(x, x') k_2(x, x') = \theta(x)^T \theta(x') (\psi(x)^T \psi(x')) \quad (829)$$

$$= \sum_{n=1}^N \theta(x_n) \theta(x'_n) \psi(x_n) \psi(x'_n) \quad (830)$$

$$= \sum_{n=1}^N \theta(x_n) \psi(x_n) \theta(x'_n) \psi(x'_n) \quad (831)$$

$$= \phi(x)^T \phi(x') \quad (832)$$

$$= k(x, x') \quad (833)$$

6.8

$$k_3(\phi(x), \phi(x')) = \psi(\phi(x))^T \psi(\phi(x')) = k(x, x') \quad (834)$$

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

$$= (Ux)^T Ux' \quad (835)$$

$$= \phi(x)^T \phi(x') \quad (836)$$

$$= k(x, x') \quad (837)$$

6.9

$$\sum_{n=1}^N \psi(x_{na})\psi(x'_{na}) + \sum_{m=1}^M \theta(x_{mb})\theta(x'_{mb}) \quad (838)$$

$$= \begin{bmatrix} \psi(x_{a1}) & \dots & \psi(x_{aN}) & \theta(x_{b1}) & \dots & \theta(x_{bM}) \end{bmatrix} \begin{bmatrix} \psi(x'_{a1}) \\ \vdots \\ \psi(x'_{aN}) \\ \theta(x'_{b1}) \\ \vdots \\ \theta(x'_{bM}) \end{bmatrix} = k(x, x') \quad (839)$$

$$x \in \mathbb{R}^n \quad (840)$$

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

$$y = f(g(x)) \quad (841)$$

$$y'(x) = \nabla_{g(x)} f(g(x)) J_x(g(x)) \quad (842)$$

$$k_a(x_a, x'_a)k(x_b, x'_b) = \sum_{n=1}^N \theta(x_{an})\theta(x'_{an}) \sum_{m=1}^M \psi(x_{bm})\psi(x'_{bm}) \quad (843)$$

$$= \sum_{n=1}^N \sum_{m=1}^M \theta(x_{an})\theta(x'_{an})\psi(x_{bm})\psi(x'_{bm}) \quad (844)$$

$$= \sum_{n=1}^N \sum_{m=1}^M \theta(x_{an})\psi(x_{bm})\theta(x'_{an})\psi(x'_{bm}) \quad (845)$$

$$= \phi(x)^T \phi(x') = k(x, x') \quad (846)$$

6.10

$$y(x) = \sum_{n=1}^N f(x_n)f(x_n)(K + \lambda I_N)^{-1}t \quad (847)$$

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 ϕ is an invertible differentiable transformation of θ .

$$g(\theta, \mathbf{x}) = J_\theta(\phi) \nabla_\phi \ln p(\mathbf{x}, \phi) = J_\theta(\phi) g(\phi, \mathbf{x}) \quad (848)$$

$$g(\phi, \mathbf{x}) = J_\theta^{-1}(\phi) g(\theta, \mathbf{x}) \quad (849)$$

$$\mathbf{F}' = J_\theta^{-1}(\phi) \mathbb{E}_x[g(\theta, \mathbf{x}) g(\theta, \mathbf{x})^T] J_\theta^{-T}(\phi) = J_\theta^{-1}(\phi) \mathbf{F} J_\theta^{-T}(\phi) \quad (850)$$

$$k'(\mathbf{x}, \mathbf{x}') = g(\phi, \mathbf{x})^T \mathbf{F}'^{-1} g(\phi, \mathbf{x}') \quad (851)$$

$$= (J_\theta^{-1}(\phi) g(\theta, \mathbf{x}))^T (J_\theta^{-1}(\phi) \mathbf{F} J_\theta^{-T}(\phi))^{-1} J_\theta^{-1}(\phi) g(\theta, \mathbf{x}') \quad (852)$$

$$= g(\theta, \mathbf{x})^T J_\theta^{-T}(\phi) J_\theta^T(\phi) \mathbf{F} J_\theta(\phi) J_\theta^{-1}(\phi) g(\theta, \mathbf{x}') \quad (853)$$

$$= g(\theta, \mathbf{x})^T \mathbf{F} g(\theta, \mathbf{x}') \quad (854)$$

$$= k(\mathbf{x}, \mathbf{x}') \quad (855)$$

Therefore, the Fisher kernel is invariant.

6.14

$$\nabla_\theta (\ln p(\mathbf{x}, \theta)) = \nabla_\mu \ln p(\mathbf{x}, \mu) \quad (856)$$

$$= S^{-1}(\mathbf{x} - \mu) \quad (857)$$

$$\mathbf{F} = \mathbb{E}[S^{-1}(\mathbf{x} - \mu)(S^{-1}(\mathbf{x} - \mu))^T] \quad (858)$$

$$= S^{-1} \mathbb{E}_x[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] S^{-1} \quad (859)$$

$$= S^{-1} S S^{-1} = S^{-1} \quad (860)$$

$$k(\mathbf{x}, \mathbf{x}') = (S^{-1}(\mathbf{x} - \mu))^T \mathbf{F}^{-1} (S^{-1}(\mathbf{x}' - \mu)) \quad (861)$$

$$= (\mathbf{x} - \mu)^T S^{-1} S S^{-1} (\mathbf{x}' - \mu) \quad (862)$$

$$= (\mathbf{x} - \mu)^T S^{-1} (\mathbf{x}' - \mu) \quad (863)$$

6.15

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

$$|K| = k(x_1, x_1)k(x_2, x_2) - k(x_2, x_1)k(x_1, x_2) \geq 0 \quad (864)$$

from which the result follows.

6.18

$$\mathbf{z} = (x, t) \quad (865)$$

$$p(\mathbf{z}) = \frac{\sum_n \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)}{\sum_m \int \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I) dt} = \frac{\sum_n \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)}{\sum_m \mathcal{N}(x - x_n, \sigma^2 I)} \quad (866)$$

$$\frac{\sum_n \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)}{\sum_m \mathcal{N}(x - x_n, \sigma^2 I)} = \frac{\sum_n \frac{1}{2\pi\sigma^2} \exp[-\frac{1}{2\sigma^2}(\mathbf{z} - \mathbf{z}_n)^T(\mathbf{z} - \mathbf{z}_n)]}{\sum_m \frac{1}{\sqrt{2\pi}\sigma} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]} \quad (867)$$

$$= \frac{\sum_n \frac{1}{2\pi\sigma^2} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2 + -\frac{1}{2\sigma^2}(t - t_n)^2]}{\sum_m \frac{1}{\sqrt{2\pi}\sigma} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]} \quad (868)$$

$$= \sum_n \frac{\frac{1}{\sqrt{2\pi}\sigma^2} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]}{\sum_m \frac{1}{\sqrt{2\pi}\sigma} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]} \frac{1}{\sqrt{2\pi}\sigma^2} \exp[-\frac{1}{2\sigma^2}(t - t_n)^2] \quad (869)$$

$$= \sum_n \pi_n \mathcal{N}(t|t_n, \sigma^2) \quad (870)$$

6.24

$$u^T W u = u^T \sqrt{W} \sqrt{W} u \quad (871)$$

$$= (\sqrt{W} u)^T \sqrt{W} u \quad (872)$$

$$> 0 \quad \forall u \in \mathbb{R} \setminus \mathbf{0} \quad (873)$$

$$u^T (W + V) u = u^T W u + u^T V u > 0 \quad (874)$$

6.25

$$a_N = a_N - \nabla \nabla \Psi(a_N) \nabla \Psi(a_N) \quad (875)$$

$$= a_N + (W_N + C^{-1})^{-1} [t_N - \sigma_N - C^{-1} a_N] \quad (876)$$

$$= (W_N + C_N^{-1}) [W_N a_N - \sigma_N + t_N] \quad (877)$$

$$= C_N ((W_N + C_N^{-1}) C_N)^{-1} [W_N a_N - \sigma_N + t_N] \quad (878)$$

$$= C_N (C_N W_N + I)^{-1} [W_N a_N - \sigma_N + t_N] \quad (879)$$

6.26

7.2

$$t_n(\mathbf{w}^T \phi(\mathbf{x}) + b) = \gamma \quad (880)$$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \gamma) \quad (881)$$

$$\nabla_{\mathbf{w}} = \mathbf{w} - \sum_{n=1}^N \nabla a_n t_n \mathbf{w}^T \phi(\mathbf{x}_n) \quad (882)$$

$$= \mathbf{w} - \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) = 0 \quad (883)$$

$$\iff \mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \quad (884)$$

$$\frac{\partial}{\partial b} = - \sum_{n=1}^N a_n t_n = 0 \quad (885)$$

$$= \sum_{n=1}^N a_n t_n \quad (886)$$

$$\tilde{L} = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n t_n a_m t_m \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) - \sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \gamma) \quad (887)$$

$$\sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \gamma) = \sum_{n=1}^N 0 - a_n \gamma \quad (888)$$

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

$$-\ln p(t_n) = - \sum_{n=1} \ln \sigma(t_n y_n) \quad (889)$$

which is the cross-entropy error function.

7.7

It is quite straightforward if you use the following steps.

$$\sum_{n=1}^N \xi_n (C - \mu_n - a_n) + \sum_{n=1}^N \hat{\xi}_n (C - \hat{\mu}_n - \hat{a}_n) = 0 \quad (890)$$

and

$$- \sum_{n=1}^N a_n (y_n) - \sum_{n=1}^N \hat{a}_n (-y_n) = - \sum_{n=1}^N (a_n - \hat{a}_n) y_n = 0 \quad (891)$$

7.8

This follows from 7.67 and 7.68.

7.9

$$p(\mathbf{w}, \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{S}_N, \mathbf{S}_N) \quad (892)$$

$$\mathbf{S}_N = \mathbf{S}_N(\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{\Phi}^T \mathbf{t}) \quad (893)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}^T \mathbf{\Phi} \quad (894)$$

We have

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}^T \mathbf{\Phi} = \mathbf{A} + \beta \mathbf{\Phi}^T \mathbf{\Phi} = \Sigma^{-1} \quad (895)$$

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

7.16

7.18

$$\nabla \ln p(\mathbf{w} | \mathbf{t}, \alpha) = \nabla \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) - \mathbf{A} \mathbf{w} \quad (896)$$

$$\nabla \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) = \sum_{n=1}^N t_n \nabla \ln y_n + (1 - t_n) \nabla \ln(1 - y_n) \quad (897)$$

$$\nabla \ln y_n = (1 - y_n) \phi(\mathbf{x}_n) \quad (898)$$

$$\nabla \ln(1 - y_n) = -y_n \phi(\mathbf{x}_n) \quad (899)$$

$$\nabla \ln p(\mathbf{w} | \mathbf{t}, \alpha) = \sum_{n=1}^N \phi(\mathbf{x}_n) [t_n - y_n t_n - y_n + y_n t_n] - \mathbf{A} \mathbf{w} \quad (900)$$

$$= \mathbf{\Phi}^T [\mathbf{t} - \mathbf{y}] - \mathbf{A} \mathbf{w} \quad (901)$$

$$\nabla \nabla \ln p(\mathbf{w} | \mathbf{t}, \alpha) = \nabla \mathbf{\Phi}^T [\mathbf{t} - \mathbf{y}] - \mathbf{A} \quad (902)$$

$$= \nabla \sum_{n=1}^N [t_n - y_n] \phi_n - \mathbf{A} \quad (903)$$

$$= - \sum_{n=1}^N \phi_n \nabla^T y_n - \mathbf{A} \quad (904)$$

$$= - \sum_{n=1}^N \phi_n y_n (1 - y_n) \phi_n^T - \mathbf{A} \quad (905)$$

$$= - \mathbf{\Phi}^T \mathbf{B} \mathbf{\Phi} - \mathbf{A} \quad (906)$$

Which gives the result.

8.1

$$\int p(\mathbf{x})d\mathbf{x} = \int \int p(x_K|pa_K) \prod_{k=1}^{K-1} p(x_k|pa_k)dx_K dx_1 \dots dx_{K-1} \quad (907)$$

$$= \int \int p(x_K|pa_K)dx_K \prod_{k=1}^{K-1} p(x_k|pa_k)dx_1 \dots dx_{K-1} \quad (908)$$

$$= \int \prod_{k=1}^{K-1} p(x_k|pa_k)dx_1 \dots dx_{K-1} \quad (909)$$

$$\vdots \quad (910)$$

$$= \int p(x_1)dx_1 = 1 \quad (911)$$

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 μ_0 ensures the OR-function. The $\mu_i :> 0$ control the increase in probability for that x_i .

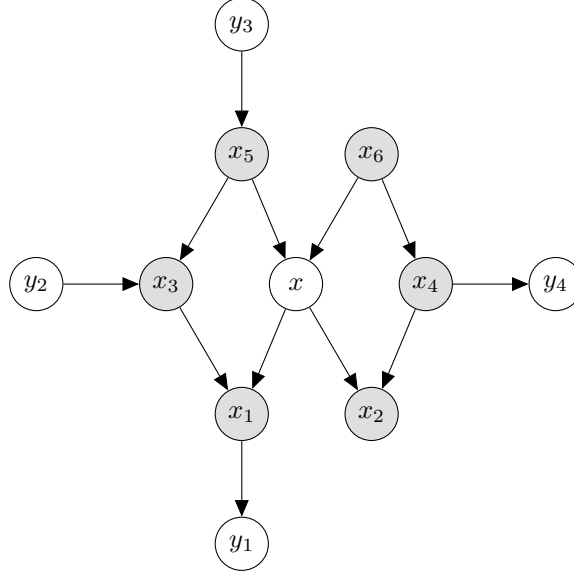
8.8

$$p(a, b, c|d) = \frac{p(a, b, c)}{p(d)} = \frac{p(a|b, c)p(b)p(c)}{p(d)} = \frac{p(a)p(b)p(c)}{p(d)} \quad (912)$$

$$\int \frac{p(a)p(b)p(c)}{p(d)}dc = p(a|d)p(b|d) \quad (913)$$

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

$$p(a, b, c, d) = p(a)p(b)p(c|a, b)p(d|c) \quad (914)$$

$$p(a, b) = \int \int p(a, b, c, d) dc dd = p(a)p(b) \quad (915)$$

Second part:

$$p(a, b, c|d) = \frac{p(a, b, c, d)}{p(d)} = \frac{p(a)p(b)p(c|a, b)p(d|c)}{p(d)} \quad (916)$$

$$p(a, b|d) = \int p(a, b, c|d) dc = \frac{p(a)p(b)p(d|c)}{p(d)} \quad (917)$$

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

$$E(x, y)_{x_k=1} - E(x, y)_{x_k=-1} = h \sum_i x_i - \beta \sum_{i,j} x_i x_j - \eta \sum_i x_i y_i + h - \beta x_k \sum_j x_j - \eta x_k \sum_i y_i \quad (918)$$

$$- h \sum_i x_i - \beta \sum_{i,j} x_i x_j - \eta \sum_i x_i y_i + h - \beta x_k \sum_j x_j - \eta x_k \sum_i y_i \quad (919)$$

$$= 2h - 2\beta x_k \sum_j x_j - \eta x_k \sum_i y_i \quad (920)$$

8.14

$$p(x, y) = \frac{1}{Z} \exp[-E(x, y)] \quad (921)$$

$$\ln p(x, y) = -\ln(Z) - E(x, y) = -\ln(Z) + \eta \sum_i x_i y_i \quad (922)$$

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 μ_k . Therefore, it will always converge.

9.2

9.3

$$p(\mathbf{x}, \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} \pi_k^{z_k} \quad (923)$$

$$p(\mathbf{z}) = \sum_{\mathbf{z}} \prod_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} \quad (924)$$

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:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{old})_{MAP} = \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\boldsymbol{\theta}, \mathbf{z}|\mathbf{x}) \quad (925)$$

$$\propto \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\boldsymbol{\theta}, \mathbf{z}|\mathbf{x}) \quad (926)$$

$$= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \quad (927)$$

$$= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) \quad (928)$$

9.5

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

9.7

$$\nabla_{\mu} \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{kn} \nabla_{\mu} \ln \mathcal{N}(\mathbf{x}_n, \mu_k, \Sigma_k) \quad (929)$$

$$\nabla \ln \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) = -\frac{1}{2} \nabla (\mathbf{x}_n - \mu_k)^T \sigma_k^{-1} (\mathbf{x}_n - \mu_k) = \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \quad (930)$$

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

$$\mathcal{L}(\pi_k, \lambda) \ln p(\mathbf{X}, \mathbf{Z}|\mu_k \Sigma_k \pi_k) + \lambda[-1 + \sum_k \pi_k] \quad (931)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \sum_{n=1}^N z_{nk} \frac{1}{\pi_k} - \lambda \quad (932)$$

$$= \sum_{n=1}^N z_{nk} - \pi_k \lambda \quad (933)$$

$$= \sum_{n=1}^N \sum_{k=1}^K z_{nk} - \pi_k \lambda \quad (934)$$

$$= N - \lambda = 0 \iff \lambda = N \quad (935)$$

Substituting back:

$$\sum_{n=1}^N z_{nk} \frac{1}{\pi_k} = N \iff \pi_k = \frac{1}{N} \sum_{n=1}^N z_{nk} = \frac{N_k}{N} \quad (936)$$

$$\nabla_{\mu_k} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) [\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \quad (937)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \nabla \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (938)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (939)$$

$$\Rightarrow \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (940)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n - \underbrace{\sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\mu}_k}_{N_k} = 0 \quad (941)$$

$$(942)$$

Which leads to the correct answer.

9.9

$$\frac{\partial}{\partial \boldsymbol{\Sigma}_k} = \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) - \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \right] \quad (943)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T - \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln |\boldsymbol{\Sigma}^{-1}|^{\frac{1}{2}} \right] \quad (944)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T - \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{2} \ln |\boldsymbol{\Sigma}^{-1}| \right] \quad (945)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T - \frac{1}{2} \boldsymbol{\Sigma} \right] = 0 \quad (946)$$

$$\Leftrightarrow \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right] = \sum_{n=1}^N \gamma(z_{nk}) \frac{1}{2} \boldsymbol{\Sigma} = N_k \frac{1}{2} \boldsymbol{\Sigma}_k \quad (947)$$

From which the answer is easily seen.

9.11

$$\mathbb{E}_{\mathbf{z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(\mathbf{z}_{nk}) [\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \quad (948)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \gamma(\mathbf{z}_{nk}) [\ln \pi_k + \ln \frac{1}{(2\pi\epsilon)^{M/2}} - \frac{1}{2\epsilon} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2] \quad (949)$$

$$\propto \sum_{n=1}^N \sum_{k=1}^K \gamma(\mathbf{z}_{nk}) [\epsilon \ln \pi_k + \epsilon \ln \frac{1}{(2\pi\epsilon)^{M/2}} - \frac{1}{2} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2] \quad (950)$$

$$= \lim_{\epsilon \rightarrow 0} -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 + C \quad (951)$$

9.12

$$\mathbb{E}[\mathbf{x}] = \int \mathbf{x} p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\pi}) d\mathbf{x} \quad (952)$$

$$= \int \mathbf{x} \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} \quad (953)$$

$$= \sum_{k=1}^K \pi_k \int \mathbf{x} p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} \quad (954)$$

$$= \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k \quad (955)$$

$$\text{cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \quad (956)$$

$$= \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T \quad (957)$$

$$(958)$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \int \mathbf{x}\mathbf{x}^T \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} \quad (959)$$

$$= \sum_{k=1}^K \pi_k \int \mathbf{x}\mathbf{x}^T p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} = \sum_{k=1}^K \pi_k \mathbb{E}_k[\mathbf{x}\mathbf{x}^T] \quad (960)$$

$$(961)$$

$$\mathbb{E}_k[\mathbf{x}\mathbf{x}^T] = \mathbb{E}_k[(\mathbf{x} - \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{x}])^T] \quad (962)$$

$$= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T + (\mathbf{x} - \mathbb{E}[\mathbf{x}]) \mathbb{E}[\mathbf{x}]^T + \mathbb{E}[\mathbf{x}] (\mathbf{x} - \mathbb{E}[\mathbf{x}])^T + \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T] \quad (963)$$

$$= \underbrace{\text{cov}_k(\mathbf{x})}_{=\boldsymbol{\Sigma}_k} + \boldsymbol{\mu}\boldsymbol{\mu}^T \quad (964)$$

9.14

$$p(\mathbf{x}, \mathbf{z}, \boldsymbol{\mu})p(\mathbf{z}|\boldsymbol{\pi}) = \prod_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \pi_k^{z_k} \quad (965)$$

$$= \sum_{\mathbf{z}} \prod_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \pi_k^{z_k} \quad (966)$$

$$(967)$$

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

$$\frac{\partial}{\partial \mu_{ki}} = \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{\partial}{\partial \mu_{ki}} x_{ni} \ln \mu_{ki} + \frac{\partial}{\partial \mu_{ki}} (1 - x_{ni}) \ln(1 - \mu_{ki}) \right] \quad (968)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{x_{ni}}{\mu_{ki}} - \frac{1 - x_{ni}}{1 - \mu_{ki}} \right] \quad (969)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{x_{ni}(1 - \mu_{ki}) - \mu_{ki}(1 - x_{ni})}{\mu_{ki}(1 - \mu_{ki})} \right] = 0 \quad (970)$$

$$\Rightarrow \sum_{n=1}^N \gamma(z_{nk}) \mu_{ki} = \sum_{n=1}^N \gamma(z_{nk}) x_n \quad (971)$$

$$= N_k \mu_{ki} = \sum_{n=1}^N \gamma(z_{nk}) x_n \quad (972)$$

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

$$\frac{\partial}{\partial \alpha} = \frac{M}{2} \frac{1}{\alpha} - \frac{1}{2} \mathbb{E}[\mathbf{w}^T \mathbf{w}] = 0 \quad (973)$$

$$\Rightarrow \frac{M}{\alpha} = \mathbb{E}[\mathbf{w}^T \mathbf{w}] \quad (974)$$

$$\alpha = \frac{M}{\mathbb{E}[\mathbf{w}^T \mathbf{w}]} \quad (975)$$

9.24

$$\ln p(\mathbf{x}|\boldsymbol{\theta}) = \sum_z q(\mathbf{z}) \ln p(\mathbf{x}|\boldsymbol{\theta}) \quad (976)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \quad (977)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \frac{q(\mathbf{z})}{q(\mathbf{z})} \quad (978)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} + \sum_z q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \quad (979)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} - \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} \quad (980)$$

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.

$$\ln p(\mathbf{x}) = \int q(\mathbf{z}) \ln p(\mathbf{x}) d\mathbf{z} \quad (981)$$

$$= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \quad (982)$$

$$= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z})q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})q(\mathbf{z})} \quad (983)$$

$$= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} - \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{x}|\mathbf{z})} \quad (984)$$

$$(985)$$

10.2

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

10.8

Filling in:

$$\mathbb{E}[\tau] = \frac{a}{b} = (a_0 + \frac{N}{2})(b_0 + \frac{1}{2} \mathbb{E}_\mu[\sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^2 + \lambda_0(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^2])^{-1} \underset{N \rightarrow \infty}{\approx} b_N^{-1} \quad (986)$$

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

$$\ln p(\mathbf{x}) = \sum_m \sum_{\mathbf{z}} q(\mathbf{z}, m) \ln p(\mathbf{x}) \quad (987)$$

$$= \sum_m \sum_{\mathbf{z}} q(\mathbf{z}, m) \ln \frac{p(\mathbf{x}, m, \mathbf{z})}{p(m, \mathbf{z}|\mathbf{x})} \frac{q(\mathbf{z}, m)}{q(\mathbf{z}, m)} \quad (988)$$

Which results in 10.35 after rearrangement.

10.15

We have:

$$\mathbb{E}[\pi_k] = \frac{\alpha_k}{\sum_{k=1}^K \alpha_k} \quad (989)$$

and $\alpha_k = \alpha_0 + N_k$.

$$\mathbb{E}[\pi_k] = \frac{\alpha_0 + N_k}{\sum_{k=1}^K \alpha_0 + N_k} \quad (990)$$

$$= K\alpha_0 + \underbrace{\sum_{k=1}^K N_k}_{=N} \quad (991)$$

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

$$\frac{dd \ln(x)}{dx} = -x^{-2} \quad (992)$$

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

Taylor approximation:

$$y(x) = \ln(\xi) + \frac{1}{\xi}(x - \xi) = \lambda x - \ln(\lambda) - 1 \quad \lambda = \frac{1}{\xi} \quad (993)$$

$$g(\lambda) = \min_x [\lambda x - f(x)] \quad (994)$$

$$\frac{d}{dx} = \lambda - \frac{1}{x} = 0 \iff x = \frac{1}{\lambda} \quad (995)$$

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

$$\frac{d}{d\lambda} \lambda x - g(\lambda) = x - \frac{1}{\lambda} \iff x = \frac{1}{\lambda} \quad (996)$$

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

10.30

$$\frac{d}{dx} = \frac{1}{1 + e^{-x}} e^{-x} = \sigma(x) e^{-x} \quad (997)$$

$$\frac{dd}{dx} = \sigma(x)(1 - \sigma(x))e^{-x} - e^{-x}\sigma(x) = -\sigma^2(x)e^{-x} \quad (998)$$

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

$$\frac{d}{dx} = \frac{1}{1 + e^{-x}} e^{-x} \quad (999)$$

$$\frac{dd}{dx} = \sigma(x)(1 - \sigma(x))e^{-x} - e^{-x}\sigma(x) = -\sigma^2(x)e^{-x} \quad (1000)$$

Taylor:

$$f(x) = f(\xi) + f'(\xi)(x - \xi) + \mathcal{O}(\xi^2) \quad (1001)$$

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

$$f(x) \leq -\ln(1 + e^{-xi}) + \sigma(\xi)e^{-\xi}(x - \xi) + \mathcal{O}(\xi^2) \quad (1002)$$

$$= -\ln(1 + e^{-xi}) + \sigma(\xi)e^{-\xi}x - \sigma(\xi)e^{-\xi}\xi + \mathcal{O}(\xi^2) \quad (1003)$$

$$= \lambda x - g(\lambda) \quad \lambda = \sigma(\xi)e^{-\xi} \quad (1004)$$

10.33

$$\frac{d}{d\xi_n} = (1 - \sigma(\xi_n)) - \frac{1}{2} - \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n \lambda'(\xi_n) + 2\xi_n \lambda(\xi_n) + \lambda'(\xi_n) \xi_n^2 \quad (1005)$$

$$= -2\xi_n \lambda(\xi_n) - \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n \lambda'(\xi_n) + 2\xi_n \lambda(\xi_n) + \lambda'(\xi_n) \xi_n^2 \quad (1006)$$

$$= -\lambda'(\xi_n)(\xi_n^2 - \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n) = 0 \quad (1007)$$

$$\iff \xi_n^2 = \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n \quad (1008)$$

10.37

$$q^{\setminus j}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_j(\boldsymbol{\theta})} \quad (1009)$$

$$Z_j = \int q^{\setminus j}(\boldsymbol{\theta}) f_j(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int q^{\setminus j}(\boldsymbol{\theta}) \tilde{f}_j(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int (q(\boldsymbol{\theta}) d\boldsymbol{\theta}) = 1 \quad (1010)$$

$$\tilde{f}_j(\boldsymbol{\theta}) = \frac{q^{new}}{q^{\setminus j}(\boldsymbol{\theta})} = \tilde{f}_j = f_j \quad (1011)$$

11.1

$$\mathbb{E}[\hat{f}] = \mathbb{E}\left[\frac{1}{L} \sum_{l=1}^L f(z^l)\right] \quad (1012)$$

$$= \frac{1}{L} \sum_p (z) f(z^l) dz \quad (1013)$$

$$= \frac{1}{L} \sum_{l=1}^L \mathbb{E}[f] = \mathbb{E}[f] \quad (1014)$$

For the variance, we have to note that

$$\mathbb{E}[f(z^l), f(z^k)] = \text{Var}(f(z)) + \mathbb{E}[f(z)]^2 \quad (1015)$$

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

$$\text{Var}[\hat{f}] = \mathbb{E}[\hat{f}^2] - \mathbb{E}[\hat{f}]^2 \quad (1016)$$

$$= \frac{1}{L^2} \left(\sum_{l=1}^L \sum_{k=1}^L \mathbb{E}[f(z^l) f(z^k)] \right) - \mathbb{E}[f]^2 \quad (1017)$$

$$= \frac{1}{L^2} \left(\sum_{l=1}^L \sum_{k=1}^L \mathbb{E}[f(z^l), f(z^k)] \right) - \mathbb{E}[f]^2 \quad (1018)$$

$$= \frac{1}{L^2} L \text{Var}(f(z)) + L^2 \mathbb{E}[f]^2 - \mathbb{E}[f]^2 \quad (1019)$$

$$= \frac{1}{L} \text{Var}(f(z)) \quad (1020)$$

11.2

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

$$p(y) = 1 \cdot \left| \frac{dz}{dy} \right| = 1 \cdot p(y) \quad (1021)$$

11.3

$$z = \left(\frac{1}{\pi} \int_{-\infty}^y \frac{1}{1 + \hat{y}^2} \right) \quad (1022)$$

$$= \left(\frac{1}{\pi} [\arctan(y) - \arctan(-\infty)] \right) \quad (1023)$$

$$= \left(\frac{1}{\pi} \arctan(y) + \frac{1}{2} \right) \quad (1024)$$

$$y = h^{-1}(z) \quad (1025)$$

$$\iff \pi z - \frac{\pi}{2} = \arctan(y) \quad (1026)$$

$$\iff h^{-1}(z) = \tan\left(\pi z - \frac{\pi}{2}\right) \quad (1027)$$

11.5

Hint 1: Show that the expectation and covariance are equal to μ and Σ .

$$\mathbb{E}[y] = \mathbb{E}_z[\mu] + \mathbb{E}_z[Lz] \quad (1028)$$

$$= \mathbb{E}_z[\mu] + L \mathbb{E}_z[z] \quad (1029)$$

$$= \mathbb{E}_z[\mu] + L\mathbf{0} \quad (1030)$$

$$= \mu \quad (1031)$$

$$\text{cov}[y] = \mathbb{E}[(y - \mu)(y - \mu)^T] \quad (1032)$$

$$= \mathbb{E}[yy^T - \mu y^T - y \mu^T + \mu \mu^T] \quad (1033)$$

$$= \mathbb{E}[(\mu + Lz)(\mu + Lz)^T - \mu(\mu + Lz)^T - (\mu + Lz)\mu^T + \mu \mu^T] \quad (1034)$$

$$= \mathbb{E}[\mu \mu^T + \mu(Lz)^T + Lz \mu^T + Lz(Lz)^T - \mu \mu^T - \mu(Lz)^T - \mu \mu^T - Lz \mu^T + \mu \mu^T] \quad (1035)$$

$$= \mathbb{E}[Lz(Lz)^T] \quad (1036)$$

$$= L \mathbb{E}[zz^T] L^T \quad (1037)$$

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

$$p(y) = p(z) \left| \frac{dz}{dy} \right| \quad (1038)$$

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

$$y = b \tan z + c \iff \frac{y - c}{b} = \tan z \iff \arctan \frac{y - c}{b} = z \quad (1039)$$

$$\frac{dz}{dy} = \frac{1}{1 + \left(\frac{y-c}{b}\right)^2} \frac{1}{b} \quad (1040)$$

$$(1041)$$

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 $E[(z'_i - \mu_i)^2]$.

$$\mathbb{E}_{z,\nu}[z'_i] = \mathbb{E}[\mu_i] + \mathbb{E}[\alpha(z_i - \mu_i)] + \mathbb{E}[\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu] \quad (1042)$$

$$= \mu_i \quad (1043)$$

$$\text{Var}_{z,\nu} = \mathbb{E}[z_i'^2] - \mathbb{E}[z'_i]^2 \quad (1044)$$

$$= \mathbb{E}[(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)] - \mu_i^2 \quad (1045)$$

$$= \mathbb{E}[(\mu_i + \alpha(z_i - \mu_i))(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu) + (\mu_i + \alpha(z_i - \mu_i))(\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)] \quad (1046)$$

$$+ (\mu_i + \alpha(z_i - \mu_i))(\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu) + (\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)(\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)] - \mu_i^2 \quad (1047)$$

$$\vdots \quad (1048)$$

$$(1049)$$

11.15

$$\frac{dH}{dr_i} = \frac{1}{Z} \frac{d}{dr_i} r_i^2 = r_i \quad (1050)$$

$$\frac{dr_i}{d\tau} = -\frac{dE(z)}{dz_i} = -\frac{dH}{dz_i} \quad (1051)$$

11.16

$$p(\mathbf{r}|\mathbf{z}) = \frac{p(\mathbf{r}, \mathbf{z})}{p(\mathbf{z})} = \frac{Z_p}{Z_H} \exp(-K(\mathbf{r})) \quad (1052)$$

11.17

12.3

$$\|u_i\|^2 = [\frac{1}{(N\lambda_i)^{1/2}} X^T v_i]^T [\frac{1}{(N\lambda_i)^{1/2}} X^T v_i] \quad (1053)$$

$$= (N\lambda_i)^{-1} [X^T v_i]^T [X^T v_i] \quad (1054)$$

$$= (N\lambda_i)^{-1} v^T X X^T v_i \quad (1055)$$

$$= \lambda_i^{-1} v^T \lambda_i v_i \quad (1056)$$

$$= 1 \quad (1057)$$

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

$$\frac{\partial \ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2)}{\partial \boldsymbol{\mu}} = -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \quad (1058)$$

$$= \sum_{n=1}^N \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = 0 \quad (1059)$$

$$\Rightarrow 0 = -N\boldsymbol{\mu} + \sum_{n=1}^N \mathbf{x}_n \quad (1060)$$

$$\Rightarrow \boldsymbol{\mu} = \bar{\mathbf{x}} \quad (1061)$$

12.14

$$D(D-1) + 1 - (D-1)(D-2)/2 = D^2 - D + 1 - \frac{1}{2}D^2 + \frac{3}{2}D - 1 = \frac{1}{2}D^2 + \frac{1}{2}D = D(D+1)/2 \quad (1062)$$

12.18

$$M \times D + 1 \quad (1063)$$

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 π_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

$$\frac{\partial Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})}{\boldsymbol{\mu}_k} = \sum_{n=1}^N \gamma(z_{nk}) \frac{1}{p(\mathbf{x}_n|\boldsymbol{\phi}_k)} \frac{\partial}{\partial \boldsymbol{\mu}_k} p(\mathbf{x}_n|\boldsymbol{\phi}_k) \quad (1064)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \frac{1}{p(\mathbf{x}_n|\boldsymbol{\phi}_k)} p(\mathbf{x}_n|\boldsymbol{\phi}_k) \frac{\partial}{\partial \boldsymbol{\mu}_k} - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad (1065)$$

$$= 0 \iff \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (1066)$$

$$\iff \boldsymbol{\mu} = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \quad (1067)$$

$$(1068)$$

$$\frac{\partial Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})}{\boldsymbol{\Sigma}_k} = \sum_{n=1}^N \gamma(z_{nk}) \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln p(\mathbf{x}|\boldsymbol{\phi}_k) \quad (1069)$$

$$(1070)$$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln p(\mathbf{x}|\boldsymbol{\phi}_k) = \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \left[\ln \frac{1}{(2\pi)^{D/2}} + \ln \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} - \frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (1071)$$

$$= \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{2} \ln |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \frac{\partial}{\partial \boldsymbol{\Sigma}_k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) \quad (1072)$$

$$= \frac{1}{2} \boldsymbol{\Sigma} - \frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \quad (1073)$$

$$\frac{\partial}{\partial \Sigma_k} = \sum_{n=1}^N \gamma(z_{nk}) \left(\frac{1}{2} \Sigma - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right) = 0 \quad (1074)$$

$$\iff \Sigma_k = \frac{\sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma(z_{nk})} \quad (1075)$$

13.17

$$h(z_1) = p(z_1|u_1)p(x_1|z_1, u_1) \quad (1076)$$

$$f_n(z_{n-1}, z_n) = p(z_n|z_{n-1}, u_n)p(x_n|z_n, u_n) \quad (1077)$$

13.19

This follows from:

$$\frac{d}{dz} \ln p(\mathbf{Z}) = \frac{d}{dz} \ln p(z_0)p(z_1|z_0) \times \cdots \times p(z_n|z_{n-1}) = \begin{bmatrix} \frac{\partial}{\partial z_0} \ln p(z_0) \\ \frac{\partial}{\partial z_1} \ln p(z_1|z_0) \\ \vdots \\ \frac{\partial}{\partial z_n} p(z_n|z_{n-1}) \end{bmatrix} \quad (1078)$$

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 Γ . I don't see how this term doesn't influence the distribution.

14.2

$$\mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M \epsilon_m(x) \right)^2 \right] = \mathbb{E} \left[\left(\frac{1}{M^2} \sum_{m=1}^M \sum_{l=1}^M \epsilon_m(x) \epsilon_l(x) \right) \right] \quad (1079)$$

$$= \mathbb{E} \left[\frac{1}{M^2} \left(\sum_{m=1}^M \sum_{\substack{l=1 \\ l \neq m}}^M \epsilon_m(x) \epsilon_l(x) + \sum_{l=1}^M \epsilon_l(x)^2 \right) \right] \quad (1080)$$

$$= \frac{1}{M^2} \sum_{l=1}^M \mathbb{E}[\epsilon_l(x)^2] \quad (1081)$$

$$= \frac{1}{M} E_{AV} \quad (1082)$$

14.3

Through Jensen's inequality:

$$\sum_{m=1}^M \frac{1}{M} \epsilon_m(x)^2 \geq \left(\sum_{m=1}^M \frac{1}{M} \epsilon_m(x) \right)^2 \quad (1083)$$

Therefore:

$$\mathbb{E}_x \left[\sum_{m=1}^M \frac{1}{M} \epsilon_m(x)^2 \right] \geq \mathbb{E} \left[\left(\sum_{m=1}^M \frac{1}{M} \epsilon_m(x) \right)^2 \right] = E_{com} \quad (1084)$$

14.6

$$\frac{d}{d\alpha_m} = (e^{\alpha_m/2} + e^{-\alpha/2}) \sum_{n=1}^N w_n^{(m)} I(y_m(x_n) = t_n) = e^{\alpha_m/2} \sum_{n=1}^N w_n^{(m)} = 0 \quad (1085)$$

$$\frac{e^{-\alpha_m/2}}{e^{\alpha_m/2} + e^{-\alpha_m/2}} = \epsilon_m \quad (1086)$$

Rearranging this results in the desired formula.

[14.7](#) [14.8](#)

14.9

$$\mathcal{L}_m(\mathbf{x}_n) = (y_n - \sum_{l=1}^M \alpha_l \hat{y}_l)^2 \quad (1087)$$

$$= (y - \underbrace{\sum_{l=1}^{M-1} \alpha_l \hat{y}_l}_{\text{residual}} + \alpha_M \hat{y}_M)^2 \quad (1088)$$

$$(1089)$$

14.10

$$L = \frac{1}{2} \sum_{n=1}^N (t_n - t)^2 \quad (1090)$$

$$\frac{dL}{dt} = - \sum_{n=1}^N (t_n - t) = 0 \iff Nt = \sum_{n=1}^N t_n \iff t = \frac{1}{N} \sum_{n=1}^N t_n \quad (1091)$$

14.13

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) = \prod_{n=1}^N \prod_{k=1}^K (\pi_k \mathcal{N}(t_n | \mathbf{w}_k^T \boldsymbol{\phi}_n, \beta^{-1})^{z_{nk}} \quad (1092)$$

Taking the log of this will yield the result.

14.14

$$\frac{dL}{d\pi_k} = \sum_{n=1}^N \gamma_{nk} / \pi_k - \lambda = 0 \quad (1093)$$

$$\iff \sum_k \sum_{n=1}^N \gamma_{nk} = \sum_k \pi_k \lambda \quad (1094)$$

$$(1095)$$

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

$$\mathbb{E}[t | \hat{\boldsymbol{\phi}}, \boldsymbol{\theta}] = \sum_{k=1}^K \pi_k \mathbb{E}[t | \hat{\boldsymbol{\phi}}, \mathbf{w}_k, \beta] \quad (1096)$$