

Bishop Questions

David Ruhe

April 23, 2020

2.1

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

$$\sum_{x=0}^1 p(x|\mu) = 1 - \mu + \mu = 1 \quad (2)$$

$$\mathbb{E}[x] = \sum_{x=0}^1 x\mu^x(1-\mu)^{1-x} \quad (3)$$

$$= \mu(1-\mu)^0 = \mu \quad (4)$$

$$\text{Var}[x] = [(x - \mathbb{E}[x])^2] \quad (5)$$

$$= [(x - \mu)^2] \quad (6)$$

$$= \sum_{i=0}^1 (x - \mu)^2 \mu^x(1-\mu)^{1-x} \quad (7)$$

$$(8)$$

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

$$H[x] = \sum i = 0^1 p(x) \ln p(x) = \mu \ln \mu + (1 - \mu) \ln(1 - \mu) \quad (9)$$

2.2

$$\sum_x p(x|\mu) = \sum_x \left(\frac{1-\mu}{2}\right)^{(1-x)/2} \left(\frac{1+\mu}{2}\right)^{(1+x)/2} \quad (10)$$

$$= \frac{1-\mu}{2} + \frac{1+\mu}{2} = 1 \quad (11)$$

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

$$\mathbb{E}[x] = \sum_x xp(x) = \frac{-1+\mu}{2} + \frac{1+\mu}{2} = \mu \quad (12)$$

$$\mathbb{E}[(x - \mathbb{E}[x])^2] = \mathbb{E}[x^2] - 2\mathbb{E}[x\mathbb{E}[x]] + \mathbb{E}[\mathbb{E}[x]^2] \quad (13)$$

$$= 1 - 2\mu\mathbb{E}[x] + \mu^2 \quad (14)$$

$$= 1 - 2\mu^2 + \mu^2 = 1 - \mu^2 \quad (15)$$

$$H[x] = -\sum_x p(x) \ln p(x) = \frac{1+\mu}{2} \ln \frac{1+\mu}{2} + \frac{1-\mu}{2} \ln \frac{1-\mu}{2} \quad (16)$$

2.6

$$p(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \quad (17)$$

$$\int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \quad (18)$$

$$\Gamma(x) = \int_0^\infty \mu^{x-1} e^{-u} du \quad (19)$$

$$\mathbb{E}[\mu] = \int_0^1 p(\mu|a, b) \mu d\mu \quad (20)$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \mu d\mu \quad (21)$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^a (1-\mu)^{b-1} d\mu \quad (22)$$

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

$$= \frac{a}{a+b} \int_0^1 \frac{\Gamma(a+1+b)}{\Gamma(a+1)\Gamma(b)} \mu^a (1-\mu)^{b-1} d\mu \quad (23)$$

$$= \frac{a}{a+b} \quad (24)$$

Where we used

$$\Gamma(a+1) = a\Gamma(a) \text{ and } \Gamma(a+b+1) = (a+b)\Gamma(a+b) \quad (25)$$

$$\text{Var}[\mu] = \mathbb{E}[\mu^2] - \mathbb{E}[\mu]^2 \quad (26)$$

$$\mathbb{E}[\mu] = \int_0^1 p(\mu|a, b) \mu d\mu^2 \quad (27)$$

$$= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a+1} (1-\mu)^{b-1} \mu d\mu \quad (28)$$

Which looks like Beta with a+2

$$= \frac{a}{a+b} \int_0^1 \frac{\Gamma(a+1+b)}{\Gamma(a+1)\Gamma(b)} \mu^{a+2} (1-\mu)^{b-1} d\mu \quad (29)$$

$$= \frac{a(a+1)}{(a+b)(a+1+b)} \int_0^1 \frac{\Gamma(a+2+b)}{\Gamma(a+2)\Gamma(b)} \mu^{a+2} (1-\mu)^{b-1} d\mu \quad (30)$$

$$= \frac{a(a+1)}{(a+b)(a+1+b)} \quad (31)$$

$$\text{Var}[\mu] = \frac{a(a+1)}{(a+b)(a+1+b)} - \frac{a^2}{(a+b)^2} \quad (32)$$

$$= \frac{(a+b)a(a+1)}{(a+b)^2(a+1+b)} - \frac{(a+b+1)a^2}{(a+b+1)(a+b)^2} \quad (33)$$

$$= \frac{(a^2+ab)(a+1)}{(a+b)^2(a+1+b)} - \frac{(a+b+1)a^2}{(a+b+1)(a+b)^2} \quad (34)$$

$$= \frac{a^3+a^2b+a^2+ab}{(a+b)^2(a+1+b)} - \frac{a^3+ba^2+a^2}{(a+b+1)(a+b)^2} \quad (35)$$

Which gives you the desired result.

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

$$\frac{\partial p(\mu|a, b)}{\partial \mu} = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} [(a-1)\mu^{a-2} - \mu^{a-1}(b-1)(1-\mu)^{b-2}] = 0 \quad (36)$$

$$(a-1)\mu^{a-2}(1-\mu)^{b-1} = \mu^{a-1}(b-1)(1-\mu)^{b-2} \quad (37)$$

$$\frac{a-1}{\mu} = \frac{b-1}{1-\mu} \quad (38)$$

$$\mu = \frac{-a+1}{-a-b+2} = \frac{a-1}{a+b-2} \quad (39)$$

2.8

$$\mathbb{E}[x] = \int p(x)xdx \quad (40)$$

$$= \int \int p(x,y)xdxdy \quad (41)$$

$$= \int \int p(x|y)p(y)xdxdy \quad (42)$$

$$= \int p(y) \int p(x|y)xdxdy \quad (43)$$

$$= \int p(y) \mathbb{E}_x[x|y]dy \quad (44)$$

$$= \mathbb{E}_y \mathbb{E}_x[x|y] \quad (45)$$

$$\mathbb{E}_y[\text{Var}_x[x|y]] = \mathbb{E}_y \mathbb{E}_{x|y}[(x - \mathbb{E}[x|y])^2] \quad (46)$$

$$= \mathbb{E}_y \mathbb{E}_{x|y}[(x^2 + \mathbb{E}[x|y]^2 - 2x \mathbb{E}[x|y])] \quad (47)$$

$$= \mathbb{E}_y \mathbb{E}_{x|y} x^2 + \mathbb{E}_y \mathbb{E}_{x|y} \mathbb{E}[x|y]^2 - 2 \mathbb{E}_y \mathbb{E}_{x|y} x \mathbb{E}[x|y]] \quad (48)$$

$$= \mathbb{E}_y \int p(x|y)x^2dx + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \int p(y) \int p(x|y)x \mathbb{E}[x|y]dxdy \quad (49)$$

$$= \int p(y) \int p(x|y)x^2dxdy + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \int \mathbb{E}[x|y] \int p(x|y)xdxdy \quad (50)$$

$$= \int \int p(x,y)x^2dxdy + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \int \mathbb{E}[x|y]^2dy \quad (51)$$

$$= \int \int p(x,y)x^2dxdy + \mathbb{E}_y \mathbb{E}[x|y]^2 - 2 \mathbb{E}_y \mathbb{E}[x|y]^2] \quad (52)$$

$$= \mathbb{E}[x^2] - \mathbb{E}_y \mathbb{E}[x|y]^2 \quad (53)$$

$$(54)$$

$$\text{Var}_y[\mathbb{E}_x[x|y]] = \mathbb{E}_y [(\mathbb{E}_x[x|y] - \mathbb{E}_y \mathbb{E}_x[x|y])^2] \quad (55)$$

$$= \mathbb{E}_y [(\mathbb{E}_x[x|y] - \mathbb{E}_x[x])^2] \quad (56)$$

$$= \mathbb{E}_y [(\mathbb{E}_x[x|y]^2 + \mathbb{E}_x[x]^2 - 2 \mathbb{E}_x[x] \mathbb{E}_x[x|y])] \quad (57)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] + \mathbb{E}_y[\mathbb{E}_x[x]^2] - 2 \mathbb{E}_y \mathbb{E}_x[x] \mathbb{E}_x[x|y] \quad (58)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] + \mathbb{E}_x[x]^2 - 2 \mathbb{E}_x[x] \mathbb{E}_y \mathbb{E}_x[x|y] \quad (59)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] + \mathbb{E}_x[x]^2 - 2 \mathbb{E}_x[x]^2 \quad (60)$$

$$= \mathbb{E}_y[\mathbb{E}_x[x|y]^2] - \mathbb{E}_x[x]^2 \quad (61)$$

$$(62)$$

Sum of these expressions gives desired result.

2.11

2.12

$$\int_a^b \frac{1}{b-a} dx = \frac{x}{b-a} - \frac{x}{b-a} \Big|_a^b = \frac{b-a}{b-a} = 1 \quad (63)$$

$$\mathbb{E}[x] = \int_a^b p(x) x dx \quad (64)$$

$$= \int_a^b \frac{x}{b-a} dx \quad (65)$$

$$= \frac{x^2}{2(b-a)} \Big|_a^b \quad (66)$$

$$= \frac{b^2 - a^2}{2(b-a)} = \frac{(b-a)(b+a)}{2(b-a)} = \frac{b+a}{2} \quad (67)$$

$$\text{Var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 \quad (68)$$

$$= \int_a^b \frac{x^2}{b-a} dx - \frac{(b+a)^2}{4} \quad (69)$$

$$= \frac{\frac{1}{3}x^3}{b-a} \Big|_a^b - \frac{(b+a)^2}{4} \quad (70)$$

$$= \frac{\frac{1}{3}b^3}{b-a} - \frac{\frac{1}{3}a^3}{b-a} - \frac{(b+a)^2}{4} \quad (71)$$

$$= \frac{(b-a)(b^2 + a^2 + ab)}{3(b-a)} - \frac{(b+a)^2}{4} \quad (72)$$

$$= \frac{b^2 + a^2 + ab}{3} - \frac{b^2 + a^2 + 2ab}{4} \quad (73)$$

$$= \frac{4b^2 + 4a^2 + 4ab}{12} - \frac{3b^2 + 3a^2 + 6ab}{12} \quad (74)$$

$$= \frac{b^2 + a^2 - 2ab}{12} \quad (75)$$

$$= \frac{(b-a)^2}{12} \quad (76)$$

$$(77)$$

2.17

$$A = \frac{1}{2} \underbrace{A + A^T}_{\text{Symmetric}} + \frac{1}{2} \underbrace{A - A^T}_{\text{Non-symmetric}} \quad (78)$$

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp -\frac{1}{2}(x - \mu)^T \frac{1}{2}(\Sigma^{-1} + \Sigma^{-T} + \Sigma^{-1} - \Sigma^{-T})(x - \mu) \quad (79)$$

$$= \exp -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M (x_n - \mu_n)(\Sigma_{nm}^{-1} + \Sigma_{nm}^{-T})(x_m - \mu_m) \quad (80)$$

$$- \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M (x_n - \mu_n) \underbrace{(\Sigma_{nm}^{-1} - \Sigma_{nm}^{-T})}_{=0} (x_m - \mu_m) \quad (81)$$

Which leaves us only with a symmetric covariance matrix.

2.21

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

$$\# \text{parameters} = D + (D - 1) + (D - 2) + \dots + D - (D - 2) + D - (D - 1) \quad (82)$$

$$= (D + 1) + (D + 1) + \dots + (D + 1) \quad (83)$$

$$= \frac{D}{2}(D + 1) \quad (84)$$

2.22

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

$$A^{-1}AA^{T^{-1}} = A^{-1}A^T A^{T^{-1}} \quad (85)$$

$$A^{-1T} = A^{-1}A^T A^{T^{-1}} \quad (86)$$

$$= A^{-1} \quad (87)$$

2.27

$$\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{x} + \mathbf{z}] \quad (88)$$

$$= \int \int p(\mathbf{x}, \mathbf{z})(\mathbf{x} + \mathbf{z}) d\mathbf{x} d\mathbf{z} \quad (89)$$

$$= \int \int p(\mathbf{x}, \mathbf{z})\mathbf{x} + p(\mathbf{x}, \mathbf{z})\mathbf{z} d\mathbf{x} d\mathbf{z} \quad (90)$$

$$= \int \int p(\mathbf{x})p(\mathbf{z})\mathbf{x} + p(\mathbf{z})p(\mathbf{x})\mathbf{z} d\mathbf{x} d\mathbf{z} \quad (91)$$

$$= \int p(\mathbf{z}) \int p(\mathbf{x})\mathbf{x} d\mathbf{x} d\mathbf{z} + \int p(\mathbf{x}) \int p(\mathbf{z})\mathbf{z} d\mathbf{x} d\mathbf{z} \quad (92)$$

$$= \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{z}] \quad (93)$$

2.30

$$\mathbb{E}[z] = R^{-1} \begin{pmatrix} \Lambda\mu - A^T Lb \\ Lb \end{pmatrix} \quad (94)$$

$$= \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} A^T \\ A\Lambda^{-1} & L^{-1} + A\Lambda^{-1} A^T \end{pmatrix} \begin{pmatrix} \Lambda\mu - A^T Lb \\ Lb \end{pmatrix} \quad (95)$$

$$= \begin{pmatrix} \Lambda^{-1}(\Lambda\mu - A^T Lb) + \Lambda^{-1} A^T Lb \\ A\Lambda^{-1}(\Lambda\mu - A^T Lb) + \mathcal{L}^{-1} Lb + A\Lambda^{-1} A^T Lb \end{pmatrix} \quad (96)$$

$$= \begin{pmatrix} \mu - \Lambda^{-1} A^T Lb + \Lambda^{-1} A^T Lb \\ A\mu - A\Lambda^{-1} A^T Lb + b + A\Lambda^{-1} A^T Lb \end{pmatrix} \quad (97)$$

Which gives the desired result.

2.41

$$\Gamma(x) = \int_0^\infty u^{x-1} \exp(-u) du \quad (98)$$

$$\int_0^\infty \text{Gam}(\lambda|a, b) d\lambda = \int_0^\infty \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) d\lambda \quad (99)$$

$$= \frac{1}{\Gamma(a)} b^a \int_0^\infty \lambda^{a-1} \exp(-b\lambda) d\lambda \quad (100)$$

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

$$= \frac{\Gamma(a)}{\Gamma(a)} = 1 \quad (101)$$

2.46

$$a = \nu/2 \quad b = \frac{\nu}{2\lambda} \quad (102)$$

$$p(x|\mu, a, b) = \frac{b^a}{\Gamma(a)} \left(\frac{1}{2a} \right)^{1/2} \left[b + \frac{1}{2}(x - \mu)^2 \right]^{-a-\frac{1}{2}} \Gamma(a + \frac{1}{2}) \quad (103)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left(\frac{1}{2\pi} \right)^{1/2} \left[\frac{\nu}{2\lambda} + \frac{1}{2}(x - \mu)^2 \right]^{-\frac{\nu}{2}-\frac{1}{2}} \left(\frac{\nu}{2\lambda} \right)^{\nu/2} \quad (104)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left[\frac{2\lambda}{\nu} \left(\frac{\nu}{2\lambda} + \frac{1}{2}(x - \mu)^2 \right) \right]^{-\frac{\nu}{2}} \left[2\pi \left(\frac{\nu}{2\lambda} + \frac{1}{2}(x - \mu)^2 \right) \right]^{-\frac{1}{2}} \quad (105)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left[1 + \frac{\lambda(x - \mu)^2}{\nu} \right]^{-\frac{\nu}{2}} \left[\frac{\pi\nu}{\lambda} + \pi(x - \mu)^2 \right]^{-\frac{1}{2}} \quad (106)$$

$$= \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left[1 + \frac{\lambda(x - \mu)^2}{\nu} \right]^{-\frac{\nu}{2}} \left[1 + \frac{\lambda}{\nu}(x - \mu)^2 \right]^{-\frac{1}{2}} \left(\frac{\lambda}{\pi\nu} \right)^{1/2} \quad (107)$$

$$(108)$$

2.47

$$\text{St}(x|\mu, \lambda, \nu) \propto \left[1 + \frac{\lambda(x - \mu)^2}{\nu}\right]^{-v/2-1/2} \quad (109)$$

$$= \exp \ln \left[1 + \frac{\lambda(x - \mu)^2}{\nu}\right]^{-v/2-1/2} \quad (110)$$

$$= \exp \left[\frac{-v-1}{2} \ln \left(1 + \frac{\lambda(x - \mu)^2}{\nu}\right) \right] \quad (111)$$

$$= \exp \left[\frac{-v-1}{2} \left(\frac{\lambda(x - \mu)^2}{\nu} - \mathcal{O} \left[\frac{\lambda(x - \mu)^2}{2\nu} \right]^2 \right) \right] \quad (112)$$

$$= \exp \left[\frac{-v-1}{2} \left(\frac{\lambda(x - \mu)^2}{\nu} - \mathcal{O} \left[\frac{\lambda^2(x - \mu)^4}{4\nu^2} \right] \right) \right] \quad (113)$$

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

2.48

$$\Gamma(x) = \int \mu^{x-1} e^{-u} du \quad (114)$$

$$|(\eta\Lambda)^{-1}|^{1/2} = [\eta^{-D}|\Lambda^{-1}|]^{1/2} = \eta^{-D/2}|\Lambda|^{-1/2} \quad (115)$$

$$\Gamma(D/2 + \nu/2) = \int \eta^{D/2+\nu/2-1} \exp(-\eta) d\eta \quad (116)$$

$$z = \eta/2(\Delta^2 + \nu) \iff \eta = 2 \frac{z}{\Delta^2 + \nu} \quad (117)$$

$$d\eta = \frac{d\eta}{dz} dz = \frac{z}{\Delta^2 + \nu} dz \quad (118)$$

$$\int \mathcal{N}(x|\mu, (\eta\Lambda)^{-1}) \text{Gam}\left(\eta\left|\frac{\nu}{2}, \frac{\nu}{2}\right.\right) d\eta \quad (119)$$

$$= \int \frac{1}{(2\pi)^{D/2}} \frac{1}{|(\eta\Lambda)^{-1}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T(\eta\Lambda)(\mathbf{x} - \mu)\right\} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (120)$$

$$= \int \frac{1}{(2\pi)^{D/2}} \frac{1}{|(\eta\Lambda)^{-1}|^{1/2}} \exp\left\{-\frac{\eta}{2}\Delta^2\right\} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (121)$$

$$= \int \frac{1}{(2\pi)^{D/2}} \eta^{D/2} |\Lambda|^{1/2} \exp\left\{-\frac{\eta}{2}\Delta^2\right\} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (122)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \int \eta^{D/2} \exp\left\{-\frac{\eta}{2}\Delta^2\right\} \eta^{\frac{\nu}{2}-1} \exp(-\nu/2\eta) d\eta \quad (123)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} \int \eta^{D/2+\nu/2-1} \exp\left\{-\frac{\eta}{2}(\Delta^2 + \nu)\right\} d\eta \quad (124)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} 2(\Delta^2 + \nu)^{-D/2-\nu/2} \int z^{1/2D+\nu/2} \exp\{-z\} dz \quad (125)$$

$$= \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \frac{1}{\Gamma(\nu/2)} (\nu/2)^{\nu/2} 2(\Delta^2 + \nu)^{-D/2-\nu/2} \Gamma(D/2 + \nu/2) \quad (126)$$

$$(127)$$

Which works out algebraically to the desired result.

2.50

$$\text{St}(\mathbf{x}|\mu, \Lambda, \nu) \propto \left[1 + \frac{1}{\nu}\Delta^2\right]^{-D/2-\nu/2} \quad (128)$$

$$= \exp \frac{-D-\nu}{2} \ln \left[1 + \frac{1}{\nu}\Delta^2\right] \quad (129)$$

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

$$= \exp \frac{-D-\nu}{2} \frac{1}{\nu}\Delta^2 - \frac{-D-\nu}{2} \mathcal{O}(1/\nu^2) \quad (130)$$

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

2.51

$$1 = [\cos(A) + i \sin(A)][\cos(A) - i \sin(A)] \quad (131)$$

$$= \cos^2(A) - \sin^2(A) \quad (132)$$

$$\cos(A - B) = \Re \exp[i(A - B)] \quad (133)$$

$$= \Re \exp(iA) \exp(-iB) \quad (134)$$

$$= \Re[\cos(A) + i \sin(A)][\cos(B) - i \sin(B)] \quad (135)$$

$$= \cos(A) \cos(B) - \sin(A) \sin(B) \quad (136)$$

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

2.54

$$0 = \sum_{n=1}^N \cos(\theta_0) \sin(\theta_n) - \cos(\theta_n) \sin(\theta_0) \quad (137)$$

$$= \cos(\theta_0) \sum_{n=1}^N \sin(\theta_n) - \sin(\theta_0) \sum_{n=1}^N \cos(\theta_n) \quad (138)$$

$$(139)$$

$$\frac{\sin(\theta_0)}{\cos(\theta_0)} = \frac{\sum_{n=1}^N \sin(\theta_n)}{\sum_{n=1}^N \cos(\theta_n)} \quad (140)$$

$$\tan(\theta_0) = \frac{\sum_{n=1}^N \sin(\theta_n)}{\sum_{n=1}^N \cos(\theta_n)} \quad (141)$$

$$\theta_0 = \arctan \left\{ \frac{\sum_{n=1}^N \sin(\theta_n)}{\sum_{n=1}^N \cos(\theta_n)} \right\} \quad (142)$$

2.54

$$\frac{\partial p}{\partial \theta} = -(2\pi I_0(m))^{-1} \exp\{m \cos(\theta - \theta_0)\} (m \sin(\theta - \theta_0)) \quad (143)$$

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

$$\frac{\partial \partial p}{\partial \theta} = \frac{1}{2\pi I_0(m)} [-\exp(m \cos(\theta - \theta_0))(m \sin(\theta - \theta_0))^2 - m \cos(\theta - \theta_0) \exp(m \cos(\theta - \theta_0))] \quad (144)$$

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

$$A(m_{ML}) = \left(\frac{1}{N} \sum_{n=1}^N \cos \theta_n \right) \cos \theta_0^{ML} - \left(\frac{1}{N} \sum_{n=1}^N \sin \theta_n \right) \sin \theta_0^{ML} \quad (145)$$

$$= \bar{r} (\cos \bar{\theta} \cos \theta_0^{ML} - \sin \bar{\theta} \sin \theta_0^{ML}) \quad (146)$$

$$= \bar{r} (\cos(\bar{\theta} - \theta_0^{ML})) \quad (147)$$

$$= \bar{r} \quad (148)$$

2.57

For this question, the following knowledge is necessary:

$$a^T B a = \underbrace{B : aa^T}_{\text{Frobenius product (Hadamard \& sum)}} = \text{vec}(B)^T \text{vec}(aa^T) \quad (149)$$

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

$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (150)$$

$$= (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp \left[-\frac{1}{2} (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right] \quad (151)$$

$$= (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp(-1/2 \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \exp \left[\begin{pmatrix} -1/2 \text{vec}(\boldsymbol{\Sigma}^{-1}) \\ -\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \end{pmatrix} \begin{pmatrix} \text{vec}(\mathbf{x} \mathbf{x}^T) & \mathbf{x} \end{pmatrix} \right] \quad (152)$$

$$= (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{1/2} \exp(-1/2 \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \exp [\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})] \quad (153)$$

$$(154)$$

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

$$g(\eta) = |\text{vec}^{-1}(-2\eta_1)|^{-1/2} \eta_2 \text{vec}^{-1}(-2\eta_1) \eta_2^T \quad (155)$$

and

$$h(x) = (2\pi)^{-D/2} \quad (156)$$

2.59

$$\int_0^\infty \frac{1}{x} f\left(\frac{1}{x}\right) dx = \int_0^\infty \frac{1}{x} f(y) x dy = \int_0^\infty f(y) dy = 1 \quad (157)$$

2.61

3.1

$$2\sigma(2a) - 1 = \frac{2}{1 + e^{-2a}} - 1 \quad (158)$$

$$= \frac{e^a}{e^a} \frac{2}{1 + e^{-2a}} - 1 \quad (159)$$

$$= \frac{2e^a}{e^a + e^{-a}} - 1 \quad (160)$$

$$= \frac{2e^a}{e^a + e^{-a}} - \frac{e^a + e^{-a}}{e^a + e^{-a}} \quad (161)$$

$$= \tanh(a) \quad (162)$$

$$a = \frac{x - \mu_j}{s} \quad (163)$$

$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M w_j \sigma(a) \quad (164)$$

$$= w_0 + \sum_{j=1}^M \frac{w_j}{2} \tanh\left(\frac{1}{2}a\right) + w_j \quad (165)$$

$$= w_0 + \sum_{j=1}^M + \sum_{j=1}^M \frac{w_j}{2} \tanh(a') \quad (166)$$

$$= u_0 + \sum_{j=1}^M u_j \tanh(a') \quad (167)$$

3.3

$$\frac{dE_D(\mathbf{w})}{d\mathbf{w}} = \sum_{n=1}^N r_n [t_n - \mathbf{w}^T \phi(\mathbf{x}_n)] \phi(\mathbf{x}_n) = 0 \quad (168)$$

$$\sum_{n=1}^N r_n t_n \phi(\mathbf{x}_n) = \sum_{n=1}^N \mathbf{w}^T \phi(\mathbf{x}_n) \phi(\mathbf{x}_n) \quad (169)$$

$$= \sum_{n=1}^N \phi(\mathbf{x}_n)^T \mathbf{w} \phi(\mathbf{x}_n) \quad (170)$$

$$= \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \mathbf{w} \quad (171)$$

$$= \Phi^T \Phi \mathbf{w} \quad (172)$$

$$([\mathbf{r} \cdot \mathbf{t}]^T \Phi)^T = \Phi^T \Phi \mathbf{w} \quad (173)$$

$$\Phi^T [\mathbf{r} \cdot \mathbf{t}] = \Phi^T \Phi \mathbf{w} \quad (174)$$

$$(\Phi^T \Phi)^{-1} \Phi^T [\mathbf{r} \cdot \mathbf{t}] = \mathbf{w} \quad (175)$$

(i) ?

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

3.4

$$\mathbb{E}[E_d(\mathbf{w})] = \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2\right] \quad (176)$$

$$= \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (w_0 + \sum_{i=1}^D w_i (x_i + \epsilon_i) - t_n)^2\right] \quad (177)$$

$$= \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) + \sum_{i=1}^D w_i \epsilon_i - t_n)^2\right] \quad (178)$$

$$= \mathbb{E}\left[\frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2 + \underbrace{\left(\sum_{i=1}^D w_i \epsilon_i\right)^2 + 2(y(\mathbf{x}_n, \mathbf{w}) - t_n) \sum_{i=1}^D w_i \epsilon_i}_{=0(\mathbb{E}[\epsilon_i]=0)}\right] \quad (179)$$

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

$$\mathbb{E}\left(\sum_{i=1}^D w_i \epsilon_i\right)^2 = \mathbb{E}\left[\sum_{i=1}^D \sum_{j=1}^D w_j \epsilon_j w_i \epsilon_i\right] = \sum_{i=1}^D w_i^2 \mathbb{E} \epsilon_i^2 = \sum_{i=1}^D w_i^2 \mathbb{E} \epsilon_i^2 = \sum_{i=1}^D w_i^2 (\sigma^2 + \underbrace{E[\epsilon_i^2]}_{=0}) \quad (180)$$

Which gives us our desired result.

3.5

$$\mathcal{L}(\mathbf{w}, \lambda) = \frac{1}{2} \sum_{n=1}^N [t_n - \mathbf{w}^T \phi(\mathbf{x}_n)]^2 + \lambda \left[\sum_{j=1}^M |w_j|^q - \eta \right] \quad (181)$$

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

3.6

$$p(\mathbf{T}|\mathbf{W}, \Sigma) = \prod_{n=1}^N p(\mathbf{t}_n|\mathbf{W}, \Sigma) \quad (182)$$

$$\ln p(\mathbf{T}|\mathbf{W}, \mathbf{\Sigma}) = \ln \prod_{n=1}^N p(\mathbf{t}_n|\mathbf{W}, \mathbf{\Sigma}) = \sum_{n=1}^N \ln p(\mathbf{t}_n|\mathbf{W}, \mathbf{\Sigma}) \quad (183)$$

$$= \sum_{n=1}^N \ln \left((2\pi)^{-D/2} |\mathbf{\Sigma}|^{-1/2} \right) + \sum_{n=1}^N \frac{1}{2} (\mathbf{y}(\mathbf{x}, \mathbf{W}) - \mathbf{t})^T \mathbf{\Sigma}^{-1} (\mathbf{y}(\mathbf{x}, \mathbf{W}) - \mathbf{t}) \quad (184)$$

$$\frac{\partial p}{\partial \mathbf{W}} = \frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial \mathbf{W}} [(\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t})^T \mathbf{\Sigma}^{-1} (\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t})] \quad (185)$$

$$= \frac{1}{2} \sum_{n=1}^N (\mathbf{\Sigma}^{-1} + \mathbf{\Sigma}^{-1T}) (\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t}) \phi(\mathbf{x})^T = 0 \quad (186)$$

$$= \sum_{n=1}^N \mathbf{\Sigma}^{-1} (\mathbf{W}^T \phi(\mathbf{x}) - \mathbf{t}) \phi(\mathbf{x})^T \quad (187)$$

$$\mathbf{\Sigma}^{-1} \sum_{n=1}^N \mathbf{W}^T \phi(\mathbf{x}) \phi(\mathbf{x})^T = \mathbf{\Sigma}^{-1} \sum_{n=1}^N \mathbf{t} \phi(\mathbf{x})^T \quad (188)$$

$$\mathbf{W}^T \mathbf{\Phi}^T \mathbf{\Phi} = \mathbf{T} \mathbf{\Phi}^T \quad (189)$$

$$\mathbf{W}^T = \mathbf{T} \mathbf{\Phi}^T (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \quad (190)$$

$$\mathbf{W} = (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \mathbf{\Phi}^T \mathbf{T} \quad (191)$$

$$\frac{d \ln p(\mathbf{T}|\mathbf{x}, \mathbf{W}, \mathbf{\Sigma})}{d \mathbf{\Sigma}} = -\frac{N}{2} \frac{d}{d \mathbf{\Sigma}} \ln |\mathbf{\Sigma}| - \frac{1}{2} \sum_{n=1}^N \frac{d}{d \mathbf{\Sigma}} (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \mathbf{\Sigma}^{-1} (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) \quad (192)$$

$$= \frac{N}{2} \frac{d}{d \mathbf{\Sigma}} \ln |\mathbf{\Sigma}|^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (193)$$

$$= \frac{N}{2} \frac{d}{d \mathbf{\Sigma}} \ln |\mathbf{\Sigma}^{-1}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (194)$$

$$= \frac{N}{2} \mathbf{\Sigma} - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (195)$$

$$(196)$$

$$\mathbf{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \phi(\mathbf{x}_n))^T \quad (197)$$

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.

$$p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta^{-1}) \quad (198)$$

$$= \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{S}_0) \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (199)$$

$$= (2\pi)^{-D/2} |\mathbf{S}_0|^{-1/2} \exp[-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0)] \prod_{n=1}^N \left(\frac{\beta}{2\pi}\right)^{\frac{1}{2}} \exp[-\frac{\beta}{2}(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2] \quad (200)$$

$$= (2\pi)^{-D/2} |\mathbf{S}_0|^{-1/2} \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \exp[-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N -\frac{\beta}{2}(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2] \quad (201)$$

$$(202)$$

Now we have to get the exponent right.

$$-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N -\frac{\beta}{2}(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 \quad (203)$$

$$= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) - \frac{\beta}{2} \sum_{n=1}^N (t_n^2 + \mathbf{w}^T \phi(\mathbf{x}_n) \mathbf{w}^T \phi(\mathbf{x}_n) - 2\mathbf{w}^T \phi(\mathbf{x}_n) t_n) \quad (204)$$

$$= -\frac{1}{2}(\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{w} - \mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{m}_0 - \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{w} + \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0) - \frac{\beta}{2}(\mathbf{t}^T \mathbf{t} + \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\mathbf{w}^T \Phi^T \mathbf{t}) \quad (205)$$

$$= -\frac{1}{2}(\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{w} - 2\mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{m}_0 + \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0 - \beta \mathbf{t}^T \mathbf{t} + \beta \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2\beta \mathbf{w}^T \Phi^T \mathbf{t}) \quad (206)$$

$$= \mathbf{w}^T (-\frac{1}{2}(\mathbf{S}_0^{-1} + \beta \Phi^T \Phi)) \mathbf{w} + \mathbf{w}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}) + c \quad (207)$$

$$= \mathbf{w}^T (-\frac{1}{2} \mathbf{S}_N^{-1}) \mathbf{w} + \mathbf{w}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}) + c \quad (208)$$

Now completing the square.

$$\mathbf{w}^T (-\frac{1}{2} \mathbf{S}_N^{-1}) \mathbf{w} + \mathbf{w}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}) + c \quad (209)$$

$$= (\mathbf{w} + \frac{1}{2}(-\frac{1}{2} \mathbf{S}_N^{-1})^{-1}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t}))^T (-\frac{1}{2} \mathbf{S}_N^{-1})(\mathbf{w} + \frac{1}{2}(-\frac{1}{2} \mathbf{S}_N^{-1})^{-1}(\mathbf{S}_0^{-1} \mathbf{m}_0 + \Phi^T \mathbf{t})) + k \quad (210)$$

$$= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^T \mathbf{S}_N^{-1}(\mathbf{w} - \mathbf{m}_N) + k \quad (211)$$

3.15

Hint 1: Use 3.95 and 3.92

$$E(\mathbf{m}_N) = \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N \quad (212)$$

$$= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\gamma(\mathbf{m}_N^T \mathbf{m}_N)^{-1}}{2} \mathbf{m}_N^T \mathbf{m}_N \quad (213)$$

$$= \frac{\beta}{2} \|\mathbf{t} - \Phi \mathbf{m}_N\|^2 + \frac{\gamma}{2} \quad (214)$$

$$(215)$$

$$\beta = (N - \gamma) \|\mathbf{t} - \Phi^T \mathbf{m}_n\|^{-2} \quad (216)$$

$$E(\mathbf{m}_N) = \frac{1}{2}(\mathcal{N} - \gamma + \gamma) = \frac{N}{2} \quad (217)$$

3.17

Evidence function:

$$p(D|M_i) = \int p(D|\mathbf{w}, M_i) p(\mathbf{w}, M_i) d\mathbf{w} \quad (218)$$

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

$$p(D|\alpha, \beta) = \int p(\mathbf{w}|\alpha^{-1}I) \prod_{n=1}^N p(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (219)$$

$$= \int (2\pi)^{-M/2} |\alpha^{-1}I|^{-1/2} \exp[-1/2(\mathbf{w}^T(\alpha^{-1}I)^{-1}\mathbf{w})] \prod_{n=1}^N \frac{1}{\sqrt{2\pi\beta^{-1}}} \exp[\beta/2(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2] \quad (220)$$

$$= \int (2\pi)^{-M/2} (\alpha^{-1})^{-M/2} \frac{\beta^{N/2}}{2\pi} \exp[-\alpha/2\mathbf{w}^T \mathbf{w} - \beta/2\|\mathbf{t} - \Phi \mathbf{w}\|^2] \quad (221)$$

4.4

$$\mathcal{L} = \mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1) + \lambda(\mathbf{w}^T \mathbf{w} - 1) \quad (222)$$

$$\nabla_{\mathbf{w}} \mathcal{L} = (\mathbf{m}_2 - \mathbf{m}_1) + 2\lambda \mathbf{w} = 0 \quad (223)$$

$$\mathbf{w} = -\frac{1}{2\lambda}(\mathbf{m}_2 - \mathbf{m}_1) \propto (\mathbf{m}_2 - \mathbf{m}_1) \quad (224)$$

4.5

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

$$(m_2 - m_1)^2 = (\mathbf{w}^T \mathbf{m}_2 - \mathbf{w}^T \mathbf{m}_1)^2 \quad (225)$$

$$= \mathbf{w}^T \mathbf{m}_2 \mathbf{w}^T \mathbf{m}_2 - \mathbf{w}^T \mathbf{m}_2 \mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{m}_2 + \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{m}_1 \quad (226)$$

$$= \mathbf{w}^T (\mathbf{m}_2 \mathbf{m}_2^T - \mathbf{m}_2 \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{m}_2^T + \mathbf{m}_1 \mathbf{m}_1^T) \mathbf{w} \quad (227)$$

$$= \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w} \quad (228)$$

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

$$s_1^2 = \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m}_1)^2 \quad (229)$$

$$= \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n \mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{x}_n \mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{x}_n + \mathbf{w}^T \mathbf{m}_1 \mathbf{w}^T \mathbf{m}_1) \quad (230)$$

$$= \sum_{n \in C_1} \mathbf{w}^T (\mathbf{x}_n \mathbf{x}_n^T - \mathbf{x}_n \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{x}_n^T + \mathbf{m}_1 \mathbf{m}_1^T) \mathbf{w} \quad (231)$$

$$= \sum_{n \in C_1} \mathbf{w}^T (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T \mathbf{w} \quad (232)$$

4.6

$$\sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n) \mathbf{x}_n = 0 \quad (233)$$

$$= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - \mathbf{w}^T \mathbf{m} - t_n) \mathbf{x}_n \quad (234)$$

$$\mathbf{w}^T \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m}) \mathbf{x}_n = \sum_{n=1}^N t_n \mathbf{x}_n \quad (235)$$

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

$$\sum_{n=1}^N t_n \mathbf{x}_n = \sum_{n \in C_1}^{N_1} t_n \mathbf{x}_n - \sum_{n \in C_2}^{N_2} t_n \mathbf{x}_n \quad (236)$$

$$= \sum_{n \in C_1}^{N_1} \frac{N}{N_1} \mathbf{x}_n - \sum_{n \in C_2}^{N_2} \frac{N}{N_2} \mathbf{x}_n \quad (237)$$

$$= N \left(\sum_{n \in C_1}^{N_1} \frac{1}{N_1} \mathbf{x}_n - \sum_{n \in C_2}^{N_2} \frac{1}{N_2} \mathbf{x}_n \right) \quad (238)$$

$$= N(\mathbf{m}_1 - \mathbf{m}_2) \quad (239)$$

$$-\mathbf{w}^T \mathbf{m} \sum_{n=1}^N \mathbf{x}_n = -\frac{1}{N} \mathbf{w}^T (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) \quad (240)$$

$$= -\frac{1}{N} (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2) (N_1 \mathbf{m}_1 + N_2 \mathbf{m}_2)^T \mathbf{w} \quad (241)$$

$$= -\frac{1}{N} (N_1^2 \mathbf{m}_1 \mathbf{m}_1^T + N_1 N_2 \mathbf{m}_1 \mathbf{m}_2^T + N_2 N_1 \mathbf{m}_2 \mathbf{m}_1^T + N_2^2 \mathbf{m}_2 \mathbf{m}_2^T) \mathbf{w} \quad (242)$$

$$= -\frac{1}{N} ((N - N_2) N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_1 N_2 \mathbf{m}_1 \mathbf{m}_2^T + N_2 N_1 \mathbf{m}_2 \mathbf{m}_1^T + (N - N_1) N_2 \mathbf{m}_2 \mathbf{m}_2^T) \mathbf{w} \quad (243)$$

$$= \left(-N_1 \mathbf{m}_1 \mathbf{m}_1^T + \frac{N_1 N_2}{N} \mathbf{m}_1 \mathbf{m}_1^T - \frac{N_1 N_2}{N} \mathbf{m}_1 \mathbf{m}_2^T - \frac{N_2 N_1}{N} \mathbf{m}_2 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \frac{N_1 N_2}{N} \mathbf{m}_2 \mathbf{m}_2^T \right) \mathbf{w} \quad (244)$$

$$= \left(-N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \frac{N_1 N_2}{N} (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \right) \mathbf{w} \quad (245)$$

$$= \left(-N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \frac{N_1 N_2}{N} \mathbf{S}_B \right) \mathbf{w} \quad (246)$$

$$(247)$$

We add the remaining terms to the final term:

$$\left(-N_1 \mathbf{m}_1 \mathbf{m}_1^T - N_2 \mathbf{m}_2 \mathbf{m}_2^T + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} \quad (248)$$

$$= \left(N_1 \mathbf{m}_1 \mathbf{m}_1^T - 2N_1 \mathbf{m}_1 \mathbf{m}_1^T + N_2 \mathbf{m}_2 \mathbf{m}_2^T - 2N_2 \mathbf{m}_2 \mathbf{m}_2^T + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} \quad (249)$$

$$= \left(N_1 \mathbf{m}_1 \mathbf{m}_1^T - \sum_{n \in C_1} \mathbf{x}_n \mathbf{m}_1^T - \mathbf{m}_1 \sum_{n \in C_1} \mathbf{x}_n^T + \dots + \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} \quad (250)$$

$$= \left(\sum_{n \in C_1} \mathbf{m}_1 \mathbf{m}_1^T - \mathbf{x}_n \mathbf{m}_1^T - \mathbf{m}_1 \mathbf{x}_n^T + \mathbf{x}_n \mathbf{x}_n^T + \dots \right) \mathbf{w} \quad (251)$$

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

4.7

$$\sigma(-a) = \frac{1}{e^x + 1} = \frac{e^{-x}}{1 + e^{-x}} = \frac{e^{-x} + 1}{e^{-x} + 1} - \frac{1}{e^{-x} + 1} \quad (252)$$

$$y = \frac{1}{1 + e^{-x}} \quad (253)$$

$$e^{-x} = \frac{1 - y}{y} \quad (254)$$

$$y = e^x (1 - y) \quad (255)$$

$$e^x = y/(1-y) \quad (256)$$

$$x = \ln[y/(1-y)] \quad (257)$$

4.8

$$p(C_1|\mathbf{x}) = \sigma(a) \quad (258)$$

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

$$a = \ln \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} + \ln \frac{p(C_1)}{p(C_2)} \quad (259)$$

$$\ln \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} = \ln \left[\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) \right) \right] \quad (260)$$

$$- \ln \left[\frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \right) \right] \quad (261)$$

$$= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \quad (262)$$

$$= \frac{1}{2} [2\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 - 2\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_2 + \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2] \quad (263)$$

$$(264)$$

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.

$$\ln p(\mathbf{X}|\mathbf{T}) = \sum_{n=1}^N \ln \prod_{j=1}^K (\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma))^{t_j} \quad (265)$$

$$= \sum_{n=1}^N \sum_{j=1}^K t_j \ln(\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \Sigma)) \quad (266)$$

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

$$\mathcal{L}(\boldsymbol{\pi}, \lambda) = \ln(p(\mathbf{X}, \mathbf{T})) - \lambda \left(\sum_{j=1}^K \pi_j - 1 \right) \quad (267)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_j} = \sum_{n=1}^N \frac{t_j}{\pi_j} - \lambda = 0 \quad (268)$$

$$\lambda = \sum_{n=1}^N \frac{t_j}{\pi_j} = N \frac{t_j}{\pi_j} = \frac{N_j}{\pi_j} \quad (269)$$

$$\pi_j = \frac{N_j}{\lambda} \quad (270)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sum_{j=1}^K \pi_j = 1 \quad (271)$$

Plugging this into 380.

$$\sum_{j=1}^K \frac{N_j}{\lambda} = 1 \iff \lambda = N \quad (272)$$

Gives us the desired result.

4.12

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

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

$$= \sigma(a) \frac{e^{-a}}{1 + e^{-a}} \quad (275)$$

$$= \sigma(a) \left[\frac{1 + e^{-a}}{1 + e^{-a}} - \frac{1}{1 + e^{-a}} \right] \quad (276)$$

4.13

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = - \sum_{n=1}^N \frac{t_n}{y_n} \nabla_{\mathbf{w}} \sigma(\mathbf{w}^T \phi_n) - \frac{1 - t_n}{1 - y_n} \nabla_{\mathbf{w}} \sigma(\mathbf{w}^T \phi_n) \quad (277)$$

$$= - \sum_{n=1}^N \frac{t_n}{y_n} \sigma(\mathbf{w}^T \phi_n) (1 - \sigma(\mathbf{w}^T \phi_n)) \phi_n - \frac{1 - t_n}{1 - y_n} \sigma(\mathbf{w}^T \phi_n) (1 - \sigma(\mathbf{w}^T \phi_n)) \phi_n \quad (278)$$

$$= - \sum_{n=1}^N t_n (1 - y_n) \phi_n - (1 - t_n) y_n \phi_n \quad (279)$$

$$= - \sum_{n=1}^N (t_n - y_n) \phi_n \quad (280)$$

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

$$p(\mathbf{t}, \mathbf{w}) = \prod_{n=1}^N y_n^{\pi_n} [1 - y_n]^{1-t_n} \quad (281)$$

$$\ln p = \sum_{n=1}^N \pi_n \ln y_n + (1 - \pi_n) \ln(1 - y_n) \quad (282)$$

4.17

$$p(C_k | \phi) = y_k = \frac{\exp a_k}{\sum_{j=1} \exp a_j} \quad (283)$$

$$\frac{\partial y_k}{\partial a_j} = -\exp a_k \left(\sum_j \exp(a_j) \right)^{-2} \exp(a_j) \quad (284)$$

$$= \begin{cases} y_k(0 - y_j) & j \neq k \\ y_k(1 - y_j) & j = k \end{cases} \quad (285)$$

$$= y_k(I_{kj} - y_j) \quad (286)$$

4.18

$$\nabla_{\mathbf{w}_j} E(\mathbf{W}) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \nabla_{\mathbf{w}_j} \ln y_{nk} \quad (287)$$

$$\nabla_{\mathbf{w}_j} \ln(y_{nk}) = -(I_{kj} - y_j) \phi_n \quad (288)$$

$$\nabla_{\mathbf{w}_j} E(\mathbf{W}) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} (I_{kj} - y_j) \phi_n \quad (289)$$

$$= \phi \sum_{n=1}^N \sum_{k=1}^K t_{nk} y_{jn} \phi - t_{nk} I_{kj} \quad (290)$$

$$= \phi \sum_{n=1}^N t_{nj} - y_{jn} \phi \underbrace{\sum_{k=1}^K t_{nk}}_{=1} \quad (291)$$

$$= \sum_{n=1}^N \phi (y_{jn} - t_{nj}) \quad (292)$$

4.19

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

$$p(\mathbf{t}, \mathbf{w}) = \sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)] \quad (293)$$

$$\nabla_{\mathbf{w}} \sum_{n=1}^N [t_n \ln \Phi(a) + (1 - t_n) \ln(1 - \Phi(a))] = \sum_{n=1}^N \left(\frac{t_n}{\Phi(a)} - \frac{1 - t_n}{1 - \Phi(a)} \right) \Phi(a) \phi_n \quad (294)$$

4.21

$$\Phi(a) = \int_0^a \mathcal{N}(0, 1) d\theta \quad (295)$$

$$= \int_0^a \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\theta^2\right) d\theta \quad (296)$$

$$= \frac{1}{2} + \int_{-\infty}^a \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\theta^2\right) d\theta \quad (297)$$

$$= \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \frac{\sqrt{\pi}}{2} \int_{-\infty}^a \frac{2}{\sqrt{\pi}} \exp\left(-\frac{1}{2}\theta^2\right) d\theta \quad (298)$$

$$= \frac{1}{2} \left(1 + \frac{1}{\sqrt{2}} \operatorname{erf}(a) \right) \quad (299)$$

4.22

$$\ln p(D) = \ln \left[f(z_0) \frac{(2\pi)^{M/2}}{|A|^{1/2}} \right] \quad (300)$$

$$= \ln f(z_0) = \frac{M}{2} \ln 2\pi - \frac{1}{2} \ln |A| \quad (301)$$

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

$$\ln f(\boldsymbol{\theta}) \Big|_{z_0} = \ln f(\boldsymbol{\theta}_{MAP}) = \ln p(D|\boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) \quad (302)$$

5.2

$$p(\mathbf{T}, \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \mathcal{N}(\mathbf{t}_n | y(\mathbf{x}_n, \mathbf{w}), \beta^{-1} I) \quad (303)$$

$$= \prod_{n=1}^N \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\beta^{-1} I|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))^T \beta^{-1} I (\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))\right) \quad (304)$$

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

$$\left(-\frac{1}{2}(\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))^T \beta^{-1} I(\mathbf{t}_n - y(\mathbf{x}_n, \mathbf{w}))\right) = -\frac{1}{2\beta} \|y(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2. \quad (305)$$

5.5

$$p(\mathbf{T}, \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \prod_{k=1}^K y_k(\mathbf{x}_n, \mathbf{w})^{t_n^k} (1 - y_k)(\mathbf{x}_n, \mathbf{w})^{1-t_n^k} \quad (306)$$

Taking the log becomes the cross-entropy function.

5.6

$$\frac{\partial E(\mathbf{w})}{\partial a_k} = \sum_{n=1}^N \left(\frac{\partial}{\partial a_k} t_n \ln y_n + \frac{\partial}{\partial a_k} (1 - t_n) \ln(1 - y_n) \right) \quad (307)$$

$$= -\frac{t_k}{y_k} \frac{\partial}{\partial a_k} y_k - \frac{1 - t_k}{1 - y_k} \frac{\partial}{\partial a_k} (1 - y_k) \quad (308)$$

$$= -t_k(1 - y_k) + (1 - t_k)y_k \quad (309)$$

$$= y_k - t_k \quad (310)$$

5.7

$$E(\mathbf{w}) = \sum_{k=1}^K t_k \ln y_k(\mathbf{x}, \mathbf{w}) \quad (311)$$

$$\frac{\partial E}{\partial a_j} = -\sum_{k=1}^K \frac{t_k}{y_k} y_k (I_{kj} - y_j) \quad (312)$$

$$= -\sum_{k=1}^K t_k (I_{kj} - y_j) \quad (313)$$

$$= -t_j + y_j \quad (314)$$

$$\frac{d \tanh}{da} = (e^a - e^{-a}) \frac{d}{da} (e^a + e^{-a})^{-1} + (e^a + e^{-a})^{-1} \frac{d}{da} (e^a - e^{-a}) \quad (315)$$

$$= -(e^a - e^{-a})(e^a + e^{-a})^{-2} \frac{d}{da} (e^a + e^{-a}) + (e^a + e^{-a})^{-1} (e^a - e^{-a}) \quad (316)$$

$$= -(e^a - e^{-a})(e^a + e^{-a})^{-2} (e^a - e^{-a}) + 1 \quad (317)$$

$$= -h^2(a) + 1 \quad (318)$$

5.9

Still Bernoulli, so

$$p(t|\mathbf{x}, \mathbf{w}) = \left(\frac{1+y}{2}\right)^{\frac{1+t}{2}} \left(\frac{1-y}{2}\right)^{\frac{1-t}{2}} \quad (319)$$

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}) \quad (320)$$

$$-\ln(p(\mathbf{t}|\mathbf{X}, \mathbf{w})) = -\sum_{n=1}^N \ln p(t_n|\mathbf{x}_n, \mathbf{w}) \quad (321)$$

$$= -\sum_{n=1}^N \left(\left(\frac{1+t}{2}\right) \ln\left(\frac{1+y}{2}\right) + \left(\frac{1-t}{2}\right) \ln\left(\frac{1-y}{2}\right) \right) \quad (322)$$

We can use tanh as activation function.

5.10

$$u_i^T H u_i = u_i^T \lambda u_i = \delta_i i \lambda_i = \lambda_i \quad (323)$$

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.

$$\frac{N(N+1)}{2} + N = \frac{N(N+3)}{2} \quad (324)$$

5.14

Hint 1: Taylor expansion on both terms in numerator.

Taylor:

$$E_n(w_j + \epsilon) = E_n(w_{ji}) + \epsilon \frac{\partial E_n}{\partial w_{ji}} + \frac{\epsilon^2}{2} \frac{\partial^2 E_n}{\partial w_{ji}^2} + O(\epsilon^3) \quad (325)$$

$$E_n(w_j - \epsilon) = E_n(w_{ji}) - \epsilon \frac{\partial E_n}{\partial w_{ji}} + \frac{\epsilon^2}{2} \frac{\partial^2 E_n}{\partial w_{ji}^2} - O(\epsilon^3) \quad (326)$$

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

5.16

$$E = \frac{1}{2} \sum_{n=1}^N \|\mathbf{y}_n - \mathbf{t}_n\|^2 = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M (\mathbf{y}_n - \mathbf{t}_n)^2 \quad (327)$$

$$\nabla E = \sum_{n=1}^N \sum_{m=1}^M (\mathbf{y}_n - \mathbf{t}_n) \nabla \mathbf{y}_n \quad (328)$$

$$\nabla \nabla E = \sum_{n=1}^N \sum_{m=1}^M \nabla \mathbf{y}_n \nabla (\mathbf{y}_n - \mathbf{t}_n) + (\mathbf{y}_n - \mathbf{t}_n) \nabla \nabla \mathbf{y}_n \approx \sum_{n=1}^N \sum_{m=1}^M \nabla \mathbf{y}_n \nabla \mathbf{y}_n^T \quad (329)$$

5.17

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

$$\frac{\partial E}{\partial w_r} = \int \int (y(x, w) - t) \frac{\partial y}{\partial w_r} p(x, t) dx dt \quad (330)$$

$$\frac{\partial^2 E}{\partial w_r \partial w_s} = \int \int \left[(y(x, w) - t) \frac{\partial^2 y(x, w)}{\partial w_r \partial w_s} + \frac{\partial y(x, w)}{\partial w_r} \frac{\partial y}{\partial w_s} \right] p(x, t) dx dt \quad (331)$$

$$\int \int (y(x, w) - t) p(x, t) dx dt = \int \int (y(x, w) - t) p(t|x) p(x) dx dt \quad (332)$$

$$= \int p(x) \left(y(x, w) - \underbrace{\int t p(t|x)}_{=y(x, w)} \right) dx dt = 0 \quad (333)$$

The remaining integral is the answer.

5.18

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left(\sum_{j=1}^M w_{kj} h \left(\sum_{i=1}^D w_{ji} x_i + w_{j0} \right) + \sum_{l=1}^D w_l x_l \right) \quad (334)$$

Finding the derivatives to these skip weights is straightforward.

5.19

$$E = - \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \quad (335)$$

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

$$\nabla \nabla E = \sum_{n=1}^N y_n (1 - y_n) \nabla a_n \nabla a_n^T + (y_n - t_n) \nabla \nabla a_n \quad (336)$$

5.20

$$\nabla_{w_j} E(\mathbf{W}) = \sum_{n=1}^N (y_{n_j} - bt_{n_j}) \nabla a_j \quad (337)$$

$$\nabla \nabla \mathbf{w}_j \approx \sum_{n=1} y_k (I - y_j) \nabla a_j \nabla a_j^T \quad (338)$$

5.24

$$\sum_i \frac{1}{a} w_{ji} (ax_i + b) + w_{j0} - \frac{b}{a} \sum_i w_{ji} = \sum_i w_{ji} x + \frac{b}{a} w_{ji} + w_{j0} - \frac{b}{a} \sum_i w_{ji} \quad (339)$$

$$= \sum_i w_{ji} x_i + w_{j0} \quad (340)$$

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

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \quad (341)$$

$$\frac{\partial \tilde{E}}{\partial w_i} = \frac{\partial E}{\partial w_i} + \frac{\partial}{\partial w_i} \lambda \Omega(\mathbf{w}) \quad (342)$$

$$\frac{\partial}{\partial w_i} \lambda \Omega(\mathbf{w}) = -\lambda \frac{\partial}{\partial w_i} \ln \left(\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right) \quad (343)$$

$$\frac{\partial}{\partial w_i} = \frac{1}{\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2)} \frac{\partial}{\partial w_i} \ln \left(\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right) \quad (344)$$

$$\frac{\partial}{\partial w_i} \ln \left(\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j^2) \right) = \sum_{j=1}^M \pi_j \frac{\partial}{\partial w_i} \mathcal{N}(w_i | \mu_j, \sigma_j^2) \quad (345)$$

$$\frac{\partial}{\partial w_i} \mathcal{N}(w_i | \mu_j, \sigma_j^2) = \frac{1}{\sqrt{2\pi\sigma_j^2} \exp(-\frac{1}{2} \frac{(w_i - \mu_j)^2}{\sigma_j^2})} \frac{\partial}{\partial w_i} \left(-\frac{1}{2} \frac{(w_i - \mu_j)^2}{\sigma_j^2} \right) \quad (346)$$

$$\frac{\partial}{\partial w_i} = -\sigma_j^{-2} (w_i - \mu_j) \quad (347)$$

Plugging everything in gives

$$\frac{\partial}{\partial w_i} = \frac{\sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j)}{\sum_{k=1}^K \mathcal{N}(w_i | \mu_j, \sigma_j^2)} \quad (348)$$

5.30

$$\frac{\partial \tilde{E}}{\partial \mu_j} = -\lambda \frac{\partial}{\partial \mu_j} \sum_i \ln \sum_{j=1}^M \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j) \quad (349)$$

$$= -\lambda \sum_i \frac{1}{\sum_j \pi_j \mathcal{N}(w_i)} \quad (350)$$

$$= -\lambda \sum_i \frac{1}{\sum_j \pi_j N(w_i | \mu_j, \sigma_j)} \sum_j \pi_j \frac{d}{d\mu_j} N(w_i | \mu_j, \sigma_j) \quad (351)$$

$$\frac{d}{d\mu_j} N(w_i | \mu_j, \sigma_j) = -\frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{1}{2\sigma_j^2}(w_i - \mu_j)^2\right) \left[\frac{1}{\sigma_j^2}(w_i - \mu_j)\right] \quad (352)$$

Plugging in gives the result.

5.31

$$\frac{\partial}{\partial \sigma_j} \pi_j \mathcal{N}(w_i | \mu_j, \sigma_j) \quad (353)$$

$$= \pi \left[\exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \frac{\partial}{\partial \sigma_j} \frac{1}{\sqrt{2\pi\sigma_j^2}} + \frac{1}{\sqrt{2\pi\sigma_j^2}} \frac{\partial}{\partial \sigma_j} \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \right] \quad (354)$$

$$\frac{\partial}{\partial \sigma_j} (2\pi\sigma_j^2)^{-\frac{1}{2}} = -\frac{1}{2} (2\pi\sigma_j^2)^{-\frac{3}{2}} \frac{\partial}{\partial \sigma_j} (2\pi\sigma_j^2) \quad (355)$$

$$= -\frac{1}{2} (2\pi\sigma_j^2)^{-\frac{3}{2}} 4\pi\sigma_j \quad (356)$$

$$= -\frac{1}{\sigma_j} \frac{1}{\sqrt{2\pi\sigma_j^2}} \quad (357)$$

$$\frac{\partial}{\partial \sigma_j} \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) = \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \frac{\partial}{\partial \sigma_j} \left(-\frac{1}{2\sigma_j^2}(\mu_j - w_i)^2\right) \quad (358)$$

$$= \exp\left(-\frac{(\mu_j - w_i)^2}{2\sigma_j^2}\right) \frac{1}{\sigma_j^3} (\mu_j - w_i)^2 = \quad (359)$$

Plugging these values in gives the result.

5.33

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

$$\mathbf{y} = \begin{pmatrix} \cos(\pi - \theta)L_1 + v_1 \\ \sin(\pi - \theta)L_1 + v_2 \end{pmatrix} \quad (360)$$

$$\mathbf{x} = \begin{pmatrix} y_1 + \cos(\theta_1 + \theta_2 - \pi)L_2 \\ y_2 + \sin(\theta_1 + \theta_2 - \pi)L_2 \end{pmatrix} \quad (361)$$

5.34

$$\frac{\partial E_n}{\partial a_{kl}^\pi} = - \frac{1}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (362)$$

$$\frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) = \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_k \quad (363)$$

$$\frac{\partial}{\partial a_k^\pi} \sum_{l=1} \pi_k = \sum_{l=1} \pi_k (I_{kl} - \pi_l) \quad (364)$$

$$= \pi_k - \pi_k \sum_{l=1} \pi_l \quad (365)$$

$$\frac{\partial E_n}{\partial a_{kl}^\pi} = - \frac{\pi_k \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) - \pi_k \sum_{l=1} \pi_l \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \quad (366)$$

$$= -\gamma_k(w) + \pi_k \quad (367)$$

5.35

$$\frac{\partial E_n}{\partial a_{kl}^\mu} = - \frac{1}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \frac{\partial}{\partial a_k^\mu} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (368)$$

$$\frac{\partial}{\partial a_k^\mu} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (369)$$

$$= \pi_k(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \frac{\partial}{\partial a_k^\mu} - \frac{1}{2} (\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1} (\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \quad (370)$$

$$(371)$$

$$\frac{\partial}{\partial a_k^\mu} - \frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \quad (372)$$

$$= -\sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \underbrace{\frac{\partial \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w})}{\partial a_k^\mu}}_{=1} \quad (373)$$

$$\frac{\partial E_n}{\partial a_{kl}^\mu} = \gamma_{nk} \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}) \quad (374)$$

5.36

$$\frac{\partial E_n}{\partial a_{kl}^\sigma} = - \frac{1}{\sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w}))} \frac{\partial}{\partial a_k^\sigma} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (375)$$

$$\frac{\partial}{\partial a_k^\sigma} \sum_{l=1} \pi_l(\mathbf{x}_n, \mathbf{w}) \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) = \pi_k(\mathbf{x}_n, \mathbf{w}) \frac{\partial}{\partial a_k^\sigma} \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) \quad (376)$$

$$\frac{\partial}{\partial a_k^\sigma} \mathcal{N}(\mathbf{t}_n | \boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}), \sigma_k^2(\mathbf{x}_n, \mathbf{w})) = \frac{\partial}{\partial a_k^\sigma} \frac{1}{(2\pi)^{\frac{L}{2}}} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) \quad (377)$$

$$= \frac{1}{(2\pi)^{\frac{L}{2}}} \left(\exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) \frac{\partial}{\partial a_k^\sigma} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \right. \quad (378)$$

$$\left. + \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \frac{\partial}{\partial a_k^\sigma} \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) \right) \quad (379)$$

$$\frac{\partial}{\partial a_k^\sigma} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} = \frac{\partial}{\partial a_k^\sigma} \frac{1}{\sigma_k(\mathbf{x}, \mathbf{w})^L} = -\frac{L}{\sigma^2(\mathbf{x}, \mathbf{w})^{L-1}} = \frac{L}{\sigma_k} \frac{1}{\sigma_k(\mathbf{x}, \mathbf{w})^L} = \frac{L}{\sigma_k} \frac{1}{|\sigma_k^2(\mathbf{x}_n, \mathbf{w})|^{\frac{1}{2}}} \quad (380)$$

$$\frac{\partial}{\partial a_k^\sigma} \left(-\frac{1}{2}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})^T \sigma_k^2(\mathbf{x}_n, \mathbf{w})^{-1}(\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t})\right) = \frac{\|\boldsymbol{\mu}_k(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}\|^2}{\sigma_k^3} \quad (381)$$

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

5.37

5.38

5.39

Using the approximation:

$$p(D|\alpha, \beta) = \int p(D|w, \beta)p(w|\alpha)dw \approx p(D, w_{map})p(w_{map})\frac{(2\pi)^{W/2}}{|A|^{\frac{1}{2}}} \quad (382)$$

Taking the log gives the result.

5.40

Use a softmax activation.

6.3

$$\|x - x_n\| = x^T x + x_n^T x - 2x_n^T x \quad (383)$$

$$= k(x, x) + k(x_n, x) - 2k(x_n, x) \quad (384)$$

6.4

$$\begin{pmatrix} 2 & 0 \\ -1 & 3 \end{pmatrix} \quad (385)$$

6.5

$$ck(\mathbf{x}, \mathbf{x}') = c\psi^T \psi = (c^{\frac{1}{2}}\psi)^T (c^{\frac{1}{2}}\psi) = \phi^T \phi \quad (386)$$

$$f(\mathbf{x})k(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') = f(\mathbf{x})\psi^T \psi f(\mathbf{x}') = (f(\mathbf{x})\psi^T)^T (\psi f(\mathbf{x}')) = \phi^T \phi \quad (387)$$

6.6

$$q(k(x, x')) = \sum_{n=1}^N \beta_n k(x, x')^n \quad (388)$$

$$= \sum_{n=1}^N \beta_0 k(x, x') \quad 6.18 \quad (389)$$

$$= \sum_{n=1}^N k(x, x') \quad 6.13 \quad (390)$$

$$= k(x, x') \quad 6.17 \quad (391)$$

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

6.7

$$\psi(x)^T \psi(x') + \theta(x)^T \theta(x') = \sum_{n=1}^N (\psi(x_n) + \theta(x_n))(\psi(x'_n) + \theta(x'_n)) - \psi(x_n)\theta(x'_n) - \theta(x_n)\psi(x'_n) \quad (392)$$

$$= \begin{bmatrix} \dots & \psi(x_n) + \theta(x_n) & -\psi(x_n) & -\theta(x_n) & \dots \end{bmatrix} \begin{bmatrix} \vdots \\ \psi(x'_n) + \theta(x'_n) \\ \theta(x'_n) \\ \psi(x'_n) \\ \vdots \end{bmatrix} \quad (393)$$

$$= \phi(x)^T \phi(x') \quad (394)$$

$$k_1(x, x') k_2(x, x') = \theta(x)^T \theta(x') (\psi(x)^T \psi(x')) \quad (395)$$

$$= \sum_{n=1}^N \theta(x_n) \theta(x'_n) \psi(x_n) \psi(x'_n) \quad (396)$$

$$= \sum_{n=1}^N \theta(x_n) \psi(x_n) \theta(x'_n) \psi(x'_n) \quad (397)$$

$$= \phi(x)^T \phi(x') \quad (398)$$

$$= k(x, x') \quad (399)$$

6.8

$$k_3(\phi(x), \phi(x')) = \psi(\phi(x))^T \psi(\phi(x')) = k(x, x') \quad (400)$$

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

$$= (Ux)^T Ux' \quad (401)$$

$$= \phi(x)^T \phi(x') \quad (402)$$

$$= k(x, x') \quad (403)$$

6.9

$$\sum_{n=1}^N \psi(x_{na})\psi(x'_{na}) + \sum_{m=1}^M \theta(x_{mb})\theta(x'_{mb}) \quad (404)$$

$$= \begin{bmatrix} \psi(x_{a1}) & \dots & \psi(x_{aN}) & \theta(x_{b1}) & \dots & \theta(x_{bM}) \end{bmatrix} \begin{bmatrix} \psi(x'_{a1}) \\ \vdots \\ \psi(x'_{aN}) \\ \theta(x'_{b1}) \\ \vdots \\ \theta(x'_{bM}) \end{bmatrix} = k(x, x') \quad (405)$$

$$x \in \mathbb{R}^n \quad (406)$$

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

$$y = f(g(x)) \quad (407)$$

$$y'(x) = \nabla_{g(x)} f(g(x)) J_x(g(x)) \quad (408)$$

$$k_a(x_a, x'_a)k(x_b, x'_b) = \sum_{n=1}^N \theta(x_{an})\theta(x'_{an}) \sum_{m=1}^M \psi(x_{bm})\psi(x'_{bm}) \quad (409)$$

$$= \sum_{n=1}^N \sum_{m=1}^M \theta(x_{an})\theta(x'_{an})\psi(x_{bm})\psi(x'_{bm}) \quad (410)$$

$$= \sum_{n=1}^N \sum_{m=1}^M \theta(x_{an})\psi(x_{bm})\theta(x'_{an})\psi(x'_{bm}) \quad (411)$$

$$= \phi(x)^T \phi(x') = k(x, x') \quad (412)$$

6.10

$$y(x) = \sum_{n=1}^N f(x_n)f(x_n)(K + \lambda I_N)^{-1}t \quad (413)$$

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 .

$$g(\theta, \mathbf{x}) = J_\theta(\phi) \nabla_\phi \ln p(\mathbf{x}, \phi) = J_\theta(\phi) g(\phi, \mathbf{x}) \quad (414)$$

$$g(\phi, \mathbf{x}) = J_\theta^{-1}(\phi) g(\theta, \mathbf{x}) \quad (415)$$

$$\mathbf{F}' = J_\theta^{-1}(\phi) \mathbb{E}_x[g(\theta, \mathbf{x}) g(\theta, \mathbf{x})^T] J_\theta^{-T}(\phi) = J_\theta^{-1}(\phi) \mathbf{F} J_\theta^{-T}(\phi) \quad (416)$$

$$k'(\mathbf{x}, \mathbf{x}') = g(\phi, \mathbf{x})^T \mathbf{F}'^{-1} g(\phi, \mathbf{x}') \quad (417)$$

$$= (J_\theta^{-1}(\phi) g(\theta, \mathbf{x}))^T (J_\theta^{-1}(\phi) \mathbf{F} J_\theta^{-T}(\phi))^{-1} J_\theta^{-1}(\phi) g(\theta, \mathbf{x}') \quad (418)$$

$$= g(\theta, \mathbf{x})^T J_\theta^{-T}(\phi) J_\theta^T(\phi) \mathbf{F} J_\theta(\phi) J_\theta^{-1}(\phi) g(\theta, \mathbf{x}') \quad (419)$$

$$= g(\theta, \mathbf{x})^T \mathbf{F} g(\theta, \mathbf{x}') \quad (420)$$

$$= k(\mathbf{x}, \mathbf{x}') \quad (421)$$

Therefore, the Fisher kernel is invariant.

6.14

$$\nabla_\theta (\ln p(\mathbf{x}, \theta)) = \nabla_\mu \ln p(\mathbf{x}, \mu) \quad (422)$$

$$= S^{-1}(\mathbf{x} - \mu) \quad (423)$$

$$\mathbf{F} = \mathbb{E}[S^{-1}(\mathbf{x} - \mu)(S^{-1}(\mathbf{x} - \mu))^T] \quad (424)$$

$$= S^{-1} \mathbb{E}_x[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] S^{-1} \quad (425)$$

$$= S^{-1} S S^{-1} = S^{-1} \quad (426)$$

$$k(\mathbf{x}, \mathbf{x}') = (S^{-1}(\mathbf{x} - \mu))^T \mathbf{F}^{-1} (S^{-1}(\mathbf{x}' - \mu)) \quad (427)$$

$$= (\mathbf{x} - \mu)^T S^{-1} S S^{-1} (\mathbf{x}' - \mu) \quad (428)$$

$$= (\mathbf{x} - \mu)^T S^{-1} (\mathbf{x}' - \mu) \quad (429)$$

6.15

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

$$|K| = k(x_1, x_1)k(x_2, x_2) - k(x_2, x_1)k(x_1, x_2) \geq 0 \quad (430)$$

from which the result follows.

6.18

$$\mathbf{z} = (x, t) \quad (431)$$

$$p(\mathbf{z}) = \frac{\sum_n \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)}{\sum_m \int \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)} dt = \frac{\sum_n \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)}{\sum_m \mathcal{N}(x - x_n, \sigma^2 I)} \quad (432)$$

$$\frac{\sum_n \mathcal{N}(\mathbf{z} - \mathbf{z}_n, \sigma^2 I)}{\sum_m \mathcal{N}(x - x_n, \sigma^2 I)} = \frac{\sum_n \frac{1}{2\pi\sigma^2} \exp[-\frac{1}{2\sigma^2}(\mathbf{z} - \mathbf{z}_n)^T(\mathbf{z} - \mathbf{z}_n)]}{\sum_m \frac{1}{\sqrt{2\pi}\sigma} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]} \quad (433)$$

$$= \frac{\sum_n \frac{1}{2\pi\sigma^2} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2 + -\frac{1}{2\sigma^2}(t - t_n)^2]}{\sum_m \frac{1}{\sqrt{2\pi}\sigma} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]} \quad (434)$$

$$= \sum_n \frac{\frac{1}{\sqrt{2\pi}\sigma^2} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]}{\sum_m \frac{1}{\sqrt{2\pi}\sigma} \exp[-\frac{1}{2\sigma^2}(x - x_n)^2]} \frac{1}{\sqrt{2\pi}\sigma^2} \exp[-\frac{1}{2\sigma^2}(t - t_n)^2] \quad (435)$$

$$= \sum_n \pi_n \mathcal{N}(t|t_n, \sigma^2) \quad (436)$$

6.24

$$u^T W u = u^T \sqrt{W} \sqrt{W} u \quad (437)$$

$$= (\sqrt{W} u)^T \sqrt{W} u \quad (438)$$

$$> 0 \quad \forall u \in \mathbb{R} \setminus \mathbf{0} \quad (439)$$

$$u^T (W + V) u = u^T W u + u^T V u > 0 \quad (440)$$

6.25

$$a_N = a_N - \nabla \nabla \Psi(a_N) \nabla \Psi(a_N) \quad (441)$$

$$= a_N + (W_N + C^{-1})^{-1} [t_N - \sigma_N - C^{-1} a_N] \quad (442)$$

$$= (W_N + C_N^{-1}) [W_N a_N - \sigma_N + t_N] \quad (443)$$

$$= C_N ((W_N + C_N^{-1}) C_N)^{-1} [W_N a_N - \sigma_N + t_N] \quad (444)$$

$$= C_N (C_N W_N + I)^{-1} [W_N a_N - \sigma_N + t_N] \quad (445)$$

6.26

7.2

$$t_n(\mathbf{w}^T \phi(\mathbf{x}) + b) = \gamma \quad (446)$$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \gamma) \quad (447)$$

$$\nabla_{\mathbf{w}} = \mathbf{w} - \sum_{n=1}^N \nabla a_n t_n \mathbf{w}^T \phi(\mathbf{x}_n) \quad (448)$$

$$= \mathbf{w} - \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) = 0 \quad (449)$$

$$\iff \mathbf{w} = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) \quad (450)$$

$$\frac{\partial}{\partial b} = - \sum_{n=1}^N a_n t_n = 0 \quad (451)$$

$$= \sum_{n=1}^N a_n t_n \quad (452)$$

$$\tilde{L} = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n t_n a_m t_m \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) - \sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \gamma) \quad (453)$$

$$\sum_{n=1}^N a_n (t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - \gamma) = \sum_{n=1}^N 0 - a_n \gamma \quad (454)$$

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

$$-\ln p(t_n) = - \sum_{n=1} \ln \sigma(t_n y_n) \quad (455)$$

which is the cross-entropy error function.

7.7

It is quite straightforward if you use the following steps.

$$\sum_{n=1}^N \xi_n (C - \mu_n - a_n) + \sum_{n=1}^N \hat{\xi}_n (C - \hat{\mu}_n - \hat{a}_n) = 0 \quad (456)$$

and

$$- \sum_{n=1}^N a_n (y_n) - \sum_{n=1}^N \hat{a}_n (-y_n) = - \sum_{n=1}^N (a_n - \hat{a}_n) y_n = 0 \quad (457)$$

7.8

This follows from 7.67 and 7.68.

7.9

$$p(\mathbf{w}, \mathbf{t}) = \mathcal{N}(\mathbf{w}|_N, \mathbf{S}_N) \quad (458)$$

$$\mathbf{S}_N = \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\mathbf{\Phi}^T\mathbf{t}) \quad (459)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta\mathbf{\Phi}^T\mathbf{\Phi} \quad (460)$$

We have

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta\mathbf{\Phi}^T\mathbf{\Phi} = \mathbf{A} + \beta\mathbf{\Phi}^T\mathbf{\Phi} = \Sigma^{-1} \quad (461)$$

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

7.16

7.18

$$\nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) = \nabla \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) - \mathbf{A}\mathbf{w} \quad (462)$$

$$\nabla \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) = \sum_{n=1}^N t_n \nabla \ln y_n + (1 - t_n) \nabla \ln(1 - y_n) \quad (463)$$

$$\nabla \ln y_n = (1 - y_n) \phi(\mathbf{x}_n) \quad (464)$$

$$\nabla \ln(1 - y_n) = -y_n \phi(\mathbf{x}_n) \quad (465)$$

$$\nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) = \sum_{n=1}^N \phi(\mathbf{x}_n) [t_n - y_n t_n - y_n + y_n t_n] - \mathbf{A}\mathbf{w} \quad (466)$$

$$= \mathbf{\Phi}^T [\mathbf{t} - \mathbf{y}] - \mathbf{A}\mathbf{w} \quad (467)$$

$$\nabla \nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) = \nabla \mathbf{\Phi}^T [\mathbf{t} - \mathbf{y}] - \mathbf{A} \quad (468)$$

$$= \nabla \sum_{n=1}^N [t_n - y_n] \phi_n - \mathbf{A} \quad (469)$$

$$= - \sum_{n=1}^N \phi_n \nabla^T y_n - \mathbf{A} \quad (470)$$

$$= - \sum_{n=1}^N \phi_n y_n (1 - y_n) \phi_n^T - \mathbf{A} \quad (471)$$

$$= -\mathbf{\Phi}^T \mathbf{B} \mathbf{\Phi} - \mathbf{A} \quad (472)$$

Which gives the result.

8.1

$$\int p(\mathbf{x})d\mathbf{x} = \int \int p(x_K|pa_K) \prod_{k=1}^{K-1} p(x_k|pa_k)dx_Kdx_1 \dots dx_{K-1} \quad (473)$$

$$= \int \int p(x_K|pa_K)dx_K \prod_{k=1}^{K-1} p(x_k|pa_k)dx_1 \dots dx_{K-1} \quad (474)$$

$$= \int \prod_{k=1}^{K-1} p(x_k|pa_k)dx_1 \dots dx_{K-1} \quad (475)$$

$$\vdots \quad (476)$$

$$= \int p(x_1)dx_1 = 1 \quad (477)$$

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 .

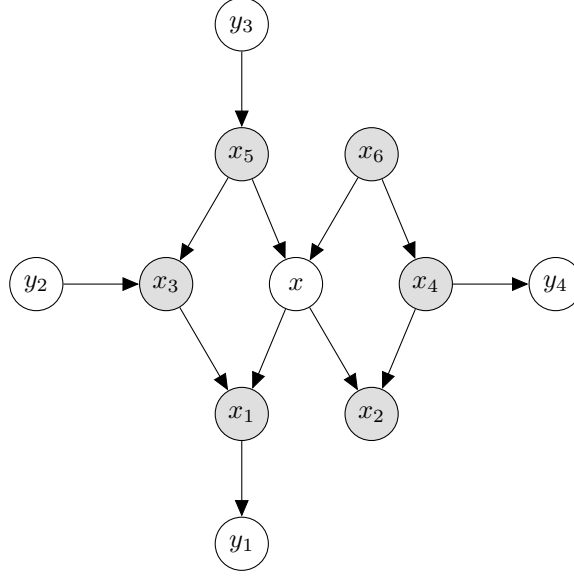
8.8

$$p(a, b, c|d) = \frac{p(a, b, c)}{p(d)} = \frac{p(a|b, c)p(b)p(c)}{p(d)} = \frac{p(a)p(b)p(c)}{p(d)} \quad (478)$$

$$\int \frac{p(a)p(b)p(c)}{p(d)}dc = p(a|d)p(b|d) \quad (479)$$

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

$$p(a, b, c, d) = p(a)p(b)p(c|a, b)p(d|c) \quad (480)$$

$$p(a, b) = \int \int p(a, b, c, d) dc dd = p(a)p(b) \quad (481)$$

Second part:

$$p(a, b, c|d) = \frac{p(a, b, c, d)}{p(d)} = \frac{p(a)p(b)p(c|a, b)p(d|c)}{p(d)} \quad (482)$$

$$p(a, b|d) = \int p(a, b, c|d) dc = \frac{p(a)p(b)p(d|c)}{p(d)} \quad (483)$$

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

$$E(x, y)_{x_k=1} - E(x, y)_{x_k=-1} = h \sum_i x_i - \beta \sum_{i,j} x_i x_j - \eta \sum_i x_i y_i + h - \beta x_k \sum_j x_j - \eta x_k \sum_i y_i \quad (484)$$

$$- h \sum_i x_i - \beta \sum_{i,j} x_i x_j - \eta \sum_i x_i y_i + h - \beta x_k \sum_j x_j - \eta x_k \sum_i y_i \quad (485)$$

$$= 2h - 2\beta x_k \sum_j x_j - \eta x_k \sum_i y_i \quad (486)$$

8.14

$$p(x, y) = \frac{1}{Z} \exp[-E(x, y)] \quad (487)$$

$$\ln p(x, y) = -\ln(Z) - E(x, y) = -\ln(Z) + \eta \sum_i x_i y_i \quad (488)$$

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

$$p(\mathbf{x}, \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} \pi_k^{z_k} \quad (489)$$

$$p(\mathbf{z}) = \sum_{\mathbf{z}} \prod_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} \quad (490)$$

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:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{old})_{MAP} = \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\boldsymbol{\theta}, \mathbf{z}|\mathbf{x}) \quad (491)$$

$$\propto \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\boldsymbol{\theta}, \mathbf{z}|\mathbf{x}) \quad (492)$$

$$= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \quad (493)$$

$$= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}_{old}) \ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) \quad (494)$$

9.5

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

9.7

$$\nabla_{\mu} \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{kn} \nabla_{\mu} \ln \mathcal{N}(\mathbf{x}_n, \mu_k, \Sigma_k) \quad (495)$$

$$\nabla \ln \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) = -\frac{1}{2} \nabla (\mathbf{x}_n - \mu_k)^T \sigma_k^{-1} (\mathbf{x}_n - \mu_k) = \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \quad (496)$$

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

$$\mathcal{L}(\pi_k, \lambda) \ln p(\mathbf{X}, \mathbf{Z}|\mu_k \Sigma_k \pi_k) + \lambda[-1 + \sum_k \pi_k] \quad (497)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \sum_{n=1}^N z_{nk} \frac{1}{\pi_k} - \lambda \quad (498)$$

$$= \sum_{n=1}^N z_{nk} - \pi_k \lambda \quad (499)$$

$$= \sum_{n=1}^N \sum_{k=1}^K z_{nk} - \pi_k \lambda \quad (500)$$

$$= N - \lambda = 0 \iff \lambda = N \quad (501)$$

Substituting back:

$$\sum_{n=1}^N z_{nk} \frac{1}{\pi_k} = N \iff \pi_k = \frac{1}{N} \sum_{n=1}^N z_{nk} = \frac{N_k}{N} \quad (502)$$

$$\nabla_{\mu_k} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) [\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \quad (503)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \nabla \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (504)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (505)$$

$$\Rightarrow \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (506)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n - \underbrace{\sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\mu}_k}_{N_k} = 0 \quad (507)$$

$$(508)$$

Which leads to the correct answer.

9.9

$$\frac{\partial}{\partial \boldsymbol{\Sigma}_k} = \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) - \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \right] \quad (509)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T - \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln |\boldsymbol{\Sigma}^{-1}|^{\frac{1}{2}} \right] \quad (510)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T - \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{2} \ln |\boldsymbol{\Sigma}^{-1}| \right] \quad (511)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T - \frac{1}{2} \boldsymbol{\Sigma} \right] = 0 \quad (512)$$

$$\Leftrightarrow \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right] = \sum_{n=1}^N \gamma(z_{nk}) \frac{1}{2} \boldsymbol{\Sigma} = N_k \frac{1}{2} \boldsymbol{\Sigma}_k \quad (513)$$

From which the answer is easily seen.

9.11

$$\mathbb{E}_{\mathbf{z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(\mathbf{z}_{nk}) [\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \quad (514)$$

$$= \sum_{n=1}^N \sum_{k=1}^K \gamma(\mathbf{z}_{nk}) [\ln \pi_k + \ln \frac{1}{(2\pi\epsilon)^{M/2}} - \frac{1}{2\epsilon} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2] \quad (515)$$

$$\propto \sum_{n=1}^N \sum_{k=1}^K \gamma(\mathbf{z}_{nk}) [\epsilon \ln \pi_k + \epsilon \ln \frac{1}{(2\pi\epsilon)^{M/2}} - \frac{1}{2} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2] \quad (516)$$

$$= \lim_{\epsilon \rightarrow 0} -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 + C \quad (517)$$

9.12

$$\mathbb{E}[\mathbf{x}] = \int \mathbf{x} p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\pi}) d\mathbf{x} \quad (518)$$

$$= \int \mathbf{x} \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} \quad (519)$$

$$= \sum_{k=1}^K \pi_k \int \mathbf{x} p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} \quad (520)$$

$$= \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k \quad (521)$$

$$\text{cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \quad (522)$$

$$= \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T \quad (523)$$

$$(524)$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \int \mathbf{x}\mathbf{x}^T \sum_{k=1}^K \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} \quad (525)$$

$$= \sum_{k=1}^K \pi_k \int \mathbf{x}\mathbf{x}^T p(\mathbf{x}|\boldsymbol{\mu}_k) d\mathbf{x} = \sum_{k=1}^K \pi_k \mathbb{E}_k[\mathbf{x}\mathbf{x}^T] \quad (526)$$

$$(527)$$

$$\mathbb{E}_k[\mathbf{x}\mathbf{x}^T] = \mathbb{E}_k[(\mathbf{x} - \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}] + \mathbb{E}[\mathbf{x}])^T] \quad (528)$$

$$= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T + (\mathbf{x} - \mathbb{E}[\mathbf{x}]) \mathbb{E}[\mathbf{x}]^T + \mathbb{E}[\mathbf{x}] (\mathbf{x} - \mathbb{E}[\mathbf{x}])^T + \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}]^T] \quad (529)$$

$$= \underbrace{\text{cov}_k(\mathbf{x})}_{=\boldsymbol{\Sigma}_k} + \boldsymbol{\mu}\boldsymbol{\mu}^T \quad (530)$$

9.14

$$p(\mathbf{x}, \mathbf{z}, \boldsymbol{\mu})p(\mathbf{z}|\boldsymbol{\pi}) = \prod_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \pi_k^{z_k} \quad (531)$$

$$= \sum_{\mathbf{z}} \prod_{k=1}^K p(\mathbf{x}|\boldsymbol{\mu}_k)^{z_k} \pi_k^{z_k} \quad (532)$$

$$(533)$$

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

$$\frac{\partial}{\partial \mu_{ki}} = \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{\partial}{\partial \mu_{ki}} x_{ni} \ln \mu_{ki} + \frac{\partial}{\partial \mu_{ki}} (1 - x_{ni}) \ln(1 - \mu_{ki}) \right] \quad (534)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{x_{ni}}{\mu_{ki}} - \frac{1 - x_{ni}}{1 - \mu_{ki}} \right] \quad (535)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \left[\frac{x_{ni}(1 - \mu_{ki}) - \mu_{ki}(1 - x_{ni})}{\mu_{ki}(1 - \mu_{ki})} \right] = 0 \quad (536)$$

$$\Rightarrow \sum_{n=1}^N \gamma(z_{nk}) \mu_{ki} = \sum_{n=1}^N \gamma(z_{nk}) x_n \quad (537)$$

$$= N_k \mu_{ki} = \sum_{n=1}^N \gamma(z_{nk}) x_n \quad (538)$$

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

$$\frac{\partial}{\partial \alpha} = \frac{M}{2} \frac{1}{\alpha} - \frac{1}{2} \mathbb{E}[\mathbf{w}^T \mathbf{w}] = 0 \quad (539)$$

$$\Rightarrow \frac{M}{\alpha} = \mathbb{E}[\mathbf{w}^T \mathbf{w}] \quad (540)$$

$$\alpha = \frac{M}{\mathbb{E}[\mathbf{w}^T \mathbf{w}]} \quad (541)$$

9.24

$$\ln p(\mathbf{x}|\boldsymbol{\theta}) = \sum_z q(\mathbf{z}) \ln p(\mathbf{x}|\boldsymbol{\theta}) \quad (542)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \quad (543)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \frac{q(\mathbf{z})}{q(\mathbf{z})} \quad (544)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} + \sum_z q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})} \quad (545)$$

$$= \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} - \sum_z q(\mathbf{z}) \ln \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} \quad (546)$$

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.

$$\ln p(\mathbf{x}) = \int q(\mathbf{z}) \ln p(\mathbf{x}) d\mathbf{z} \quad (547)$$

$$= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \quad (548)$$

$$= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z})q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})q(\mathbf{z})} \quad (549)$$

$$= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} - \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{x}|\mathbf{z})} \quad (550)$$

$$(551)$$

10.2

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

10.8

Filling in:

$$\mathbb{E}[\tau] = \frac{a}{b} = (a_0 + \frac{N}{2})(b_0 + \frac{1}{2} \mathbb{E}_\mu[\sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^2 + \lambda_0(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^2])^{-1} \underset{N \rightarrow \infty}{\approx} b_N^{-1} \quad (552)$$

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

$$\ln p(\mathbf{x}) = \sum_m \sum_z q(\mathbf{z}, m) \ln p(\mathbf{x}) \quad (553)$$

$$= \sum_m \sum_z q(\mathbf{z}, m) \ln \frac{p(\mathbf{x}, m, \mathbf{z})}{p(m, \mathbf{z}|\mathbf{x})} \frac{q(\mathbf{z}, m)}{q(\mathbf{z}, m)} \quad (554)$$

Which results in 10.35 after rearrangement.

10.15

We have:

$$\mathbb{E}[\pi_k] = \frac{\alpha_k}{\sum_{k=1}^K \alpha_k} \quad (555)$$

and $\alpha_k = \alpha_0 + N_k$.

$$\mathbb{E}[\pi_k] = \frac{\alpha_0 + N_k}{\sum_{k=1}^K \alpha_0 + N_k} \quad (556)$$

$$= K\alpha_0 + \underbrace{\sum_{k=1}^K N_k}_{=N} \quad (557)$$

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

$$\frac{dd \ln(x)}{dx} = -x^{-2} \quad (558)$$

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

Taylor approximation:

$$y(x) = \ln(\xi) + \frac{1}{\xi}(x - \xi) = \lambda x - \ln(\lambda) - 1 \quad \lambda = \frac{1}{\xi} \quad (559)$$

$$g(\lambda) = \min_x [\lambda x - f(x)] \quad (560)$$

$$\frac{d}{dx} = \lambda - \frac{1}{x} = 0 \iff x = \frac{1}{\lambda} \quad (561)$$

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

$$\frac{d}{d\lambda} \lambda x - g(\lambda) = x - \frac{1}{\lambda} \iff x = \frac{1}{\lambda} \quad (562)$$

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

10.30

$$\frac{d}{dx} = \frac{1}{1 + e^{-x}} e^{-x} = \sigma(x) e^{-x} \quad (563)$$

$$\frac{dd}{dx} = \sigma(x)(1 - \sigma(x))e^{-x} - e^{-x}\sigma(x) = -\sigma^2(x)e^{-x} \quad (564)$$

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

$$\frac{d}{dx} = \frac{1}{1 + e^{-x}} e^{-x} \quad (565)$$

$$\frac{dd}{dx} = \sigma(x)(1 - \sigma(x))e^{-x} - e^{-x}\sigma(x) = -\sigma^2(x)e^{-x} \quad (566)$$

Taylor:

$$f(x) = f(\xi) + f'(\xi)(x - \xi) + \mathcal{O}(\xi^2) \quad (567)$$

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

$$f(x) \leq -\ln(1 + e^{-xi}) + \sigma(\xi)e^{-\xi}(x - \xi) + \mathcal{O}(\xi^2) \quad (568)$$

$$= -\ln(1 + e^{-xi}) + \sigma(\xi)e^{-\xi}x - \sigma(\xi)e^{-\xi}\xi + \mathcal{O}(\xi^2) \quad (569)$$

$$= \lambda x - g(\lambda) \quad \lambda = \sigma(\xi)e^{-\xi} \quad (570)$$

10.33

$$\frac{d}{d\xi_n} = (1 - \sigma(\xi_n)) - \frac{1}{2} - \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n \lambda'(\xi_n) + 2\xi_n \lambda(\xi_n) + \lambda'(\xi_n) \xi_n^2 \quad (571)$$

$$= -2\xi_n \lambda(\xi_n) - \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n \lambda'(\xi_n) + 2\xi_n \lambda(\xi_n) + \lambda'(\xi_n) \xi_n^2 \quad (572)$$

$$= -\lambda'(\xi_n)(\xi_n^2 - \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n) = 0 \quad (573)$$

$$\iff \xi_n^2 = \phi_n^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi_n \quad (574)$$

10.37

$$q^{\setminus j}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{\tilde{f}_j(\boldsymbol{\theta})} \quad (575)$$

$$Z_j = \int q^{\setminus j}(\boldsymbol{\theta}) f_j(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int q^{\setminus j}(\boldsymbol{\theta}) \tilde{f}_j(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int (q(\boldsymbol{\theta}) d\boldsymbol{\theta}) = 1 \quad (576)$$

$$\tilde{f}_j(\boldsymbol{\theta}) = \frac{q^{\text{new}}}{q^{\setminus j}(\boldsymbol{\theta})} = \tilde{f}_j = f_j \quad (577)$$

11.1

$$\mathbb{E}[\hat{f}] = \mathbb{E}\left[\frac{1}{L} \sum_{l=1}^L f(z^l)\right] \quad (578)$$

$$= \frac{1}{L} \sum_p (z) f(z^l) dz \quad (579)$$

$$= \frac{1}{L} \sum_{l=1}^L \mathbb{E}[f] = \mathbb{E}[f] \quad (580)$$

For the variance, we have to note that

$$\mathbb{E}[f(z^l), f(z^k)] = \text{Var}(f(z)) + \mathbb{E}[f(z)]^2 \quad (581)$$

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

$$\text{Var}[\hat{f}] = \mathbb{E}[\hat{f}^2] - \mathbb{E}[\hat{f}]^2 \quad (582)$$

$$= \frac{1}{L^2} \left(\sum_{l=1}^L \sum_{k=1}^L \mathbb{E}[f(z^l) f(z^k)] \right) - \mathbb{E}[f]^2 \quad (583)$$

$$= \frac{1}{L^2} \left(\sum_{l=1}^L \sum_{k=1}^L \mathbb{E}[f(z^l), f(z^k)] \right) - \mathbb{E}[f]^2 \quad (584)$$

$$= \frac{1}{L^2} L \text{Var}(f(z)) + L^2 \mathbb{E}[f]^2 - \mathbb{E}[f]^2 \quad (585)$$

$$= \frac{1}{L} \text{Var}(f(z)) \quad (586)$$

11.2

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

$$p(y) = 1 \cdot \left| \frac{dz}{dy} \right| = 1 \cdot p(y) \quad (587)$$

11.3

$$z = \left(\frac{1}{\pi} \int_{-\infty}^y \frac{1}{1 + \hat{y}^2} \right) \quad (588)$$

$$= \left(\frac{1}{\pi} [\arctan(y) - \arctan(-\infty)] \right) \quad (589)$$

$$= \left(\frac{1}{\pi} \arctan(y) + \frac{1}{2} \right) \quad (590)$$

$$y = h^{-1}(z) \quad (591)$$

$$\iff \pi z - \frac{\pi}{2} = \arctan(y) \quad (592)$$

$$\iff h^{-1}(z) = \tan\left(\pi z - \frac{\pi}{2}\right) \quad (593)$$

11.5

Hint 1: Show that the expectation and covariance are equal to μ and Σ .

$$\mathbb{E}[y] = \mathbb{E}_z[\mu] + \mathbb{E}_z[\mathbf{L}z] \quad (594)$$

$$= \mathbb{E}_z[\mu] + \mathbf{L} \mathbb{E}_z[z] \quad (595)$$

$$= \mathbb{E}_z[\mu] + \mathbf{L}\mathbf{0} \quad (596)$$

$$= \mu \quad (597)$$

$$\text{cov}[y] = \mathbb{E}[(y - \mu)(y - \mu)^T] \quad (598)$$

$$= \mathbb{E}[yy^T - \mu y^T - y \mu^T + \mu \mu^T] \quad (599)$$

$$= \mathbb{E}[(\mu + \mathbf{L}z)(\mu + \mathbf{L}z)^T - \mu(\mu + \mathbf{L}z)^T - (\mu + \mathbf{L}z)\mu^T + \mu \mu^T] \quad (600)$$

$$= \mathbb{E}[\mu \mu^T + \mu(\mathbf{L}z)^T + \mathbf{L}z \mu^T + \mathbf{L}z(\mathbf{L}z)^T - \mu \mu^T - \mu(\mathbf{L}z)^T - \mu \mu^T - \mathbf{L}z \mu^T + \mu \mu^T] \quad (601)$$

$$= \mathbb{E}[\mathbf{L}z(\mathbf{L}z)^T] \quad (602)$$

$$= \mathbf{L} \mathbb{E}[zz^T] \mathbf{L}^T \quad (603)$$

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

$$p(y) = p(z) \left| \frac{dz}{dy} \right| \quad (604)$$

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

$$y = b \tan z + c \iff \frac{y - c}{b} = \tan z \iff \arctan \frac{y - c}{b} = z \quad (605)$$

$$\frac{dz}{dy} = \frac{1}{1 + \left(\frac{y-c}{b}\right)^2} \frac{1}{b} \quad (606)$$

$$(607)$$

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 $E[(z'_i - \mu_i)^2]$.

$$\mathbb{E}_{z,\nu}[z'_i] = \mathbb{E}[\mu_i] + \mathbb{E}[\alpha(z_i - \mu_i)] + \mathbb{E}[\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu] \quad (608)$$

$$= \mu_i \quad (609)$$

$$\text{Var}_{z,\nu} = \mathbb{E}[z'^2_i] - \mathbb{E}[z'_i]^2 \quad (610)$$

$$= \mathbb{E}[(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)] - \mu_i^2 \quad (611)$$

$$= \mathbb{E}[(\mu_i + \alpha(z_i - \mu_i))(\mu_i + \alpha(z_i - \mu_i) + \sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu) + (\mu_i + \alpha(z_i - \mu_i))(\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)] - \mu_i^2 \quad (612)$$

$$+ (\mu_i + \alpha(z_i - \mu_i))(\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu) + (\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)(\sigma_i(1 - \alpha_i^2)^{\frac{1}{2}}\nu)] - \mu_i^2 \quad (613)$$

$$\vdots \quad (614)$$

$$(615)$$

11.15

$$\frac{dH}{dr_i} = \frac{1}{Z} \frac{d}{dr_i} r_i^2 = r_i \quad (616)$$

$$\frac{dr_i}{d\tau} = -\frac{dE(z)}{dz_i} = -\frac{dH}{dz_i} \quad (617)$$

11.16

$$p(\mathbf{r}|\mathbf{z}) = \frac{p(\mathbf{r}, \mathbf{z})}{p(\mathbf{z})} = \frac{Z_p}{Z_H} \exp(-K(\mathbf{r})) \quad (618)$$

11.17

12.3

$$\|u_i\|^2 = [\frac{1}{(N\lambda_i)^{1/2}} X^T v_i]^T [\frac{1}{(N\lambda_i)^{1/2}} X^T v_i] \quad (619)$$

$$= (N\lambda_i)^{-1} [X^T v_i]^T [X^T v_i] \quad (620)$$

$$= (N\lambda_i)^{-1} v^T X X^T v_i \quad (621)$$

$$= \lambda_i^{-1} v^T \lambda_i v_i \quad (622)$$

$$= 1 \quad (623)$$

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

$$\frac{\partial \ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2)}{\partial \boldsymbol{\mu}} = -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \quad (624)$$

$$= \sum_{n=1}^N \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = 0 \quad (625)$$

$$\Rightarrow 0 = -N\boldsymbol{\mu} + \sum_{n=1}^N \mathbf{x}_n \quad (626)$$

$$\Rightarrow \boldsymbol{\mu} = \bar{\mathbf{x}} \quad (627)$$

12.14

$$D(D-1) + 1 - (D-1)(D-2)/2 = D^2 - D + 1 - \frac{1}{2}D^2 + \frac{3}{2}D - 1 = \frac{1}{2}D^2 + \frac{1}{2}D = D(D+1)/2 \quad (628)$$

12.18

$$M \times D + 1 \quad (629)$$

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

$$\frac{\partial Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})}{\boldsymbol{\mu}_k} = \sum_{n=1}^N \gamma(z_{nk}) \frac{1}{p(\mathbf{x}_n|\boldsymbol{\phi}_k)} \frac{\partial}{\partial \boldsymbol{\mu}_k} p(\mathbf{x}_n|\boldsymbol{\phi}_k) \quad (630)$$

$$= \sum_{n=1}^N \gamma(z_{nk}) \frac{1}{p(\mathbf{x}_n|\boldsymbol{\phi}_k)} p(\mathbf{x}_n|\boldsymbol{\phi}_k) \frac{\partial}{\partial \boldsymbol{\mu}_k} - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad (631)$$

$$= 0 \iff \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (632)$$

$$\iff \boldsymbol{\mu} = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \quad (633)$$

$$(634)$$

$$\frac{\partial Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})}{\boldsymbol{\Sigma}_k} = \sum_{n=1}^N \gamma(z_{nk}) \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln p(\mathbf{x}|\boldsymbol{\phi}_k) \quad (635)$$

$$(636)$$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}_k} \ln p(\mathbf{x}|\boldsymbol{\phi}_k) = \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \left[\ln \frac{1}{(2\pi)^{D/2}} + \ln \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} - \frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (637)$$

$$= \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{2} \ln |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \frac{\partial}{\partial \boldsymbol{\Sigma}_k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) \quad (638)$$

$$= \frac{1}{2} \boldsymbol{\Sigma} - \frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \quad (639)$$

$$\frac{\partial}{\partial \Sigma_k} = \sum_{n=1}^N \gamma(z_{nk}) \left(\frac{1}{2} \Sigma - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \right) = 0 \quad (640)$$

$$\iff \Sigma_k = \frac{\sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma(z_{nk})} \quad (641)$$

13.17

$$h(z_1) = p(z_1|u_1)p(x_1|z_1, u_1) \quad (642)$$

$$f_n(z_{n-1}, z_n) = p(z_n|z_{n-1}, u_n)p(x_n|z_n, u_n) \quad (643)$$

13.19

This follows from:

$$\frac{d}{dz} \ln p(\mathbf{Z}) = \frac{d}{dz} \ln p(z_0)p(z_1|z_0) \times \cdots \times p(z_n|z_{n-1}) = \begin{bmatrix} \frac{\partial}{\partial z_0} \ln p(z_0) \\ \frac{\partial}{\partial z_1} \ln p(z_1|z_0) \\ \vdots \\ \frac{\partial}{\partial z_n} p(z_n|z_{n-1}) \end{bmatrix} \quad (644)$$

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

$$\mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M \epsilon_m(x) \right)^2 \right] = \mathbb{E} \left[\left(\frac{1}{M^2} \sum_{m=1}^M \sum_{l=1}^M \epsilon_m(x) \epsilon_l(x) \right) \right] \quad (645)$$

$$= \mathbb{E} \left[\frac{1}{M^2} \left(\sum_{m=1}^M \sum_{\substack{l=1 \\ l \neq m}}^M \epsilon_m(x) \epsilon_l(x) + \sum_{l=1}^M \epsilon_l(x)^2 \right) \right] \quad (646)$$

$$= \frac{1}{M^2} \sum_{l=1}^M \mathbb{E}[\epsilon_l(x)^2] \quad (647)$$

$$= \frac{1}{M} E_{AV} \quad (648)$$

14.3

Through Jensen's inequality:

$$\sum_{m=1}^M \frac{1}{M} \epsilon_m(x)^2 \geq \left(\sum_{m=1}^M \frac{1}{M} \epsilon_m(x) \right)^2 \quad (649)$$

Therefore:

$$\mathbb{E}_x \left[\sum_{m=1}^M \frac{1}{M} \epsilon_m(x)^2 \right] \geq \mathbb{E} \left[\left(\sum_{m=1}^M \frac{1}{M} \epsilon_m(x) \right)^2 \right] = E_{com} \quad (650)$$

14.6

$$\frac{d}{d\alpha_m} = (e^{\alpha_m/2} + e^{-\alpha/2}) \sum_{n=1}^N w_n^{(m)} I(y_m(x_n) = t_n) = e^{\alpha_m/2} \sum_{n=1}^N w_n^{(m)} = 0 \quad (651)$$

$$\frac{e^{-\alpha_m/2}}{e^{\alpha_m/2} + e^{-\alpha_m/2}} = \epsilon_m \quad (652)$$

Rearranging this results in the desired formula.

[14.7](#) [14.8](#)

14.9

$$\mathcal{L}_m(\mathbf{x}_n) = (y_n - \sum_{l=1}^M \alpha_l \hat{y}_l)^2 \quad (653)$$

$$= (y - \underbrace{\sum_{l=1}^{M-1} \alpha_l \hat{y}_l}_{\text{residual}} + \alpha_M \hat{y}_M)^2 \quad (654)$$

$$(655)$$

14.10

$$L = \frac{1}{2} \sum_{n=1}^N (t_n - t)^2 \quad (656)$$

$$\frac{dL}{dt} = - \sum_{n=1}^N (t_n - t) = 0 \iff Nt = \sum_{n=1}^N t_n \iff t = \frac{1}{N} \sum_{n=1}^N t_n \quad (657)$$

14.13

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) = \prod_{n=1}^N \prod_{k=1}^K (\pi_k \mathcal{N}(t_n | \mathbf{w}_k^T \boldsymbol{\phi}_n, \beta^{-1})^{z_{nk}} \quad (658)$$

Taking the log of this will yield the result.

14.14

$$\frac{dL}{d\pi_k} = \sum_{n=1}^N \gamma_{nk} / \pi_k - \lambda = 0 \quad (659)$$

$$\iff \sum_k \sum_{n=1}^N \gamma_{nk} = \sum_k \pi_k \lambda \quad (660)$$

$$(661)$$

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

$$\mathbb{E}[t | \hat{\boldsymbol{\phi}}, \boldsymbol{\theta}] = \sum_{k=1}^K \pi_k \mathbb{E}[t | \hat{\boldsymbol{\phi}}, \mathbf{w}_k, \beta] \quad (662)$$