

1 Introduction

Problem 1.1

Given that:

$$\sum_{j=0}^M A_{ij} w_j = T_i \quad (1)$$

Where:

$$\begin{aligned} A_{ij} &= \sum_{n=1}^N (x_n)^{i+j} \\ T_i &= \sum_{n=1}^N (x_n)^i t_n \\ y(x, w) &= w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M \\ E(w) &= \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 \end{aligned}$$

Show that the coefficients $\mathbf{w} = \{w_i\}$ that minimizes (1) error equation. To solve this problem, we will derivative the sum-of-squares error function (5).

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{1}{2} \sum_{i=0}^N 2\{y(x_n, w) - t_n\} x_i = 0 \\ \iff \sum_{i=0}^N \{y(x_n, w) x^i - t_n x^i\} &= 0 \\ \iff \sum_{i=0}^N \{y(x_n, w) x^i\} &= \sum_{i=0}^N t_n x^i \\ \iff \sum_{i=0}^N \sum_{j=0}^M w_j x_i^j x^i &= \sum_{i=0}^N t_n x_i \\ \iff \sum_{j=0}^M \sum_{i=0}^N w_j x_i^{i+j} &= \sum_{i=0}^N x_n^i t_n \end{aligned}$$

If we denote the above equation by (2) and (3), we get the answer. The problem is solved.

Problem 1.2

Given the regularized sum-of-squares equation:

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

Write down the set of coupled linear equations, analogous to equation from previous exercise, satisfied by the coefficients $\{w_i\}$. To solve this problem, we derivative the given equation like ex 1.1:

$$\begin{aligned}\frac{\partial \tilde{E}}{\partial w_i} &= \sum_{n=1}^N \{y(x_n, w) - t_n\} x_n^i + \lambda \|w\| = 0 \\ \sum_{n=1}^N \sum_{j=1}^M (w_j x_n^j x_n^i - t_n x_n^i) + \lambda \|w\| &= 0 \\ \sum_{j=1}^M (A_{ij} - T_i) + \lambda \|w\| &= 0\end{aligned}$$

*Note: Incorrect, the theta variable called kronecker delta and the derivative of equation is incorrect

Problem 1.3

Based on Bayes's Theorem:

$$P(X) = \sum_Y P(X|Y)P(Y)$$

The probability of selecting an apple P(a) from all boxes:

$$\begin{aligned}P(a) &= P(a|r) * P(r) + P(a|b) * P(b) + P(a|g) * P(g) = 0.3 * 0.2 + 0.5 * 0.2 + 0.3 * 0.6 \\ &= 0.34\end{aligned}$$

The probabiltiy that orange came from green box P(g|o) and based on Bayes's Theorem, we have:

$$P(g|o) = \frac{P(o|g) * P(g)}{P(o)}$$

We will calculate P(o) like P(a):

$$\begin{aligned}P(o) &= P(o|r) * P(r) + P(o|b) * P(b) + P(o|g) * P(g) = 0.4 * 0.2 + 0.5 * 0.2 + 0.3 * 0.6 \\ &= 0.36 \\ P(g|o) &= \frac{0.3 * 0.6}{0.36} = 0.5\end{aligned}$$

Problem 1.4

Equation (1.27) from the change of variable theorem:

$$\begin{aligned}p_y(y) &= p_x(x) \left| \frac{dx}{dy} \right| \\ &= p_x(g(y)) \|g'(y)\|\end{aligned}$$

We will calculate the derivative of equation (1.27) with respect to y, we got:

$$\frac{p_y(y)}{dy} = \frac{dp_x(x)|g'(y)|}{dy} = \frac{dp_x(x)}{dy} |g'(y)| + p_x(x) \frac{g'(y)}{dy} \quad (3)$$

Applying Chain Rule into first term in equation (1):

$$\frac{dp_x(x)}{dy}|g'(y)| = \frac{p_x(g(y))}{dg(y)} \frac{dg(y)}{dy}|g'(y)| \quad (4)$$

if \hat{x} is maximum of the density over x , we obtain:

$$\frac{dp_x(x)}{dx}|_{\hat{x}} = 0$$

Therefore, the equation (2) will equal 0, leading to the first term of equation (1) equal 0. In case of linear transformation, the second term of equation (1) will be vanish (ex: $y = ax + b$) so equation(1) equal 0. In conclusion, in case of linear transformation, the location of the maximum transforms in the same way as the variable itself.

Problem 1.5

$$\begin{aligned} E[E[f(x)]] &= E[f(x)] && \text{cause } E[f(x)] \text{ is constant} \\ \text{var}[f(x)] &= E[(f(x) - E[f(x)])^2] \\ &= E[f(x)^2 - 2f(x)E[f(x)] + E[f(x)]^2] \\ &= E[f(x)^2] - 2E[f(x)]^2 + E[f(x)]^2 \\ &= E[f(x)^2] - E[f(x)]^2 \end{aligned}$$

Problem 1.6

If x and y are independent, we obtain:

$$\begin{aligned} p_{xy}(x, y) &= p_x(x)p_y(y) \\ \int_{xy} p_{xy}(x, y) &= \int_x \int_y xyp_x(x)p_y(y) \\ &= \int_x xp_x(x) \int_y yp_y(y) \\ &= E[x]E[y] \end{aligned}$$

The result is leading to make the covariance of x and y equal 0. The problem is solved.

Problem 1.7

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

Transform Cartesian coordination $(x, y) \rightarrow$ polar coordination (r, θ) :

$$\begin{aligned} x &= r\cos(\theta), y = r\sin(\theta) \\ 0 &\leq r \leq +\infty, -\infty \leq \theta \leq 2\pi \\ I^2 &= \int_{-\infty}^{+\infty} \int_0^{+\infty} \exp\left(\frac{-1}{2\sigma^2}(r\cos\theta)^2 - \frac{1}{2\sigma^2}(r\sin\theta)^2\right) r dr d\theta \\ I^2 &= \int_{-\infty}^{+\infty} \int_0^{+\infty} \exp\left(\frac{-1}{2\sigma^2}r^2\right) r dr d\theta \end{aligned}$$

Substituting $u = r^2$ and calculate integral over u firstly:

$$\begin{aligned} & \int_0^{+\infty} \frac{1}{2} \exp\left(\frac{-1}{2\sigma^2}u\right) du \\ &= -\sigma^2 \exp\left(\frac{-1}{2\sigma^2}u\right) \Big|_0^{+\infty} \\ &= -\sigma^2(0 - 1) = \sigma^2 \end{aligned}$$

Integral over θ :

$$\begin{aligned} I^2 &= \int_{-\infty}^{2\pi} \sigma^2 d\theta = \sigma^2 \theta \Big|_{-\infty}^{2\pi} \\ I &= \sqrt{2\pi\sigma^2} \end{aligned}$$

A probability distribution is called "Normalized" if the sum of all possible result equal one.

$$\begin{aligned} & \int_{-\infty}^{+\infty} \mathcal{N}(x|\mu, \sigma) \\ &= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-1}{2\sigma^2}(x - \mu)^2\right) dx \\ &= \frac{\sqrt{2\pi\sigma^2}}{\sqrt{2\pi\sigma^2}} = 1 \end{aligned} \quad (\text{Solved})$$

Problem 1.8

$$\begin{aligned} E[x] &= \int_{-\infty}^{+\infty} \mathcal{N}(x|\mu, \sigma^2) x dx \\ &= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-1}{2\sigma^2}(x - \mu)^2\right) x dx \end{aligned}$$

Substituting $y = x - \mu$:

$$\begin{aligned} E[x] &= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-1}{2\sigma^2}y^2\right) (y + \mu) dy \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \left(\int_{-\infty}^{+\infty} \exp\left(\frac{-1}{2\sigma^2}y^2\right) y dy + \int_{-\infty}^{+\infty} \exp\left(\frac{-1}{2\sigma^2}y^2\right) \mu dy \right) \end{aligned}$$

Get the result from problem (1.7), we obtain:

$$E[x] = 0 + \mu = \mu$$

So the univariate Gauss Distribution given from (1.46) satisfies (1.49). Then we differentiate the equation (1.127) by using product rule.

$$\begin{aligned} \mathcal{N}(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{(x - \mu)^2}{2\sigma^2}\right) \\ \frac{d\mathcal{N}(x|\mu, \sigma^2)}{d\sigma^2} &= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)' \exp\left(\frac{(x - \mu)^2}{2\sigma^2}\right) + \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{(x - \mu)^2}{2\sigma^2}\right)' \end{aligned}$$

The derivative of First term:

$$\frac{-1}{2\sqrt{2\pi}\sigma^3} \exp\left(\frac{(x-\mu)^2}{2\sigma^2}\right)$$

The derivative of Second term:

$$\frac{(x-\mu)^2}{2\sqrt{2\pi}\sigma^5} \exp\left(\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Simplify the equation, we will obtain:

$$\begin{aligned} \int_{-\infty}^{+\infty} \left(\frac{-1}{2\sigma^2} + \frac{(x-\mu)^2}{2\sigma^4}\right) \mathcal{N}(x|\mu, \sigma^2) dx &= 0 \\ \int_{-\infty}^{+\infty} (x-\mu)^2 \mathcal{N}(x|\mu, \sigma^2) dx &= \sigma^2 \end{aligned}$$

Based on definition of variance, we proved:

$$\text{var}[x] = \sigma^2$$

Also we have:

$$E[x^2] = \text{var}[x] + E[x]^2 = \sigma^2 + \mu^2$$

Problem 1.9

The formula (1.52) is the general form of (1.46) so we will only prove this formula:

$$\begin{aligned} \frac{\partial \left(\frac{1}{2\pi^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp \frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right)}{\partial x} &= 0 \\ \frac{1}{2} (\Sigma^{-1} + (\Sigma^{-1})^T) (x-\mu) \mathcal{N}(x|\mu, \Sigma) &= 0 \\ \Sigma^{-1} (x-\mu) \mathcal{N}(x|\mu, \Sigma) &= 0 \\ x &= \mu \end{aligned}$$

Using these formulas to get above result:

$$\begin{aligned} x^T A &= A^T x \\ \frac{\partial x^T A x}{\partial x} &= (A + A^T) x \end{aligned}$$

The Second Derivative of formula (1.52):

$$\begin{aligned} \Sigma^{-1} \mathcal{N}(x|\mu, \Sigma) + \Sigma^{-1} (x-\mu) \frac{1}{2} (\Sigma^{-1} + (\Sigma^{-1})^T) (x-\mu) \mathcal{N}(x|\mu, \Sigma) \\ = \Sigma^{-1} \mathcal{N}(x|\mu, \Sigma) (1 + \Sigma^{-1} (x-\mu)^2) \end{aligned}$$

As we can see, Σ is positive semi-definite, the distribution is always positive so the the above equation is always positive. Applying it to Hessian matrix, all the eigenvalues of matrix are positive. It mean the point $x = \mu$ is the maximum point.

Problem 1.10

Proof of First Term:

$$E[x + z] = E[x] + E[z]$$

$$\int p(x)dx = 1$$

$$p(x, z) = p(x)p(z) \quad (\text{cause } x, z \text{ are independent})$$

$$\begin{aligned} E[x + z] &= \int \int (x + z)p(x, z)dx dz \\ &= \int \int xp(x)p(z)dx dz + \int \int zp(x)p(z)dx dz \\ &= E[x] + E[z] \end{aligned}$$

Proof of Second Term:

$$\begin{aligned} var[x + z] &= E[(x + z)^2] - E[x + z]^2 \\ &= \int \int (x + z)^2 p(x, z) - (\int \int (x + z)p(x, z))^2 \\ &= E[x^2] + E[z^2] + \int \int 2xz p(x)p(z)dx dz - (E[x]^2 + E[z]^2 + 2 \int \int 2xz p(x)p(z)dx dz) \\ &= E[x^2] - E[x]^2 + E[z^2] - E[z]^2 \\ &= var[x] + var[z] \end{aligned}$$

Problem 1.11

$$\ln \prod_{n=1}^N f(x) = \sum_{n=1}^N \ln(f(x))$$

From above equation, we obtain:

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

Setting the derivative of the log likelihood function with respect to μ :

$$\begin{aligned} \frac{d \ln(p(x|\mu, \sigma^2))}{d\mu} &= \frac{1}{\sigma^2} \sum_{n=1}^N (x_n - \mu) = 0 \\ \sum_{n=1}^N x_n - N\mu &= 0 \\ \mu_{ML} &= \frac{1}{N} \sum_{n=1}^N x_n \end{aligned}$$

Setting the derivative of the log likelihood function with respect to σ^2 :

$$\frac{dp(x|\mu, \sigma^2)}{d\sigma^2} = \frac{1}{2\sigma^4} \sum_{n=1}^N N(x_n - \mu)^2 - \frac{N}{2\sigma^2} = 0$$

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

Problem 1.12

If $n \neq m$ then x_n and x_m are independent. Hence $E[x_n x_m] = E[x_n][x_m] = \mu^2 + \sigma^2$. In case of $n = m$, $E[x_n x_m] = E[x] = \mu^2$. Combine those results, we obtain the equation (1.130).

Easy to prove the equation (1.49):

$$E[\mu_{ML}] = E\left[\frac{1}{N} \sum_{n=1}^N x_n\right] = \frac{1}{N} \sum_{n=1}^N E[x] = \frac{1}{N} N\mu = \mu$$

With equation (1.50):

$$\begin{aligned} E[\sigma_{ML}^2] &= E\left[\frac{1}{N} \sum_{n=1}^N (x - \mu_{ML})^2\right] = \frac{1}{N} \sum_{n=1}^N E[(x^2 - 2x\mu_{ML} + \mu_{ML}^2)] \\ &= \frac{1}{N} \sum_{n=1}^N (E[x^2] - 2E[x\mu_{ML}] + E[\mu_{ML}^2]) \\ &= \mu^2 + \sigma^2 - \frac{2}{N^2} E\left[\sum_{n=1}^N \sum_{m=1}^M x_{nm}^2\right] + \frac{1}{N^2} E\left[\left(\sum_{p=1}^P x_p\right)^2\right] \end{aligned}$$

Because $x = m = p$, hence:

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

Based on equation (1.130), we will obtain

$$\begin{aligned} : E[\mu_{ML}] &= \mu^2 + \sigma^2 - \mu^2 - \frac{1}{N} \sigma^2 \\ &= \frac{N-1}{N} \sigma^2 \end{aligned}$$

Problem 1.13

We will solve similarly like problem (1.12):

$$\begin{aligned}
E[\sigma_{ML}^2] &= E\left[\frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2\right] \\
&= \frac{1}{N} \sum_{n=1}^N (E[x^2] - 2E[x\mu] + E[\mu^2]) \\
&= \frac{1}{N} \sum_{n=1}^N (\mu^2 + \sigma^2 - 2\mu^2 + \mu^2) \\
&= \sigma^2
\end{aligned}$$

In case of Maximum Likelihood Mean:

$$\begin{aligned}
E[\mu_{ML}^2] &= E\left[\left(\frac{1}{N} \sum_{n=1}^N x_n\right)^2\right] \\
&= \frac{1}{N^2} (N(N\mu^2 + \sigma^2)) = \mu^2 + \frac{\sigma^2}{N}
\end{aligned}$$

Problem 1.17

The derivative of $\Gamma(x+1)$:

$$\begin{aligned}
\Gamma(x+1) &= \int_0^\infty u^x e^{-u} du \\
&= -u^x e^{-u} \Big|_0^\infty + x\Gamma(x)
\end{aligned}$$

Through L'hospital rule, we obtain:

$$\lim_{u \rightarrow \infty} \frac{u^x}{e^u} = \lim_{u \rightarrow \infty} \frac{x!}{e^\infty} = 0$$

So we have proved $\Gamma(x+1) = x\Gamma(x)$.

$$\Gamma(1) = \int_0^\infty e^{-u} du = e^{-u} \Big|_0^\infty = (1 - 0) = 1$$

Based on two above proved equation, we will obtain:

$$\Gamma(x+1) = x\Gamma(x) = x!\Gamma(1) = x!$$

Problem 1.35

The Differential Entropy of Gaussian Distribution:

$$\begin{aligned}
H(x) &= - \int p(x) \ln p(x) dx \\
&= - \int p(x) \ln \left\{ \frac{1}{(2\pi\sigma^2)^{1/2}} \exp -\frac{(x-\mu)^2}{2\sigma^2} \right\} dx \\
&= - \int p(x) \ln \left\{ \frac{1}{(2\pi\sigma^2)^{1/2}} \right\} dx + \int p(x) \left(\frac{(x-\mu)^2}{2\sigma^2} \right) dx \\
&= \frac{1}{2} \ln(2\pi\sigma^2) + \frac{1}{2} = \frac{1}{2} (1 + \ln(2\pi\sigma^2))
\end{aligned}$$

Problem 1.37

We already have:

$$\begin{aligned}
H[x] &= - \int p(x) \ln p(x) dx \\
H[y|x] &= - \int \int p(x, y) \ln p(y|x) dx dy
\end{aligned}$$

It's straightforward by adding two above equations:

$$\begin{aligned}
H[x] + H[y|x] &= - \left(\int p(x) \ln p(x) dx + \int \int p(x, y) \ln p(y|x) dx dy \right) \\
&= - \left(\int \int p(x, y) \ln p(x) dx dy + \int \int p(x, y) \ln p(y|x) dx dy \right) \\
&= - \left(\int \int p(x, y) \ln (p(x)p(y|x)) dx dy \right) \\
&= - \int \int p(x, y) \ln p(x, y) dx dy = H[x, y] \quad (\text{Solved})
\end{aligned}$$

2 Probability Distributions

Problem 2.1

Bernoulli Distribution:

$$\begin{aligned}
Bern(x|\mu) &= \mu^x (1-\mu)^{1-x} \sum_{x=0,1} Bern(x|\mu) = \sum_{x=0,1} \mu^x (1-\mu)^{1-x} = (1-\mu) + \mu = 1 \\
E[x] &= \sum_{x=0,1} Bern(x|\mu) x = \sum_{x=0,1} x \mu^x (1-\mu)^{1-x} = 0 + \mu = \mu \\
var[x] &= \sum_{x=0,1} (E[x^2] - E[x]^2) = \sum_{x=0,1} \left(\sum_{x=0,1} \mu^x (1-\mu)^{1-x} x^2 \right) + \mu^2 \\
&= \mu + \mu^2 = \mu(1 + \mu)
\end{aligned}$$

The entropy $H(x)$ of Bernoulli Distribution:

$$\begin{aligned} H(x) &= - \sum_{x=0,1} \text{Bern}(x|\mu) \ln \text{Bern}(x|\mu) = - \sum_{x=0,1} \mu^x (1-\mu)^{1-x} \ln(\mu^x (1-\mu)^{1-x}) \\ &= -(1-\mu) \ln 1-\mu - \mu \ln \mu \end{aligned}$$

Problem 2.2

This problem is the same as problem (2.1), so we will solve similarly:

$$\begin{aligned} \sum_{x=-1,1} p(x|\mu) &= \sum_{x=-1,1} \left(\frac{1-\mu}{2}\right)^{(1-x)/2} \left(\frac{1+\mu}{2}\right)^{(1+x)/2} \\ &= \frac{1-\mu}{2} + \frac{1+\mu}{2} = 1 \end{aligned}$$

So the distribution (2.261) is normalized

$$\begin{aligned} E[x] &= \sum_{x=-1,1} \left(\frac{1-\mu}{2}\right)^{(1-x)/2} \left(\frac{1+\mu}{2}\right)^{(1+x)/2} x \\ &= -\frac{1-\mu}{2} + \frac{1+\mu}{2} = \frac{2\mu}{2} = \mu \text{var}[x] = E[x^2] - E[x]^2 = \sum_{x=-1,1} \left(\frac{1-\mu}{2}\right)^{(1-x)/2} \left(\frac{1+\mu}{2}\right)^{(1+x)/2} x^2 - \mu^2 \\ &= \frac{1-\mu}{2} + \frac{1+\mu}{2} - \mu^2 = \mu - \mu^2 = \mu(1-\mu) \\ H[x] &= - \sum_{x=-1,1} p(x|\mu) \ln(p(x|\mu)) = \sum_{x=-1,1} \left(\frac{1-\mu}{2}\right)^{(1-x)/2} \left(\frac{1+\mu}{2}\right)^{(1+x)/2} \ln \left(\frac{1-\mu}{2}\right)^{(1-x)/2} \left(\frac{1+\mu}{2}\right)^{(1+x)/2} \\ &= -\left(\frac{1-\mu}{2}\right) \ln \frac{1-\mu}{2} - \left(\frac{1+\mu}{2}\right) \ln \frac{1+\mu}{2} \end{aligned}$$

Problem 2.3

Transforming the equation (2.262) based on the definition (2.10) of the number of combinations, we will obtain:

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

So we will solve the problem with above equation:

$$\begin{aligned} \frac{N!}{(N-m)!m!} + \frac{N!}{(N-m+1)!(m-1)!} &= \frac{N!(N-m+1)}{(N-m+1)!m!} + \frac{N!m}{(N-m+1)!m!} \\ &= \frac{N!(N+1)}{(N-m+1)!m!} = \frac{(N+1)!}{(N-m+1)!m!} \quad (\text{Solved}) \end{aligned}$$

We will prove equation (2.263) by induction and easy to see that the equation is true if $N = 1$. Assuming the equation (2.263) holds true on N then we will prove that equation also holds:

$$\begin{aligned}
(1+x)^{N+1} &= (1+x) \sum_{m=0}^N C_N^m x^m = x \sum_{m=0}^N C_N^m x^m + \sum_{m=0}^N C_N^m x^m \\
&= \sum_{m=0}^N C_N^m x^{m+1} + \sum_{m=0}^N C_N^m x^m = \sum_{m=1}^{N+1} C_N^{m-1} x^m + \sum_{m=0}^N C_N^m x^m \\
&= \sum_{m=1}^N (C_N^m + C_N^{m-1}) x^m + x^{N+1} + x^0 \\
&= \sum_{m=1}^N C_{N+1}^m x^m + x^{N+1} + x^0 \\
&= \sum_{m=0}^{N+1} C_{N+1}^m x^m \quad (Solved)
\end{aligned}$$

Then we will show that the equation (2.264) holds true by induction like above. Substituting $y = (1-x)$, we obtain:

$$\begin{aligned}
(x+y)^{N+1} &= (x+y) \sum_{m=0}^N C_N^m x^m y^{N-m} \\
&= x \sum_{m=0}^N C_N^m x^m y^{N-m} + y \sum_{m=0}^N C_N^m x^m y^{N-m} \\
&= \sum_{m=0}^N C_N^m x^{m+1} y^{N-m} + y \sum_{m=0}^N C_N^m x^m y^{N-m+1} \\
&= \sum_{m=1}^{N+1} C_{N+1}^{m-1} x^m y^{N-m} + \sum_{m=0}^N C_N^m x^m y^{N-m+1} \\
&= \sum_{m=1}^N (C_{N+1}^{m-1} + C_N^m) x^m y^{N-m} + x^{N+1} + y^{N+1} \\
&= \sum_{m=1}^N C_{N+1}^m x^m y^{N-m} + x^{N+1} + y^{N+1} \\
&= \sum_{m=0}^{N+1} C_{N+1}^m x^m y^{N-m+1}
\end{aligned}$$

We have solved the equation by induction then substitute again $y = (1-x)$, we will prove the binomial distribution is normalized.

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

Problem 2.5

let $t = x + y$ and $x = t\mu$, we obtain:

$$x = t\mu \quad y = t(1 - \mu) \quad t = x + y \quad \mu = \frac{x}{x + y}$$

Using Jacobian Determinant for multiple integral:

$$\begin{aligned} \frac{\partial(x, y)}{\partial(\mu, t)} &= \begin{bmatrix} \frac{\partial x}{\partial \mu} & \frac{\partial x}{\partial t} \\ \frac{\partial y}{\partial \mu} & \frac{\partial y}{\partial t} \end{bmatrix} = \begin{bmatrix} t & \mu \\ -t & 1 - \mu \end{bmatrix} = t \\ \Gamma(a)\Gamma(b) &= \int_0^\infty \exp(-x)x^{a-1}dx \int_0^\infty \exp(-y)y^{b-1}dy \\ &= \int_0^\infty \int_0^\infty \exp(-x-y)x^{a-1}y^{b-1}dxdy \\ &= \int_0^\infty \int_0^1 \exp(-t)(t\mu)^{a-1}(t(1-\mu))^{b-1}tdtd\mu \\ &= \int_0^\infty \exp(-t)t^{a+b-1}dt \int_0^1 \mu^{a-1}(1-\mu)^{b-1}d\mu \\ &= \Gamma(a+b) \int_0^1 \mu^{a-1}(1-\mu)^{b-1}d\mu \\ &\rightarrow \int_0^1 \mu^{a-1}(1-\mu)^{b-1}d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \end{aligned}$$

Problem 2.6

$$\begin{aligned} E[\mu] &= \int_0^1 \mu^{a-1}(1-\mu)^{b-1} \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \mu d\mu \\ &= \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^a(1-\mu)^{b-1}d\mu \end{aligned}$$

We have equation $\Gamma(x+1) = x\Gamma(x)$ that have proved in problem (1.17). Hence we obtain:

$$\begin{aligned} &= \frac{\Gamma(a+b+1)a}{\Gamma(a)\Gamma(b)(a+b)} \int_0^1 \mu^a(1-\mu)^{b-1}d\mu \\ &= \frac{a}{a+b} \frac{\Gamma(a+b+1)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} \\ &= \frac{a}{a+b} \end{aligned}$$

In case of variance of beta distribution, we will solve it similarly like $E[\mu]$:

$$\begin{aligned} var[\mu] &= E[x^2] - E[x]^2 = \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a+1}(1-\mu)^{b-1} d\mu - \left(\frac{a}{a+b}\right)^2 \\ &= \frac{(a+1)a}{(a+b)(a+b+1)} \frac{\Gamma(a+b+2)}{\Gamma(a+2)\Gamma(b)} \int_0^1 \mu^{a+1}(1-\mu)^{b-1} d\mu - \left(\frac{a}{a+b}\right)^2 \\ &= \frac{(a+1)a}{(a+b)(a+b+1)} - \left(\frac{a}{a+b}\right)^2 = \frac{ab}{(a+b)^2(a+b+1)} \end{aligned}$$

Differentiating the beta distribution and give it equal zero:

$$\begin{aligned} \frac{\partial Beta(\mu|a,b)}{\partial \mu} &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} ((a-1)\mu^{a-2}(1-\mu)^{b-1} - (b-1)\mu^{a-1}(1-\mu)^{b-2}) = 0 \\ \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} (\mu^{a-2}(1-\mu)^{b-2}((1-\mu)(a-1) - (b-1)\mu)) &= 0 \\ (1-\mu)(a-1) - (b-1)\mu &= 0 \\ \mu &= \frac{a-1}{a+b-2} \end{aligned}$$

So the mode of the beta distribution is $\mu = \frac{a-1}{a+b-2}$ when $a > 1$ and $b > 1$.

Problem 2.8

Proof of equation (2.270):

$$\begin{aligned} E_y[E_x[x|y]] &= E_y\left[\int xp(x|y)dx\right] \\ &= \int \int xp(y)p(x|y)dxdy \\ &= \int \int xp(x,y)dxdy \\ &= \int xp(x)dx = E[x] \end{aligned} \quad (\text{Solved})$$

Proof of equation (2.271):

$$var_x[x] = E_x[x^2] - E_x[x]^2$$

Separating the above equation into two terms. The First Term:

$$\begin{aligned} E_x[x^2] &= E_y[E_x[x^2|y]] \\ &= E_y[var_x[x|y] + E_x[x|y]^2] \\ &= E_y[var_x[x|y]] + E_y[E_x[x|y]^2] \\ &= E_y[var_x[x|y]] + var_y[E_x[x|y]] + E_y[E_x[x|y]]^2 \end{aligned}$$

The Second Term:

$$E_x[x]^2 = E_y[E_x[x|y]]^2 \quad (\text{Based on equation 2.270})$$

Combining two terms, we will obtain:

$$\begin{aligned} \text{var}[x] &= E_y[\text{var}_x[x|y]] + \text{var}_y[E_x[x|y]] + E_y[E_x[x|y]]^2 - E_y[E_x[x|y]]^2 \\ &= E_y[\text{var}_x[x|y]] + \text{var}_y[E_x[x|y]] \end{aligned} \quad (\text{Solved})$$

Problem 2.10

Based on Dirichlet Distribution given by 2.38, we have:

$$\begin{aligned} E[\mu_j] &= \int \text{Dir}(\mu|\alpha) \mu_j d\mu \\ &= \int \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k-1} \mu_j du \\ &= \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \int \prod_{k=1}^K \mu_k^{\alpha_k-1} \mu_j du \\ &= \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \frac{\Gamma(\alpha_1) \cdots \Gamma(\alpha_j + 1) \cdots \Gamma(\alpha_K)}{\Gamma(\alpha_0 + 1)} \\ &= \frac{\Gamma(\alpha_0) \Gamma(\alpha_j + 1)}{\Gamma(\alpha_0 + 1) \Gamma(\alpha_j)} \end{aligned}$$

We also have $\Gamma(\alpha + 1) = \alpha \Gamma(\alpha)$. Hence:

$$E[\mu_j] = \frac{\alpha_j}{\alpha_0}$$

We will solve as above with $E[\mu_j^2]$:

$$\begin{aligned} E[\mu_j^2] &= \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \frac{\Gamma(\alpha_1) \cdots \Gamma(\alpha_j + 2) \cdots \Gamma(\alpha_K)}{\Gamma(\alpha_0 + 2)} \\ &= \frac{\Gamma(\alpha_0) \Gamma(\alpha_j + 2)}{\Gamma(\alpha_0 + 2) \Gamma(\alpha_j)} \\ &= \frac{\alpha_j(\alpha_j + 1)}{\alpha_0(\alpha_0 + 1)} \end{aligned}$$

Bring the above result into equation $\text{var}[\mu_j]$:

$$\begin{aligned} \text{var}[\mu_j] &= E[\mu_j^2] - E[\mu_j]^2 \\ &= \frac{\alpha_j(\alpha_j + 1)}{\alpha_0(\alpha_0 + 1)} - \frac{\alpha_j^2}{\alpha_0^2} \\ &= \frac{\alpha_j(\alpha_0 + 1)}{\alpha_0^2(\alpha_0 + 1)} \end{aligned}$$

In case of Covariance:

$$\begin{aligned}
cov[\mu_j \mu_l] &= E[(\mu_j - E[\mu_j])(\mu_l - E[\mu_l])] \\
&= \int (\mu_j - E[\mu_j])(\mu_l - E[\mu_l]) Dir(x|\mu) d\mu \\
&= \int (\mu_j \mu_l - E[\mu_l] \mu_j - E[\mu_j] \mu_l + E[\mu_j] E[\mu_l]) Dir(x|\mu) d\mu \\
&= \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \frac{\Gamma(\alpha_1) \cdots \Gamma(\alpha_j + 1) \Gamma(\alpha_l + 1) \cdots \Gamma(\alpha_K)}{\Gamma(\alpha_0 + 2)} - E[\mu_l] E[\mu_j] - E[\mu_j] E[\mu_l] + E[\mu_j] E[\mu_l] \\
&= \frac{\alpha_j \alpha_l}{\alpha_0(\alpha_0 + 1)} - E[\mu_j] E[\mu_l] \\
&= \frac{\alpha_j \alpha_l}{\alpha_0(\alpha_0 + 1)} - \frac{\alpha_j \alpha_l}{\alpha_0^2} \\
&= -\frac{\alpha_j \alpha_l}{\alpha_0^2(\alpha_0 + 1)}
\end{aligned}$$

Problem 2.11

Differentiating the Dirichlet Distribution:

$$\begin{aligned}
\frac{\partial Dir(\mu|\alpha)}{\partial \alpha_j} &= \partial \left(\frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k - 1} \right) / \partial \alpha_j \\
&= \frac{\partial \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)}}{\partial \alpha_j} \prod_{k=1}^K \mu_k^{\alpha_k - 1} + \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \frac{\partial \prod_{k=1}^K \mu_k^{\alpha_k - 1}}{\partial \alpha_j}
\end{aligned}$$

We will handle with second term:

$$\frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \frac{\partial \prod_{k=1}^K \mu_k^{\alpha_k - 1}}{\partial \alpha_j} = \ln(\alpha_j) Dir(\mu|\alpha)$$

With above result, we integrate both sides of equation and handle the left side firstly:

$$\int \frac{\partial Dir(\mu|\alpha)}{\partial \alpha_j} d\mu = \int \frac{\partial 1}{\partial \alpha_j} d\mu = 0$$

The Right Side:

$$\begin{aligned}
&\int \frac{\partial \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)}}{\partial \alpha_j} \prod_{k=1}^K \mu_k^{\alpha_k - 1} d\mu + \int \ln(\alpha_j) Dir(\mu|\alpha) d\mu \\
&= \frac{\partial \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)}}{\partial \alpha_j} \frac{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)}{\Gamma(\alpha_0)} + E[\ln(\alpha_j)] \\
&= \frac{\partial \ln \left(\frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \right)}{\partial \alpha_j} + E[\ln(\alpha_j)]
\end{aligned}$$

Combining both sides, we obtain:

$$\begin{aligned}
E[\ln(\mu)] &= -\frac{\partial \ln\left(\frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_K)}\right)}{\partial \alpha_j} \\
&= -\frac{\partial \ln(\Gamma(\alpha_0)) - \ln(\Gamma(\alpha_1)\cdots\Gamma(\alpha_K))}{\partial \alpha_j} \\
&= \frac{\partial \ln(\Gamma(\alpha_j))}{\partial \alpha_j} - \frac{\partial \ln(\Gamma(\alpha_0))}{\partial \alpha_j} \\
&= \frac{\partial \ln(\Gamma(\alpha_j))}{\partial \alpha_j} - \frac{\partial \ln(\Gamma(\alpha_0))}{\partial \alpha_0} \frac{\partial \alpha_0}{\partial \alpha_j} = \psi(\alpha_j) - \psi(\alpha_0)
\end{aligned}$$

Problem 2.12

Based on equation (2.278):

$$\int_a^b U(x|a, b)dx = \int_a^b \frac{1}{b-a}dx = \frac{b}{b-a} - \frac{a}{b-a} = 1$$

Hence the Uniform Distribution is normalized. Next is mean of the distribution:

$$E[x] = \int_a^b xU(x|a, b)dx = \frac{1}{2(b-a)}x^2\Big|_a^b = \frac{b^2}{2(b-a)} - \frac{a^2}{2(b-a)} = \frac{(a+b)}{2}$$

In case of variance of the Uniform Distribution:

$$\begin{aligned}
var[x] &= E[x^2] - E[x]^2 = \int_a^b x^2U(x|a, b)dx - \left(\frac{(a+b)}{2}\right)^2 \\
&= \frac{b^3 - a^3}{3(b-a)} - \frac{(a+b)^2}{4} \\
&= \frac{a^2 + b^2 - 2ab}{12} = \frac{(a+b)^2}{12}
\end{aligned}$$

Problem 2.13

The Gaussian Distribution:

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

KL Divergence:

$$KL(p(x)||q(x)) = \int p(x) \ln\left\{\frac{p(x)}{q(x)}\right\} dx$$

With two above equations, we obtain:

$$KL(p(x)||q(x)) = \int p(x) \frac{1}{2} \left(\ln\left(\frac{|L|}{|\Sigma|}\right) - (x-\mu)^T \Sigma^{-1}(x-\mu) + (x-m)^T L^{-1}(x-m) \right) dx$$

We have property of expectation:

$$\begin{aligned} E[(x - \mu)^T A(x - \mu)] &= \text{Tr}(A\Sigma) + (x - \mu)^T A(x - \mu) \\ E[x] &= \int p(x)x dx \\ \int p(x)dx &= 1 \end{aligned}$$

Hence:

$$\begin{aligned} KL(p(x)||q(x)) &= \frac{1}{2}(\ln(\frac{|L|}{|\Sigma|}) + E[(x - \mu)^T \Sigma^{-1}(x - \mu)] - E[(x - m)^T L^{-1}(x - m)]) \\ &= \frac{1}{2}(\ln(\frac{|L|}{|\Sigma|}) + \text{Tr}(\Sigma^{-1}\Sigma) + (x - \mu)^T \Sigma^{-1}(x - \mu) - \text{Tr}(L^{-1}\Sigma) - (x - m)^T L^{-1}(x - m)) \\ &= \frac{1}{2}(\ln(\frac{|L|}{|\Sigma|}) + \text{Tr}(I_D) + (\mu - \mu)^T \Sigma^{-1}(\mu - \mu) - \text{Tr}(L^{-1}\Sigma) - (\mu - m)^T L^{-1}(\mu - m)) \\ &= \frac{1}{2}(\ln(\frac{|L|}{|\Sigma|}) + D - \text{Tr}(L^{-1}\Sigma) - (\mu - m)^T L^{-1}(\mu - m)) \end{aligned}$$

Problem 2.30

Mean of z:

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

Covariance of z:

$$\text{cov}[z] = R^{-1} = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1}A^T \\ A\Lambda^{-1} & L^{-1} + A\Lambda^{-1}A^T \end{pmatrix}$$

Hence:

$$E[z] = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1}A^T \\ A\Lambda^{-1} & L^{-1} + A\Lambda^{-1}A^T \end{pmatrix} \begin{pmatrix} \Lambda\mu - A^T Lb \\ Lb \end{pmatrix} = \begin{pmatrix} \mu \\ A\mu + b \end{pmatrix}$$

Problem 2.35

The proof of 2.62 is in book so we don't prove it. Consequently, we will prove the equation (2.124)

$$\begin{aligned} E[\Sigma_{ML}] &= E[\frac{1}{N} \sum_{n=1}^N (x - \mu_{ML})^T (x - \mu_{ML})] \\ &= \frac{1}{N} E[\sum_{n=1}^N (x_n x_n^T - 2x_n \mu_{ML} + \mu_{ML}^T \mu_{ML})] \\ &= \frac{1}{N} \sum_{n=1}^N E[x_n x_n^T] - \frac{2}{N} \sum_{n=1}^N E[x_n \mu_{ML}] + \frac{1}{N} \sum_{n=1}^N E[\mu_{ML}^T \mu_{ML}] \end{aligned}$$

For the first term:

$$\frac{1}{N} \sum_{n=1}^N E[x_n x_n^T] = \frac{1}{N} (N(\mu\mu^T + \Sigma)) = \mu\mu^T + \Sigma$$

For the second term:

$$\begin{aligned} -\frac{2}{N} \sum_{n=1}^N E[x_n \mu_{ML}] &= -\frac{2}{N} \sum_{n=1}^N E[x_n (\frac{1}{N} \sum_{m=1}^M x_m)] \\ &= -\frac{2}{N^2} \sum_{n=1}^N \sum_{m=1}^M E[x_n x_m] \\ &= -\frac{2}{N^2} \sum_{n=1}^N \sum_{m=1}^M (\mu\mu^T + I_{nm}\Sigma) \\ &= -\frac{2}{N^2} N^2 (\mu\mu^T + \frac{1}{N} \Sigma) = -2(\mu\mu^T + \frac{1}{N} \Sigma) \end{aligned}$$

For the third term:

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^N E[\mu_{ML} \mu_{ML}^T] &= \frac{1}{N} E[\sum_{n=1}^N \sum_{m=1}^M x_n x_m] \\ &= \frac{1}{N^2} N^2 (\mu\mu^T + \frac{1}{N} \Sigma) \end{aligned}$$

Combine all above results, we obtain:

$$\begin{aligned} E[\Sigma_{ML}] &= \mu\mu^T + \Sigma - 2(\mu\mu^T + \frac{1}{N} \Sigma) + \mu\mu^T + \frac{1}{N} \Sigma \\ &= \Sigma - \frac{1}{N} \Sigma = \frac{N-1}{N} \Sigma \end{aligned}$$

Problem 2.41

$$\begin{aligned} \int_0^\infty \text{Gamma}(\lambda|a, b) d\lambda &= 1 \int_0^\infty \text{Gamma}(\lambda|a, b) d\lambda = \int_0^\infty \frac{1}{\Gamma(a)} b^a \lambda^{a-1} e^{-b\lambda} d\lambda \\ &= \frac{b^a}{\Gamma(a)} \int_0^\infty \lambda^{a-1} e^{-b\lambda} d\lambda \end{aligned}$$

Substituting $u = b\lambda$, we obtain:

$$\begin{aligned} &= \frac{b^a}{\Gamma(a)} \int_0^\infty \frac{1}{b} \left(\frac{u}{b}\right)^{a-1} e^{-u} du \\ &= \frac{b^a}{\Gamma(a)} \int_0^\infty \left(\frac{1}{b}\right)^a u^{a-1} e^{-u} du \\ &= \frac{1}{\Gamma(a)} \Gamma(a) = 1 \end{aligned}$$

Solved

Problem 2.42

I will get result from problem (2.41):

$$\begin{aligned} E[\lambda] &= \int_0^\infty \text{Gamma}(\lambda|a, b) \lambda d\lambda \\ &= \frac{b^a}{\Gamma(a)} \int_0^\infty \int_0^\infty \lambda^a e^{-b\lambda} d\lambda \\ &= \frac{b^a}{\Gamma(a)} \frac{1}{b^{a+1}} \Gamma(a+1) \end{aligned}$$

We have proved $\Gamma(x+1) = x\Gamma(x)$. Hence:

$$\begin{aligned} E[\lambda] &= b^a \frac{1}{b^{a+1}} \frac{1}{\Gamma(a)} a\Gamma(a) = \frac{a}{b} \text{(Solved)} \\ \text{var}[\lambda] &= E[\lambda^2] - E[\lambda]^2 \\ &= \int_0^\infty \text{Gamma}(\lambda|a, b) \lambda^2 d\lambda - \left(\frac{a}{b}\right)^2 \\ &= \frac{b^a}{\Gamma(a)} \frac{1}{b^{a+2}} \Gamma(a+2) - \frac{a^2}{b^2} \\ &= \frac{a(a+1)}{b^2} - \frac{a^2}{b^2} = \frac{a}{b^2} \text{(Solved)} \\ \text{mode}[\lambda] &= \frac{\partial}{\partial \lambda} \frac{1}{\Gamma(a)} b^a \lambda^{a-1} e^{-b\lambda} \\ &= \frac{1}{\Gamma(a)} b^a ((a-1)\lambda^{a-2} e^{-b\lambda} - b\lambda^{a-1} e^{-b\lambda}) \\ &= \frac{1}{\Gamma(a)} b^a \lambda^{a-2} e^{-b\lambda} ((a-1) - b\lambda) \end{aligned}$$

Give it equal zero so we obtain:

$$\text{mode}[\lambda] = \frac{a-1}{b}$$

When $b \neq 0$

Linear Models for Regression

Problem 3.1

We have formula of $\tanh(a)$:

$$\begin{aligned}\tanh(a) &= \frac{e^a - e^{-a}}{e^a + e^{-a}} = \frac{e^a(1 - e^{-2a})}{e^a(1 + e^{-2a})} \\ &= \frac{(1 - e^{-2a})}{(1 + e^{-2a})} = \frac{-1 - e^{-2a} + 2}{1 + e^{-2a}} \\ &= -1 + 2\sigma(-2a)\end{aligned}$$

Linear combination of logistic sigmoid function:

$$\begin{aligned}y(x, w) &= w_0 + \sum_{j=1}^M w_j \sigma\left(\frac{x - \mu_j}{s}\right) \\ &= w_0 + \sum_{j=1}^M w_j \left(\frac{\tanh\left(\frac{x - \mu_j}{s}\right) + 1}{2}\right) \\ &= w_0 + \frac{1}{2} \sum_{j=1}^M w_j + \sum_{j=1}^M \frac{w_j}{2} \tanh\left(\frac{x - \mu_j}{s}\right)\end{aligned}$$

So the linear combination of sigmoid function will be equivalent to a linear combination of \tanh function when:

$$\mu_0 = w_0 + \frac{1}{2} \sum_{j=1}^M w_j \qquad \mu_j = \frac{w_j}{2}$$