



Student Name: Saeedreza Zouashkiani

Student ID: 400206262

Deep Learning Homework 1

1)

1-a)

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \sqrt{\text{rank}(\mathbf{A})} \|\mathbf{A}\|_2$$

$$\|\mathbf{A}\|_2 = \sigma_{\max}, \|\mathbf{A}\|_F = \sqrt{\text{trace}(\mathbf{A}^H \mathbf{A})} = \sqrt{\sum_{i=1}^{\text{rank}(\mathbf{A})} \sigma_i^2} = \sqrt{\sigma_{\max}^2 + \dots + \sigma_{\min}^2}$$

$$\begin{aligned} \|\mathbf{A}\|_2 = \sigma_{\max} &= \sqrt{\sigma_{\max}^2} \leq \sqrt{\sigma_{\max}^2 + \dots + \sigma_{\min}^2} = \sqrt{\sum_{i=1}^{\text{rank}(\mathbf{A})} \sigma_i^2} = \|\mathbf{A}\|_F \leq \sqrt{\text{rank}(\mathbf{A}) \cdot \sigma_{\max}^2} \\ &= \sqrt{\text{rank}(\mathbf{A})} \|\mathbf{A}\|_2 \blacksquare \end{aligned}$$

1-b)

1-b-i)

$$P(X \geq a) \leq \frac{E(X)}{a}$$

$$E(X) = P(X < a) \cdot E(X|X < a) + P(X \geq a) \cdot E(X|X \geq a)$$

Since X is non-negative then $E(X|X < a)$ is positive, and $E(X|X \geq a)$ is larger than a . Thus:

$$E(X) \geq P(X \geq a) \cdot E(X|X \geq a) \geq P(X \geq a) \cdot a$$

Therefore

$$P(X \geq a) \leq \frac{E(X)}{a} \blacksquare$$

1-b-ii)

$$P(|Z - \mu| \geq \varepsilon) = P((Z - \mu)^2 \geq \varepsilon^2) \leq \frac{E((Z - \mu)^2)}{\varepsilon^2} = \frac{\sigma^2}{\varepsilon^2} \blacksquare$$

1-b-iii)

Let Z_i denote the random variable (indicator) that determines whether the random point falls within the circle. Z_i can take a value of 1 with probability of $\frac{\pi}{4}$ and 0 with a probability of $1 - \frac{\pi}{4}$. Therefore

$$E(Z_i) = 1 \cdot \frac{\pi}{4} + 0 \cdot \left(1 - \frac{\pi}{4}\right) = \frac{\pi}{4}$$

$$\text{Var}(Z_i) = \frac{\pi}{4} \cdot \left(1 - \frac{\pi}{4}\right)$$

The estimator then is $\hat{\pi} = \frac{4}{n} \sum Z_i$. We check whether $\hat{\pi}$ is an unbiased estimator of π .

$$E(\hat{\pi}) = E\left(\frac{4}{n} \sum Z_i\right) = \frac{4}{n} \cdot \frac{n\pi}{4} = \pi$$

$$Var(\hat{\pi}) = \frac{16}{n^2} Var\left(\sum Z_i\right) = \frac{\pi(4-\pi)}{n}$$

1-b-iv)

Using Chebyshev's inequality, we have:

$$P(|\hat{\pi} - \pi| \geq 0.01) \leq \frac{Var(\hat{\pi})}{(0.01)^2} = \frac{\pi(4-\pi)}{n(0.01)^2} \leq 1 - 0.95 = 0.05$$

Then if we solve for n , we get $n \geq 539353.24$. Therefore $n = 539354$

2)

2-i)

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial (a_1 x_1 + \dots + a_n x_n)}{\partial \mathbf{x}} = [a_1 \dots a_n] = \mathbf{a}^T$$

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \frac{\partial (\mathbf{x}^T \mathbf{A}) \mathbf{x}}{\partial \mathbf{x}} + \frac{\partial \mathbf{x}^T (\mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = \mathbf{x}^T \mathbf{A} + (\mathbf{A} \mathbf{x})^T = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T)$$

2-ii)

$$\frac{\partial \mathbf{A} \mathbf{A}^{-1}}{\partial \beta} = 0 = \frac{\partial \mathbf{A}}{\partial \beta} \mathbf{A}^{-1} + \mathbf{A} \frac{\partial \mathbf{A}^{-1}}{\partial \beta}$$

By rearranging the terms and multiplying by \mathbf{A}^{-1} from left we get

$$\frac{\partial \mathbf{A}^{-1}}{\partial \beta} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \beta} \mathbf{A}^{-1}$$

2-iii)

Let M_{ij} , C_{ij} , $adj(\mathbf{A})$ be the (i, j) minor of \mathbf{A} , (i, j) element of cofactor matrix of \mathbf{A} , and adjugate matrix of \mathbf{A} respectively.

$$adj(\mathbf{A}) = \mathbf{C}^T$$

$$\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} adj(\mathbf{A})$$

$$(\mathbf{A}^{-1})_{ij}^T = \frac{1}{\det(\mathbf{A})} C_{ij}$$

By cofactor expansion of \mathbf{A}

$$\det(\mathbf{A}) = \sum_{k=1}^n A_{ik} C_{ik}$$

$$\frac{\partial \det(\mathbf{A})}{\partial A_{ij}} = \sum_{k=1}^n \left(\frac{\partial A_{ik}}{\partial A_{ij}} C_{ik} + A_{ik} \frac{\partial C_{ik}}{\partial A_{ij}} \right) = C_{ij}$$

$$\frac{\partial \det(\mathbf{A})}{\partial \mathbf{A}} = \mathbf{C} = \text{adj}(\mathbf{A}^T) = \det(\mathbf{A}) \mathbf{A}^{-T} = |\mathbf{A}| \mathbf{A}^{-T} = \nabla_{\mathbf{A}} |\mathbf{A}|$$

$$\nabla_{\mathbf{A}} \log |\mathbf{A}| = \frac{\partial \log(\det(\mathbf{A}))}{\partial \mathbf{A}} = \frac{1}{\det(\mathbf{A})} \frac{\partial \det(\mathbf{A})}{\partial \mathbf{A}} = |\mathbf{A}|^{-1} |\mathbf{A}| \mathbf{A}^{-T} = \mathbf{A}^{-T}$$

3)

Let \mathbf{A} be an $n \times n$ matrix.

3-i) $\text{tr}(\mathbf{A}) = \lambda_1 + \dots + \lambda_n$

The characteristic polynomial of \mathbf{A} is defined as

$$p_{\mathbf{A}}(t) = \det(t\mathbf{I} - \mathbf{A}) = (-1)^n (t^n - (\text{tr}(\mathbf{A})t^{n-1} + \dots + (-1)^n \det(\mathbf{A}))$$

Also, the characteristic polynomial can be factorized as

$$p_{\mathbf{A}}(t) = (-1)^n (t - \lambda_1) \dots (t - \lambda_n)$$

So, by comparing terms we get

$$\text{tr}(\mathbf{A}) = \lambda_1 + \dots + \lambda_n$$

3-ii) $\det(\mathbf{A}) = \lambda_1 \dots \lambda_n$

By comparing terms from last part, it is easily derived.

4)

4-i)

If \mathbf{A} has full column rank, the inverse of $\mathbf{A}^T \mathbf{A}$ exists

To check whether \mathbf{A}^\dagger is a pseudoinverse, we should have:

$$\mathbf{A} = \mathbf{A} \mathbf{A}^\dagger \mathbf{A}$$

$$\mathbf{A} \mathbf{A}^\dagger \mathbf{A} = \mathbf{A} \cdot (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{A} = \mathbf{A}$$

Or if we have $\mathbf{A}^\dagger = \mathbf{V} \mathbf{\Sigma}^\dagger \mathbf{U}^T$

$$\begin{aligned} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T &= (\mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T = (\mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T = \mathbf{V} (\mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{V}^T \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \\ &= \mathbf{V} \mathbf{\Sigma}^\dagger \mathbf{\Sigma}^T \mathbf{\Sigma}^T \mathbf{U}^T = \mathbf{V} \mathbf{\Sigma}^\dagger \mathbf{U}^T = \mathbf{A}^\dagger \end{aligned}$$

4-ii)

If \mathbf{A} has full row rank, the inverse of $\mathbf{A}\mathbf{A}^T$ exists

$$\mathbf{A}.\mathbf{A}^\dagger.\mathbf{A} = \mathbf{A}.\mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}.\mathbf{A} = \mathbf{A}$$

Or

$$\begin{aligned}\mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1} &= \mathbf{V}\Sigma^T\mathbf{U}^T(\mathbf{U}\Sigma\mathbf{V}^T\mathbf{V}\Sigma^T\mathbf{U}^T)^{-1} = \mathbf{V}\Sigma^T\mathbf{U}^T\mathbf{U}(\Sigma\Sigma^T)^{-1}\mathbf{U}^T = \mathbf{V}\Sigma^T(\Sigma\Sigma^T)^{-1}\mathbf{U}^T = \mathbf{V}\Sigma^T\Sigma^{T\dagger}\Sigma^\dagger\mathbf{U}^T \\ &= \mathbf{V}\Sigma^\dagger\mathbf{U}^T = \mathbf{A}^\dagger\end{aligned}$$

5)

5-i)

We start by eliminating the block matrix under A

$$\begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ -\mathbf{C}\mathbf{A}^{-1} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{0} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix}$$

Then by eliminating the element above D

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I}_n & -\mathbf{A}^{-1}\mathbf{B} \\ \mathbf{0} & \mathbf{I}_k \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{C} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix}$$

Therefore, by combining the two

$$\begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ -\mathbf{C}\mathbf{A}^{-1} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I}_n & -\mathbf{A}^{-1}\mathbf{B} \\ \mathbf{0} & \mathbf{I}_k \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix}$$

Thus

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{C}\mathbf{A}^{-1} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{A}^{-1}\mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$$

Which is the LDU decomposition of M.

Therefore

$$\begin{aligned}\det(\mathbf{M}) &= \det\left(\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{C}\mathbf{A}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{A}^{-1}\mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}\right) \\ &= \det\left(\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{C}\mathbf{A}^{-1} & \mathbf{I} \end{bmatrix}\right) \det\left(\begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix}\right) \det\left(\begin{bmatrix} \mathbf{I} & \mathbf{A}^{-1}\mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}\right) \\ &= 1. \det(\mathbf{A}) \det(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B}). 1 = \det(\mathbf{A}) \det(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})\end{aligned}$$

5-ii)

Like the last part, using block LDU decomposition when D is invertible, we get

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{D}^{-1}\mathbf{C} & \mathbf{I} \end{bmatrix}$$

$$\det(\mathbf{M}) = \det\left(\begin{bmatrix} \mathbf{I} & \mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}\right) \det\left(\begin{bmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}\right) \det\left(\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{D}^{-1}\mathbf{C} & \mathbf{I} \end{bmatrix}\right) = \det(\mathbf{D}) \det(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})$$

5-iii)

Assume that \mathbf{A} and \mathbf{D} are both invertible, therefore by inverting the LDU decomposition of part i and ii we get. First invert using part i.

$$\begin{aligned} M^{-1} &= \begin{bmatrix} I & A^{-1}B \\ 0 & I \end{bmatrix}^{-1} \begin{bmatrix} A & 0 \\ 0 & D - CA^{-1}B \end{bmatrix}^{-1} \begin{bmatrix} I_n & 0 \\ CA^{-1} & I_k \end{bmatrix}^{-1} \\ &= \begin{bmatrix} I_n & -A^{-1}B \\ 0 & I_k \end{bmatrix} \begin{bmatrix} A^{-1} & 0 \\ 0 & (D - CA^{-1}B)^{-1} \end{bmatrix} \begin{bmatrix} I_n & 0 \\ -CA^{-1} & I_k \end{bmatrix} \\ &= \begin{bmatrix} A^{-1} + A^{-1}B(D - CA^{-1}B)^{-1}CA^{-1} & -A^{-1}B(D - CA^{-1}B)^{-1} \\ -(D - CA^{-1}B)^{-1}CA^{-1} & (D - CA^{-1}B)^{-1} \end{bmatrix} \end{aligned}$$

Now inverting part ii, results in:

$$\begin{aligned} M^{-1} &= \begin{bmatrix} I & 0 \\ D^{-1}C & I \end{bmatrix}^{-1} \begin{bmatrix} A - BD^{-1}C & 0 \\ 0 & D \end{bmatrix}^{-1} \begin{bmatrix} I & BD^{-1} \\ 0 & I \end{bmatrix}^{-1} \\ &= \begin{bmatrix} I & 0 \\ -D^{-1}C & I \end{bmatrix} \begin{bmatrix} (A - BD^{-1}C)^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix} \begin{bmatrix} I & -BD^{-1} \\ 0 & I \end{bmatrix} \\ &= \begin{bmatrix} (A - BD^{-1}C)^{-1} & -(A - BD^{-1}C)^{-1}BD^{-1} \\ -D^{-1}C(A - BD^{-1}C)^{-1} & D^{-1} + D^{-1}C(A - BD^{-1}C)^{-1}BD^{-1} \end{bmatrix} \end{aligned}$$

By comparing the first block matrix we get

$$(A - BD^{-1}C)^{-1} = A^{-1} + A^{-1}B(D - CA^{-1}B)^{-1}CA^{-1}$$

If we substitute \mathbf{D} with $-\mathbf{D}$

$$(A + BD^{-1}C)^{-1} = A^{-1} - A^{-1}B(D + CA^{-1}B)^{-1}CA^{-1}$$

5-iv)

Using part I, with $D = 1, \mathbf{u} = B, \mathbf{v}^T = C$:

$$\det \left(\begin{bmatrix} I_n & 0 \\ \mathbf{v}^T A^{-1} & 1 \end{bmatrix} \begin{bmatrix} A & 0 \\ 0 & 1 - \mathbf{v}^T A^{-1} \mathbf{u} \end{bmatrix} \begin{bmatrix} I & A^{-1} \mathbf{u} \\ 0 & 1 \end{bmatrix} \right) = \det(A)(1 - \mathbf{v}^T A^{-1} \mathbf{u})$$

6)

6-i)

Using question 5

$$\begin{aligned} \det \begin{pmatrix} tI - B & -x \\ \mathbf{y}^* & t - a \end{pmatrix} &= (t - a) \det \left((tI - B) - \frac{x\mathbf{y}^*}{t - a} \right) = (t - a) \cdot \det(tI - B) \left(1 - \frac{\mathbf{y}^*(tI - B)^{-1}x}{t - a} \right) \\ &= \det(tI - B) (t - a - \mathbf{y}^*(tI - B)^{-1}x) = (t - a) \cdot p_B(t) - \mathbf{y}^* \text{adj}(tI - B)x \end{aligned}$$

6-ii)

Using Courant-Fischer's theorem

$$\lambda_i(\mathbf{M}) = \min_{\dim V=i} \max_{\substack{x \in V, \\ \|x\|=1}} \langle \mathbf{M}x, x \rangle$$

$$\lambda_i(\mathbf{A} + \mathbf{B}) = \min_{\dim V=i} \max_{\substack{x \in V, \\ \|x\|=1}} \langle \mathbf{A}x, x \rangle + \langle \mathbf{B}x, x \rangle$$

To give an upper bound on a minimum value of a function, we just need to give an upper bound on some value it takes.

Let V_A, V_B be subspaces of \mathbb{R}^n with dimensions of $i + j$ and $n - j$ respectively which achieve the minimum values of $\max_{\substack{x \in V_A, \\ \|x\|=1}} \langle \mathbf{A}x, x \rangle$, $\max_{\substack{x \in V_B, \\ \|x\|=1}} \langle \mathbf{B}x, x \rangle$ and let $W = V_A \cap V_B$ be their intersection. W has dimension of at least i .

$$\max_{\substack{x \in W, \\ \|x\|=1}} \langle (\mathbf{A} + \mathbf{B})x, x \rangle \leq \max_{\substack{x \in W, \\ \|x\|=1}} \langle \mathbf{A}x, x \rangle + \max_{\substack{x \in W, \\ \|x\|=1}} \langle \mathbf{B}x, x \rangle \leq \lambda_{i+j}(\mathbf{A}) + \lambda_{n-j}(\mathbf{B})$$

Since W has dimension of at least i , the above is an upper bound on the value of $\max_{\substack{x \in V, \\ \|x\|=1}} \langle (\mathbf{A} + \mathbf{B})x, x \rangle$

for any i 'th dimensional subspace $V \subseteq W$

6-iii)

Let $\mathbf{M} = \begin{pmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & 0 \end{pmatrix}$ and $\mathbf{N} = \begin{pmatrix} \mathbf{0} & \mathbf{y} \\ \mathbf{y}^* & a \end{pmatrix}$. Another inequality from Weyl states that:

$$\lambda_i(\mathbf{M} + \mathbf{N}) = \lambda_i(\mathbf{A}) \geq \lambda_j(\mathbf{N}) + \lambda_{i-j+1}(\mathbf{M}) \quad (1)$$

Using part i:

$$p_{\mathbf{M}}(t) = t p_{\mathbf{B}}(t) \quad (2)$$

$$p_{\mathbf{N}}(t) = (t - a) \det(t\mathbf{I}) - \mathbf{y}^* \det(t\mathbf{I})(t\mathbf{I})^{-1} \mathbf{y} = t^{n-1}(t(t - a) - \mathbf{y}^* \mathbf{y}) \quad (3)$$

Therefore $\lambda_{i+1}(\mathbf{M}) = \lambda_i(\mathbf{B})$ and $\lambda_1(\mathbf{M}) = 0$. Also $\lambda_i(\mathbf{N}) = 0, \forall i \neq n + 1, i \neq 1$ and it has a negative and positive eigenvalue which correspond to $\lambda_1(\mathbf{N})$ and $\lambda_{n+1}(\mathbf{N})$ respectively.

Using Weyl inequalities and (2), (3):

$$\begin{aligned} \lambda_i(\mathbf{A}) &\leq \lambda_{i+j}(\mathbf{M}) + \lambda_{n+1-j}(\mathbf{N}) \\ \lambda_i(\mathbf{A}) &\leq \lambda_{i+1}(\mathbf{M}) + \lambda_{n+1-1}(\mathbf{N}) = \lambda_{i+1}(\mathbf{M}) = \lambda_i(\mathbf{B}) \end{aligned} \quad (4)$$

And:

$$\lambda_i(\mathbf{A}) \geq \lambda_1(\mathbf{N}) + \lambda_i(\mathbf{M}) = \lambda_i(\mathbf{M}) = \lambda_{i-1}(\mathbf{B}) \quad (5)$$

Using (4), (5):

$$\lambda_1(\mathbf{A}) \leq \lambda_1(\mathbf{B}) \leq \lambda_2(\mathbf{A}) \leq \lambda_2(\mathbf{B}) \leq \dots \leq \lambda_n(\mathbf{A}) \leq \lambda_n(\mathbf{B}) \leq \lambda_{n+1}(\mathbf{A})$$

7)

$$L(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f_{x_i}(x_i; \theta) = \prod_{i=1}^n \frac{1}{\theta^2} x_i e^{-\frac{x_i}{\theta}} = \frac{1}{\theta^{2n}} e^{-\frac{\sum_{i=1}^n x_i}{\theta}} \prod_{i=1}^n x_i$$

$$\ln(L(x_1, \dots, x_n; \theta)) = -\frac{\sum_{i=1}^n x_i}{\theta} - 2n \ln(\theta) + \sum_{i=1}^n \ln(x_i)$$

Differentiating with respect to theta and setting to zero, yields:

$$\frac{\partial \ln(L(x_1, \dots, x_n; \theta))}{\partial \theta} = \frac{\sum_{i=1}^n x_i}{\theta^2} - \frac{2n}{\theta} = 0$$

$$\widehat{\theta}_{ML} = \frac{\sum_{i=1}^n x_i}{2n}$$

8)

8-i) ML

$$L(x_1, \dots, x_n; \mu) = \prod_{i=1}^n f(x_i; \mu) = \frac{1}{(\sqrt{2\pi}\sigma^2)^n} e^{-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}}$$

$$\ln(L(x_1, \dots, x_n; \mu)) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}$$

Taking derivative of log likelihood function with respect to μ and setting it to zero

$$\frac{\partial \ln(L(x_1, \dots, x_n; \mu))}{\partial \mu} = 0 \rightarrow \hat{\mu}_{ML} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\mathbb{E}[\hat{\mu}_{ML}] = \mathbb{E}\left[\frac{\sum_{i=1}^n x_i}{n}\right] = \mu$$

Therefore $\hat{\mu}_{ML}$ is an unbiased estimator of μ

$$\text{Var}(\hat{\mu}_{ML}) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n x_i\right) = \frac{\sigma^2}{n}$$

Using Chebyshev's inequality:

$$P[|\mu - \hat{\mu}_{ML}| < \epsilon] \leq \frac{\sigma^2}{n\epsilon^2} \xrightarrow{n \rightarrow \infty} 0$$

8-ii) MAP

We want to maximize $f(\mu)f(x|\mu)$

$$\frac{1}{(\sqrt{2\pi}\beta^2)} e^{-\frac{(\mu-\gamma)^2}{2\beta^2}} \frac{1}{(\sqrt{2\pi}\sigma^2)^n} e^{-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}}$$

Taking the log

$$-\frac{1}{2} \ln(2\pi\beta^2) - \frac{(\mu - \gamma)^2}{2\beta^2} - \frac{n}{2} \ln(2\pi\sigma^2) - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}$$

Taking derivative with respect to μ and setting it to zero

$$-\frac{(\mu - \gamma)}{\beta^2} + \frac{\sum_{i=1}^n (x_i - \mu)}{\sigma^2} = -\frac{(\mu - \gamma)}{\beta^2} + \frac{\sum_{i=1}^n (x_i)}{\sigma^2} - \frac{n\mu}{\sigma^2} \rightarrow \hat{\mu}_{MAP} = \frac{\beta^2 \sum_{i=1}^n (x_i) + \sigma^2 \gamma}{n\beta^2 + \sigma^2}$$

$$\mathbb{E}[\hat{\mu}_{MAP}] = \mathbb{E}\left[\frac{\beta^2 \sum_{i=1}^n (x_i) + \sigma^2 \gamma}{n\beta^2 + \sigma^2}\right] = \lim_{n \rightarrow \infty} \frac{n\beta^2 \mu + \sigma^2 \gamma}{n\beta^2 + \sigma^2} = \mu$$

Therefore $\hat{\mu}_{MAP}$ is an unbiased estimator of μ

$$Var(\hat{\mu}_{MAP}) = \left(\frac{\beta^2}{n\beta^2 + \sigma^2}\right)^2 Var\left(\sum_{i=1}^n (x_i)\right) = \left(\frac{\beta^2}{n\beta^2 + \sigma^2}\right)^2 n\sigma^2$$

Using Chebyshev's inequality:

$$P[|\mu - \hat{\mu}_{ML}| < \epsilon] \leq \frac{Var(\hat{\mu}_{MAP})}{\epsilon^2} \xrightarrow{n \rightarrow \infty} 0$$

9)

Let $\mathbf{x} \in \mathbb{R}^n, \mathbf{x}_a \in \mathbb{R}^{n_a}, \mathbf{x}_b \in \mathbb{R}^{n_b}$ so that $n_a + n_b = n$

We will at first prove part ii

ii)

Let \mathbf{S}_a be a subset matrix of dimension $n_a \times n$ such that $s_{ij} = 1$, if the j' th element in \mathbf{x}_a corresponds to i' th element in \mathbf{x} , and zero otherwise.

$$\mathbf{S}_a = \begin{pmatrix} \mathbf{I}_{n_a} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{n_b} \end{pmatrix}$$

$$\mathbf{x}_a = \mathbf{S}_a \mathbf{x}$$

Therefore, by applying linear transformation

$$\mathbf{x}_a \sim \mathcal{N}(\mathbf{S}_a \boldsymbol{\mu}, \mathbf{S}_a \boldsymbol{\Sigma} \mathbf{S}_a^T)$$

Therefore \mathbf{x}_a is a normal distribution with

$$E(\mathbf{x}_a) = \boldsymbol{\mu}_a$$

$$Cov(\mathbf{x}_a) = \boldsymbol{\Sigma}_{aa}$$

i)

The joint distribution of \mathbf{x}_a and \mathbf{x}_b is $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and from part ii, the marginal distribution of \mathbf{x}_b is $\mathcal{N}(\boldsymbol{\mu}_b, \boldsymbol{\Sigma}_{bb})$. According to Bayes' law

$$\begin{aligned} p(\mathbf{x}_a|\mathbf{x}_b) &= \frac{p(\mathbf{x}_a, \mathbf{x}_b)}{p(\mathbf{x}_b)} = \frac{\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\mathcal{N}(\mathbf{x}_b; \boldsymbol{\mu}_b, \boldsymbol{\Sigma}_{bb})} = \frac{\frac{1}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)}{\frac{1}{\sqrt{(2\pi)^{n_b} |\boldsymbol{\Sigma}_{bb}|}} \exp\left(-\frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)\right)} \\ &= \frac{1}{\sqrt{(2\pi)^{n-n_b}}} \sqrt{\frac{|\boldsymbol{\Sigma}_{bb}|}{|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) + \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)\right) \end{aligned}$$

Denote the inverse $\boldsymbol{\Sigma}^{-1}$ as $\boldsymbol{\Sigma}^{-1} = \begin{pmatrix} \boldsymbol{\Sigma}'_{aa} & \boldsymbol{\Sigma}'_{ab} \\ \boldsymbol{\Sigma}'_{ba} & \boldsymbol{\Sigma}'_{bb} \end{pmatrix}$

$$\begin{aligned} p(\mathbf{x}_a|\mathbf{x}_b) &= \frac{1}{\sqrt{(2\pi)^{n-n_b}}} \sqrt{\frac{|\boldsymbol{\Sigma}_{bb}|}{|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Sigma}'_{aa} & \boldsymbol{\Sigma}'_{ab} \\ \boldsymbol{\Sigma}'_{ba} & \boldsymbol{\Sigma}'_{bb} \end{pmatrix} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix}\right) \\ &\quad + \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \boldsymbol{\mu}_b) \stackrel{(\boldsymbol{\Sigma}'_{ba})^T = \boldsymbol{\Sigma}'_{ab})}{=} \frac{1}{\sqrt{(2\pi)^{n-n_b}}} \sqrt{\frac{|\boldsymbol{\Sigma}_{bb}|}{|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}((\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Sigma}'_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) \right. \\ &\quad + 2(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Sigma}'_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) + (\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Sigma}'_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b)) \\ &\quad \left. + \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)\right) \quad (1) \end{aligned}$$

According to question 5

$$\begin{aligned} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} &= \begin{bmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & -(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \end{bmatrix} \\ \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}^{-1} &= \begin{pmatrix} (\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1} & -(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1}\boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1} \\ -\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1} & \boldsymbol{\Sigma}_{bb}^{-1} + \boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1}\boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1} \end{pmatrix} \end{aligned}$$

Plugging the inverse into (1):

$$\begin{aligned}
p(x_a|x_b) &= \frac{1}{\sqrt{(2\pi)^{n-n_b}}} \sqrt{\frac{|\Sigma_{bb}|}{|\Sigma|}} \exp\left(-\frac{1}{2}((x_a - \mu_a)^T(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}(x_a - \mu_a) \right. \\
&\quad - 2(x_a - \mu_a)^T(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b) + (x_b - \mu_b)^T(\Sigma_{bb}^{-1} \\
&\quad \left. + \Sigma_{bb}^{-1}\Sigma_{ba}(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1})(x_b - \mu_b)) + \frac{1}{2}(x_b - \mu_b)^T\Sigma_{bb}^{-1}(x_b - \mu_b)\right) \\
&= \frac{1}{\sqrt{(2\pi)^{n-n_b}}} \sqrt{\frac{|\Sigma_{bb}|}{|\Sigma|}} \exp\left(-\frac{1}{2}((x_a - \mu_a)^T(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}(x_a - \mu_a) \right. \\
&\quad - 2(x_a - \mu_a)^T(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b) \\
&\quad \left. + (x_b - \mu_b)^T\Sigma_{bb}^{-1}\Sigma_{ba}(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b))\right) \\
&= \frac{1}{\sqrt{(2\pi)^{n-n_b}}} \sqrt{\frac{|\Sigma_{bb}|}{|\Sigma|}} \exp\left(-\frac{1}{2}[(x_a - \mu_a) - \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b)]^T(\Sigma_{aa} \right. \\
&\quad \left. - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}[(x_a - \mu_a) - \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b)])\right)
\end{aligned}$$

Also using the determinant derived in question 5

$$\det\begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} = |\Sigma_{bb}| |\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}|$$

Thus

$$\begin{aligned}
p(x_a|x_b) &= \frac{1}{\sqrt{(2\pi)^{n_a} |\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}|}} \exp\left(-\frac{1}{2}[(x_a - \mu_a) - \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b)]^T(\Sigma_{aa} \right. \\
&\quad \left. - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}[(x_a - \mu_a) - \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b)])\right)
\end{aligned}$$

Which is the probability distribution of multivariate normal distribution

$$\mu_{a|b} = \mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b)$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}$$

10)

10-i)

$$\min_x \|Ax - b\| = \min_x (Ax - b)^T(Ax - b) = \min_x (x^T A^T Ax - x^T A^T b - b^T Ax + b^T b)$$

The last term does not affect the minimization over x. By differentiating with respect to x and setting to zero:

$$\frac{\partial x^T A^T Ax - x^T A^T b - b^T Ax}{\partial x} \Big|_{x=x^*} = 0 = 2x^{*T} A^T A - 2b^T A$$

Therefore

$$x^* = (A^T A)^{-1} A^T b = A^\dagger b$$

To prove that \mathbf{x}^* has the smallest norm 2 among all solutions, suppose that there exists \mathbf{x} such that $\mathbf{Ax} = \mathbf{b}$. Therefore $\mathbf{A}(\mathbf{x} - \mathbf{x}^*) = \mathbf{0}$. This leads to the conclusion that $(\mathbf{x} - \mathbf{x}^*) \in \mathbf{K}(\mathbf{A})$, but the least square solution $\mathbf{x}^* \in \mathbf{K}(\mathbf{A})^\perp$. Therefore:

$$\|\mathbf{x}\|_2^2 = \|\mathbf{x} - \mathbf{x}^* + \mathbf{x}^*\|_2^2 = \|\mathbf{x} - \mathbf{x}^*\|_2^2 + \|\mathbf{x}^*\|_2^2 \geq \|\mathbf{x}^*\|_2^2$$

10-ii)

Assume that the algorithm has converged. We shall have $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} = \mathbf{x}$

$$\mathbf{x} = \mathbf{x} - \nu(\mathbf{A}^T \mathbf{Ax} - \mathbf{A}^T \mathbf{b}) \rightarrow \nu(\mathbf{A}^T \mathbf{Ax} - \mathbf{A}^T \mathbf{b}) = \mathbf{0} \rightarrow \mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$$

10-iii)

Let $\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$ be the optimum solution. By adding and subtracting \mathbf{x}^* from both sides (Using t'th iteration notation by subscript)

$$\mathbf{x}_{t+1} - \mathbf{x}^* = \mathbf{x}_t - \mathbf{x}^* - \nu(\mathbf{A}^T \mathbf{A})(\mathbf{x}_t - (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}) = \mathbf{x}_t - \mathbf{x}^* - \nu(\mathbf{A}^T \mathbf{A})(\mathbf{x}_t - \mathbf{x}^*)$$

Define $\mathbf{y}_t = \mathbf{x}_t - \mathbf{x}^*$

$$\mathbf{y}_{t+1} = \mathbf{y}_t - \nu(\mathbf{A}^T \mathbf{A})\mathbf{y}_t = (\mathbf{I} - \nu \mathbf{A}^T \mathbf{A})\mathbf{y}_t$$

We know that $\mathbf{A}^T \mathbf{A}$ can be diagonalized by $\mathbf{A}^T \mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T$. Multiply both sides by \mathbf{Q}^T and define $\mathbf{z}_t = \mathbf{Q}^T \mathbf{y}_t$.

$$\mathbf{Q}^T \mathbf{y}_{t+1} = \mathbf{z}_{t+1} = \mathbf{Q}^T (\mathbf{I} - \nu \mathbf{A}^T \mathbf{A}) \mathbf{y}_t = \mathbf{Q}^T (\mathbf{Q} \mathbf{Q}^T - \nu \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T) \mathbf{y}_t = \mathbf{Q}^T \mathbf{Q} (\mathbf{I} - \nu \mathbf{\Lambda}) \mathbf{Q}^T \mathbf{y}_t = (\mathbf{I} - \nu \mathbf{\Lambda}) \mathbf{z}_t$$

$(\mathbf{I} - \nu \mathbf{\Lambda})$ is a diagonal matrix. Therefore, for the i'th element of \mathbf{z}_k after k iterations

$$z_k^i = (1 - \nu \lambda_i)^k z_0^i$$

Therefore, for the i'th mode to converge ($k \rightarrow \infty$) we shall have

$$|1 - \nu \lambda_i| < 1 \rightarrow 0 < \nu < \frac{2}{\lambda_i}$$

And to satisfy convergence of all modes we shall have

$$0 < \nu < \frac{2}{\lambda_{\max}}$$