



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.

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

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

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

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.

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

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

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

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

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

4)

4-i)

If A has full column rank, the inverse of $A^T A$ exists

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

$$A = AA^\dagger A$$

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

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

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

4-ii)

If A has full row rank, the inverse of AA^T exists

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

Or

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

5)

5-i)

We start by eliminating the block matrix under A

$$\begin{bmatrix} I_n & 0 \\ -CA^{-1} & I_k \end{bmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} A & B \\ 0 & D - CA^{-1}B \end{bmatrix}$$

Then by eliminating the element above D

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} I_n & -A^{-1}B \\ 0 & I_k \end{bmatrix} = \begin{bmatrix} A & 0 \\ C & D - CA^{-1}B \end{bmatrix}$$

Therefore, by combining the two

$$\begin{bmatrix} I_n & 0 \\ -CA^{-1} & I_k \end{bmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} I_n & -A^{-1}B \\ 0 & I_k \end{bmatrix} = \begin{bmatrix} A & 0 \\ 0 & D - CA^{-1}B \end{bmatrix}$$

Thus

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} I_n & 0 \\ CA^{-1} & I_k \end{bmatrix} \begin{bmatrix} A & 0 \\ 0 & D - CA^{-1}B \end{bmatrix} \begin{bmatrix} I & A^{-1}B \\ 0 & I \end{bmatrix}$$

Which is the LDU decomposition of M.

Therefore

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

5-ii)

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

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

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

5-iii)

Assume that **A** and **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 **D** with **-D**

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

5-iv)

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

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

6)

6-i)

Using question 5

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

6-ii)

Using Courant-Fischer's theorem

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

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

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 $\max_{\substack{x \in V_A, \\ \|x\|=1}} \langle Ax, x \rangle$, $\max_{\substack{x \in V_B, \\ \|x\|=1}} \langle Bx, x \rangle$ and let $W = V_A \cap V_B$ be their intersection. W has dimension

of at least i .

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

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

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

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

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

8-ii) MAP

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

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

Taking the log

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

Taking derivative with respect to μ and setting it to zero

10)

10-i)

$$\min_x \|Ax - b\|^2 = \min_x (Ax - b)^T (Ax - b) = \min_x (x^T A^T A x - x^T A^T b - b^T A x + b^T b)$$

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

$$\frac{\partial x^T A^T A x - x^T A^T b - b^T A x}{\partial x} = 0 = 2x^T A^T A - 2b^T A$$

Therefore

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

10-ii)

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

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

10-iii)

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

$$x_{t+1} - x^* = x_t - x^* - \nu(A^T A)(x_t - (A^T A)^{-1} A^T b) = x_t - x^* - \nu(A^T A)(x_t - x^*)$$

Define $y_t = x_t - x^*$

$$y_{t+1} = y_t - \nu(A^T A)y_t = (I - \nu A^T A)y_t$$

We know that $A^T A$ can be diagonalized by $A^T A = Q\Lambda Q^T$. Multiply both sides by Q^T and define $z_t = Q^T y_t$.

$$Q^T y_{t+1} = z_{t+1} = Q^T (I - \nu A^T A)y_t = Q^T (Q Q^T - \nu Q \Lambda Q^T)y_t = Q^T Q (I - \nu \Lambda) Q^T y_t = (I - \nu \Lambda) z_t$$

$(I - \nu \Lambda)$ is a diagonal matrix. Therefore, for the i 'th element of 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

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