

Solutions to Assignment 1 of CPSC 368/512 (Spring'23)

February 8, 2023

1 Problem 1

Problems 1.1. Which of the following is True?

1. $n^{1/\sqrt{\log n}} \leq (\log n)^{O(1)}$
2. $n^{\sqrt{n}} \leq 2^{O(n)}$
3. $2^n \leq n^{O(\log n)}$

Part 1 ($n^{1/\sqrt{\log n}} \leq (\log n)^{O(1)}$): False

The goal is to check whether there exist constants $k > 0$ and $n_0 \in \mathbb{N}$, such that, for all $n \geq n_0$, the following inequality holds

$$n^{1/\sqrt{\log n}} \leq (\log n)^k. \quad (1)$$

After taking log on both sides we get:

$$\sqrt{\log n} \leq k \log \log n$$

Since

$$\lim_{n \rightarrow \infty} \frac{\sqrt{\log n}}{\log \log n} = \infty,$$

Equation (1) can hold only for finitely many n 's (for any fixed k). Hence, the statement is false.

Part 2 ($n^{\sqrt{n}} \leq 2^{O(n)}$): True

We want to check if there exist constants $k > 0$ and $n_0 \in \mathbb{N}$, such that, for all $n \geq n_0$, the following inequality holds

$$n^{\sqrt{n}} \leq 2^{kn}.$$

Taking log on both sides, we get:

$$\sqrt{n} \log n \leq kn \cdot \log 2$$

On rearranging, the above gives the following inequality

$$\log n \leq k\sqrt{n} \cdot \log 2$$

Since

$$\lim_{n \rightarrow \infty} \frac{\log n}{\sqrt{n}} = 0,$$

to prove the statement it is enough to set $k = 1$ and n_0 to be large enough (e.g., $n_0 = 17$ suffices).

Part 3 ($2^n \leq n^{O(\log n)}$): False

Similar to previous exercises, we want to check if there exist constants $k > 0$ and $n_0 \in \mathbb{N}$, such that, for all $n \geq n_0$, the following inequality holds

$$2^n \leq n^{k \log n}. \quad (2)$$

Taking log on both sides, gives that

$$n \log 2 \leq k(\log n)^2.$$

Since

$$\lim_{n \rightarrow \infty} \frac{n}{(\log n)^2} = \infty,$$

Equation (2) can only hold for finitely many values of n (for any fixed k). Hence, the statement is false.

2 Problem 2

Problems 2.1. For each of the following functions, compute the gradient and the Hessian, and write the second-order Taylor approximation.

1. $f(x) = \sum_{i=1}^m (a_i^\top x - b_i)^2$, for $x \in \mathbb{R}^n$ where $a_1, \dots, a_m \in \mathbb{Q}^n$, and $b_1, \dots, b_m \in \mathbb{Q}$.
2. $f(x) = \log \left(\sum_{j=1}^m e^{\langle x, v_j \rangle} \right)$ where $v_1, \dots, v_m \in \mathbb{Q}^n$.
3. $f(X) = \text{Tr}(AX)$ where A is a symmetric $n \times n$ matrix and X runs over symmetric matrices.
4. $f(X) = -\log \det X$, where X runs over positive definite matrices.

2.1 Part 1

This is an exercise in differentiation. Fix an $i \in [m]$. Using the chain-rule for the i -th term, we get that

$$\begin{aligned} \nabla (a_i^\top x - b_i)^2 &= 2 \cdot (a_i^\top x - b_i) \cdot a_i, \\ \nabla^2 (a_i^\top x - b_i)^2 &= 2 \cdot a_i a_i^\top. \end{aligned}$$

Using the linearity of differentiation, we get that

$$\begin{aligned} \nabla f(x) &= \sum_{i=1}^m 2 \cdot (a_i^\top x - b_i) \cdot a_i, \\ \nabla^2 f(x) &= \sum_{i=1}^m 2 \cdot a_i a_i^\top. \end{aligned}$$

We can obtain the second-order Taylor approximation at $y \in \mathbb{R}^n$ by substituting the gradient and Hessian in the following equation

$$f(y) + (x - y)^\top \nabla f(y) + \frac{1}{2} (x - y)^\top \nabla^2 f(y) (x - y).$$

Since $f(x)$ is quadratic, its second order Taylor expansion is $f(x)$ itself.

2.2 Part 2

Using the chain rule and standard results, for any $i \in [n]$, we have that

$$\begin{aligned}\frac{\partial f}{\partial x_i}(x) &= \frac{\partial \log \left(\sum_{j=1}^m e^{\langle x, v_j \rangle} \right)}{\partial x_i} \\ &= \frac{\partial \log(y)}{\partial y} \Big|_{y=\sum_{j=1}^m e^{\langle x, v_j \rangle}} \cdot \frac{\partial \sum_{j=1}^m e^{\langle x, v_j \rangle}}{\partial x_i} \\ &= \frac{\sum_{j=1}^m e^{\langle x, v_j \rangle} v_{ji}}{\sum_{j=1}^m e^{\langle x, v_j \rangle}},\end{aligned}$$

where v_{ji} is the i -th coordinate of v_j . Thus, the gradient is

$$\nabla f(x) = \left[\frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_n}(x) \right]^\top = \frac{\sum_{j=1}^m e^{\langle x, v_j \rangle} v_j}{\sum_{j=1}^m e^{\langle x, v_j \rangle}}.$$

For any $i \in [n]$ and $\ell \in [n]$, we can compute $\frac{\partial^2 f}{\partial x_i \partial x_\ell}(x)$ as follows

$$\begin{aligned}\frac{\partial^2 f}{\partial x_i \partial x_\ell}(x) &= \frac{\partial}{\partial x_\ell} \frac{\sum_{j=1}^m e^{\langle x, v_j \rangle} v_{ji}}{\sum_{j=1}^m e^{\langle x, v_j \rangle}} \\ &= \frac{1}{\sum_{j=1}^m e^{\langle x, v_j \rangle}} \cdot \frac{\partial}{\partial x_\ell} \sum_{j=1}^m e^{\langle x, v_j \rangle} v_{ji} - \frac{\sum_{j=1}^m e^{\langle x, v_j \rangle} v_{ji}}{\left(\sum_{j=1}^m e^{\langle x, v_j \rangle} \right)^2} \cdot \frac{\partial}{\partial x_\ell} \sum_{j=1}^m e^{\langle x, v_j \rangle} \quad (\text{Using the product rule}) \\ &= \frac{1}{\sum_{j=1}^m e^{\langle x, v_j \rangle}} \cdot \sum_{j=1}^m e^{\langle x, v_j \rangle} v_{ji} v_{j\ell} - \frac{\sum_{j=1}^m e^{\langle x, v_j \rangle} v_{ji}}{\left(\sum_{j=1}^m e^{\langle x, v_j \rangle} \right)^2} \cdot \frac{\partial}{\partial x_\ell} \sum_{j=1}^m e^{\langle x, v_j \rangle} \\ &= \frac{1}{\sum_{j=1}^m e^{\langle x, v_j \rangle}} \cdot \sum_{j=1}^m e^{\langle x, v_j \rangle} v_{ji} v_{j\ell} - \frac{\left(\sum_{j=1}^m e^{\langle x, v_j \rangle} v_{ji} \right) \cdot \left(\sum_{j=1}^m e^{\langle x, v_j \rangle} v_{j\ell} \right)}{\left(\sum_{j=1}^m e^{\langle x, v_j \rangle} \right)^2}.\end{aligned}$$

Thus, the Hessian is

$$\nabla^2 f(x) = \frac{\sum_{j=1}^m e^{\langle x, v_j \rangle} v_j v_j^\top}{\sum_{j=1}^m e^{\langle x, v_j \rangle}} - \frac{\left(\sum_{j=1}^m e^{\langle x, v_j \rangle} v_j \right) \cdot \left(\sum_{j=1}^m e^{\langle x, v_j \rangle} v_j^\top \right)}{\left(\sum_{j=1}^m e^{\langle x, v_j \rangle} \right)^2}$$

We can obtain the second-order Taylor approximation at $y \in \mathbb{R}^n$ by substituting the gradient and Hessian in the following equation

$$f(y) + (x - y)^\top \nabla f(y) + \frac{1}{2} (x - y)^\top \nabla^2 f(y) (x - y).$$

2.3 Part 3

Unlike the first two parts, here we have to be careful because X lies in the linear-space of symmetric matrices. First, we have to identify the “directions” along which X can move in while staying in the linear-space of symmetric matrices, i.e., the tangent space at X . Then, we compute the first-order and second-order directional derivatives of $f(X)$ along these directions. These directional derivatives suffice to give the second-order Taylor approximation of $f(X)$. Finally, we can identify the gradient and Hessian from the directional derivatives.

Since the set of symmetric matrices is embedded in $\mathbb{R}^{n \times n}$ (in fact, it is a linear subspace), any “direction” in the tangent space must be a matrix in $\mathbb{R}^{n \times n}$. In Lemma 2.2 we establish that X can only move along the directions defined by symmetric matrices, and hence, the tangent space at X is exactly the set of all symmetric matrices in $\mathbb{R}^{n \times n}$.

Lemma 2.2 (Tangent space). For any symmetric matrix $X \in \mathbb{R}^{n \times n}$, nonzero constant $t \neq 0$, and matrix $S \in \mathbb{R}^{n \times n}$, $X + tS$ is symmetric iff S is symmetric.

Proof. $X + tS$ is symmetric if for all $i, j \in [n]$, $(X + tS)_{ij} = (X + tS)_{ji}$. Fix any $i, j \in [n]$. Since X is symmetric

$$\begin{aligned}(X + tS)_{ij} &= X_{ij} + tS_{ij} = X_{ji} + tS_{ij}, \\ (X + tS)_{ji} &= X_{ji} + tS_{ji}.\end{aligned}$$

Thus, $(X + tS)_{ij} = (X + tS)_{ji}$ iff $tS_{ij} = tS_{ji}$. Because t is nonzero, this holds iff $S_{ij} = S_{ji}$. Since the choice of $i, j \in [n]$ was arbitrary, it follows that $X + tS$ is symmetric iff S is symmetric. \square

Fix any symmetric matrices $S, T \in \mathbb{R}^{n \times n}$. We can compute the first-order directional derivatives of $f(X)$ along the direction S as follows

$$\begin{aligned}Df(X)[S] &:= \lim_{t \rightarrow 0} \frac{f(X + tS) - f(X)}{t} \\ &= \lim_{t \rightarrow 0} \frac{\text{Tr}(A(X + tS)) - \text{Tr}(AX)}{t} \\ &= \lim_{t \rightarrow 0} \frac{\text{Tr}(tAS)}{t} && \text{(Using that } \text{Tr}(X + Y) = \text{Tr}(X) + \text{Tr}(Y)\text{)} \\ &= \text{Tr}(AS). && \text{(Using that } \text{Tr}(tX) = t \cdot \text{Tr}(X)\text{)}\end{aligned}$$

The second-order directional derivative can be computed as follows

$$\begin{aligned}D^2f(X)[S, T] &= \lim_{t \rightarrow 0} \frac{D^2f(X + tT)[S] - D^2f(X)[S]}{t} \\ &= \lim_{t \rightarrow 0} \frac{\text{Tr}(AS) - \text{Tr}(AS)}{t} \\ &= 0.\end{aligned}$$

From this, we get the second-order Taylor approximation of f at $B \in \mathbb{R}^{n \times n}$ is

$$f(B) + Df(B)[X - B] + \frac{1}{2}D^2f(B)[X - B, X - B] = \text{Tr}(AB) + \text{Tr}(A(X - B)) = \text{Tr}(AX).$$

Note that the second-order Taylor approximation of $f(X)$ at $B \in \mathbb{R}^{n \times n}$ is $f(X)$ itself – this is true because $f(X)$ is a linear function in X .

2.4 Part 4

Like the previous part, we have to be careful about the directions along which X can move while staying in the set (to be precise, cone) of positive definite $n \times n$ matrices.

Since the set of positive definite (PD) matrices is embedded in $\mathbb{R}^{n \times n}$, any “direction” in the tangent space must be a matrix in $\mathbb{R}^{n \times n}$. First, in Lemma 2.3, we prove that X can only move along directions defined by symmetric matrices, i.e., the tangent space at X is exactly the set of symmetric matrices in $\mathbb{R}^{n \times n}$.

Lemma 2.3 (Tangent space). For any PD matrix $X \in \mathbb{R}^{n \times n}$ and direction $S \in \mathbb{R}^{n \times n}$, it holds that:

1. **(All directions defined by symmetric matrices are valid).** If S is symmetric, then there exists a positive $t_0 > 0$, such that for all $0 < t < t_0$, $X + tS$ is a PD matrix.
2. **(No directions defined by unsymmetric matrices are valid).** Otherwise, if S is not symmetric, then for all $t > 0$, $X + tS$ is not symmetric, and hence, not a PD matrix.

Proof.

Part 1. If S is symmetric, then by Lemma 2.2 and the fact that $t > 0$, it follows that $X + tS$ is also symmetric. It remains to show that for all nonzero $x \in \mathbb{R}^n$, $x^\top(X + tS)x > 0$. Since S and X are a real-symmetric, they have real eigenvalues. Let the minimum eigenvalue of X and S be $\lambda_{\min}(X)$ and $\lambda_{\min}(S)$ respectively. Because X is PD, we know that $\lambda_{\min}(X) > 0$. We consider two cases:

Case A ($\lambda_{\min}(S) \geq 0$): Set $t_0 := 1$. For any $0 < t < t_0$ and nonzero $x \in \mathbb{R}^n$, we have that

$$\begin{aligned} x^\top(X + tS)x &= x^\top Xx + tx^\top Sx \\ &\geq (\lambda_{\min}(X) + t\lambda_{\min}(S)) \cdot \|x\|_2^2 \\ &\geq \lambda_{\min}(X) \cdot \|x\|_2^2 && \text{(Using that } t > 0 \text{ and } \lambda_{\min}(S) \geq 0) \\ &> 0. \end{aligned}$$

Case B ($\lambda_{\min}(S) < 0$): Set $t_0 := -\frac{\lambda_{\min}(X)}{\lambda_{\min}(S)}$. Note that $t_0 > 0$. For any $0 < t < t_0$ and nonzero $x \in \mathbb{R}^n$, we have that

$$\begin{aligned} x^\top(X + tS)x &= x^\top Xx + tx^\top Sx \\ &\geq (\lambda_{\min}(X) + t\lambda_{\min}(S)) \cdot \|x\|_2^2 \\ &> (\lambda_{\min}(X) + t_0\lambda_{\min}(S)) \cdot \|x\|_2^2 && \text{(Using that } t < t_0, \lambda_{\min}(S) < 0, \text{ and } \|x\|_2 > 0) \\ &> 0. && \text{(By the choice of } t_0, \lambda_{\min}(X) + t_0\lambda_{\min}(S) > 0) \end{aligned}$$

Part 2. If S is not symmetric, then by Lemma 2.2 and the fact that $t > 0$, it follows that $X + tS$ is not symmetric. \square

(Alternatively, one can also prove Lemma 2.3 by using the fact that PD matrices form an open subset in the linear space of symmetric matrices.)

Fix any symmetric matrices $S, T \in \mathbb{R}^{n \times n}$. We can compute the first-order directional derivatives of $f(X)$ along direction S as follows

$$\begin{aligned} Df(X)[S] &:= \lim_{t \rightarrow 0} \frac{f(X + tS) - f(X)}{t} \\ &= \lim_{t \rightarrow 0} \frac{-\log |X + tS| + \log |X|}{t} \\ &= \lim_{t \rightarrow 0} \frac{-\log (|X|^{-1} \cdot |X + tS|)}{t} \\ &= \lim_{t \rightarrow 0} \frac{-\log |I + tX^{-1}S|}{t}. \end{aligned} \tag{3}$$

Let the eigenvalues of $X^{-1}S$ be

$$\lambda_1, \lambda_2, \dots, \lambda_n \in \mathbb{R}.$$

Then the eigenvalues of $I + tX^{-1}S$ are

$$1 + t\lambda_1, 1 + t\lambda_2, \dots, 1 + t\lambda_n \in \mathbb{R}.$$

Substituting $|I + tX^{-1}S| = \prod_{i=1}^n (1 + t\lambda_i)$ in Equation (3), we get that

$$\begin{aligned}
Df(X)[S] &= -\lim_{t \rightarrow 0} \frac{1}{t} \cdot \log \left(\prod_{i=1}^n (1 + t\lambda_i) \right) \\
&= -\lim_{t \rightarrow 0} \frac{1}{t} \cdot \sum_{i=1}^n \log(1 + t\lambda_i) \\
&= -\lim_{t \rightarrow 0} \frac{1}{t} \cdot \sum_{i=1}^n (t\lambda_i \pm O(t^2)) \quad (\text{Using that, say, for all } x \in (-0.9, \infty), \log(1+x) = x \pm O(x^2)) \\
&= -\sum_{i=1}^n \lambda_i \\
&= -\text{Tr}(X^{-1}S). \quad (\text{Using that } \lambda_1, \lambda_2, \dots, \lambda_n \text{ are eigenvalues of } X^{-1}S) \quad (4)
\end{aligned}$$

Computing the directional derivative was sufficient to get full credit for the gradient. However, the expression of $Df(X)[S]$, one can identify that $-X^{-1}$ satisfies that for all symmetric matrices $S \in \mathbb{R}^{n \times n}$, $\langle X^{-1}, S \rangle = Df(X)[S]$. Thus, $\nabla f(X) = -X^{-1}$ is a feasible gradient along all directions defined by symmetric matrices $S \in \mathbb{R}^{n \times n}$.

Using Equation (4), we can compute the second-order directional derivative as follows

$$\begin{aligned}
D^2f(X)[S, T] &= \lim_{t \rightarrow 0} \frac{Df(X + tT)[S] - Df(X)[S]}{t} \\
&= \lim_{t \rightarrow 0} \frac{-\text{Tr}((X + tT)^{-1}S) + \text{Tr}(X^{-1}S)}{t} \\
&= \lim_{t \rightarrow 0} \frac{-\text{Tr}((I + tX^{-1}T)^{-1}X^{-1}S) + \text{Tr}(X^{-1}S)}{t}. \quad (\text{Using that for any } A, B \in \mathbb{R}^{n \times n}, (AB)^{-1} = B^{-1}A^{-1}) \quad (5)
\end{aligned}$$

Recall that for any matrix $M \in \mathbb{R}^{n \times n}$, whose largest absolute eigenvalue is strictly smaller than 1, it holds that

$$(I + M)^{-1} = I + \sum_{i=1}^{\infty} (-M)^i.$$

Let the maximum absolute eigenvalue of $X^{-1}T$ be γ_T and that of $X^{-1}S$ be γ_S . Then for any $0 < t < 1/\gamma_T$, the largest absolute eigenvalue of $tX^{-1}T$ is strictly smaller than 1. Thus, for $0 < t < 1/\gamma_T$ it follows that

$$(I + tX^{-1}T)^{-1} = I + \sum_{i=1}^{\infty} (-tX^{-1}T)^i.$$

Substituting this in Equation (5), we get that

$$\begin{aligned}
D^2f(X)[S, T] &= \lim_{t \rightarrow 0} \frac{-\text{Tr}((I + \sum_{i=1}^{\infty} (-tX^{-1}T)^i) X^{-1}S) + \text{Tr}(X^{-1}S)}{t} \\
&= \lim_{t \rightarrow 0} \frac{-\sum_{i=1}^{\infty} \text{Tr}((-tX^{-1}T)^i X^{-1}S)}{t} \quad (\text{Using that } \text{Tr}(X + Y) = \text{Tr}(X) + \text{Tr}(Y)) \\
&= \lim_{t \rightarrow 0} -\sum_{i=1}^{\infty} t^{i-1} \cdot \text{Tr}((-X^{-1}T)^i X^{-1}S) \quad (\text{Using that } \text{Tr}(tX) = t \cdot \text{Tr}(X)) \\
&= \lim_{t \rightarrow 0} \text{Tr}(X^{-1}T X^{-1}S) - \sum_{i=2}^{\infty} t^{i-1} \cdot \text{Tr}((-X^{-1}T)^i X^{-1}S).
\end{aligned}$$

For a matrix $M \in \mathbb{R}^{n \times n}$, let $\rho(M) \in \mathbb{R}$ denote its largest absolute eigenvalue. Then, we can bound the each term of the sum in the above equation as follows

$$\begin{aligned}
D^2 f(X)[S, T] &= \lim_{t \rightarrow 0} \text{Tr} \left(X^{-1} T X^{-1} S \right) \pm \sum_{i=2}^{\infty} t^{i-1} n \cdot \rho \left((-X^{-1} T)^i X^{-1} S \right) \\
&\text{(Using that for any } M \in \mathbb{R}^{n \times n}, \text{ with eigenvalues } \lambda_1(M), \dots, \lambda_n(M), \text{Tr}(M) = \sum_{i=1}^n \lambda_i(M) \leq n \cdot \rho(M)) \\
&= \lim_{t \rightarrow 0} \text{Tr} \left(X^{-1} T X^{-1} S \right) \pm \sum_{i=2}^{\infty} t^{i-1} n \gamma_T^i \gamma_S \\
&\quad \text{(Using that for any two } U, V \in \mathbb{R}^{n \times n}, \rho(MN) \leq \rho(M) \cdot \rho(N)) \\
&= \lim_{t \rightarrow 0} \text{Tr} \left(X^{-1} T X^{-1} S \right) \pm O(t) \\
&= \text{Tr} \left(X^{-1} T X^{-1} S \right).
\end{aligned}$$

Computing the second-order directional derivative was sufficient to get full credit for the Hessian.

Consider the tensor $X^{-1} \otimes X^{-1}$. One can observe that for all symmetric matrices $S, T \in \mathbb{R}^{n \times n}$ this tensor satisfies that $T^\top X^{-1} \otimes X^{-1} S = D^2 f(X)[S, T]$. Thus, $\nabla^2 f(X) = X^{-1} \otimes X^{-1}$ is a feasible Hessian along all directions defined by symmetric matrices $S, T \in \mathbb{R}^{n \times n}$.

From the first-order and second-order directional derivatives, we get the second-order Taylor approximation of f at $A \in \mathbb{R}^{n \times n}$ as

$$\begin{aligned}
f(A) + Df(A)[X - A] + \frac{1}{2} D^2 f(A)[X - A, X - A] \\
= -\log |A| - \text{Tr}(A^{-1}(X - A)) + \frac{1}{2} \text{Tr} \left(A^{-1}(X - A) A^{-1}(X - A) \right).
\end{aligned}$$

3 Problem 3

Problems 3.1. Consider a real $m \times n$ matrix A with $n \leq m$ and a vector $b \in \mathbb{R}^m$. Let

$$p(x) := \arg \min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2.$$

Assuming that A is of full rank, derive a formula for $p(x)$ in terms of A and b .

Let b_A be the projection of b on the range of A . Then, for some b_N in the null space of A , we can write $b = b_A + b_N$. We have

$$\|Ax - b\|^2 = \|(Ax - b_A) - b_N\|^2 = \|Ax - b_A\|^2 + \|b_N\|^2 + \langle Ax - b_A, b_N \rangle.$$

Since, $Ax - b_A$ is in the range space of A and b_N is in the null space of A , $\langle Ax - b_A, b_N \rangle = 0$. Thus, we get that

$$\|Ax - b\|^2 = \|Ax - b_A\|^2 + \|b_N\|^2.$$

Since b_N does not depend on x , it follows that $\arg \min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 = \arg \min_{x \in \mathbb{R}^n} \|Ax - b_A\|_2^2$. Let us define

$$f(x) := \|Ax - b_A\|^2 = (Ax - b_A)^\top (Ax - b_A) = x^\top A^\top Ax - 2b_A^\top Ax + b_A^\top b_A.$$

To minimize f , we will first compute its critical point. For this we calculate its gradient:

$$\nabla f(x) = 2A^\top Ax - 2A^\top b_A.$$

We know that $A^\top A$ is invertible because A has fullrank (exercise: prove this), hence $\nabla f(x) = 0$ iff $x = (A^\top A)^{-1} A^\top b_A$. Moreover, at $x_0 = (A^\top A)^{-1} A^\top b$ it holds that $\|Ax_0 - b_A\| = 0$. Since the norm is nonnegative, x_0 is a minimizer. Further, since x_0 is the unique critical point of $f(x)$, it is the unique minimizer of $f(x)$. Hence, (by our previous discussion) it is also of unique minimizer of $\|Ax - b\|^2$.

4 Problem 4

Problems 4.1. Given an $n \times n$ real PD matrix A and $u \in \mathbb{R}^n$, show that

$$\|u\|_A := \sqrt{u^\top A u}$$

is a norm. What aspect of being a norm for $\|u\|_A$ breaks down when A is just guaranteed to be PSD (and not PD)? And when A has negative eigenvalues?

Recall that a real valued function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is defined to be a norm if it satisfies the following three properties:

1. **(Triangle inequality)** for all $x, y \in \mathbb{R}^n$, $f(x + y) \leq f(x) + f(y)$,
2. **(Absolute homogeneity)** for all $x \in \mathbb{R}^n$ and $s \in \mathbb{R}$, $f(sx) = |s| f(x)$,
3. **(Positive definiteness)** for all $x \in \mathbb{R}^n$, $f(x) = 0$ iff $x = 0$.

Proposition 4.2. If A is a PD matrix then $\|\cdot\|_A$ is a norm.

Proof. The main part of the proof is proving that $\|\cdot\|_A$ satisfies the triangle inequality. To prove the triangle inequality, we use the following fact

Fact 4.3. If A is an $n \times n$ PSD matrix, then for any $x, y \in \mathbb{R}^n$, it holds that

$$\sqrt{x^\top A x} \cdot \sqrt{y^\top A y} \geq x^\top A y.$$

Proof. Let $t \in \mathbb{R}^n$ be any real number. Then, because A is PD, it holds that

$$(tx + y)^\top A (tx + y) \geq 0, \quad x^\top A x \geq 0, \quad \text{and} \quad y^\top A y \geq 0. \quad (6)$$

(Note that $tx + y$, x , and y could be zero, so we do not have a strict inequalities.) Consider the following quadratic function in t

$$(tx + y)^\top A (tx + y) = (x^\top A x) t^2 + (2x^\top A y) t + y^\top A y. \quad (\text{Using that } A \text{ is symmetric})$$

Because $(tx + y)^\top A (tx + y)$ is nonnegative, the discriminant of the above quadratic must be at most 0, i.e.,

$$4(x^\top A y)^2 - 4(x^\top A x) \cdot (y^\top A y) \leq 0.$$

Using Equation (6), this implies that

$$x^\top A y \leq |x^\top A y| \leq \sqrt{x^\top A x} \cdot \sqrt{y^\top A y}.$$

□

In the rest of the proof assume that A is an $n \times n$ PD matrix.

Triangle inequality. Fix any $x, y \in \mathbb{R}^n$. Using the above fact, we can prove the triangle inequality as follows

$$\begin{aligned} \left(\sqrt{x^\top A x} + \sqrt{y^\top A y} \right)^2 &= x^\top A x + y^\top A y + 2\sqrt{(x^\top A x)(y^\top A y)} \\ &= 2(x + y)^\top A (x + y) + 2\sqrt{(x^\top A x)(y^\top A y)} - 2x^\top A y \\ &\stackrel{\text{Fact 4.3}}{\geq} 2(x + y)^\top A (x + y). \end{aligned} \quad (\text{Using that } A \text{ is symmetric})$$

Absolute homogeneity. We can observe that for all $s \in \mathbb{R}$ and $u \in \mathbb{R}^n$, $\|su\|_A := \sqrt{(su)^\top A (su)} = |s| \sqrt{u^\top A u} = |s| \|u\|_A$.

Positive definiteness. Since A is PD, we know that for all nonzero $u \in \mathbb{R}^n$, i.e., $u \neq 0$, it holds that $u^\top A u > 0$. Further, we can observe that $[0, \dots, 0]^\top A [0, \dots, 0] = 0$. Thus, $\|\cdot\|_A$ satisfies positive definiteness. \square

Extension 1: A is only guaranteed to be PSD. If A is only guaranteed to be a PSD matrix, the positive definiteness breaks down. The triangle inequality and absolute homogeneity still hold; this follows because the above proof of triangle inequality only uses the fact that A is a PSD matrix and the above proof of absolute homogeneity holds for any $n \times n$ matrix $A \in \mathbb{R}^{n \times n}$.

Extension 2: A can have negative eigenvalues. If A can have negative eigenvalues, then $\|\cdot\|_A$ can take complex values. Thus, the inequality in the definition of the triangle inequality is not well defined. Since A can be a PSD matrix, by the previous paragraph it follows that positive definiteness can break down. One can check that absolute homogeneity holds. To see this, note that the above proof of absolute homogeneity holds for any $n \times n$ matrix $A \in \mathbb{R}^{n \times n}$.