

ML2025 Fall Homework Assignment 1

Handwritten

Yu-Cheng Lin
b11901152@ntu.edu.tw

Lu-Fang Chiang
b12901140@ntu.edu.tw

September 2025

Tools You Need to Know/Learn

- Rank, Nullity, Vector Space, Inner Product Space
- Determinant, Trace, Inverse Matrix, Adjoint Matrix
- Symmetric Matrix, Orthogonal Matrix, Diagonalization
- Positive Definite, Positive Semidefinite
- Eigenvalue Decomposition, Singular Value Decomposition (SVD)

You are expected to know (or learn through this homework) the definitions and usages of the above concepts, and are encouraged to discuss with TAs during TA hour if you encounter difficulties.

Homework Policy

- The homework is graded out of 100 points, with up to 15 additional bonus points available (only awarded for completely correct solutions).
- The official submission deadline will follow the schedule announced on NTU COOL.
- Homework may be handwritten or typed (e.g. using \LaTeX), but must be submitted in **PDF format**.
- If you discuss the homework with classmates, you should state their student IDs in your submission.
- Plagiarism is strictly prohibited. Serious violations will be dealt with according to NTU regulations.

Problem 1 (Linear Algebra Recap) (30 pts)

In the following problems, we begin by recapping key topics from undergraduate linear algebra. We then proceed to derive several fundamental theorems and lemmas, which will play a central role in the forthcoming lectures. In the following problem, we assume V is an inner product space over \mathbb{R} with $\dim V = n$.

Definition 1.1. (Symmetric transformation) A linear transformation $\mathbf{T} : V \rightarrow V$ is called symmetric if

$$\langle \mathbf{T}\mathbf{u}, \mathbf{v} \rangle = \langle \mathbf{u}, \mathbf{T}\mathbf{v} \rangle, \quad \forall \mathbf{u}, \mathbf{v} \in V.$$

If \mathbf{A} is the matrix representation of \mathbf{T} and $V \cong \mathbb{R}^n$, this condition is equivalent to

$$\mathbf{u}^\top \mathbf{A} \mathbf{v} = \mathbf{u}^\top \mathbf{A}^\top \mathbf{v}, \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^n.$$

Definition 1.2. (positive semidefinite) A symmetric linear transformation $\mathbf{T} : V \rightarrow V$ is called positive semidefinite if

$$\langle \mathbf{T}\mathbf{u}, \mathbf{u} \rangle \geq 0, \quad \forall \mathbf{u} \in V.$$

Equivalently, for the matrix representation \mathbf{A} of \mathbf{T} , this condition is equivalent to \mathbf{A} is symmetric and

$$\mathbf{u}^\top \mathbf{A} \mathbf{u} \geq 0, \quad \forall \mathbf{u} \in \mathbb{R}^n.$$

(a) (10 pts) Prove the following statements are equivalent.

- (i) \mathbf{A} is positive semidefinite.
- (ii) \mathbf{A} is symmetric and all of the eigenvalues of \mathbf{A} are all non-negative.
- (iii) There exists some square matrix \mathbf{B} such that $\mathbf{A} = \mathbf{B}^\top \mathbf{B}$.

(Hint: The symmetric matrix is orthogonally diagonalizable).

In the following section, we will prove an important theorem in machine learning, namely the Singular Value Decomposition (SVD).

Theorem 1 (Singular Value Decomposition). For any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, there exist orthogonal matrices $\mathbf{U} \in \mathbb{R}^{m \times m}$ and $\mathbf{V} \in \mathbb{R}^{n \times n}$, and a diagonal matrix $\Sigma \in \mathbb{R}^{m \times n}$ with nonnegative diagonal entries such that

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^\top.$$

Consider $\mathbf{A}^\top \mathbf{A}$. It is symmetric and hence orthogonally diagonalizable. Suppose $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is an orthonormal set of eigenvectors of $\mathbf{A}^\top \mathbf{A}$ with $\mathbf{A}^\top \mathbf{A} \mathbf{v}_i = \lambda_i \mathbf{v}_i$

(b) (3 pts) Prove that λ_i are nonnegative $\forall i = 1, \dots, n$.

(c) (5 pts) Prove that $\text{Null}(\mathbf{A}^\top \mathbf{A}) = \text{Null}(\mathbf{A})$

Suppose $\sigma_i := \sqrt{\lambda_i}$ for each i , and let $r = \text{rank}(\mathbf{A})$. By reordering the set $\{\sigma_1, \dots, \sigma_n\}$, we may assume that

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0.$$

By problem (c), we know $\sigma_r > 0$ and $\sigma_k = 0 \quad \forall k \in \{r+1, \dots, n\}$.

(d) (5 pts) Prove that $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\}$ is an orthogonal set.

(e) (7 pts) Let $\mathbf{u}_i = \frac{\mathbf{A}\mathbf{v}_i}{\sigma_i}$, complete the proof the SVD.

(Hint: You may start with $\mathbf{A}(\mathbf{v}_1, \dots, \mathbf{v}_n) = (\sigma_1 \mathbf{u}_1, \dots, \sigma_r \mathbf{u}_r, 0, \dots, 0)$ and using Gram-Schmidt process. Don't forget to show that \mathbf{U}, \mathbf{V} are orthogonal matrices.)

Problem 2 (Definition of Derivative as Linear Operator) (10 pts)

We now introduce the rigorous definition of the derivative in finite-dimensional vector spaces. The key idea is that the derivative at a point is a linear operator that best approximates the change of the function.

General definition (Fréchet derivative). Let V, W be finite-dimensional real inner product spaces, and let $f : V \rightarrow W$. We say that f is differentiable at $\mathbf{x} \in V$ if there exists a linear operator

$$Df(\mathbf{x}) : V \rightarrow W$$

such that

$$\lim_{\|\mathbf{h}\| \rightarrow 0} \frac{\|f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - Df(\mathbf{x})[\mathbf{h}]\|_W}{\|\mathbf{h}\|_V} = 0,$$

where $\|\cdot\|_V$ and $\|\cdot\|_W$ denote the norms on V and W respectively. The operator $Df(\mathbf{x})$ is called the derivative (or Jacobian operator) of f at \mathbf{x} .

Remark (scalar-valued functions). If $W = \mathbb{R}$, then $Df(\mathbf{x}) : V \rightarrow \mathbb{R}$ is a linear transformation. In this case the definition is equivalent to

$$f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) = Df(\mathbf{x})[\mathbf{h}] + o(\|\mathbf{h}\|_V), \quad \|\mathbf{h}\| \rightarrow 0.$$

If there exists a vector $g \in V$ such that

$$Df(\mathbf{x})[\mathbf{h}] = \langle g, \mathbf{h} \rangle, \quad \forall \mathbf{h} \in V,$$

we call this vector g the gradient of f at \mathbf{x} , denoted $\nabla f(\mathbf{x})$. It turns out that a fundamental result in mathematics (the Riesz representation theorem) guarantees that such a vector always exists and is unique (you don't need to know the details in this course, but you can look it up if you're interested).

Coordinate form.

- If $V = \mathbb{R}^n$, $W = \mathbb{R}$ with the standard Euclidean inner product, then

$$Df(\mathbf{x})[d\mathbf{x}] = \nabla f(\mathbf{x})^\top d\mathbf{x}.$$

- If $V = \mathbb{R}^{m \times n}$, $W = \mathbb{R}$ with the Frobenius inner product $\langle \mathbf{A}, \mathbf{B} \rangle_F = \text{tr}(\mathbf{A}^\top \mathbf{B})$, then

$$Df(\mathbf{X})[d\mathbf{X}] = \langle \nabla f(\mathbf{X}), d\mathbf{X} \rangle_F = \text{tr}(\nabla f(\mathbf{X})^\top d\mathbf{X}).$$

Chain rule. If $f : V \rightarrow W$ and $g : U \rightarrow V$ are differentiable at $\mathbf{x} \in U$, then their composition $h = f \circ g : U \rightarrow W$ is differentiable at \mathbf{x} and satisfies

$$Dh(\mathbf{x}) = Df(g(\mathbf{x})) \circ Dg(\mathbf{x}).$$

In the scalar-valued case $f : \mathbb{R} \rightarrow \mathbb{R}$ and $g : V \rightarrow \mathbb{R}$, this reduces to the familiar chain rule

$$\nabla(f \circ g)(\mathbf{x}) = f'(g(\mathbf{x})) \nabla g(\mathbf{x}).$$

All the definitions above are just to make sure our calculus is mathematically solid. For the following problems, please use this rigorous way of thinking — that is, start from the small-increment approximation when writing your proofs (which means you can still take derivatives/gradients in the usual intuitive way).

- (a) (3 pts) Given $\mathbf{x}, \mathbf{a} \in \mathbb{R}^n$, show that

$$\nabla \|\mathbf{x} - \mathbf{a}\|_2 = \frac{\partial \|\mathbf{x} - \mathbf{a}\|_2}{\partial \mathbf{x}} = \frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|_2}.$$

- (b) (3 pts) Given $\mathbf{a} \in \mathbb{R}^m$, $\mathbf{X} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^n$, show that

$$\nabla(\mathbf{a}^\top \mathbf{X} \mathbf{b}) = \frac{\partial \mathbf{a}^\top \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^\top.$$

- (c) (4 pts) Given $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$, show that

$$\nabla(\mathbf{x}^\top \mathbf{A} \mathbf{x}) = \frac{\partial(\mathbf{x}^\top \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}.$$

Problem 3 (Matrix Calculus) (20 pts)

In this problem, you need to compute the derivative of a scalar function f with respect to a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$. For convenience, we often denote this differential simply as

$$df = Df(\mathbf{A})[d\mathbf{A}] = \sum_{i,j} \frac{\partial f}{\partial a_{ij}} da_{ij} = \text{tr}(\nabla f(\mathbf{A})^\top d\mathbf{A}).$$

There are several conventions for the matrix gradient (mainly differing by a transpose); we follow the previous convention, consistent with the differential identity and the “summation over i, j ” view of df , and define

$$\nabla f(\mathbf{A}) = \frac{\partial f(\mathbf{A})}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial f}{\partial a_{11}} & \cdots & \frac{\partial f}{\partial a_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial a_{n1}} & \cdots & \frac{\partial f}{\partial a_{nn}} \end{bmatrix}$$

(a) (10 pts) Let $\mathbf{A} \in \mathbb{R}^{n \times n}$. Show that

$$\frac{\partial \det(\mathbf{A})}{\partial \mathbf{A}} = \det(\mathbf{A}) (\mathbf{A}^{-1})^\top.$$

Hint: Recall from high school that for a 2×2 matrix

$$\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}, \quad \mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix},$$

where the matrix on the right (before dividing by $\det(\mathbf{A})$) is constructed by taking cofactors and then transposing — this is the adjugate matrix. Similarly, for 3×3 matrices, you may have seen the cofactor expansion formula for the determinant and the adjugate-based formula for the inverse.

In general, for an $n \times n$ matrix \mathbf{A} , the cofactor matrix is

$$\mathbf{C} = \begin{bmatrix} C_{11} & \cdots & C_{1n} \\ \vdots & \ddots & \vdots \\ C_{n1} & \cdots & C_{nn} \end{bmatrix}, \quad C_{ij} = (-1)^{i+j} M_{ij},$$

where M_{ij} is the determinant of the $(n-1) \times (n-1)$ minor (delete row i , column j). The transpose of \mathbf{C} is called the adjugate matrix:

$$\text{adj}(\mathbf{A}) = \mathbf{C}^\top.$$

One always has the identity

$$\mathbf{A} \text{adj}(\mathbf{A}) = \det(\mathbf{A}) \mathbf{I}.$$

Use this identity and the definition of the determinant to derive the derivative formula.

(b) (10 pts) Prove that

$$\frac{\partial \log(\det(\mathbf{A}))}{\partial a_{ij}} = \mathbf{e}_j^\top \mathbf{A}^{-1} \mathbf{e}_i,$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a (non-singular) matrix, and \mathbf{e}_j is the unit vector along the j -th axis (e.g. $\mathbf{e}_3 = [0, 0, 1, 0, \dots, 0]^\top$). It is common to write the formula as

$$\frac{\partial \log(\det(\mathbf{A}))}{\partial \mathbf{A}} = (\mathbf{A}^{-1})^\top.$$

Hint: You can continue directly from part (a), apply the chain rule to $\log(\det(\mathbf{A}))$.

Alternative hint (prove (b) \Rightarrow (a)). You may also start from Cramer's rule: for $\mathbf{A}\mathbf{x} = \mathbf{e}_i$,

$$x_j = \frac{\det(\mathbf{A}_j)}{\det(\mathbf{A})},$$

where \mathbf{A}_j is obtained from \mathbf{A} by replacing its j -th column with \mathbf{e}_i . By differentiating $\log \det(\mathbf{A})$ along coordinate directions and applying Cramer's rule, you can derive part (b). If you proceed this way and then recover part (a) from (b) via the chain rule, you will still receive full credit.

Problem 4 (Closed-Form Linear Regression Solution) (20 + 15 pts)

Consider the linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\epsilon},$$

where $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\boldsymbol{\theta} \in \mathbb{R}^d$ and $\boldsymbol{\epsilon} \in \mathbb{R}^n$. Denote $\mathbf{X}_i \in \mathbb{R}^{1 \times d}$ as the i -th row of \mathbf{X} , with the following interpretations:

- If the linear model has the bias term, then write $\boldsymbol{\theta} = [w_1, \dots, w_m, b]^\top$ and $\mathbf{X}_i = [x_{i,1}, x_{i,2}, \dots, x_{i,m}, 1]$, namely $d = m + 1$.
- If the linear model has no bias term, then write $\boldsymbol{\theta} = [w_1, \dots, w_d]^\top$ and $\mathbf{X}_i = [x_{i,1}, x_{i,2}, \dots, x_{i,m}]$, namely $d = m$.

(a) (20 pts) Without the bias term, consider the L^2 -regularized loss function:

$$\sum_i \kappa_i (y_i - \mathbf{X}_i \boldsymbol{\theta})^2 + \lambda \sum_j w_j^2, \quad \lambda > 0, \kappa_i > 0 \text{ for all } i.$$

Show that the optimal solution that minimizes the loss function is $\boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{K} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{K} \mathbf{y}$, where

$$\mathbf{K} = \begin{bmatrix} \kappa_1 & & 0 \\ & \ddots & \\ 0 & & \kappa_n \end{bmatrix}$$

is a diagonal matrix and \mathbf{I} is the $d \times d$ identical matrix.

(b) (Bonus) (15 pts) With the bias term, the L^2 -regularized loss function becomes

$$\sum_i \kappa_i (y_i - \mathbf{X}_i \boldsymbol{\theta})^2 + \lambda \sum_j w_j^2, \quad \lambda > 0, \kappa_i > 0 \text{ for all } i.$$

Show that the optimal solution that minimizes the loss function is $\boldsymbol{\theta}^* = [\mathbf{w}^{\star\top}, b^{\star\top}]^\top$, where

$$\begin{aligned} \mathbf{w}^* &= \left(\tilde{\mathbf{X}}^\top \mathbf{K} \tilde{\mathbf{X}} + \lambda \mathbf{I} - \frac{1}{\text{Tr}(\mathbf{K})} \tilde{\mathbf{X}}^\top \mathbf{K} \mathbf{e} \mathbf{e}^\top \mathbf{K} \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^\top \mathbf{K} \left(\mathbf{y} - \frac{1}{\text{Tr}(\mathbf{K})} \mathbf{e} \mathbf{e}^\top \mathbf{K} \mathbf{y} \right), \\ b^* &= \frac{1}{\text{Tr}(\mathbf{K})} (\mathbf{e}^\top \mathbf{K} \mathbf{y} - \mathbf{e}^\top \mathbf{K} \tilde{\mathbf{X}} \mathbf{w}^*) \end{aligned}$$

for which $\mathbf{e} = [1 \dots 1]^\top$ denotes the all one vector, $\mathbf{X} = [\tilde{\mathbf{X}} \mathbf{e}]$, where $\tilde{\mathbf{X}}$ is the dataset of the no bias term. $\text{Tr}(\mathbf{K})$ is the trace of the matrix \mathbf{K} , and that \mathbf{K} and \mathbf{I} are defined as in (a).

Problem 5 (Noise and Regularization) (20 pts)

Consider the linear model $f_{\mathbf{w},b} : \mathbb{R}^k \rightarrow \mathbb{R}$, where $\mathbf{w} \in \mathbb{R}^k$ and $b \in \mathbb{R}$, defined as

$$f_{\mathbf{w},b}(x) = \mathbf{w}^\top \mathbf{x} + b$$

Given dataset $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, if the inputs $\mathbf{x}_i \in \mathbb{R}^k$ are contaminated with input noise $\boldsymbol{\eta}_i \in \mathbb{R}^k$, we may consider the expected sum-of-squares loss in the presence of input noise as

$$\tilde{L}_{ss}(\mathbf{w}, b) = \mathbb{E} \left[\frac{1}{2N} \sum_{i=1}^N (f_{\mathbf{w},b}(\mathbf{x}_i + \boldsymbol{\eta}_i) - y_i)^2 \right]$$

where the expectation is taken over the randomness of input noises $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_N$. Additionally, the inputs (\mathbf{x}_i) and the input noise $(\boldsymbol{\eta}_i)$ are independent.

Now assume the input noises $\boldsymbol{\eta}_i = [\eta_{i,1}, \eta_{i,2}, \dots, \eta_{i,k}]^\top$ are random vectors with zero mean $\mathbb{E}[\eta_{i,j}] = 0$, and the covariance between components is given by

$$\mathbb{E}[\eta_{i,j} \eta_{i',j'}] = \delta_{i,i'} \delta_{j,j'} \sigma^2$$

where $\delta_{i,i'} = \begin{cases} 1 & , \text{ if } i = i' \\ 0 & , \text{ otherwise.} \end{cases}$ denotes the Kronecker delta.

(a) (20 pts) Please show that

$$\tilde{L}_{ss}(\mathbf{w}, b) = \frac{1}{2N} \sum_{i=1}^N (f_{\mathbf{w},b}(\mathbf{x}_i) - y_i)^2 + \frac{\sigma^2}{2} \|\mathbf{w}\|^2$$

That is, minimizing the expected sum-of-squares loss in the presence of input noise is equivalent to minimizing noise-free sum-of-squares loss with the addition of a L^2 -regularization term on the weights.

(Hint: $\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x} = \text{tr}(\mathbf{x}\mathbf{x}^\top)$ and the square of a vector is dot product with itself)