

DDA3020 Homework 1

Due date: Oct 16, 2023

Instructions

- The **deadline** is **23:59, Oct 16, 2023**.
- The weight of this assignment in the final grade is 15%.
- **Electronic submission:** Turn in solutions electronically via Blackboard. Be sure to submit your homework as a single file. Please name your solution file as *A1_studentID_name*
- Note that **late submissions** will result in discounted scores: 0-24 hours \rightarrow 80%, 24-120 hours \rightarrow 50%, 120 or more hours \rightarrow 0%.
- Answer the questions in English. Otherwise, you'll lose half of the points.
- Collaboration policy: You need to solve all questions independently and collaboration between students is **NOT** allowed.
- If you have any questions concerning this homework, feel free to reach out to TA Dong QIAO (dongqiao@link.cuhk.edu.cn) or Wenrang ZHANG (223040237@link.cuhk.edu.cn). You're also welcome to physically visit them during their office hours with your questions.

1 Written Problems (50 points)

1.1. (Matrix Differential, 20 points) It is a fundamental ability to derive the derivatives of some functions in machine learning. Nine basic cases can be summarized as in Table 1. Nevertheless, only three of them (*i.e.*, cells in gray background) are common in our course. In this problem, you will follow the next steps to derive the derivatives for these cases.

(1) Layout of Derivatives

- **Differentiation of a scalar function *w.r.t.* a vector:** If $f : \mathbb{R}^n \mapsto \mathbb{R}$ is a scalar function of n variables, $\mathbf{w} \in \mathbb{R}^n$ is a $n \times 1$ vector, then differentiation of $f(\mathbf{w})$ *w.r.t.* \mathbf{w}

Table 1: Vector/Matrix derivatives

	f	\mathbf{f}	\mathbf{F}
w	$\frac{f}{w}$	$\frac{d\mathbf{f}}{dw}$	$\frac{d\mathbf{F}}{dw}$
\mathbf{w}	$\frac{df}{d\mathbf{w}}$	$\frac{d\mathbf{f}}{d\mathbf{w}}$	$\frac{d\mathbf{F}}{d\mathbf{w}}$
\mathbf{W}	$\frac{df}{d\mathbf{W}}$	$\frac{d\mathbf{f}}{d\mathbf{W}}$	$\frac{d\mathbf{F}}{d\mathbf{W}}$

results in a $n \times 1$ vector;

$$\frac{df}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \vdots \\ \frac{\partial f}{\partial w_n} \end{bmatrix}$$

- **Differentiation of a scalar function *w.r.t.* a matrix:** If $f : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$ is a scalar function of mn variables, $\mathbf{w} \in \mathbb{R}^{m \times n}$ is a $m \times n$ matrix, then differentiation of $f(\mathbf{w})$ *w.r.t.* \mathbf{w} results in a $m \times n$ matrix;

$$\frac{df}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial f}{\partial w_{11}} & \cdots & \frac{\partial f}{\partial w_{1n}} \\ \vdots & \cdots & \vdots \\ \frac{\partial f}{\partial w_{m1}} & \cdots & \frac{\partial f}{\partial w_{mn}} \end{bmatrix}$$

- **Differentiation of a vector function *w.r.t.* a vector:** If $\mathbf{f} : \mathbb{R}^n \mapsto \mathbb{R}^m$ is a vector function of size $m \times 1$ and \mathbf{w} is a $n \times 1$ vectors, then differentiation of $\mathbf{f}(\mathbf{w})$ *w.r.t.* \mathbf{w} results in a $n \times m$ matrix.

$$\frac{d\mathbf{f}^T}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial f_1}{\partial w_1} & \cdots & \frac{\partial f_h}{\partial w_1} \\ \vdots & \cdots & \vdots \\ \frac{\partial f_1}{\partial w_d} & \cdots & \frac{\partial f_h}{\partial w_d} \end{bmatrix}$$

We also call this matrix the gradient matrix of \mathbf{f} *w.r.t.* \mathbf{w} . The tranpose of this matrix is exactly the Jacobian matrix $\frac{d\mathbf{f}}{d\mathbf{w}^T}$ (*i.e.*, $\frac{d\mathbf{f}}{d\mathbf{w}}$ in convention) as

$$\frac{d\mathbf{f}}{d\mathbf{w}^T} = \begin{bmatrix} \frac{\partial f_1}{\partial w_1} & \cdots & \frac{\partial f_1}{\partial w_d} \\ \vdots & \cdots & \vdots \\ \frac{\partial f_h}{\partial w_1} & \cdots & \frac{\partial f_h}{\partial w_d} \end{bmatrix}$$

- (2) **(10 points) Derivation by definition:** Use definition of multivariate calculus to derive the derivatives

Example 1: For $f(\mathbf{w}) = \mathbf{a}^T \mathbf{w}$ with $\mathbf{w} \in \mathbb{R}^n$, we have $\frac{df}{d\mathbf{w}} = \mathbf{a}$ since

$$\frac{df}{d\mathbf{w}} = \frac{\partial \sum_{i=1}^n a_i x_i}{\partial x_i} = \frac{\partial a_i x_i}{\partial x_i} = a_i$$

Derive the following derivatives by definition.

- (5 pts) $f(\mathbf{w}) = \mathbf{w}^T \mathbf{A} \mathbf{w}$ with $\mathbf{w} \in \mathbb{R}^n$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$;
- (5 pts) $f(\mathbf{w}) = \mathbf{A} \mathbf{w}$ with $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{w} \in \mathbb{R}^n$.

(3) **(Bonus task: 5 points) Derivation by differentiation:** derive the derivatives based on the relation between differentiation and gradient as

$$\begin{aligned} f : \mathbb{R}^n &\mapsto \mathbb{R} : & df &= \sum_{i=1}^n \frac{\partial f}{\partial w_i} dw_i = \left(\frac{df}{d\mathbf{w}}\right)^T d\mathbf{w} = \text{tr}\left(\left(\frac{df}{d\mathbf{w}}\right)^T d\mathbf{w}\right) \\ f : \mathbb{R}^{m \times n} &\mapsto \mathbb{R} : & df &= \sum_{i,j} \frac{\partial f}{\partial W_{ij}} dW_{ij} = \text{tr}\left(\left(\frac{df}{d\mathbf{W}}\right)^T d\mathbf{W}\right) \end{aligned}$$

Example 2: For $f(\mathbf{W}) = \mathbf{a}^T \mathbf{W} \mathbf{b}$ with $\mathbf{W} \in \mathbb{R}^{m \times n}$, $\mathbf{a} \in \mathbb{R}^m$, and $\mathbf{b} \in \mathbb{R}^n$, we have $\frac{df}{d\mathbf{W}} = \mathbf{a} \mathbf{b}^T$ since

$$\begin{aligned} df &= d(\mathbf{a}^T \mathbf{W} \mathbf{b}) = d(\mathbf{a}^T) \mathbf{W} \mathbf{b} + \mathbf{a}^T d(\mathbf{W}) \mathbf{b} + \mathbf{a}^T \mathbf{W} d(\mathbf{b}) = \mathbf{a}^T d(\mathbf{W}) \mathbf{b} \\ &= \text{tr}(\mathbf{a}^T d(\mathbf{W}) \mathbf{b}) = \text{tr}(\mathbf{b} \mathbf{a}^T d\mathbf{W}) = \text{tr}((\mathbf{a} \mathbf{b}^T)^T d\mathbf{W}) \end{aligned}$$

Hint: Some useful formulas are here.

$$\begin{aligned} d(X + Y) &= dX + dY & d(XY) &= (dX)Y + Xd(Y) & dX^T &= (dX)^T \\ \text{tr}(X^T) &= \text{tr}(X) & \text{tr}(XY) &= \text{tr}(YX) & d(\text{tr}(X)) &= \text{tr}(dX) \end{aligned}$$

Derive the following derivatives by differentiation.

- (2 pts) $f(\mathbf{w}) = \mathbf{w}^T \mathbf{A} \mathbf{w}$ with $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{w} \in \mathbb{R}^n$;
- (3 pts) $f(\mathbf{W}) = \text{tr}(\mathbf{W}^T \mathbf{A} \mathbf{W})$ with $\mathbf{A} \in \mathbb{R}^{m \times m}$ and $\mathbf{W} \in \mathbb{R}^{m \times n}$.

(4) **(10 points) Derivation by chain rule:** In the realm of machine learning, the input is commonly a feature vector while the loss is almost always a scalar objective function. Therefore, chain rule is useful for gradient update.

(a) (10 pts) Vector chain ending with a scalar: For $\mathbf{x} \rightarrow \mathbf{y}_1 \rightarrow \mathbf{y}_2 \rightarrow \cdots \rightarrow \mathbf{y}_n \rightarrow z$, the chain rule is

$$\frac{dz}{d\mathbf{x}^T} = \frac{dz}{d\mathbf{y}_n^T} \frac{\partial \mathbf{y}_n}{\partial \mathbf{y}_{n-1}^T} \cdots \frac{\partial \mathbf{y}_2}{\partial \mathbf{y}_1^T} \frac{\partial \mathbf{y}_1}{\partial \mathbf{x}^T}$$

Assume $\ell = \|\mathbf{X} \mathbf{w} - \mathbf{y}\|_2^2$, let $\mathbf{z} = \mathbf{X} \mathbf{w} - \mathbf{y}$, then $\ell = \|\mathbf{z}\|_2^2 = \mathbf{z}^T \mathbf{z}$. Consequently, here is a chain like $\mathbf{w} \rightarrow \mathbf{z} \rightarrow \ell$. Follow the chain rule, the derivative $\frac{\partial \ell}{\partial \mathbf{w}^T} = \frac{\partial \ell}{\partial \mathbf{z}^T} \frac{\partial \mathbf{z}}{\partial \mathbf{w}^T}$ which implies that

$$\frac{\partial \ell}{\partial \mathbf{w}} = \left(\frac{\partial \mathbf{z}}{\partial \mathbf{w}^T}\right)^T \frac{\partial \ell}{\partial \mathbf{z}}$$

Determine the closed-form solution of $\frac{\partial \ell}{\partial \mathbf{w}}$.

- (b) (Bonus Task, 5 points) Matrix chain ending with a scalar: For $\mathbf{X} \rightarrow \mathbf{Y} \rightarrow \ell$, the chain rule is

$$\frac{\partial \ell}{\partial X_{ij}} = \sum_{kl} \frac{\partial \ell}{\partial Y_{kl}} \frac{\partial Y_{kl}}{\partial X_{ij}}$$

Assume $\ell = f(\mathbf{Y})$ where $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{B}$, derive the derivative of $\frac{\partial \ell}{\partial \mathbf{X}}$.

Hint: Firstly, it has

$$\frac{\partial Y_{kl}}{\partial X_{ij}} = \frac{\partial \sum_s A_{ks} X_{sl}}{\partial X_{ij}} = \frac{\partial A_{ki} X_{il}}{\partial X_{ij}} = A_{ki} \delta_{lj}$$

where $\delta_{lj} = 1$ for $l = j$; otherwise, $\delta_{lj} = 0$.

Secondly, we have

$$\frac{\partial \ell}{\partial X_{ij}} = \sum_{kl} \frac{\partial \ell}{\partial Y_{kl}} A_{ki} \delta_{lj} = \sum_k \frac{\partial \ell}{\partial Y_{kj}} A_{ki} = \mathbf{A}_{:,i}^T \left(\frac{\partial \ell}{\partial \mathbf{Y}} \right)_{:,j}$$

Lastly, you can derive the form of $\frac{\partial \ell}{\partial \mathbf{X}}$.

1.2. (Convexity of functions, 10 points) Prove that: (1) $f(x) = \text{ReLU}(x) = (x)^+ = \max(0, x)$ is convex for $x \in \mathbb{R}$; (2) $f(x) = |x|$ is convex for $x \in \mathbb{R}$; (3) $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$ is convex where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{x} \in \mathbb{R}^n$.

1.3. (Gradient Descent, 10 points) Suppose we have training data $\{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$, where $\mathbf{x}_i \in \mathbb{R}^d$ and $\mathbf{y}_i \in \mathbb{R}^k$, $i = 1, 2, \dots, N$. (1) Find the closed-form solution of the following problem.

$$\min_{\mathbf{W}, \mathbf{b}} \sum_{i=1}^N \alpha_i \|\mathbf{y}_i - \mathbf{W}\mathbf{x}_i - \mathbf{b}\|_2^2,$$

where the diagonal of diagonal matrix $\text{diag}(\mathbf{A}) = (\alpha_1, \alpha_2, \dots, \alpha_N)^T$ are weights for different sample; (2) Show how to use gradient descent to solve the problem.

Hint: You can use either definition or differentiation method to derive the derivatives. If you use differentiation method, please note that

$$\sum_{i=1}^N \alpha_i \|\mathbf{y}_i - \mathbf{W}\mathbf{x}_i - \mathbf{b}\|_2^2 = \text{tr}[(\mathbf{Y} - \mathbf{X}\mathbf{W})^T \mathbf{A}(\mathbf{Y} - \mathbf{X}\mathbf{W})]$$

where $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N)^T \in \mathbb{R}^{N \times k}$, $\mathbf{X} = [(\mathbf{x}_1^T, 1)^T, (\mathbf{x}_2^T, 1)^T, \dots, (\mathbf{x}_N^T, 1)^T]^T \in \mathbb{R}^{N \times (d+1)}$, $\mathbf{W} = (\mathbf{W}, \mathbf{b})^T \in \mathbb{R}^{(d+1) \times k}$, and $\mathbf{A} = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_N)$.

1.4. (Maximum Likelihood Estimation, 10 points) Suppose x_1, x_2, \dots, x_N are drawn from $\mathcal{N}(\mu, \sigma^2)$. Show that the maximum likelihood estimation (MLE) of σ^2 is $\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{MLE})^2$.

2 Programming (50 points)

2.1. (Polynomial regression, 20 points) In this exercise, we will try to fit a non-linear function g with polynomial regression on the feasible space $\mathbf{X} = [0, 11]$:

Unknown $g(x) = ?$

$$\text{Construct } f(x) = \sum_{i=0}^n \alpha_i x^i \iff f(x) = w^T x', \quad x' = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^n \end{bmatrix}, \quad \text{s.t. } \forall x \in \mathbf{X}, \quad f(x) \approx g(x)$$

Where n is the polynomial degree of freedom and is manually chosen.

Follow the instructions given in the jupyter notebook. At the end of the exercise, you will retrieve an estimation of the desired function and make some comment on this method.

2.2. (Linear regression, 30 points) The CSV or XLS file contains a dataset for regression. There are 7750 samples with 25 features (described in the doc file). This data is for the purpose of bias correction of next-day maximum and minimum air temperatures forecast of the LDAPS model operated by the Korea Meteorological Administration over Seoul, South Korea.

This data consists of summer data from 2013 to 2017. The input data is largely composed of the LDAPS model's next-day forecast data, in-situ maximum and minimum temperatures of present-day, and geographic auxiliary variables. There are two outputs (i.e. next-day maximum and minimum air temperatures) in this data. Hindcast validation was conducted for the period from 2015 to 2017.

You need to delete the first two attributes (station and date), and use attributes 3-23 to predict attributes 24 and 25. Randomly split the data into two parts, one contains 80% of the samples and the other contains 20% of the samples. Use the first part as training data and train a linear regression model and make prediction on the second part. Report the training error and testing error in terms of RMSE.

Repeat the splitting, training, and testing for 10 times. Use a loop and print the RMSEs in each trial.

Note that you need to write the codes of learning the parameters by yourself. Do not use the classification or regression packages of Sklearn. You can use their tools to shuffle the data randomly for splitting.