

DDA5001 Machine Learning

Notes: Matrix Differentiation

Teaching Team

The Chinese University of Hong Kong, Shenzhen



Matrix Differentiation

We have learned some commonly used results in lecture:

- ▶ $\frac{\partial \mathbf{c}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{c},$
- ▶ $\frac{\partial \|\mathbf{x}\|^2}{\partial \mathbf{x}} = \frac{\partial \mathbf{x}^\top \mathbf{x}}{\partial \mathbf{x}} = 2\mathbf{x},$
- ▶ $\frac{\partial \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}^\top,$
- ▶ $\frac{\partial \mathbf{x}^\top \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A}^\top + \mathbf{A})\mathbf{x},$
- ▶ $\frac{\partial \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2}{\partial \mathbf{x}} = 2\mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}).$

How about more complex functions? What if we want to calculate the gradient with respect to a matrix?

- ▶ We need a general rule to do matrix differentiation.

Matrix Differentiation

Consider the function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and gradient $\nabla f(\mathbf{x}) := \left[\frac{\partial f(\mathbf{x})}{\partial x_i} \right] \in \mathbb{R}^d$.

Basic rule:

$$df = \sum_{i=1}^d \frac{\partial f}{\partial x_i} dx_i = \nabla f(\mathbf{x})^\top d\mathbf{x}.$$

More general case where \mathbf{X} is a $\mathbb{R}^{m \times n}$ matrix:

$$df = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial \mathbf{X}_{ij}} d\mathbf{X}_{ij} = \text{Tr} \left(\frac{\partial f}{\partial \mathbf{X}}^\top d\mathbf{X} \right), \quad (1)$$

where $\frac{\partial f}{\partial \mathbf{X}} := \left[\frac{\partial f}{\partial \mathbf{X}_{ij}} \right] \in \mathbb{R}^{m \times n}$ is the gradient.

► The equation (1) can be used to calculate the gradient $\frac{\partial f}{\partial \mathbf{X}}$.

Matrix Differentiation

$$df = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial X_{ij}} dX_{ij} = \text{Tr}\left(\frac{\partial f}{\partial \mathbf{x}}^\top d\mathbf{X}\right).$$

Example: Calculate the gradient of function $f(\mathbf{X}) := \mathbf{a}^\top \mathbf{X} \mathbf{b}$.

$$\begin{aligned} df &= \mathbf{a}^\top d\mathbf{X} \mathbf{b} \quad (d(\mathbf{X}\mathbf{Y}) = (d\mathbf{X})\mathbf{Y} + \mathbf{X}(d\mathbf{Y})) \\ \implies \text{Tr}(df) &= \text{Tr}(\mathbf{a}^\top d\mathbf{X} \mathbf{b}) \\ &= \text{Tr}(\mathbf{b} \mathbf{a}^\top d\mathbf{X}) \quad (\text{Tr}(\mathbf{X}\mathbf{Y}) = \text{Tr}(\mathbf{Y}\mathbf{X})). \end{aligned}$$

$$\text{Hence, } \frac{\partial f}{\partial \mathbf{x}}^\top = \mathbf{b} \mathbf{a}^\top \implies \frac{\partial f}{\partial \mathbf{x}} = \mathbf{a} \mathbf{b}^\top.$$

Matrix Differential Step: (scalar to matrix/vector)

- ▶ Calculate df and apply trace operator to df .
- ▶ Establish the relation in (1) based on **basic rules**.

Basic Rules

Differential Rules:

- (1) $d(\mathbf{X} \pm \mathbf{Y}) = d\mathbf{X} \pm d\mathbf{Y}$.
- (2) $d(\mathbf{X}\mathbf{Y}) = (d\mathbf{X})\mathbf{Y} + \mathbf{X}(d\mathbf{Y})$.
- (3) $d\mathbf{X}^{-1} = -\mathbf{X}^{-1}d\mathbf{X}\mathbf{X}^{-1}$.
- (4) $d(\mathbf{X} \odot \mathbf{Y}) = d\mathbf{X} \odot \mathbf{Y} + \mathbf{X} \odot d\mathbf{Y}$, where \odot is element-wise multiplication.
- (5) $d\sigma(\mathbf{X}) = \sigma'(\mathbf{X}) \odot d\mathbf{X}$, where $\sigma(\mathbf{X}) = [\sigma(\mathbf{X}_{ij})]$ is element-wise function mapping.

Trace Rules:

- (1) $a = \text{Tr}(a)$, where a is scalar.
- (2) $\text{Tr}(\mathbf{X}\mathbf{Y}) = \text{Tr}(\mathbf{Y}\mathbf{X})$, given \mathbf{X} and \mathbf{Y}^\top have the same size.
- (3) $\text{Tr}(\mathbf{X}^\top) = \text{Tr}(\mathbf{X})$.
- (4) $\text{Tr}(\mathbf{X} \pm \mathbf{Y}) = \text{Tr}(\mathbf{X}) \pm \text{Tr}(\mathbf{Y})$.
- (5) $\text{Tr}(\mathbf{X}^\top (\mathbf{Y} \odot \mathbf{C})) = \text{Tr}((\mathbf{X} \odot \mathbf{Y})^\top \mathbf{C})$, given $\mathbf{X}, \mathbf{Y}, \mathbf{C}$ have the same size.

Exercise 1: Least Square

Problem 1: calculate the gradient of the least square problem

$$f(\boldsymbol{\theta}) = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2.$$

Solution.

$$\begin{aligned} df &= d[(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})] \\ &= [d(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})]^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) + (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top [d(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})] \\ &= (\mathbf{X}d\boldsymbol{\theta})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) + (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top \mathbf{X}d\boldsymbol{\theta} \\ &= 2(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top \mathbf{X}d\boldsymbol{\theta}. \end{aligned}$$

Thus,

$$\begin{aligned} df &= \text{Tr}(df) = \text{Tr}(2(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top \mathbf{X}d\boldsymbol{\theta}) \\ &\implies \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = 2\mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}). \end{aligned}$$

Exercise 2: Multi-Class Logistic Regression

Problem 2: Let $\Theta \in \mathbb{R}^{K \times d}$ be the parameter. The objective of one-sample logistic regression is

$$\min_{\Theta} f(\Theta) := -\mathbf{y}^\top \log \sigma(\Theta \mathbf{x}),$$

where $\mathbf{x} \in \mathbb{R}^d$ is the input sample, $\mathbf{y} \in \mathbb{R}^K$ is one-hot vector whose value is 1 for the corresponding class and 0 otherwise. $\sigma(\mathbf{z}) := \frac{\exp(\mathbf{z})}{\mathbf{1}^\top \exp(\mathbf{z})} \in \mathbb{R}^K$ is the softmax function, where $\exp(\mathbf{z})$ is element-wise exponential function.

Derive the gradient $\frac{\partial f(\Theta)}{\partial \Theta}$.

Exercise 2: Solution

Let $\mathbf{1}$ denotes all-one vector, we have

$$\begin{aligned} f &= -\mathbf{y}^\top (\log \exp(\Theta \mathbf{x}) - \mathbf{1} \log(\mathbf{1}^\top \exp(\Theta \mathbf{x}))) \\ &= -\mathbf{y}^\top \Theta \mathbf{x} + \log(\mathbf{1}^\top \exp(\Theta \mathbf{x})). \end{aligned}$$

Differentiating f :

$$\begin{aligned} df &= -\mathbf{y}^\top (d\Theta) \mathbf{x} + \frac{\mathbf{1}^\top \exp(\Theta \mathbf{x}) \odot (d\Theta) \mathbf{x}}{\mathbf{1}^\top \exp(\Theta \mathbf{x})} \\ &= -\mathbf{y}^\top (d\Theta) \mathbf{x} + \frac{\exp(\Theta \mathbf{x})^\top (d\Theta) \mathbf{x}}{\mathbf{1}^\top \exp(\Theta \mathbf{x})}. \quad (\mathbf{1}^\top \mathbf{u} \odot \mathbf{v} = \mathbf{u}^\top \mathbf{v}) \end{aligned}$$

Take trace operator and rearrange using exchange property:

$$\text{Tr}(df) = \text{Tr} \left(\mathbf{x} \left(-\mathbf{y}^\top + \frac{\exp(\Theta \mathbf{x})^\top}{\mathbf{1}^\top \exp(\Theta \mathbf{x})} \right) d\Theta \right).$$

Hence, the gradient is $\frac{\partial f}{\partial \Theta} = (-\mathbf{y} + \sigma(\Theta \mathbf{x})) \mathbf{x}^\top$.

(The loss is only for one-sample. How about multi-sample case?)