

Matrix differentiation

Hung-Hsuan Chen

What is matrix differentiation

- A specialized notation for doing multivariable differentiation, especially over spaces of matrices

Notations

- x (lower case) \Rightarrow a scalar
- \mathbf{x} (lower case, bold face) \Rightarrow a column vector
 - $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$
- \mathbf{X} (upper case, bold face) \Rightarrow a matrix
 - $\mathbf{X} = (x_{ij}) = (\mathbf{x}_1^T, \mathbf{x}_2^T, \dots, \mathbf{x}_n^T)$

Example

```
4.6,3.2,1.4,0.2,Iris-setosa
5.3,3.7,1.5,0.2,Iris-setosa
5.0,3.3,1.4,0.2,Iris-setosa
7.0,3.2,4.7,1.4,Iris-versicolor
6.4,3.2,4.5,1.5,Iris-versicolor
6.9,3.1,4.9,1.5,Iris-versicolor
5.5,2.3,4.0,1.3,Iris-versicolor
6.5,2.8,4.6,1.5,Iris-versicolor
```

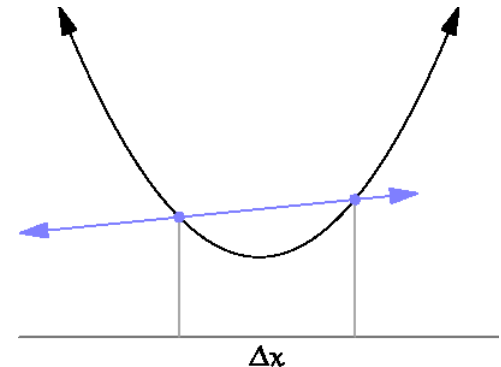
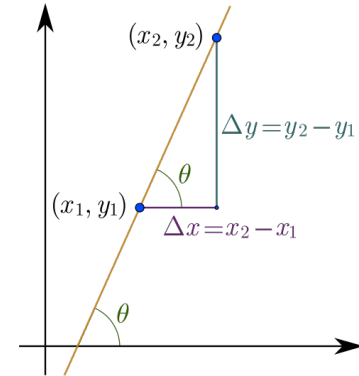
$$x_{1,2} = 3.2, x_{1,3} = 1.4$$

$$X = \begin{bmatrix} 4.6 & \cdots & 0.2 \\ 5.3 & & 0.2 \\ 5.0 & & 0.2 \\ 7.0 & \ddots & 1.4 \\ 6.4 & & 1.5 \\ 6.9 & & 1.5 \\ 5.5 & & 1.3 \\ 6.5 & \cdots & 1.5 \end{bmatrix} \in R^{8 \times 4}, y = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \in R^{8 \times 1}$$

$$\mathbf{x}_{1*} = \mathbf{x}_1 = \begin{bmatrix} 4.6 \\ 3.2 \\ 1.4 \\ 0.2 \end{bmatrix} \in R^{4 \times 1}, \quad \mathbf{x}_{*1} = \begin{bmatrix} 4.6 \\ 5.3 \\ \vdots \\ 6.5 \end{bmatrix} \in R^{8 \times 1},$$

Derivative (scalar)

- An ***infinitesimal*** change in x is denoted by dx , and the derivative of y with respect to x is written dy/dx or $f'(x)$
 - If $f(x) = x^k$, $f'(x) = kx^{k-1}$
 - If $f(x) = e^x$, $f'(x) = e^x$
 - If $f(x) = a^x$, $f'(x) = a^x \ln(a)$
 - If $f(x) = \ln x$, $f'(x) = \frac{1}{x}$
 - If $f(x) = \log_a x$, $f'(x) = \frac{1}{x \ln a}$



Combined function and chain rule

- Given f and g are functions, a and b are real numbers:
 - Sum rule: $(af + bg)' = af' + bg'$
 - Product rule: $(fg)' = f'g + fg'$
 - Quotient rule: $\left(\frac{f}{g}\right)' = \frac{f'g - fg'}{g^2}$
- If $f(x) = g(h(x))$
 - Chain rule: $\frac{df(x)}{dx} = \frac{dg(h(x))}{dh(x)} \cdot \frac{dh(x)}{dx}$
 - E.g., $\frac{de^{7x}}{dx} = \frac{de^{7x}}{d7x} \cdot \frac{d7x}{dx} = e^{7x} \cdot 7$

Partial derivative (scalar)

- A partial derivative of a function of several variables is its derivative with respect to one of those variables, with the others held constant

– E.g., $f(x_1, x_2) = x_1^2 + 3x_1x_2 + x_2^3$

- $\frac{\partial f(x_1, x_2)}{\partial x_1} = 2x_1 + 3x_2$

- $\frac{\partial f(x_1, x_2)}{\partial x_2} = 3x_1 + 3x_2^2$

Types of matrix derivatives

Types	Scalar	Vector	Matrix
Scalar	$\frac{\partial y}{\partial x}$	$\frac{\partial \mathbf{y}}{\partial x}$	$\frac{\partial \mathbf{Y}}{\partial x}$
Vector	$\frac{\partial y}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$	
Matrix	$\frac{\partial y}{\partial \mathbf{X}}$		

- The three types of derivatives that have not been considered are not as widely considered and a notation is not widely agreed upon

Layout conventions

- There are mainly two types of layout conventions in matrix calculus
 - Numerator Layout Notation
 - Denominator Layout Notation
- Most books and papers don't state which convention they use
- Even worse, sometimes the two conventions are mixed in the equations
- This confuses the beginners
- We will mostly follow the **Numerator Layout Notation** unless otherwise mentioned



Types of matrix derivatives of different layout conventions

~~note~~

	Scalar y		Vector y (size m)		Matrix Y (size $m \times n$)	
	Notation	Type	Notation	Type	Notation	Type
Scalar x	$\frac{\partial y}{\partial x}$	scalar	$\frac{\partial \mathbf{y}}{\partial x}$	(numerator layout) size- m column vector (denominator layout) size- m row vector	$\frac{\partial \mathbf{Y}}{\partial x}$	(numerator layout) $m \times n$ matrix
Vector \mathbf{x} (size n)	$\frac{\partial y}{\partial \mathbf{x}}$	(numerator layout) size- n row vector (denominator layout) size- n column vector	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$	(numerator layout) $m \times n$ matrix (denominator layout) $n \times m$ matrix	$\frac{\partial \mathbf{Y}}{\partial \mathbf{x}}$	
Matrix \mathbf{X} (size $p \times q$)	$\frac{\partial y}{\partial \mathbf{X}}$	(numerator layout) $q \times p$ matrix (denominator layout) $p \times q$ matrix	$\frac{\partial \mathbf{y}}{\partial \mathbf{X}}$		$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$	

ex $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} \end{bmatrix}$

$$\frac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \end{bmatrix}$$

Derivative by scalar

- $\frac{\partial y}{\partial x}$

- $\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix}$

- $\frac{\partial \mathbf{Y}}{\partial x} = \begin{bmatrix} \frac{\partial y_{11}}{\partial x} & \dots & \frac{\partial y_{1n}}{\partial x} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_{m1}}{\partial x} & \dots & \frac{\partial y_{mn}}{\partial x} \end{bmatrix}$

Derivative by vector

- $\frac{\partial y}{\partial \mathbf{x}} = \left[\frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$

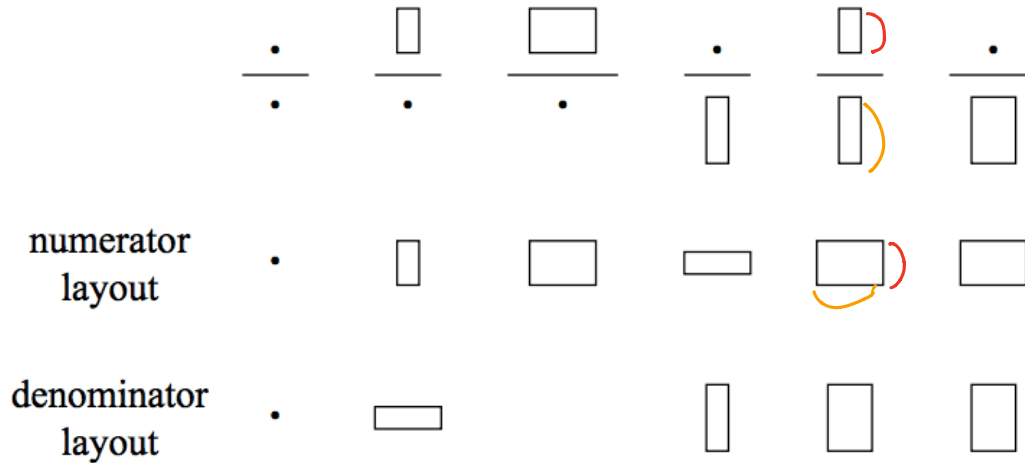
- $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$

Derivative by matrix

$$\bullet \frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \cdots & \frac{\partial y}{\partial x_{m1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{1n}} & \cdots & \frac{\partial y}{\partial x_{mn}} \end{bmatrix}$$



Pictorial representation



Commonly used derivatives

Scalar a , vector \mathbf{a} , and matrix \mathbf{A} are not functions of x , \mathbf{x} , and \mathbf{X}

- $\frac{da}{dx} = \mathbf{0}$ (column vector)
- $\frac{da}{dx} = \mathbf{0}^T$ (row vector)
- $\frac{da}{d\mathbf{X}} = \mathbf{0}^T$ (shape is the same as \mathbf{X}^T)
- $\frac{da}{dx} = \mathbf{0}$ (shape is $\text{len}(\mathbf{a}) * \text{len}(\mathbf{x})$)
- $\frac{d\mathbf{a}^T \mathbf{x}}{dx} = \frac{d\mathbf{x}^T \mathbf{a}}{dx} = \mathbf{a}^T$
- $\frac{d\mathbf{x}^T \mathbf{x}}{d\mathbf{x}} = 2\mathbf{x}^T$ (Note: \mathbf{x}^T is a scalar, \mathbf{x} is a vector)
- $\frac{d\mathbf{A}\mathbf{x}}{d\mathbf{x}} = \mathbf{A}$

Exercise

$$\bullet \frac{d(x^T a)^2}{dx} = \frac{d(x^T a)^2}{d(x^T a)} \cdot \frac{d(x^T a)}{dx} = (2x^T a) a^T$$

Multiple linear regression

- n : the number of training instances
- d : the number of features
- Training instances:

$$- \mathbf{X} = \begin{bmatrix} x_{1,1} & \dots & x_{1,d} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \dots & x_{n,d} \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

– (We assume no coefficient parameter here)

- Find $\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$ such that $(\hat{\mathbf{y}} - \mathbf{y})^T (\hat{\mathbf{y}} - \mathbf{y})$ is minimized, where

$$- \hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\theta}$$

- The solution is $\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} \quad \hat{y}_i = \mathbf{x}_i^T \boldsymbol{\theta}$$

$$\sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$\hat{\mathbf{y}} - \mathbf{y} = \begin{bmatrix} \hat{y}_1 - y_1 \\ \vdots \\ \hat{y}_n - y_n \end{bmatrix}$$

Solving θ

$$J(\theta) = \frac{1}{2}(\hat{y} - y)^T(\hat{y} - y) = \frac{1}{2}(X\theta - y)^T(X\theta - y)$$

$$\rightarrow \frac{\partial J(\theta)}{\partial \theta} = \frac{\partial (X\theta - y)^T(X\theta - y)}{\partial (X\theta - y)} \cdot \frac{\partial (X\theta - y)}{\partial \theta} = (X\theta - y)^T X := 0$$

$$\rightarrow (\theta^T X^T - y^T)X = 0$$

$$\rightarrow \theta^T X^T X = y^T X$$

$$\rightarrow (\theta^T X^T X)^T = (y^T X)^T$$

$$\rightarrow X^T X \theta = X^T y$$

$$\rightarrow \theta = (X^T X)^{-1} X^T y$$

$$\bullet \frac{dx^T x}{dx} = 2x^T$$

$$\bullet \frac{dAx}{dx} = A$$

Summary

- Matrix and vector are compact ways to denote set of variables
- Matrix and vector differentiation may be confusing sometimes, mostly because of inconsistent notations
 - Numerator vs denominator layouts
 - $\frac{da^T x}{dx} = \mathbf{a}^T$ or \mathbf{a} ?
- Sometimes write out the full matrix or vector is helpful

Quiz

- Given
 - Random variables: x is a scalar, $\mathbf{x} \in R^{n \times 1}$, $\mathbf{X} \in R^{n \times m}$
 - Functions of $x, \mathbf{x}, \mathbf{X}$: y is a scalar, $\mathbf{y} \in R^{n \times 1}$, $\mathbf{Y} \in R^{n \times m}$
- What are the shapes of the followings?
 - $\frac{\partial y}{\partial x}; \frac{\partial \mathbf{y}}{\partial \mathbf{x}}; \frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$ $\{ \text{ca. } n \times 1, n \times m$
 - $\frac{\partial y}{\partial x}; \frac{\partial \mathbf{y}}{\partial \mathbf{x}}; \frac{\partial \mathbf{Y}}{\partial x}$ $1 \times n, n \times m,$
 - $\frac{\partial y}{\partial \mathbf{X}}; \frac{\partial \mathbf{y}}{\partial \mathbf{X}}; \frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$ $m \times n,$

$$\otimes \quad \frac{\partial \vec{y}}{\partial \vec{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \end{bmatrix}$$

$$\otimes \quad \frac{\partial y}{\partial \vec{X}} \quad \vec{X} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} \quad y = f(x_1 \dots x_4)$$

$$\frac{\partial y}{\partial \vec{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{12}} \\ \frac{\partial y}{\partial x_{21}} & \frac{\partial y}{\partial x_{22}} \end{bmatrix}$$

$$\otimes \quad y_1 = x \quad \vec{y} = \begin{bmatrix} x \\ 2x \end{bmatrix} \quad \frac{\partial \vec{y}}{\partial x} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$\otimes \quad \frac{\partial \vec{y}}{\partial x}$$

$$y_2 = 2x$$

$$\otimes \quad y_1 = x_1 + x_2, \quad y_2 = 2x_1 - x_2$$

$$\frac{\partial \vec{y}}{\partial \vec{X}} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \quad \frac{\partial \vec{y}}{\partial \vec{X}} = \begin{bmatrix} 1 & 1 \\ 2 & -1 \end{bmatrix}$$

$$\vec{X} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

☆

$$\frac{\partial Y}{\partial X}$$

$$Y = \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{bmatrix} = \begin{bmatrix} x & 2x \\ -x & 3x \end{bmatrix}$$

$$\frac{\partial Y}{\partial X} = \begin{bmatrix} 1 & 2 \\ -1 & 3 \end{bmatrix}$$