

COS 302 Precept 6

Spring 2020

Princeton University

Outline

Overview of Matrix Gradients

Review of Differentiation Basics

Examples of Various Gradients

Outline

Overview of Matrix Gradients

Review of Differentiation Basics

Examples of Various Gradients

Different Flavors of Gradients

Type	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	$\frac{\partial y}{\partial \mathbf{X}}$		

¹ source: <https://www.comp.nus.edu.sg/cs5240/lecture/matrix-differentiation.pdf>

Challenges of Vector/Matrix Calculus

- Good news: most of the rules you know and love from single variable calculus generalize well (not in all cases but in some).
- Bad news: confusing notation (many more variables lead to very tedious book-keeping) and more identities to memorize.

Numerator vs Denominator Layout

- Two main conventions used in vector/matrix calculus: numerator and denominator layouts.
- Numerator layout makes the dimension of the derivative be the numerator dimension by denominator dimension.
- For example, if y is a scalar and $\mathbf{x} \in \mathbb{R}^N$ then

$$\frac{\partial y}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times N}$$

How to Become a Differentiation Master

1. Identify what is the flavor of your derivative.
2. What is the dimension of the derivative?
3. Identify any differentiation rules you will need for this case.
4. Can any part of the derivative be reduced to a particular identity?
5. Identify the partial derivatives.

Outline

Overview of Matrix Gradients

Review of Differentiation Basics

Examples of Various Gradients

Definitions

Names	Notation & Expression
Difference Quotient	$\frac{\delta y}{\delta x} = \frac{f(x+\delta x) - f(x)}{\delta x}$
Derivative	$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$
Partial Derivative	$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_i+h, \dots, x_n) - f(x)}{h}$
Gradient	$\nabla_{\mathbf{x}} f = \frac{df}{d\mathbf{x}} = \left[\frac{\partial f}{\partial x_1} \cdots \frac{\partial f}{\partial x_n} \right]$
Jacobian	$\nabla_{\mathbf{x}} \mathbf{f} = \frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \left[\frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} \cdots \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \right]$

Differentiation Rules (Scalar-Scalar)

- Sum Rule

$$(f(x) + g(x))' = f'(x) + g'(x)$$

- Product Rule

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$$

- Quotient Rule

$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$$

- Chain Rule

$$(g(f(x)))' = (g \circ f)'(x) = g'(f(x))f'(x)$$

Differentiation Rules (Scalar-Vector)

- Sum Rule

$$\frac{\partial}{\partial \mathbf{x}}(f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}}$$

- Product Rule

$$\frac{\partial}{\partial \mathbf{x}}(f(\mathbf{x})g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}}g(\mathbf{x}) + f(\mathbf{x})\frac{\partial g}{\partial \mathbf{x}}$$

- Chain Rule

$$\frac{\partial}{\partial \mathbf{x}}(g(f(\mathbf{x}))) = \frac{\partial}{\partial \mathbf{x}}(g \circ f)(\mathbf{x}) = \frac{\partial g}{\partial f} \frac{\partial f}{\partial \mathbf{x}}$$

Outline

Overview of Matrix Gradients

Review of Differentiation Basics

Examples of Various Gradients

Gradient of Matrix Multiplication

Consider the matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ and the vector $\mathbf{x} \in \mathbb{R}^N$. Define the vector function $\mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x}$ where $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^M$, what is $\frac{d\mathbf{f}}{d\mathbf{x}}$?

Gradient of Matrix Multiplication cont.

1. Dimension of gradient: $\frac{df}{dx} \in \mathbb{R}^{M \times N}$
2. One of these $M \times N$ partial derivatives will look like:

$$f_i(\mathbf{x}) = \sum_{j=1}^N A_{ij}x_j \implies \frac{\partial f_i}{\partial x_j} = A_{ij}$$

3. Collecting all of these partial derivatives:

$$\frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix} = \mathbf{A} \in \mathbb{R}^{M \times N}$$

Chain Rule Example

Consider the scalar function $h : \mathbb{R} \rightarrow \mathbb{R}$ where $h(t) = f(g(t))$ with $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}^2$ such that

$$f(\mathbf{x}) = \exp(x_1 x_2^2) \ ,$$
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = g(t) = \begin{bmatrix} t \cos t \\ t \sin t \end{bmatrix} .$$

What is $\frac{dh}{dt}$?

Chain Rule Example cont.

Even though the gradient is a scalar we need to compute two vector gradients (gradients) because of the chain rule:

$$\begin{aligned}\frac{dh}{dt} &= \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial t} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{bmatrix} \\ &= \begin{bmatrix} \exp(x_1 x_2^2) x_2^2 & 2 \exp(x_1 x_2^2) x_1 x_2 \end{bmatrix} \begin{bmatrix} \cos t - t \sin t \\ \sin t + t \cos t \end{bmatrix} \\ &= \exp(x_1 x_2^2) (x_2^2 (\cos t - x_2) + 2x_1 x_2 (\sin t + x_1)),\end{aligned}$$

where $x_1 = t \cos t$ and $x_2 = t \sin t$

Least Squares by Chain Rule

Consider the matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ where $N > D$ and the two vectors $\mathbf{y} \in \mathbb{R}^N$ and $\boldsymbol{\beta} \in \mathbb{R}^D$. Before we saw in class that the over-determined system of linear equations:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}$$

does not always have a solution. Instead of solving the problem directly we can try to find an approximate solution $\hat{\boldsymbol{\beta}}$.

Least Squares by Chain Rule cont.

If we picked a random β , and multiplied it by \mathbf{X} , we probably wouldn't get a vector that was very close to \mathbf{y} . Specifically, the error vector

$$\mathbf{e}(\beta) = \mathbf{y} - \mathbf{X}\beta$$

would probably not be close to the zero vector. A good choice of β is one that minimizes the Euclidean distance between \mathbf{y} and $\mathbf{X}\beta$. Specifically, one that minimizes the function $L(\mathbf{e}) = \|\mathbf{e}\|^2$.

Least Squares by Chain Rule cont.

To find the best β , let's take the gradient of L with respect to β and set it equal to zero.

$$1. \frac{\partial L}{\partial \beta} \in \mathbb{R}^{1 \times D}$$

$$2. \frac{\partial L}{\partial \beta} = \frac{\partial L}{\partial \mathbf{e}} \frac{\partial \mathbf{e}}{\partial \beta}$$

$$3. \frac{\partial L}{\partial e_i} = \frac{\partial}{\partial e_i} \sum_{i=1}^N e_i^2 = 2e_i. \text{ Since, } \frac{\partial L}{\partial \mathbf{e}} \in \mathbb{R}^{1 \times N} \text{ we have } \frac{\partial L}{\partial \mathbf{e}} = 2\mathbf{e}^T$$

$$4. \frac{\partial \mathbf{e}}{\partial \beta} = -\mathbf{X} \in \mathbb{R}^{N \times D}$$

$$5. \frac{\partial L}{\partial \beta} = -2\mathbf{e}^T \mathbf{X} = -2(\mathbf{y}^T - \beta^T \mathbf{X}^T) \mathbf{X} = 0 \implies \beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$