

Review of High-Dimensional Calculus

Ju Sun*

January 23, 2025

High-dimensional calculus is typically not covered in basic calculus course series, but it is the language of modern machine learning—we almost always express quantities of interest as vectors, matrices, or even tensors. In this set of notes, we quickly go over the basics of high-dimensional calculus that are most useful for machine learning and highlight certain computational techniques that are not often taught elsewhere. Two recommended sources for learning high-dimensional calculus are [Mun97, Col12]. [Zor15, Zor16] is another set of useful references. There is also an interesting online tool for computing symbolic matrix derivatives <https://www.matrixcalculus.org/>.

1 Our notations

- scalars: x (small letters) vectors: \mathbf{x} (bold small) matrices: \mathbf{X} (bold capital) tensors: \mathcal{X} (script capital) sets: S (capital)
- vectors are always **column vectors**, unless stated otherwise
- x_i : i -th element of \mathbf{x} x_{ij} : (i, j) -th element of \mathbf{X} \mathbf{x}^i : i -th row of \mathbf{X} as a **row vector** \mathbf{x}_j : j -th column of \mathbf{X} as a **column vector**
- \mathbb{R} : real numbers \mathbb{R}_+ : positive reals \mathbb{R}^n : space of n -dimensional vectors $\mathbb{R}^{m \times n}$: space of $m \times n$ matrices $\mathbb{R}^{m \times n \times k}$: space of $m \times n \times k$ tensors, etc
- $[n] \doteq \{1, \dots, n\}$ (a notation often used by theoretical computer scientists; \doteq means “defined as”)

2 Differentiability

2.1 First-order differentiability

Definition 2.1 (First-order derivative or Jacobian). Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$. f is (Fréchet) differentiable at a point \mathbf{x} if there exists a matrix $\mathbf{B} \in \mathbb{R}^{m \times n}$ such that

$$\frac{f(\mathbf{x} + \boldsymbol{\delta}) - f(\mathbf{x}) - \mathbf{B}\boldsymbol{\delta}}{\|\boldsymbol{\delta}\|_2} \rightarrow \mathbf{0} \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0}, \quad (2.1)$$

*Department of Computer Science and Engineering, University of Minnesota at Twin Cities. Email: jusun@umn.edu.

or equivalently,

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \mathbf{B}\boldsymbol{\delta} + \underbrace{o(\|\boldsymbol{\delta}\|_2)}_{\text{lower-order terms than } \boldsymbol{\delta}} \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0}. \quad (2.2)$$

Here \mathbf{B} is called the (Fréchet) derivative, or the Jacobian of f at \mathbf{x} , denoted as $\mathbf{J}_f(\mathbf{x})$.

Here, for a vector-valued function $h(\boldsymbol{\delta}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$,

$$h(\boldsymbol{\delta}) \in o(\|\boldsymbol{\delta}\|_2) \text{ as } \boldsymbol{\delta} \rightarrow \mathbf{0} \iff \lim_{\boldsymbol{\delta} \rightarrow \mathbf{0}} \frac{h(\boldsymbol{\delta})}{\|\boldsymbol{\delta}\|_2} = \mathbf{0} \iff \lim_{\boldsymbol{\delta} \rightarrow \mathbf{0}} \frac{\|h(\boldsymbol{\delta})\|_2}{\|\boldsymbol{\delta}\|_2} = 0. \quad (2.3)$$

Definition 2.2 (Gradient). Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$. The gradient $\nabla f(\mathbf{x}) \in \mathbb{R}^{n \times 1}$, which is a column vector, is the transpose of the Jacobian $\mathbf{J}_f(\mathbf{x}) \in \mathbb{R}^{1 \times n}$.

According to this convention, for $f(\mathbf{X}) : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, the gradient should be $\mathbb{R}^{(m \times n) \times 1}$, which is a length- $(m \times n)$ vector. This is *inconvenient* for many purposes. So in practice, the vector is shaped into an $m \times n$ matrix, i.e., having the same shape as \mathbf{X} . We will use this convention, i.e., $\nabla f(\mathbf{X}) = \left[\frac{\partial f}{\partial x_{i,j}}(\mathbf{X}) \right]_{i,j} \in \mathbb{R}^{m \times n}$.¹

For computation of Jacobian $\mathbf{J}_f(\mathbf{x})$ for $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$, we have

$$\mathbf{J}_f(\mathbf{x}) = \left[\frac{\partial f_i}{\partial x_j}(\mathbf{x}) \right]_{i,j}, \quad (2.4)$$

i.e., $\mathbf{J}_f(\mathbf{x})$ is the collection of all the first-order partial derivatives $\frac{\partial f_i}{\partial x_j}(\mathbf{x})$ for all i, j .

Theorem 2.3 (A sufficient condition for first-order differentiability). If all first-order partial derivatives exist and are continuous at \mathbf{x} , then $f(\mathbf{x})$ is first-order differentiable at \mathbf{x} .

However, the condition is not necessary for first-order differentiability.

Theorem 2.4 (A sufficient and necessary condition for first-order differentiability). Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and let $f_i(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ be the i -th component function of f , i.e.,

$$f(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix}. \quad (2.5)$$

Then f is differentiable at a point \mathbf{x} if and only if each component function f_i is differentiable at \mathbf{x} . If so, its derivative is the $m \times n$ matrix whose i -th row is the derivative of f_i .

Theorem 2.5 (Calculus rules of Jacobian). Assume $f, g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are differentiable at a point $\mathbf{x} \in \mathbb{R}^n$.

- **linearity:** $\lambda_1 f + \lambda_2 g$ is differentiable at \mathbf{x} and $\mathbf{J}_{\lambda_1 f + \lambda_2 g}(\mathbf{x}) = \lambda_1 \mathbf{J}_f(\mathbf{x}) + \lambda_2 \mathbf{J}_g(\mathbf{x})$.
- **product:** assume $m = 1$, f, g is differentiable at \mathbf{x} and $\nabla[f g](\mathbf{x}) = f(\mathbf{x})\nabla g(\mathbf{x}) + g(\mathbf{x})\nabla f(\mathbf{x})$.
- **quotient:** assume $m = 1$ and $g(\mathbf{x}) \neq 0$, $\frac{f}{g}$ is differentiable at \mathbf{x} and $\nabla \left[\frac{f}{g} \right](\mathbf{x}) = \frac{g(\mathbf{x})\nabla f(\mathbf{x}) - f(\mathbf{x})\nabla g(\mathbf{x})}{g^2(\mathbf{x})}$.
- **Chain rule:** Let $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $h : \mathbb{R}^n \rightarrow \mathbb{R}^k$. If f is differentiable at \mathbf{x} , and h is differentiable at \mathbf{y} where $\mathbf{y} = f(\mathbf{x})$. Then, $h \circ f : \mathbb{R}^m \rightarrow \mathbb{R}^k$ is differentiable at \mathbf{x} , and

$$\mathbf{J}_{[h \circ f]}(\mathbf{x}) = \mathbf{J}_h(f(\mathbf{x}))\mathbf{J}_f(\mathbf{x}). \quad (2.6)$$

When $k = 1$,

$$\nabla[h \circ f](\mathbf{x}) = \mathbf{J}_f^\top(\mathbf{x})\nabla h(f(\mathbf{x})). \quad (2.7)$$

¹We write $\mathbf{M} = [m_{ij}]_{i,j}$ to mean \mathbf{M} is a matrix in which the (i, j) -th entry takes the form m_{ij} .

2.2 Deriving the Jacobian

There are two main methods. One is by invoking Eq. (2.4) and the calculus rules in Theorem 2.5. The other is by the perturbation-expansion method based on Definition 2.1, as explained below.

Let's take an example $f(\mathbf{x}) = \mathbf{Ax} - \mathbf{b}$ for $\mathbf{A} \in \mathbb{R}^{n \times m}$ and $\mathbf{b} \in \mathbb{R}^n$ and try to derive the Jacobian.

- **Elementwise calculation and calculus rules.** Obviously $\mathbf{J}_b(\mathbf{x}) = \mathbf{0}$ so we can focus on the \mathbf{Ax} term. Now

$$\frac{\partial(\mathbf{Ax})_i}{\partial x_j} = \frac{\partial(\mathbf{a}^i \mathbf{x})}{\partial x_j} = \frac{\partial \sum_k a_{ik} x_k}{\partial x_j} = a_{ij}. \quad (2.8)$$

So $\mathbf{J}_{\mathbf{Ax}}(\mathbf{x}) = \mathbf{A}$, and by linearity $\mathbf{J}_f(\mathbf{x}) = \mathbf{A}$.

- **Perturbation-expansion method.** We make a sufficiently small perturbation δ to \mathbf{x} , so

$$f(\mathbf{x} + \delta) = \mathbf{A}(\mathbf{x} + \delta) - \mathbf{b} = (\mathbf{Ax} - \mathbf{b}) + \mathbf{A}\delta = f(\mathbf{x}) + \mathbf{A}\delta. \quad (2.9)$$

Comparing this to Eq. (2.2) in Definition 2.1, we easily obtain that $\mathbf{J}_f(\mathbf{x}) = \mathbf{A}$.

In the perturbation-expansion method, after the infinitesimal perturbation δ , we rearrange the terms to match the form of Eq. (2.2), i.e.,

$$f(\mathbf{x}) + \text{linear term in } \delta + \text{lower-order term in } \|\delta\|_2 \quad (2.10)$$

so that we can read off the Jacobian from the linear term.

Example 2.6. Consider the least-squares objective $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$. We will derive the Jacobian, which is the transpose of the gradient.

- **Chain rule.** We can view f as composition of $g(\mathbf{x}) = \mathbf{Ax} - \mathbf{b}$ and $h(\mathbf{y}) = \|\mathbf{y}\|_2^2$ so that $f = h \circ g(\mathbf{x})$. From our last example, $\mathbf{J}_g(\mathbf{x}) = \mathbf{A}$. For h , it is easy to check that $\mathbf{J}_h(\mathbf{y}) = 2\mathbf{y}^\top$. Applying the chain rule, we obtain that

$$\mathbf{J}_f(\mathbf{x}) = \mathbf{J}_h(\mathbf{Ax} - \mathbf{b})\mathbf{J}_g(\mathbf{x}) = 2(\mathbf{Ax} - \mathbf{b})^\top \mathbf{A}. \quad (2.11)$$

- **Perturbation-expansion method.** By making an infinitesimal perturbation δ to \mathbf{x} , we obtain

$$f(\mathbf{x} + \delta) = \|\mathbf{A}(\mathbf{x} + \delta) - \mathbf{b}\|_2^2 \quad (2.12)$$

$$= \|(\mathbf{Ax} - \mathbf{b}) + \mathbf{A}\delta\|_2^2 \quad (2.13)$$

$$= \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \|\mathbf{A}\delta\|_2^2 + 2\langle \mathbf{Ax} - \mathbf{b}, \mathbf{A}\delta \rangle. \quad (2.14)$$

Let us make some clarification before proceeding. We use $\langle \cdot, \cdot \rangle$ to mean the inner product for vectors, i.e., for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$, $\langle \mathbf{u}, \mathbf{v} \rangle \doteq \sum_i u_i v_i$. For all $p \geq 1$, the ℓ_p norm of the vector $\mathbf{u} \in \mathbb{R}^n$ is defined as $\|\mathbf{u}\|_p \doteq (\sum_i |u_i|^p)^{1/p}$. For $p = 2$, the norm is also called the Euclidean norm and it can be easily verified that $\|\mathbf{u}\|_2 = \sqrt{\langle \mathbf{u}, \mathbf{u} \rangle}$. So for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$,

$$\|\mathbf{u} + \mathbf{v}\|_2^2 = \langle \mathbf{u} + \mathbf{v}, \mathbf{u} + \mathbf{v} \rangle = \langle \mathbf{u}, \mathbf{u} \rangle + \langle \mathbf{v}, \mathbf{v} \rangle + 2\langle \mathbf{u}, \mathbf{v} \rangle = \|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 + 2\langle \mathbf{u}, \mathbf{v} \rangle. \quad (2.15)$$

We have used this identity to arrive at Eq. (2.14). In Eq. (2.14), $\|\mathbf{Ax} - \mathbf{b}\|_2^2 = f(\mathbf{x})$, and $\|\mathbf{A}\delta\|_2^2 \in O(\|\delta\|_2^2) \implies \|\mathbf{A}\delta\|_2^2 \in o(\|\delta\|_2)$ which we do not care. The linear term is $2\langle \mathbf{Ax} - \mathbf{b}, \mathbf{A}\delta \rangle$. We now invoke another identity $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top \mathbf{v}$ to obtain that

$$2\langle \mathbf{Ax} - \mathbf{b}, \mathbf{A}\delta \rangle = 2(\mathbf{Ax} - \mathbf{b})^\top \mathbf{A}\delta. \quad (2.16)$$

Comparing this with Eq. (2.2), we conclude that

$$\mathbf{J}_f(\mathbf{x}) = 2(\mathbf{Ax} - \mathbf{b})^\top \mathbf{A}. \quad (2.17)$$

The chain rule is fine for simple compositions. But it quickly leads to fatigue when there are many compositions. Moreover, when intermediate variables involve matrices, tensors will often be involved. An example is when deriving gradients for functions involving deep neural networks, e.g.,

$$f(\mathbf{W}) = \sum_i \|\mathbf{y}_i - \mathbf{W}_k \sigma(\mathbf{W}_{k-1} \dots \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}_i))\|_2^2. \quad (2.18)$$

So, the perturbation-expansion method becomes handy when we deal with matrix/vector variables. They allow us to derive Jacobian, gradient, and Hessian directly in matrix form, *without taking partial derivatives*.

Example 2.7 (Chain rule follows from the perturbation-expansion method). *The chain rule in Theorem 2.5 can be easily derived from the perturbation-expansion method. Consider an infinitesimal perturbation δ to \mathbf{x} in $h \circ f$:*

$$h \circ f(\mathbf{x} + \delta) = h(f(\mathbf{x} + \delta)) = h(f(\mathbf{x}) + \mathbf{J}_f(\mathbf{x})\delta + o(\|\delta\|_2)), \quad (2.19)$$

where we expand $f(\mathbf{x} + \delta)$ by the definition of first-order differentiability, because f is differentiable at \mathbf{x} . Now $h(f(\mathbf{x}) + \mathbf{J}_f(\mathbf{x})\delta + o(\|\delta\|_2))$ is h at the point $f(\mathbf{x})$ perturbed by the infinitesimal quantity $\mathbf{J}_f(\mathbf{x})\delta + o(\|\delta\|_2)$. Since h is differentiable at the point $f(\mathbf{x})$, we can invoke Eq. (2.2) again and obtain that

$$\begin{aligned} & h(f(\mathbf{x}) + \mathbf{J}_f(\mathbf{x})\delta + o(\|\delta\|_2)) \\ &= h(f(\mathbf{x})) + \mathbf{J}_h(f(\mathbf{x}))(\mathbf{J}_f(\mathbf{x})\delta + o(\|\delta\|_2)) + \underbrace{o(\|\mathbf{J}_f(\mathbf{x})\delta + o(\|\delta\|_2)\|_2)}_{o(\|\delta\|_2)} \end{aligned} \quad (2.20)$$

$$= h(f(\mathbf{x})) + \mathbf{J}_h(f(\mathbf{x}))\mathbf{J}_f(\mathbf{x})\delta + \underbrace{\mathbf{J}_h(f(\mathbf{x}))o(\|\delta\|_2)}_{o(\|\delta\|_2)} + o(\|\delta\|_2) \quad (2.21)$$

$$= h(f(\mathbf{x})) + \mathbf{J}_h(f(\mathbf{x}))\mathbf{J}_f(\mathbf{x})\delta + o(\|\delta\|_2). \quad (2.22)$$

So $h \circ f$ is differentiable at \mathbf{x} , and $\mathbf{J}_{h \circ f}(\mathbf{x}) = \mathbf{J}_h(f(\mathbf{x}))\mathbf{J}_f(\mathbf{x})$.

In the above derivation, we have used several basic properties of the small- o notation:

Lemma 2.8. *As $\delta \rightarrow \mathbf{0}$, (1) $o(\|f(\delta)\|) \in o(\|h(\delta)\|)$ if $\|f(\delta)\| \leq \|h(\delta)\|$ everywhere; (2) $o(\|\delta\|) + o(\|\delta\|) \in o(\|\delta\|)$; Moreover, for any matrix \mathbf{A} , (3) $o(\|\mathbf{A}\delta\|) \in o(\|\delta\|)$; (4) $\mathbf{A}o(\|\delta\|) \in o(\|\delta\|)$.*

Proof. (1) Suppose that $g(\delta) \in o(\|f(\delta)\|)$, i.e., $\lim_{\delta \rightarrow \mathbf{0}} \|g(\delta)\|/\|f(\delta)\| = 0$, then

$$\lim_{\delta \rightarrow \mathbf{0}} \|g(\delta)\|/\|h(\delta)\| \leq \lim_{\delta \rightarrow \mathbf{0}} \|g(\delta)\|/\|f(\delta)\| = 0.$$

(2) Suppose that $g_1(\delta) \in o(\|\delta\|)$ and $g_2(\delta) \in o(\|\delta\|)$, then

$$\lim_{\delta \rightarrow \mathbf{0}} \|g_1(\delta) + g_2(\delta)\|/\|\delta\| \leq \lim_{\delta \rightarrow \mathbf{0}} \|g_1(\delta)\|/\|\delta\| + \lim_{\delta \rightarrow \mathbf{0}} \|g_2(\delta)\|/\|\delta\| = 0.$$

(3) Since $\|\mathbf{A}\delta\| \leq \|\mathbf{A}\|_{\text{op}}\|\delta\|$ (here $\|\mathbf{A}\|_{\text{op}}$ denotes the operator norm of \mathbf{A} , i.e., the largest singular value of \mathbf{A}), we have $o(\|\mathbf{A}\delta\|) \in o(\|\mathbf{A}\|_{\text{op}}\|\delta\|)$ by (1). Since $\|\mathbf{A}\|_{\text{op}}$ is a constant independent of δ , we can omit it as we only worry about the order of the term with respect to δ .

(4) For any $g(\delta) \in o(\|\delta\|)$, i.e., $\lim_{\delta \rightarrow \mathbf{0}} \|g(\delta)\|/\|\delta\| = 0$, we have

$$\lim_{\delta \rightarrow \mathbf{0}} \|\mathbf{A}g(\delta)\|/\|\delta\| \leq \lim_{\delta \rightarrow \mathbf{0}} \|\mathbf{A}\|_{\text{op}}\|g(\delta)\|/\|\delta\| = \|\mathbf{A}\|_{\text{op}} \lim_{\delta \rightarrow \mathbf{0}} \|g(\delta)\|/\|\delta\| = 0,$$

implying the claimed results. ■

2.3 Second-order differentiability

It is possible to define second- or even higher-order differentiability for general $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$. For our purposes, it suffices to consider real-valued functions $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$, on which we focus exclusively here. Assume f is first-order differentiable in a small ball around \mathbf{x} .

- Write $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \doteq \left[\frac{\partial}{\partial x_j} \left(\frac{\partial f}{\partial x_i} \right) \right](\mathbf{x})$ provided the right side is well defined.
- **Symmetry:** If both $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$ and $\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x})$ exist and both are continuous at \mathbf{x} , then they are equal.
- **Hessian (matrix):**

$$\nabla^2 f(\mathbf{x}) \doteq \left[\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x}) \right]_{j,i}.$$

$\nabla^2 f$ is symmetric due to the symmetry property above.

- **Sufficient condition:** if all $\frac{\partial^2 f}{\partial x_j \partial x_i}(\mathbf{x})$ exist and are **continuous** near \mathbf{x} , f is 2nd-order differentiable at \mathbf{x} (the converse is not true; we omit the precise definition of 2nd-order differentiability due to its technicality).

3 Taylor's theorems

Taylor's theorems take several forms. Here we focus on the form useful for gradient and Hessian derivation.

Theorem 3.1 (Taylor's theorem: scalar-variable version). Consider $f(x) : \mathbb{R} \rightarrow \mathbb{R}$.

- If f is 1st-order differentiable at x , then

$$f(x + \delta) = \underbrace{f(x) + \delta f'(x)}_{\text{first-order Taylor expansion}} + o(|\delta|) \quad \text{as } \delta \rightarrow 0. \quad (3.1)$$

- If f is 2nd-order differentiable at x , then

$$f(x + \delta) = \underbrace{f(x) + \delta f'(x) + \frac{1}{2} \delta^2 f''(x)}_{\text{second-order Taylor expansion}} + o(|\delta|^2) \quad \text{as } \delta \rightarrow 0. \quad (3.2)$$

The result can be easily generalized to real-valued vector- and matrix-variable functions.

Theorem 3.2 (Taylor's theorem: vector-variable version). Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$.

- If f is 1st-order differentiable at \mathbf{x} , then

$$f(\mathbf{x} + \boldsymbol{\delta}) = \underbrace{f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle}_{\text{first-order Taylor expansion}} + o(\|\boldsymbol{\delta}\|_2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0}. \quad (3.3)$$

- If f is 2nd-order differentiable at \mathbf{x} , then

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \underbrace{\langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}) \boldsymbol{\delta} \rangle}_{\text{second-order Taylor expansion}} + o(\|\boldsymbol{\delta}\|_2^2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0}. \quad (3.4)$$

To present the matrix version, we need to clarify the definitions of inner product and Euclidean norm for matrices, both natural generalizations of those for vectors. For $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{m \times n}$,

$$\langle \mathbf{U}, \mathbf{V} \rangle = \sum_{i,j} u_{ij} v_{ij} \quad \text{and} \quad \|\mathbf{U}\|_F = \sqrt{\sum_{i,j} u_{ij}^2} = \sqrt{\langle \mathbf{U}, \mathbf{U} \rangle}. \quad (3.5)$$

In other words, let $\text{vec}(\mathbf{U})$ be the vectorized version of \mathbf{U} by sequentially stacking its columns into a long vector. We have

$$\langle \mathbf{U}, \mathbf{V} \rangle = \langle \text{vec}(\mathbf{U}), \text{vec}(\mathbf{V}) \rangle \quad \text{and} \quad \|\mathbf{U}\|_F = \|\text{vec}(\mathbf{U})\|_2. \quad (3.6)$$

Theorem 3.3 (Taylor's theorem: matrix-variable version). Consider $f(\mathbf{X}) : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$.

- If f is 1st-order differentiable at \mathbf{X} , then

$$f(\mathbf{X} + \boldsymbol{\Delta}) = \underbrace{f(\mathbf{X}) + \langle \nabla f(\mathbf{X}), \boldsymbol{\Delta} \rangle}_{\text{first-order Taylor expansion}} + o(\|\boldsymbol{\Delta}\|_F) \quad \text{as } \boldsymbol{\Delta} \rightarrow \mathbf{0}. \quad (3.7)$$

- If f is 2nd-order differentiable at \mathbf{X} , then

$$f(\mathbf{X} + \boldsymbol{\Delta}) = \underbrace{f(\mathbf{X}) + \langle \nabla f(\mathbf{X}), \boldsymbol{\Delta} \rangle + \frac{1}{2} \langle \boldsymbol{\Delta}, \nabla^2 f(\mathbf{X})[\boldsymbol{\Delta}] \rangle}_{\text{second-order Taylor expansion}} + o(\|\boldsymbol{\Delta}\|_F^2) \quad \text{as } \boldsymbol{\Delta} \rightarrow \mathbf{0}. \quad (3.8)$$

Here, $\nabla^2 f(\mathbf{X})$ is a 4-dimensional tensor (i.e., 4-dimensional array) that collects all the second-order derivatives. It is possible to define tensor-matrix product, but it can quickly get technical, and there are several nonequivalent ways of doing it. For our purposes, we only need to know that: (1) $\nabla^2 f(\mathbf{X})[\boldsymbol{\Delta}]$ is an $m \times n$ matrix here; (2) this form is consistent with the vector-variable version, if we vectorize the variable: let $F(\text{vec}(\mathbf{X})) \doteq f(\mathbf{X})$, we have

$$F(\text{vec}(\mathbf{X} + \boldsymbol{\Delta})) = F(\text{vec}(\mathbf{X})) + \langle \nabla F(\text{vec}(\mathbf{X})), \text{vec}(\boldsymbol{\Delta}) \rangle + \frac{1}{2} \langle \text{vec}(\boldsymbol{\Delta}), \nabla^2 F(\text{vec}(\mathbf{X})) \text{vec}(\boldsymbol{\Delta}) \rangle + o(\|\text{vec}(\boldsymbol{\Delta})\|_2^2). \quad (3.9)$$

So, if we have to implement the second-order expansion for matrix-variable functions, a safe way is to vectorize the variable and then call the vector-variable version.

Now we want to put Taylor's theorems to good use. But before that, we need another important property of the Taylor expansion. In short, *Taylor expansion is unique*.

Theorem 3.4 (Asymptotic uniqueness of Taylor expansion—scalar version). Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be k ($k \geq 1$ integer) times differentiable at a point x . If $P(\delta)$ is a k -th order polynomial satisfying $f(x + \delta) - P(\delta) = o(\delta^k)$ as $\delta \rightarrow 0$, then $P(\delta) = f(x) + \sum_{i=1}^k \frac{1}{i!} f^{(i)}(x) \delta^i$, i.e., k -th order Taylor expansion.

Why is this useful? Typically, we calculate derivatives to obtain the Taylor expansion. This theorem enables the reverse path. Suppose we somehow obtain a k -th order polynomial $P(\delta)$ satisfying $f(x + \delta) = P(\delta) + o(|\delta|)$ —e.g., by the perturbation-expansion technique described above, we can compare it to the standard Taylor expansion form and read off the derivatives.

Example 3.5. Consider $f(x) = x^3$ and let us calculate $\nabla f(x)$ and $\nabla^2 f(x)$. For any infinitesimal perturbation δ ,

$$f(x + \delta) = (x + \delta)^3 = x^3 + 3x^2\delta + 3x\delta^2 + \delta^3. \quad (3.10)$$

First-order term in δ is $3x^2\delta$ and so $\nabla f(x) = 3x^2$. Second-order term in δ is $3x\delta^2$ and so $\nabla^2 f(x) = 6x$.

The uniqueness property also holds for the vector and matrix versions.

Theorem 3.6 (Asymptotic uniqueness of Taylor expansion—vector version). Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$.

- Assume $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is 1st-order differentiable at \mathbf{x} . If $P(\boldsymbol{\delta}) \doteq f(\mathbf{x}) + \langle \mathbf{v}, \boldsymbol{\delta} \rangle$ satisfies that

$$f(\mathbf{x} + \boldsymbol{\delta}) - P(\boldsymbol{\delta}) = o(\|\boldsymbol{\delta}\|_2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0},$$

then $P(\boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle$, i.e., the 1st-order Taylor expansion, and $\mathbf{v} = \nabla f(\mathbf{x})$.

- Assume $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is 2nd-order differentiable at \mathbf{x} . If $P(\boldsymbol{\delta}) \doteq f(\mathbf{x}) + \langle \mathbf{v}, \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \mathbf{H}\boldsymbol{\delta} \rangle$ with \mathbf{H} symmetric satisfies that

$$f(\mathbf{x} + \boldsymbol{\delta}) - P(\boldsymbol{\delta}) = o(\|\boldsymbol{\delta}\|_2^2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0},$$

then $P(\boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x})\boldsymbol{\delta} \rangle$, i.e., the 2nd-order Taylor expansion, and $\mathbf{v} = \nabla f(\mathbf{x})$, $\mathbf{H} = \nabla^2 f(\mathbf{x})$.

The matrix version, as well as proofs of the asymptotic uniqueness properties and other forms of Taylor's theorems, can be found in Chapter 5 of [Col12].

Now we provide a couple of examples to show how the perturbation-expansion technique can help us to move from Taylor expansions to gradient and Hessian.

Example 3.7. Let's consider $f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$ again and try to derive $\nabla f(\mathbf{x})$ and $\nabla^2 f(\mathbf{x})$. From Example 2.6, we know that

$$f(\mathbf{x} + \boldsymbol{\delta}) = \underbrace{\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2}_{1\text{-st order in } \boldsymbol{\delta}} + 2 \underbrace{\langle \mathbf{A}\mathbf{x} - \mathbf{b}, \mathbf{A}\boldsymbol{\delta} \rangle}_{1\text{-st order in } \boldsymbol{\delta}} + \underbrace{\|\mathbf{A}\boldsymbol{\delta}\|_2^2}_{2\text{-nd order in } \boldsymbol{\delta}}. \quad (3.11)$$

To read off the gradient, we need to rearrange the 1-st order term into the form $\langle \clubsuit, \boldsymbol{\delta} \rangle$ for some \clubsuit . Now we need a useful rule for manipulating vector/matrix inner products.

Any leading matrix can be transposed and moved to the leading position of the other side of the inner product; similarly, any trailing matrix can be transposed and moved to the trailing position of the other side of the inner product. For example, consider matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$ with compatible dimensions so that $\langle \mathbf{AB}, \mathbf{CD} \rangle$ is well defined. Then

$$\langle \mathbf{AB}, \mathbf{CD} \rangle = \langle \mathbf{B}, \mathbf{A}^\top \mathbf{CD} \rangle = \langle \mathbf{C}^\top \mathbf{AB}, \mathbf{D} \rangle = \langle \mathbf{A}, \mathbf{CDB}^\top \rangle = \langle \mathbf{ABD}^\top, \mathbf{C} \rangle. \quad (3.12)$$

This property can be derived from the cyclic property of matrix traces.

So we can rearrange the 1st order term as

$$2 \langle \mathbf{A}\mathbf{x} - \mathbf{b}, \mathbf{A}\boldsymbol{\delta} \rangle = \langle 2\mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}), \boldsymbol{\delta} \rangle, \quad (3.13)$$

implying that $\nabla f(\mathbf{x}) = 2\mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{b})$. For the 2nd order term,

$$\|\mathbf{A}\boldsymbol{\delta}\|_2^2 = \langle \mathbf{A}\boldsymbol{\delta}, \mathbf{A}\boldsymbol{\delta} \rangle = \langle \boldsymbol{\delta}, \mathbf{A}^\top \mathbf{A}\boldsymbol{\delta} \rangle, \quad (3.14)$$

which is to be compared to $\frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x})\boldsymbol{\delta} \rangle$, implying that $\nabla^2 f(\mathbf{x}) = 2\mathbf{A}^\top \mathbf{A}$.

Example 3.8. We now consider a matrix-variable problem with two blocks of variables

$$f(\mathbf{W}_1, \mathbf{W}_2) = \sum_i \|\mathbf{y}_i - \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_i\|_2^2 = \|\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}\|_F^2 \quad (3.15)$$

and try to derive the gradient. This is an objective corresponding to a two-layer linear neural network. Making infinitesimal perturbation to $\mathbf{W}_1, \mathbf{W}_2$, we obtain

$$f(\mathbf{W}_1 + \Delta_1, \mathbf{W}_2 + \Delta_2) = \|\mathbf{Y} - (\mathbf{W}_2 + \Delta_2)(\mathbf{W}_1 + \Delta_1)\mathbf{X}\|_F^2 \quad (3.16)$$

$$= \|(\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}) - \mathbf{W}_2 \Delta_1 \mathbf{X} - \Delta_2 \mathbf{W}_1 \mathbf{X} - \Delta_2 \Delta_1 \mathbf{X}\|_F^2. \quad (3.17)$$

Now we need the identity: for matrices \mathbf{U}, \mathbf{V} of the same size, $\|\mathbf{U} + \mathbf{V}\|_F^2 = \|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2 + 2\langle \mathbf{U}, \mathbf{V} \rangle$, which can be derived similarly to the vector version. First,

$$\|(\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}) - \mathbf{W}_2 \Delta_1 \mathbf{X} - \Delta_2 \mathbf{W}_1 \mathbf{X} - \Delta_2 \Delta_1 \mathbf{X}\|_F^2 \quad (3.18)$$

$$= \|(\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}) - \mathbf{W}_2 \Delta_1 \mathbf{X} - \Delta_2 \mathbf{W}_1 \mathbf{X}\|_F^2 + \underbrace{\|\Delta_2 \Delta_1 \mathbf{X}\|_F^2}_{o(\|\Delta\|_F)} - 2 \underbrace{\langle (\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}) - \mathbf{W}_2 \Delta_1 \mathbf{X} - \Delta_2 \mathbf{W}_1 \mathbf{X}, \Delta_2 \Delta_1 \mathbf{X} \rangle}_{o(\|\Delta\|_F)}. \quad (3.19)$$

So we only need to focus on

$$\|(\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}) - \mathbf{W}_2 \Delta_1 \mathbf{X} - \Delta_2 \mathbf{W}_1 \mathbf{X}\|_F^2 = \underbrace{\|\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}\|_F^2}_{f(\mathbf{W}_1, \mathbf{W}_2)} + \underbrace{\|\mathbf{W}_2 \Delta_1 \mathbf{X} + \Delta_2 \mathbf{W}_1 \mathbf{X}\|_F^2}_{o(\|\Delta\|_F)} - 2 \langle \mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}, \mathbf{W}_2 \Delta_1 \mathbf{X} + \Delta_2 \mathbf{W}_1 \mathbf{X} \rangle. \quad (3.20)$$

We now only need to compare the linear term $-2 \langle \mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}, \mathbf{W}_2 \Delta_1 \mathbf{X} + \Delta_2 \mathbf{W}_1 \mathbf{X} \rangle$ with

$$\left\langle \begin{bmatrix} \nabla_{\mathbf{W}_1} f \\ \nabla_{\mathbf{W}_2} f \end{bmatrix}, \begin{bmatrix} \Delta_1 \\ \Delta_2 \end{bmatrix} \right\rangle = \langle \nabla_{\mathbf{W}_1} f, \Delta_1 \rangle + \langle \nabla_{\mathbf{W}_2} f, \Delta_2 \rangle. \quad (3.21)$$

We have that

$$-2 \langle \mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}, \mathbf{W}_2 \Delta_1 \mathbf{X} \rangle = -2 \langle \mathbf{W}_2^\top (\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}) \mathbf{X}^\top, \Delta_1 \rangle, \quad (3.22)$$

$$-2 \langle \mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}, \Delta_2 \mathbf{W}_1 \mathbf{X} \rangle = -2 \langle (\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}) \mathbf{X}^\top \mathbf{W}_1^\top, \Delta_2 \rangle, \quad (3.23)$$

implying that

$$\nabla_{\mathbf{W}_1} f = -2 \mathbf{W}_2^\top (\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}) \mathbf{X}^\top, \quad (3.24)$$

$$\nabla_{\mathbf{W}_2} f = -2 (\mathbf{Y} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{X}) \mathbf{X}^\top \mathbf{W}_1^\top. \quad (3.25)$$

Note that in this example, due to Eq. (3.21), one can also perturb \mathbf{W}_1 only (i.e., $\Delta_2 = \mathbf{0}$) to obtain $\nabla_{\mathbf{W}_1} f$ and similarly perturb \mathbf{W}_2 only to obtain $\nabla_{\mathbf{W}_2} f$. This tends to make the process less messy. Similarly, for functions with multiple blocks, one can take turns to perturb one block each time to derive block-wise gradients (in general, this does NOT work for higher-order derivatives!).

Final words on this: we have discussed two or three techniques for deriving derivatives. For practical problems, often a mixture of these techniques works the best. So, be flexible!

4 Directional derivatives and curvatures

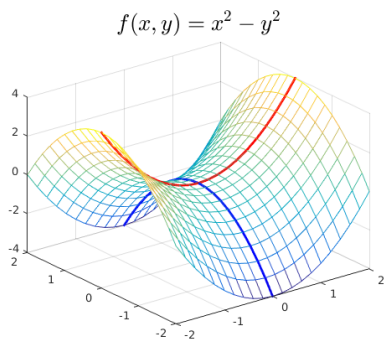


Figure 1: Blue: negative curvature (bending down); Red: positive curvature (bending up)

Consider $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$.

- **directional derivative:** $D_v f(\mathbf{x}) \doteq \left. \frac{d}{dt} f(\mathbf{x} + t\mathbf{v}) \right|_{t=0}$, i.e., rate of change at \mathbf{x} along \mathbf{v}

- When f is 1-st order differentiable at \mathbf{x} ,

$$D_v f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle.$$

- Now $D_v f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is another function. What is $D_u(D_v f)(\mathbf{x})$? If f is 2nd-order differentiable at \mathbf{x} ,

$$D_u(D_v f)(\mathbf{x}) = \langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{v} \rangle.$$

When $\mathbf{u} = \mathbf{v}$,

$$D_u(D_u f)(\mathbf{x}) = \langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle = \left. \frac{d^2}{dt^2} f(\mathbf{x} + t\mathbf{u}) \right|_{t=0},$$

which is the **directional curvature** along \mathbf{u} and grows quadratically with respect to $\|\mathbf{u}\|_2$. To make it independent of the norm $\|\mathbf{u}\|_2$, one can consider $\frac{\langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle}{\|\mathbf{u}\|_2^2}$.

The spectral property (i.e., the distribution of eigenvalues and eigenvectors) of $\nabla^2 f(\mathbf{x})$ determines the directional curvatures. In particular, eigenvector directions corresponding to negative (positive) eigenvalues of $\nabla^2 f(\mathbf{x})$ have negative (positive) curvatures.

Further reading

Chapters 3 & 5 of [DFO20] are particularly relevant, and you are encouraged to review the materials there.

Disclaimer

This set of notes is preliminary and has not been thoroughly proofread. Typos and factual errors are well expected, and hence use it with caution. Bug reports are very welcome and should be sent to Prof. Ju Sun via jusun@umn.edu.

References

- [Col12] Rodney Coleman, *Calculus on normed vector spaces*, Springer New York, 2012.
- [DFO20] Marc Peter Deisenroth, A. Aldo Faisal, and Cheng Soon Ong, *Mathematics for machine learning*, Cambridge University Press, 2020.
- [Mun97] James R. Munkres, *Analysis on manifolds*, Taylor & Francis Inc, 1997.
- [Zor15] Vladimir A. Zorich, *Mathematical analysis I*, Springer Berlin Heidelberg, 2015.
- [Zor16] ———, *Mathematical analysis II*, Springer Berlin Heidelberg, 2016.