

Lecture 2 Part 1: Derivatives in Higher Dimensions: Jacobians and Matrix Functions

MIT 18.S096 Matrix Calculus For Machine Learning And Beyond

March 20, 2024

All definitions and theorems follow from Chapter 5 of Magnus and Neudecker (1988).

1 Vector-to-vector functions

Differentiability. Let f be a function mapping from \mathbb{R}^m to \mathbb{R}^n . Let $\mathbf{x}, d\mathbf{x} \in \mathbb{R}^m$. If there exists a real $m \times n$ matrix A , depending on \mathbf{x} but not $d\mathbf{x}$, such that

$$f(\mathbf{x} + d\mathbf{x}) = f(\mathbf{x}) + A(\mathbf{x})(d\mathbf{x}) + r_{\mathbf{x}}(d\mathbf{x}) \quad \text{with} \quad \lim_{\|d\mathbf{x}\| \rightarrow 0} \frac{r_{\mathbf{x}}(d\mathbf{x})}{\|d\mathbf{x}\|} = 0,$$

then f is said to be *differentiable* at \mathbf{x} .

Derivative. The matrix $A(\mathbf{x}) \in \mathbb{R}^{n \times m}$ is called the (first) *derivative* of f at \mathbf{x} . The *gradient* is simply the transpose of the derivative.

Magnus and Neudecker (1988) goes on to prove several important facts about differentiable functions:

- The derivative $A(\mathbf{x})$ is unique.
- If f is differentiable at \mathbf{x} , then f is continuous at \mathbf{x} .
- If f is differentiable at \mathbf{x} , then all partial derivatives

$$D_j f_i(\mathbf{x}) \triangleq \lim_{t \rightarrow 0} \frac{f_i(\mathbf{x} + t e_j) - f_i(\mathbf{x})}{t}$$

exist at \mathbf{x} ¹. Let $Df(\mathbf{x})$ be a matrix where the ij -entry is $D_j f_i(\mathbf{x})$, i.e., the *Jacobian*. One can show that $A(\mathbf{x}) = Df(\mathbf{x})$. This “reveals that the elements $a_{ij}(\mathbf{x})$ of the matrix $A(\mathbf{x})$ are, in fact, precisely the partial derivatives $D_j f_i(\mathbf{x})$ ”.

2 A very simple example from the lecture

Too simple to be included.

3 My first example

Function:

$$f(\mathbf{x}) = \mathbf{x} \odot \mathbf{x}$$

1. The converse is not true. For the converse to be true, we must add the assumption that the partial derivatives are continuous at \mathbf{x} . This result can be found in college-level calculus textbooks. See, for example, Chapter 14 Partial Derivatives in *Calculus – Early Transcendentals* (8th edition) by James Stewart. I must admit that I didn’t pay attention to this in college.

Difference:

$$\begin{aligned}
 f(\mathbf{x} + d\mathbf{x}) - f(\mathbf{x}) &= (\mathbf{x} + d\mathbf{x}) \odot (\mathbf{x} + d\mathbf{x}) - \mathbf{x} \odot \mathbf{x} \\
 &= \mathbf{x} \odot \mathbf{x} + 2\mathbf{x} \odot (d\mathbf{x}) + (d\mathbf{x}) \odot (d\mathbf{x}) - \mathbf{x} \odot \mathbf{x} \\
 &= 2\mathbf{x} \odot (d\mathbf{x}) + (d\mathbf{x}) \odot (d\mathbf{x}) \\
 &= \text{diag}(2\mathbf{x})(d\mathbf{x}) + (d\mathbf{x}) \odot (d\mathbf{x})
 \end{aligned}$$

Differential:

$$df(\mathbf{x}; d\mathbf{x}) = \text{diag}(2\mathbf{x})(d\mathbf{x})$$

Gradient:

$$\nabla f(\mathbf{x}) = \text{diag}(2\mathbf{x})$$

Verification:

derivative of w.r.t.

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x} \odot \mathbf{x}) = 2 \cdot \text{diag}(\mathbf{x})$$

4 My second example

Function:

$$f(\mathbf{x}) = \mathbf{x}\mathbf{x}^T\mathbf{x}$$

Difference:

$$\begin{aligned}
 &f(\mathbf{x} + d\mathbf{x}) - f(\mathbf{x}) \\
 &= (\mathbf{x} + d\mathbf{x})(\mathbf{x} + d\mathbf{x})^T(\mathbf{x} + d\mathbf{x}) - \mathbf{x}\mathbf{x}^T\mathbf{x} \\
 &= (\mathbf{x} + d\mathbf{x})(\mathbf{x}^T + (d\mathbf{x})^T)(\mathbf{x} + d\mathbf{x}) - \mathbf{x}\mathbf{x}^T\mathbf{x} \\
 &= (\mathbf{x}\mathbf{x}^T + \mathbf{x}(d\mathbf{x})^T + (d\mathbf{x})\mathbf{x}^T + (d\mathbf{x})(d\mathbf{x})^T)(\mathbf{x} + d\mathbf{x}) - \mathbf{x}\mathbf{x}^T\mathbf{x} \\
 &= \mathbf{x}(d\mathbf{x})^T\mathbf{x} + (d\mathbf{x})\mathbf{x}^T\mathbf{x} + (d\mathbf{x})(d\mathbf{x})^T\mathbf{x} + \mathbf{x}\mathbf{x}^T(d\mathbf{x}) + \mathbf{x}(d\mathbf{x})^T(d\mathbf{x}) + (d\mathbf{x})\mathbf{x}^T(d\mathbf{x}) + (d\mathbf{x})(d\mathbf{x})^T(d\mathbf{x})
 \end{aligned}$$

Differential (the first line is actually what the product rule says):

$$\begin{aligned}
 df &= \mathbf{x}(d\mathbf{x})^T\mathbf{x} + (d\mathbf{x})\mathbf{x}^T\mathbf{x} + \mathbf{x}\mathbf{x}^T(d\mathbf{x}) \\
 &= \mathbf{x}\mathbf{x}^T(d\mathbf{x}) + \mathbb{I}(d\mathbf{x})\mathbf{x}^T\mathbf{x} + \mathbf{x}\mathbf{x}^T(d\mathbf{x}) \\
 &= (2\mathbf{x}\mathbf{x}^T + (\mathbf{x}^T\mathbf{x})\mathbb{I})(d\mathbf{x})
 \end{aligned}$$

Gradient:

$$\nabla f = 2\mathbf{x}\mathbf{x}^T + (\mathbf{x}^T\mathbf{x})\mathbb{I}$$

Verification:

derivative of

$\mathbf{x} * \mathbf{x}' * \mathbf{x}$

w.r.t.

\mathbf{x}



$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x} \cdot \mathbf{x}^\top \cdot \mathbf{x}) = \mathbf{x}^\top \cdot \mathbf{x} \cdot \mathbb{I} + 2 \cdot \mathbf{x} \cdot \mathbf{x}^\top$$

5 Chain rule: derivative of a composition of vector-to-vector functions

Theorem (chain rule). Let $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $g: \mathbb{R}^n \rightarrow \mathbb{R}^p$ be differentiable. Then the composition function $h: \mathbb{R}^m \rightarrow \mathbb{R}^p$ defined by $h(\mathbf{x}) = g(f(\mathbf{x}))$ is differentiable. Further, the Jacobian of h is

$$Dh(\mathbf{x}) = (Dg(\mathbf{y}))(Df(\mathbf{x})),$$

where $\mathbf{x} \in \mathbb{R}^m$ and $\mathbf{y} = f(\mathbf{x}) \in \mathbb{R}^n$.

Practical note. Suppose we need to multiply two Jacobians together. Their shapes are $(m \times q)$ and $(q \times p)$. There would be mp dot products of length q , which involves about mpq scalar operations in total. Now consider multiplying three Jacobians together. Their shapes are $(1 \times n)$, $(n \times n)$ and $(n \times n)$. This can be done in two opposite orders, each offering a different cost:

$$\underbrace{\underbrace{(1 \times n), (n \times n)}_{(1,n) \text{ with cost } n^2}, (n \times n)}_{(1,p) \text{ with cost } n^2} \quad \text{vs.} \quad (1 \times n), \underbrace{\underbrace{(n \times n), (n \times n)}_{(n,n) \text{ with cost } n^3}}_{(1,p) \text{ with cost } n^2}.$$

Clearly, the first way is cheaper, so order does matter here.

Theorem (Cauchy's rule of invariance). (Using the setup from the theorem above)

$$dh(\mathbf{x}; d\mathbf{x}) = dg(\mathbf{y}; df(\mathbf{x}; d\mathbf{x})).$$

Proof. From the theorem above, it follows that

$$\begin{aligned} dh(\mathbf{x}; d\mathbf{x}) &= Dh(\mathbf{x}) d\mathbf{x} \\ &= (Dg(\mathbf{y}))(Df(\mathbf{x})) d\mathbf{x} \\ &= (Dg(\mathbf{y})) \underbrace{df(\mathbf{x}; d\mathbf{x})}_{d\mathbf{y}} \\ &= dg(\mathbf{y}; df(\mathbf{x}; d\mathbf{x})). \end{aligned}$$