



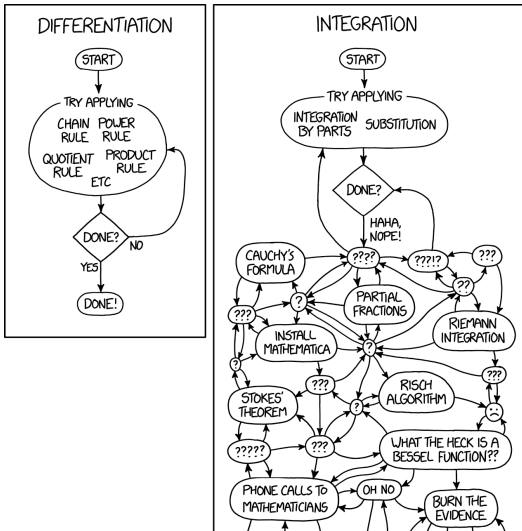
# Matrix Differentiation

Wen Perng

NTUEE

2024 September 18

## Derivatives are easy! (Source: XKCD)



In machine learning and many other studies, we often want to optimize functions:

$$\min_x f(x).$$

We know that **extremas** (max,min,saddle points) occur at places with zero derivatives:

$$\left. \frac{df(x)}{dx} \right|_{x^*} = 0,$$

whether it is a minima or not can be determined by a second order derivative test:

$$\left. \frac{d^2 f(x)}{dx^2} \right|_{x^*} > 0.$$

We denote the point of minima as

$$x^* = \arg \min_x f(x).$$

What if we want to minimize a function of two variables?

$$\min_{x_1, x_2} f(x_1, x_2)$$

We do the same test for finding extremas: by **partial derivatives** of first and second order:

$$\frac{\partial f}{\partial x_1} = 0, \quad \frac{\partial f}{\partial x_2} = 0, \quad \det \left( \frac{\partial^2 f}{\partial x_i \partial x_j} \right) > 0.$$

This requires the notions of **gradient** and **Hessian**, which will be discussed later.

But instead of writing the derivatives out term-by-term, is there a

1. coordinate free,
2. compact way

of writing out the derivatives? I.e., we would want to justify the notations of

$$\frac{df(A)}{dA} \text{ and } \frac{df(\boldsymbol{x})}{d\boldsymbol{x}},$$

where  $A$  is a matrix and  $\boldsymbol{x}$  is a vector.

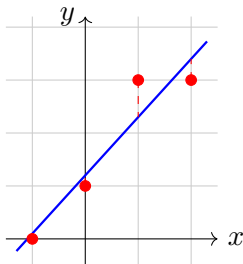
What do the derivatives mean and how do they behave? Note that we will work in cases where matrices are treated as *matrices* instead of *arrays*.

So... what kind of functions will we encounter? A prime example is **regression** by minimizing the **squared error**:

$$\mathcal{E} = ||X\mathbf{w} - \mathbf{y}||^2,$$

where  $X \in \mathbb{R}^{m \times n}$ ,  $\mathbf{w} \in \mathbb{R}^n$ ,  $\mathbf{y} \in \mathbb{R}^m$ . We would like to find a  $\mathbf{w}$  that *best estimates*  $\mathbf{y}$  by the approximation:

$$\hat{\mathbf{y}} = X\mathbf{w}.$$

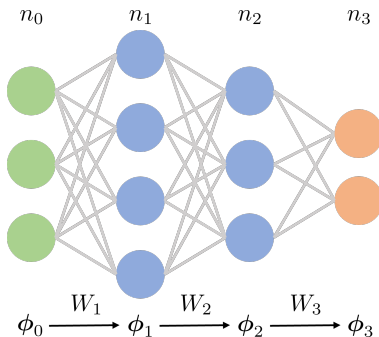


$$X = [1, \mathbf{x}]$$

$$\mathbf{w} = [a_0, a_1]^T$$

$$\hat{y} = a_0 + a_1x$$

In machine learning, when we want to calculate the backpropagation of a multilayer perceptron (MLP), we also need to learn how to differentiate with respect to matrix variables.



$$f(W, \phi) := \sigma(W^T \phi) \quad (\text{perceptron})$$

$$\phi_3 = f(W_3, \cdot) \circ f(W_2, \cdot) \circ f(W_1, \phi_0)$$

Lastly, others might want to differentiate functions such as

$$f = \mathbf{x}^T A \mathbf{x}, \det(A), \operatorname{tr}(A), \lambda_i(A), \dots$$

Or even crazier, the QR-decomposition:

$$\begin{aligned} X(t) &= Q(t)R(t) \\ \Rightarrow \frac{dQ(t)}{dt} &\stackrel{?}{=} \frac{\partial Q}{\partial X} \cdot \frac{dX}{dt}. \end{aligned}$$

This last part will be introduced in the next course. So stay tuned.



## Def. (Transpose)

Flip the matrix along its diagonal:

$$\begin{bmatrix} \textcolor{red}{a}_{11} & a_{12} & a_{13} \\ a_{21} & \textcolor{red}{a}_{22} & a_{23} \end{bmatrix}^T = \begin{bmatrix} \textcolor{red}{a}_{11} & a_{21} \\ a_{12} & \textcolor{red}{a}_{22} \\ a_{13} & a_{23} \end{bmatrix}. \quad (1)$$

What number when transposed doubles?

$$\text{午}^T = \text{什}$$

## Def. (Inner Product)

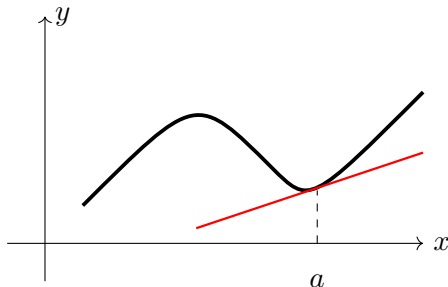
We denote the *Euclidean* inner product of two column vectors  $\mathbf{x}$  and  $\mathbf{y}$  as:

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = \langle \mathbf{x}, \mathbf{y} \rangle.$$

- 1 Differentiation Revisited
- 2 Matrix Differentiation
- 3 Second Derivative
- 4 Examples of Matrix Derivatives

- 1 Differentiation Revisited
  - What is Derivative
  - Derivative as Linear Operator
- 2 Matrix Differentiation
- 3 Second Derivative
- 4 Examples of Matrix Derivatives

Derivatives are linearizations. Consider the easiest case below.



We know that the first-order approximation will be

$$y(x) = \left. \frac{dy}{dx} \right|_a (x - a) + y(a).$$

We can rewrite the previous equation into the form below:

$$y(a + dx) = \left. \frac{dy}{dx} \right|_a dx + y(a)$$

$$dy = y(a + dx) - y(a) = \left. \frac{dy}{dx} \right|_a dx$$

Thus obtaining:

$$dy = \left( \frac{dy}{dx} \right) dx \quad (2)$$

We don't need to "divide" the  $dx$  over to the other side. The equation above will generalize to future cases.

## Prop. (Derivatives as Linear Approximation)

For  $y = f(x)$ ,

$$dy = \frac{dy}{dx} dx.$$

The above definition holds for  $x$  and  $y$  being scalars, vectors or matrices so long as the dimensions are compatible.

e.g.

$$f(\mathbf{x}) = A\mathbf{x}, \quad A \in \mathbb{R}^{m \times n}, \quad \mathbf{x} \in \mathbb{R}^{n \times 1}$$

By plugging in the definition above:

$$d\mathbf{f} = \mathbf{f}(\mathbf{x} + d\mathbf{x}) - \mathbf{f}(\mathbf{x}) = A(\mathbf{x} + d\mathbf{x}) - A\mathbf{x} = A d\mathbf{x}$$

$$\frac{d\mathbf{f}}{d\mathbf{x}} = A$$

## Prop. (Derivatives as Linear Approximation)

For  $y = f(x)$ ,

$$dy = \frac{dy}{dx} dx.$$

The above definition holds for  $x$  and  $y$  being scalars, vectors or matrices so long as the dimensions are compatible.

e.g.

$$\mathbf{f}(\mathbf{A}) = \mathbf{A}\mathbf{x}, \quad \mathbf{A} \in \mathbb{R}^{m \times n}, \quad \mathbf{x} \in \mathbb{R}^{n \times 1}$$

By plugging in the definition above:

$$d\mathbf{f} = \mathbf{f}(\mathbf{A} + d\mathbf{A}) - \mathbf{f}(\mathbf{A}) = (\mathbf{A} + d\mathbf{A})\mathbf{x} - \mathbf{A}\mathbf{x} = d\mathbf{A} \mathbf{x}$$

$$\frac{d\mathbf{f}}{d\mathbf{A}} \stackrel{?}{=} \mathbf{x}$$

Perhaps we shouldn't stick to the *left-multiplications*, but simply view derivatives as a **linear operator**:

$$dy = \mathcal{L}\{dx\}. \quad (3)$$

For  $y = f(x)$ , the differentiation linear operator can have many forms:

$$\begin{aligned}\mathcal{L}\{dx\} &= f'dx \\ \mathcal{L}\{dx\} &= dx f' \\ \mathcal{L}\{dx\} &= \sqrt{f'}dx\sqrt{f'}.\end{aligned}$$

All works as long as it is **linear**:

$$\mathcal{L}\{adx_1 + dx_2\} = a\mathcal{L}\{dx_1\} + \mathcal{L}\{dx_2\}.$$



### Prop. (Derivatives as Linear Operators (Fréchet Derivative))

For  $y = f(x)$ ,

$$dy = \mathcal{L}\{dx\}.$$

Where  $\mathcal{L}\{\cdot\}$  is a linear operator, denoting the notion of a derivative.

The derivative is a **linear operator (function)** that maps the changes in the variable to the first-order changes in the function. <sup>1</sup>

---

<sup>1</sup>Since all linear operators can be represented by matrices, we can, still, represent derivatives as matrices and vectors alike. But it is often not needed and might further complicate the results.

- 1 Differentiation Revisited
- 2 Matrix Differentiation
  - How to Calculate ■ First Derivative
- 3 Second Derivative
- 4 Examples of Matrix Derivatives

Two methods are possible:

1. Term-wise differentiation,
2. **Directional derivatives** (if it exists): the directional derivative of  $f(x)$  by  $x$  at  $x_0$  in the direction of  $h$  will be

$$df = \frac{\partial f(x)}{\partial x}(x_0)[h] := \left. \frac{d}{dt} f(x_0 + th) \right|_{t=0}, \quad (4)$$

or equivalently, find a linear  $\mathcal{L}$  such that

$$f(x) = f(x_0) + \mathcal{L}\{x - x_0\} + o(\|x - x_0\|). \quad (5)$$

The first method is quite boring, we will dive into the second method and some examples in the next page.

(1) Differentiate scalar function  $f$  by vector  $\mathbf{x} \in \mathbb{R}^n$ :

Let us suppose

$$f(\mathbf{x}) = \mathbf{x}^\top A \mathbf{x}, \text{ where } A \text{ is symmetric.}$$

[Sol.]

$$\begin{aligned} \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x})[\mathbf{h}] &= \left. \frac{d}{dt} ((\mathbf{x} + t\mathbf{h})^\top A (\mathbf{x} + t\mathbf{h})) \right|_{t=0} \\ &= \mathbf{h}^\top A \mathbf{x} + \mathbf{x}^\top A \mathbf{h} \\ &= \mathbf{2x}^\top A \mathbf{h}. \end{aligned}$$

Since the final result is of the **matrix-vector multiplication form**, we can have the following statement:

$$\frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}) \equiv \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \equiv \frac{\partial f}{\partial \mathbf{x}} = \mathbf{2x}^\top A. \quad (6)$$

Further, given a function  $f(\mathbf{x}) = f(x_1, \dots, x_n)$  with

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = [x_1 \quad \cdots \quad x_n]^T \in \mathbb{R}^{n \times 1},$$

we know from calculus that its **infinitesimal** change is:

$$df = \frac{\partial f}{\partial x_1} dx_1 + \cdots + \frac{\partial f}{\partial x_n} dx_n. \quad (7)$$

This result is obtained by term-wise differentiation. Intimidating?  
But this is just the dot product of two vectors!

Notice that

$$df = \frac{\partial f}{\partial x_1} dx_1 + \cdots + \frac{\partial f}{\partial x_n} dx_n$$

is just the dot product of  $\nabla f$  and  $d\mathbf{x}$ :

$$df = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix} \begin{bmatrix} dx_1 \\ \vdots \\ dx_n \end{bmatrix} = (\nabla f)^\top d\mathbf{x}.$$

### Def. (Gradient)

Given an inner product  $\langle \cdot, \cdot \rangle$ , the **gradient** of a function  $f(\mathbf{x})$  ( $\mathbf{x} \in \mathbb{R}^n$ ) is defined as

$$\nabla f(\mathbf{x}) \in \mathbb{R}^n \text{ s.t. } df = \langle \nabla f(\mathbf{x}), \mathbf{h} \rangle. \quad (8)$$

For  $f = f(\mathbf{x})$ , if we have

$$(\nabla f)^\top d\mathbf{x} = df = \frac{\partial f}{\partial \mathbf{x}} d\mathbf{x}$$

as operation by *left-multiplication*, then

$$\frac{\partial f}{\partial \mathbf{x}} = \left[ \frac{\partial f}{\partial x_1} \quad \cdots \quad \frac{\partial f}{\partial x_n} \right] = (\nabla f)^\top \equiv \nabla^\top f.$$

And if we define

$$df = d\mathbf{x}^\top \frac{\partial f}{\partial \mathbf{x}}$$

by *right-multiplication*, then

$$\frac{\partial f}{\partial \mathbf{x}^\top} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \nabla f.$$

Back to our example above:

$$f(\mathbf{x}) = \mathbf{x}^\top A \mathbf{x}$$

[Sol.]

$$df = \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x})[\mathbf{h}] = 2\mathbf{x}^\top A \mathbf{h} = \langle 2A\mathbf{x}, \mathbf{h} \rangle$$

$$\nabla f(\mathbf{x}) = 2A\mathbf{x}.$$



Notice that:

$$\begin{aligned} 1 \times n &\rightarrow \frac{\partial f}{\partial \mathbf{x}} \leftarrow 1 \times 1 \\ &\quad \quad \quad \leftarrow n \times 1 \\ n \times 1 &\rightarrow \frac{\partial f}{\partial \mathbf{x}^\top} \leftarrow 1 \times 1 \\ &\quad \quad \quad \leftarrow 1 \times n \end{aligned}$$

We can even have:

$$m \times n \rightarrow \frac{\partial f}{\partial \mathbf{x}} \leftarrow m \times 1 \\ \quad \quad \quad \leftarrow n \times 1$$

But remember, this is just a convention. And the main takeaway is that **derivatives are linear operators**.

(2) Differentiate vector function  $\mathbf{f}$  by another vector  $\mathbf{x}$ :  
If  $\mathbf{x}$  is a column vector and  $\mathbf{f}$  is of the form:

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

then we have:

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial f_m}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n},$$

Some terms and notations to introduce: for a scalar function  $f$  and position vector  $\mathbf{x} = [x_1, \dots, x_n]^\top$ ,

1. Gradient (transpose):

$$\nabla f := \frac{\partial f}{\partial \mathbf{x}^\top} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}, \quad \nabla^\top f := \frac{\partial f}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_n} \end{bmatrix} \quad (9)$$

e.g.  $f(\mathbf{x}) = f(x_1, x_2) = x_1^2 x_2$

Some terms to introduce: for a scalar function  $f$  and position vector  $\mathbf{x} = [x_1, \dots, x_n]^T$ ,

2. Hessian:

$$\mathbf{H}_f := \nabla \nabla^T f = \frac{\partial^2 f}{\partial \mathbf{x}^T \partial \mathbf{x}} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{bmatrix} \quad (10)$$

e.g.  $f(\mathbf{x}) = f(x_1, x_2) = x_1^2 x_2$

For a vector function  $\mathbf{y}(\mathbf{x}) = [y_1(\mathbf{x}), \dots, y_n(\mathbf{x})]^\top$  and position vector  $\mathbf{x} = [x_1, \dots, x_n]^\top$ ,

3. Jacobian:

$$J := \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial y_n}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_n} \end{bmatrix}, \quad (11)$$

The Jacobian is useful in coordinate transformations, like the change of variables in integration.

$$\int \dots \int_{\mathbf{y}(\mathcal{X})} f(\mathbf{y}) dy_1 \dots dy_n = \int \dots \int_{\mathcal{X}} f(\mathbf{y}(\mathbf{x})) |J| dx_1 \dots dx_n.$$

(3) Differentiate scalar function  $f(X)$  by a matrix  $X \in \mathbb{R}^{m \times n}$ :

$$df = \frac{\partial f}{\partial X}(X_0)[H] = \left. \frac{d}{dt} f(X_0 + tH) \right|_{t=0}.$$

Some examples are as follows:

1.  $f(A) = Ax$

$$df = (dA)x$$

2.  $f(A) = A^T A$

$$df = (A + dA)^T (A + dA) - A^T A = (dA)^T A + A^T dA$$

3.  $f(X) = \text{Tr} \{X^T A X B\}$ , with symmetric  $A$  and  $B$ .

### Def. (Trace)

The trace of a square matrix is the sum of its diagonals.

$$\text{Tr} \{A\} = \sum_{i=1}^n a_{ii}$$

[Sol.]

$$\begin{aligned} \frac{\partial f}{\partial X}[H] &= \left. \frac{d}{dt} \text{Tr} \{ (X + tH)^T A (X + tH) B \} \right|_{t=0} \\ &= 2 \text{Tr} \{ H^T A X B \} \end{aligned}$$

For  $X \in \mathbb{R}^{m \times n}$ ,  $f(X) = f(x_{11}, \dots, x_{mn})$ , we have its **infinitesimal** change as:

$$\begin{aligned}
 df &= \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial x_{ij}} dx_{ij} \\
 &= \text{Tr} \left\{ \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \frac{\partial f}{\partial x_{21}} & \cdots & \frac{\partial f}{\partial x_{m1}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{1n}} & \frac{\partial f}{\partial x_{2n}} & \cdots & \frac{\partial f}{\partial x_{mn}} \end{bmatrix} \begin{bmatrix} dx_{11} & \cdots & dx_{1n} \\ dx_{21} & \cdots & dx_{2n} \\ \vdots & \ddots & \vdots \\ dx_{m1} & \cdots & dx_{mn} \end{bmatrix} \right\} \\
 &= \frac{\partial f}{\partial X} [dX].
 \end{aligned} \tag{12}$$

Hence, it is often denoted that for  $X \in \mathbb{R}^{m \times n}$ ,  $\frac{\partial f}{\partial X} \in \mathbb{R}^{n \times m}$ , acting on the differential change by the **Frobenius inner product**.



### Def. (Inner Product)

An inner product  $\langle \cdot, \cdot \rangle$  on a vector space  $V$  over the field  $\mathbb{R}$  is a map  $V \times V \rightarrow \mathbb{R}$  that satisfies: for  $x, y, z \in V$  and  $\lambda, \mu \in \mathbb{R}$

1. Symmetry:  $\langle x, y \rangle = \langle y, x \rangle$ ,
2. Linearity:  $\langle \lambda x + \mu y, z \rangle = \lambda \langle x, z \rangle + \mu \langle y, z \rangle$ ,
3. Positive-definiteness:  $\langle x, x \rangle > 0$  for all  $x \neq 0$

### Def. (Frobenius Inner Product)

For real matrices in  $\mathbb{R}^{m \times n}$ , we have the following inner product on them:

$$\langle X, Y \rangle = \text{Tr} \{ X^T Y \}. \quad (13)$$

$$\begin{aligned}
 df &= \text{Tr} \left\{ \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \frac{\partial f}{\partial x_{21}} & \cdots & \frac{\partial f}{\partial x_{m1}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{1n}} & \frac{\partial f}{\partial x_{2n}} & \cdots & \frac{\partial f}{\partial x_{mn}} \end{bmatrix} \begin{bmatrix} dx_{11} & \cdots & dx_{1n} \\ dx_{21} & \cdots & dx_{2n} \\ \vdots & \ddots & \vdots \\ dx_{m1} & \cdots & dx_{mn} \end{bmatrix} \right\} \\
 &= \text{Tr} \{ (\nabla f)^\top dX \} = \langle \nabla f, X \rangle \\
 \rightarrow \nabla f(X) &= \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \cdots & \frac{\partial f}{\partial x_{1n}} \\ \frac{\partial f}{\partial x_{21}} & \cdots & \frac{\partial f}{\partial x_{2n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{m1}} & \cdots & \frac{\partial f}{\partial x_{mn}} \end{bmatrix}
 \end{aligned} \tag{14}$$

Let us look back at the example: for  $f(X) = \text{Tr} \{ X^\top A X B \}$ ,

$$\frac{\partial f}{\partial X} [H] = 2 \text{Tr} \{ H^\top A X B \}, \quad \nabla f(X) = 2 A X B.$$

Remember that the following properties of derivatives still applies:

1. Chain Rule: for  $f = f(y)$  and  $y = g(x)$

$$\begin{aligned}\frac{\partial}{\partial x}(f \circ g)[h] &= \frac{\partial f}{\partial y}(g(x)) \left[ \frac{\partial g}{\partial x}[h] \right] \equiv \frac{\partial f}{\partial y} \Big|_{y=g(x)} \circ \frac{\partial g}{\partial x}[h] \\ &\equiv \frac{\partial f}{\partial y} \circ \frac{\partial g}{\partial x}[h]\end{aligned}\tag{15}$$

2. Product Rule:

$$\frac{\partial}{\partial x}(f(x) \cdot g(x)) [h] = \frac{\partial f}{\partial x}[h] \cdot g(x) + f(x) \cdot \frac{\partial g}{\partial x}[h] \tag{16}$$

- 1 Differentiation Revisited
- 2 Matrix Differentiation
- 3 Second Derivative
  - Hessian ■ LLMSE Problem
- 4 Examples of Matrix Derivatives

For a given function  $f : \mathbb{R} \rightarrow \mathbb{R}$ , we can find its second-order approximation as

$$f(x+h) = f(x) + \frac{\partial f}{\partial x}(x) \cdot h + \frac{1}{2} \frac{\partial^2 f}{\partial x^2} h^2 + \dots$$

But if we have  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , then

$$\begin{aligned} f(\mathbf{x} + \mathbf{h}) &= f(\mathbf{x}) + \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x})[\mathbf{h}] + \frac{1}{2} \frac{\partial^2 f}{\partial \mathbf{x}^2}(\mathbf{x})[\mathbf{h}] + \dots \\ &= f(\mathbf{x}) + (\nabla f(\mathbf{x}))^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \mathbf{H}_f(\mathbf{x}) \mathbf{h} + \dots, \end{aligned} \tag{17}$$

where the matrix  $\mathbf{H}_f(\mathbf{x})$  is called the **Hessian matrix** of  $f$  at  $\mathbf{x}$ .

We have seen that the Hessian matrix  $\mathbf{H}_f = \frac{\partial^2 f}{\partial \mathbf{x}^\top \partial \mathbf{x}}$  determines the behavior of the extrema.

### Def. (Definiteness)

A symmetric matrix  $A$  is called **positive definite** if for any non-zero vector  $x$  it satisfies:

$$x^T A x > 0 \iff A \succ 0.$$

Moreover, we have:

$$(\text{positive semi-definite}) \quad x^T A x \geq 0 \iff A \succeq 0$$

$$(\text{negative definite}) \quad x^T A x < 0 \iff A \prec 0$$

$$(\text{negative semi-definite}) \quad x^T A x \leq 0 \iff A \preceq 0$$

If none of the above are satisfied, then the matrix is **indefinite**.

Positive definite  $\Rightarrow$  minima; negative definite  $\Rightarrow$  maxima.

And we can make the series expansion even more general to have  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ , then

$$\begin{aligned} f(X + H) &= f(X) + \frac{\partial f}{\partial X}(X)[H] + \frac{1}{2} \frac{\partial^2 f}{\partial X^2}(X)[H] + \dots \\ &= f(X) + \langle \nabla f(X), H \rangle + \frac{1}{2} \langle \mathbf{H}_f(X)[H], H \rangle + \dots \end{aligned} \tag{18}$$

The linear function  $\mathbf{H}_f(X)$  that maps  $H \in \mathbb{R}^{m \times n}$  to  $\mathbb{R}^{m \times n}$  is called the **Hessian operator**.

Like the gradient, the definition of the Hessian also relies on the inner product.

### Def. (Gradient)

Given an inner product  $\langle \cdot, \cdot \rangle$ , the **gradient** of a function  $f(x)$  ( $x \in \mathbb{R}, \mathbb{R}^n$  or  $\mathbb{R}^{m \times n}$ ) is defined as

$$\nabla f(x) \text{ s.t. } \frac{\partial f}{\partial x}(x)[h] = \langle \nabla f(x), h \rangle. \quad (19)$$

### Def. (Hessian)

Given an inner product  $\langle \cdot, \cdot \rangle$ , the **Hessian** of a function  $f(x)$  ( $x \in \mathbb{R}, \mathbb{R}^n$ , or  $\mathbb{R}^{m \times n}$ ) is defined as

$$\mathbf{H}_f(x) \text{ s.t. } \frac{\partial^2 f}{\partial x^2}(x)[h] = \langle \mathbf{H}_f(x)[h], h \rangle. \quad (20)$$



So how do we calculate the second derivative?

The second derivative can be derived via directional derivatives:

$$\frac{\partial^2 f}{\partial x^2}[h, k] = \left. \frac{d}{dt} \frac{d}{ds} f(x + th + sk) \right|_{t=0, s=0}. \quad (21)$$

Or gaining idea from Taylor series, one can first calculate

$$\frac{\partial^2 f}{\partial x^2}[h] = \left. \frac{d^2}{dt^2} f(x + th) \right|_{t=0}, \quad (22)$$

then utilize the **polarization identity**:

$$\frac{\partial^2 f}{\partial x^2}[h, k] = \frac{1}{4} \left( \frac{\partial^2 f}{\partial x^2}[h + k] - \frac{\partial^2 f}{\partial x^2}[h - k] \right). \quad (23)$$

An example is as follows:  $f(X) = \text{Tr} \{X^\top A X B\}$ , with  $A, B$  symmetric.

[Sol.]

$$\frac{\partial f}{\partial X}[H] = 2\text{Tr} \{H^\top A X B\}.$$

$$\begin{aligned}\frac{\partial^2 f}{\partial X^2}[H, K] &= \left. \frac{d}{dt} \frac{d}{ds} f(X + tH + sK) \right|_{t=0, s=0} \\ &= \left. \frac{d}{ds} \left( \frac{\partial f}{\partial X}(X + sK)[H] \right) \right|_{s=0} \\ &= \left. \frac{d}{ds} 2\text{Tr} \{H^\top A(X + sK)B\} \right|_{s=0} \\ &= \text{Tr} \{H^\top A K B\} = \text{Tr} \{K^\top A H B\} \\ \Rightarrow H_f(X)[H] &= AHB.\end{aligned}$$

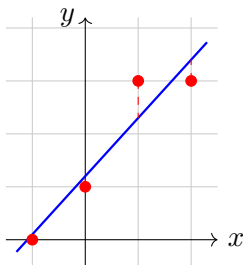
Let us also try using the polarization identity.

[Sol.]

$$\begin{aligned}\frac{\partial^2 f}{\partial X^2}[H] &= \frac{d^2}{dt^2} \text{Tr} \{ (X + tH)^\top A (X + tH) B \} \Big|_{t=0} \\ &= 2 \text{Tr} \{ H^\top A H B \} \\ \frac{\partial^2 f}{\partial X^2}[H, K] &= \frac{1}{4} (2 \text{Tr} \{ (H + K)^\top A (H + K) B \} \\ &\quad - 2 \text{Tr} \{ (H - K)^\top A (H - K) B \}) \\ &= 2 \text{Tr} \{ H^\top A K B \}\end{aligned}$$

Thus, we can see that the results are identical.

Now we turn to the least mean square problem as another example.



Suppose the regression line for the data points is

$$\hat{y} = a_0 + a_1x,$$

we can record the relationship as follows

$$\underbrace{\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ 1 & x_4 \end{bmatrix}}_X \underbrace{\begin{bmatrix} a_0 \\ a_1 \end{bmatrix}}_w \approx \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}}_y$$

**Def. (LLMSE problem)**

The problem of finding the *linear least mean square estimate* is stated as below: given measurements  $\mathbf{y}$  over sample points  $X$ , find the optimal coefficients (weights)  $\mathbf{w}$  that gives the estimate

$$\hat{\mathbf{y}} = X\mathbf{w},$$

such that the mean square error (variance)

$$\mathcal{E} = \|\hat{\mathbf{y}} - \mathbf{y}\|^2$$

is minimized.

It is often used as a cost function in filtering and machine learning. For our talk, we will be focusing on the optimization problem of:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|X\mathbf{w} - \mathbf{y}\|^2.$$

By our knowledge of extrema occurs at stationary points, we know that

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} ||X\mathbf{w} - \mathbf{y}||^2$$

occurs when the directional derivative is always zero at that point, i.e.

$$\left( \frac{\partial}{\partial \mathbf{w}} ||X\mathbf{w} - \mathbf{y}||^2 \right) \bigg|_{\mathbf{w}^*} [\mathbf{h}] = 0 \quad (\forall \mathbf{h}),$$

or simply, the derivative is zero:

$$\left( \frac{\partial}{\partial \mathbf{w}} ||X\mathbf{w} - \mathbf{y}||^2 \right) \bigg|_{\mathbf{w}^*} = 0.$$

[Sol.]

$$\begin{aligned}\frac{\partial}{\partial \mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 &= \frac{\partial}{\partial \mathbf{w}} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= 2(\mathbf{X}\mathbf{w} - \mathbf{y})^\top \mathbf{X}\end{aligned}$$

$$(\mathbf{X}\mathbf{w}^* - \mathbf{y})^\top \mathbf{X} = 0$$

$$\mathbf{X}^\top \mathbf{X}\mathbf{w}^* = \mathbf{X}^\top \mathbf{y}$$

$$\rightarrow \mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad \blacksquare$$

The term  $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  is coined the **matrix pseudo-inverse**.

Remark: What is the criterion for the existence of the pseudo-inverse?

We've only checked that  $\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}$  is an extrema, but we have yet to check whether its a maxima or a minima. A second derivative test is needed:

$$\frac{\partial}{\partial \mathbf{w}} \|X\mathbf{w} - \mathbf{y}\|^2 = 2(X\mathbf{w} - \mathbf{y})^T X$$

$$\frac{\partial^2}{\partial \mathbf{w}^T \partial \mathbf{w}} \|X\mathbf{w} - \mathbf{y}\|^2 = \frac{\partial}{\partial \mathbf{w}^T} 2(X\mathbf{w} - \mathbf{y})^T X = 2X^T X$$

Hence we know that the second derivative is **positive definite**, i.e. for all  $d\mathbf{w} \neq 0$ ,

$$\|X(\mathbf{w}^* + d\mathbf{w}) - \mathbf{y}\|^2 - \|X\mathbf{w}^* - \mathbf{y}\|^2 = (d\mathbf{w})^T (2X^T X) d\mathbf{w} > 0,$$

it is therefore a minima.



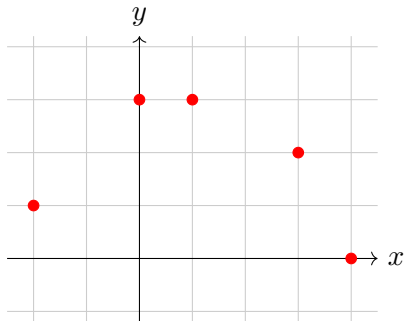
Given a set of data points  $(x, y)$ , find a function that interpolates them with the least mean square error.

$x$	$y$
-2	1
0	3
1	3
3	2
4	0

Find a quadratic:

$$\hat{y} = a_0 + a_1x + a_2x^2$$

such that the mean square error is minimized.

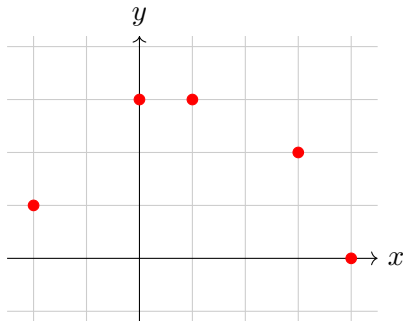


Given a set of data points  $(x, y)$ , find a function that interpolates them with the least mean square error.

We can rewrite the estimation equation as:

$$\hat{\mathbf{y}} = \begin{bmatrix} 1 & -2 & 4 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = X\mathbf{w}$$

$$\mathbf{y} = [1, 3, 3, 2, 0]^T$$

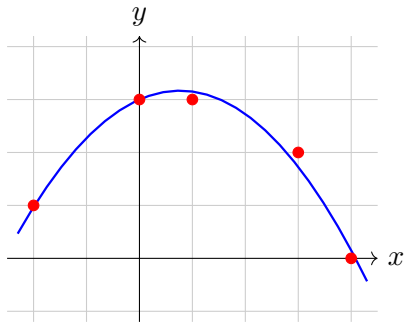


Given a set of data points  $(x, y)$ , find a function that interpolates them with the least mean square error.

$$(X^T X)^{-1} X^T = \begin{bmatrix} 0.14 & 0.43 & 0.43 & 0.14 & -0.14 \\ -0.28 & 0.047 & 0.12 & 0.10 & 0.00 \\ 0.06 & -0.05 & -0.06 & -0.01 & 0.06 \end{bmatrix}$$

$$w^* = \begin{bmatrix} 3.0000 \\ 0.4394 \\ -0.2879 \end{bmatrix}$$

$$\hat{y} = 3 + 0.44x - 0.29x^2$$



- 1 Differentiation Revisited
- 2 Matrix Differentiation
- 3 Second Derivative
- 4 Examples of Matrix Derivatives
  - Derivatives of Matrix Functions ■ MLP

Given  $A = A(t)$ , what is

$$\frac{dA^{-1}}{dt}$$

expressed using  $\frac{dA}{dt}$ ?

[Sol.] Since we have

$$AA^{-1} = \mathbb{1},$$

by product rule:

$$\begin{aligned} \frac{\partial}{\partial A} (AA^{-1}) [H] &= \cancel{\frac{\partial A}{\partial A}} [H] \cdot A^{-1} + A \cdot \frac{\partial A^{-1}}{\partial A} [H] = 0 \\ \frac{\partial A^{-1}}{\partial A} [H] &= -A^{-1} H A^{-1} \\ \Rightarrow \frac{dA^{-1}}{dt} &= \frac{\partial A^{-1}}{\partial A} \left[ \frac{dA}{dt} \right] = -A^{-1} \frac{dA}{dt} A^{-1} \quad \blacksquare \end{aligned}$$

### Lemma

The following identity holds for all square matrices  $A$ :

$$\det(e^A) = e^{\text{Tr}\{A\}}. \quad (24)$$

It can be immediately proven by Jordan canonical form of matrices.

e.g.  $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$

Hence if the matrix  $A$  is expressible as an exponential:

$$A = e^B,$$

i.e., it is invertible, then a solution is as below:

[Sol.]

$$\begin{aligned}\frac{d}{dt} \det(A(t)) &= \frac{d}{dt} e^{\text{Tr}\{B(t)\}} \\ &= e^{\text{Tr}\{B(t)\}} \text{Tr} \left\{ \frac{d}{dt} B(t) \right\} \\ &= \det(A(t)) \text{Tr} \left\{ A^{-1} \frac{dA(t)}{dt} \right\}. \quad \blacksquare\end{aligned}$$

### Def. (Eigenvalues and Eigenvectors)

For a given square matrix  $A$ , it has **eigenvalues**  $\{\lambda_i\}$  such that they satisfy:

$$A\mathbf{v}_i = \lambda_i\mathbf{v}_i,$$

where  $\mathbf{v}_i$  is the associated (right) **eigenvector**. And if  $\mathbf{u}_i$  satisfies:

$$\mathbf{u}_i^T A = \lambda_i \mathbf{u}_i^T,$$

then  $\mathbf{u}_i$  is the associated left eigenvector.



For a time varying  $A(t)$ , what is  $\frac{d\lambda_i(t)}{dt}$ ?

[Sol.]

Let us consider  $\lambda_i(t)$  with its associated right and left eigenvectors:  $\mathbf{v}_i(t)$  and  $\mathbf{u}_i(t)$  that has length satisfying:

$$\mathbf{u}_i^T \mathbf{v}_i = 1.$$

Then,

$$\mathbf{u}_i^T(t) A(t) \mathbf{v}_i(t) = \lambda_i(t)$$

$$\frac{d\lambda_i}{dt} = \frac{d\mathbf{u}_i^T}{dt} A \mathbf{v}_i + \mathbf{u}_i^T \frac{dA}{dt} \mathbf{v}_i + \mathbf{u}_i^T A \frac{d\mathbf{v}_i}{dt}$$

For a time varying  $A(t)$ , what is  $\frac{d\lambda_i(t)}{dt}$ ?

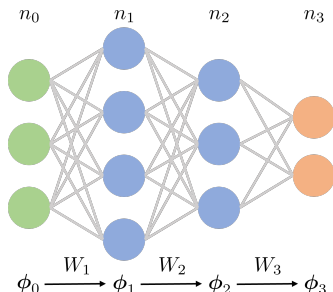
[Sol.]

$$\frac{d\lambda_i}{dt} = \frac{d\mathbf{u}_i^\top}{dt} A \mathbf{v}_i + \mathbf{u}_i^\top \frac{dA}{dt} \mathbf{v}_i + \mathbf{u}_i^\top A \frac{d\mathbf{v}_i}{dt}$$

But we also have

$$\begin{aligned} \frac{d\mathbf{u}_i^\top}{dt} A \mathbf{v}_i + \mathbf{u}_i^\top A \frac{d\mathbf{v}_i}{dt} &= \lambda_i \left( \frac{d\mathbf{u}_i^\top}{dt} \mathbf{v}_i + \mathbf{u}_i^\top \frac{d\mathbf{v}_i}{dt} \right) \\ &= \lambda_i \frac{d}{dt} (\mathbf{u}_i^\top \mathbf{v}_i) = 0 \end{aligned}$$

$$\rightarrow \frac{d\lambda_i}{dt} = \mathbf{u}_i^\top \frac{dA}{dt} \mathbf{v}_i$$



$$\phi = \begin{bmatrix} \phi_1 \\ \vdots \\ \phi_n \end{bmatrix} \rightarrow \left( \begin{array}{c|c} \mathbf{w}_i & \sigma \end{array} \right) \rightarrow \phi_i = \sigma(\mathbf{w}_i^T \phi)$$

For layer- $\ell$  with  $n_\ell$  perceptrons,

$$\mathbb{R}^{n_\ell} \ni \phi_\ell = \sigma(W_\ell^T \phi_{\ell-1}),$$

the weight matrix

$$W_\ell = [\mathbf{w}_{\ell,1} \cdots \mathbf{w}_{\ell,n_\ell}] \in \mathbb{R}^{n_\ell \times n_{\ell+1}}.$$

Hence if we denote  $\sigma(W^T \phi) = f(W, \phi)$ , the network above can be reduced to

$$\phi_3 = f(W_3, \cdot) \circ f(W_2, \cdot) \circ f(W_1, \phi_0).$$

Our goal is to do backpropagation: given a cost function  $\mathcal{E}(\phi_3)$ , find the optimal weights  $W_1, W_2, W_3$  such that the cost is minimized.

Hence, we need to first calculate the derivatives of  $f(W, \phi)$ .

- Differentiate w.r.t.  $W$ :

$$\begin{aligned} \frac{\partial f}{\partial W}(W_\ell, \phi_{\ell-1})[H] &= \begin{bmatrix} \frac{\partial}{\partial W} \sigma(W_{\ell,1}^\top \phi_{\ell-1})[H] \\ \vdots \end{bmatrix} = \begin{bmatrix} \sigma'(W_\ell^\top \phi_{\ell-1}) H_1^\top \phi \\ \vdots \end{bmatrix} \\ &= \begin{bmatrix} \sigma'(W_{\ell,1}^\top \phi_{\ell-1}) & 0 \\ 0 & \ddots \end{bmatrix} \begin{bmatrix} H_1^\top \phi \\ \vdots \end{bmatrix} = \text{diag}(\sigma'(W_\ell^\top \phi_{\ell-1})) H^\top \phi_{\ell-1}. \end{aligned}$$

- Differentiate w.r.t.  $\phi$ :

$$\frac{\partial f}{\partial \phi}(W_\ell, \phi_{\ell-1})[h] = \text{diag}(\sigma'(W_\ell^\top \phi_{\ell-1})) W_\ell^\top h.$$

Suppose we want to alter  $W_1$  by the methods of **gradient descent**. The gradient to  $W_1$  can be calculated:

$$\begin{aligned}
 \frac{\partial \mathcal{E}}{\partial W_1}(\phi_3)[H] &= \frac{\partial \mathcal{E}}{\partial \phi_3} \circ \left. \frac{\partial f}{\partial \phi_2} \right|_{(W_3, \phi_2)} \circ \left. \frac{\partial f}{\partial \phi_1} \right|_{(W_2, \phi_1)} \circ \left. \frac{\partial f}{\partial W_1} \right|_{(W_1, \phi_0)} [H] \\
 &= \frac{\partial \mathcal{E}}{\partial \phi_3} \left[ \text{diag}(\sigma'(W_3^T \phi_2)) W_3^T \cdot \text{diag}(\sigma'(W_2^T \phi_1)) W_2^T \cdot \text{diag}(\sigma'(W_1^T \phi_0)) H^T \phi_0 \right] \\
 &= \text{Tr} \left\{ H^T \phi_0 \frac{\partial \mathcal{E}}{\partial \phi_3} \text{diag}(\sigma'(W_3^T \phi_2)) W_3^T \text{diag}(\sigma'(W_2^T \phi_1)) W_2^T \text{diag}(\sigma'(W_1^T \phi_0)) \right\} \\
 \Rightarrow \nabla_{W_1} \mathcal{E} &= \phi_0 \frac{\partial \mathcal{E}}{\partial \phi_3} \text{diag}(\sigma'(W_3^T \phi_2)) W_3^T \text{diag}(\sigma'(W_2^T \phi_1)) W_2^T \text{diag}(\sigma'(W_1^T \phi_0))
 \end{aligned}$$

After each iteration,  $W_1$  is updated via

$$W_1 \mapsto W_1 - \alpha \nabla_{W_1} \mathcal{E},$$

with  $\alpha$  a suitable decreasing step size.

Notes: