

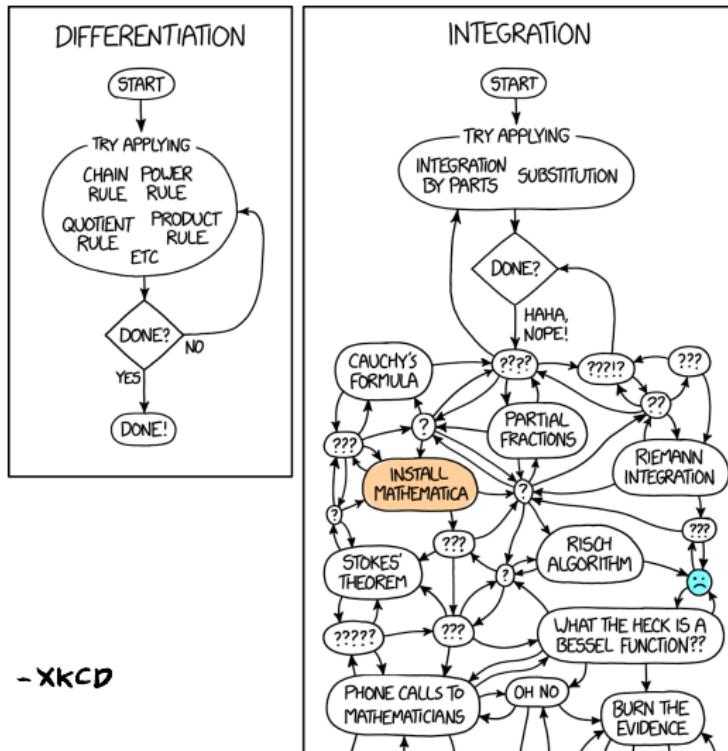
Matrix Differentiation

Wen Perng

Electrical Engineering, NTU

2023/11/02

Derivatives are easy!



- XKCD

Figure: XKCD: differentiation and integration

What do we want?

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:

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

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

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

We denote the point of minima as:

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

What do we want?

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

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

What do we want?

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 extrema: by **partial derivatives**

$$\frac{\partial f}{\partial x_1} = 0, \frac{\partial f}{\partial x_2} = 0.$$

What do we want?

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 extrema: by **partial derivatives**

$$\frac{\partial f}{\partial x_1} = 0, \frac{\partial f}{\partial x_2} = 0.$$

But if the function to minimize is a function of a **matrix** or a **vector**:

$$\min_{A \in \mathbb{R}^{m \times n}} f(A) \text{ or } \min_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x}),$$

What do we want?

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 extrema: by **partial derivatives**

$$\frac{\partial f}{\partial x_1} = 0, \frac{\partial f}{\partial x_2} = 0.$$

But if the function to minimize is a function of a **matrix** or a **vector**:

$$\min_{A \in \mathbb{R}^{m \times n}} f(A) \text{ or } \min_{x \in \mathbb{R}^n} f(x),$$

how do we find the **derivative of a matrix function**?

What do we want?

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 extrema: by **partial derivatives**

$$\frac{\partial f}{\partial x_1} = 0, \frac{\partial f}{\partial x_2} = 0.$$

But if the function to minimize is a function of a **matrix** or a **vector**:

$$\min_{A \in \mathbb{R}^{m \times n}} f(A) \text{ or } \min_{x \in \mathbb{R}^n} f(x),$$

how do we find the **derivative of a matrix function**? I.e., how do we justify the notation of:

$$\frac{df(A)}{dA} = 0 \text{ or } \frac{df(x)}{dx}?$$

A glimpse of what is to come

So... what kind of functions will we encounter? A prime example is 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:

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

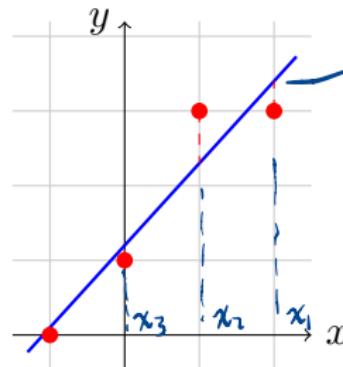
A glimpse of what is to come

So... what kind of functions will we encounter? A prime example is 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:

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



$$\hat{\mathbf{y}} = a_0 + a_1 x$$

$$\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \end{bmatrix} = \begin{bmatrix} a_0 + a_1 x_1 \\ a_0 + a_1 x_2 \\ a_0 + a_1 x_3 \end{bmatrix}$$

$$X = [1, x]$$

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

$$\hat{\mathbf{y}} = a_0 + a_1 x$$

A glimpse of what is to come

Others might want to differentiate functions such as

$$f = \mathbf{x}^T A \mathbf{x}, \det(A), \text{tr}(A), \lambda(A), \dots.$$

A glimpse of what is to come

Others might want to differentiate functions such as

$$f = \mathbf{x}^T A \mathbf{x}, \det(A), \text{tr}(A), \lambda(A), \dots.$$

轉置

Def. (Transpose)

Flip the matrix along its diagonal:

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

A glimpse of what is to come

Others might want to differentiate functions such as

$$f = \mathbf{x}^T A \mathbf{x}, \det(A), \text{tr}(A), \lambda(A), \dots.$$

Def. (Transpose)

Flip the matrix along its diagonal:

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

What number when transposed doubles?

A glimpse of what is to come

Others might want to differentiate functions such as

$$f = \mathbf{x}^T A \mathbf{x}, \det(A), \text{tr}(A), \lambda(A), \dots.$$

Def. (Transpose)

Flip the matrix along its diagonal:

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

What number when transposed doubles?

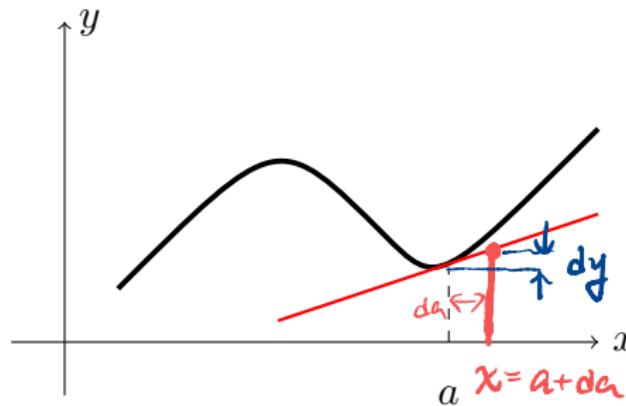
$$\text{午}^T = \text{什}$$

Table of Contents

- 1 Differentiation Revisited
- 2 Matrix Differentiation
- 3 Examples
- 4 Derivatives of “Matrix Derivatives”

What is a derivative?

Derivatives are linearizations. Consider the easiest case below.



We know that the first order approximation will be

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

What is a derivative?

We can rewrite the previous equation into the form below:

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

What is a derivative?

We can rewrite the previous equation into the form below:

$$y(a + dx) = \frac{dy}{dx} \Big|_a dx + y(a)$$

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

What is a derivative?

We can rewrite the previous equation into the form below:

$$y(a + dx) = \frac{dy}{dx} \Big|_a dx + y(a)$$

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

Thus obtaining:

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

What is a derivative?

We can rewrite the previous equation into the form below:

$$y(a + dx) = \frac{dy}{dx} \Big|_a dx + y(a)$$

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

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 works for matrices too.

Table of Contents

- 1 Differentiation Revisited
- 2 Matrix Differentiation
- 3 Examples
- 4 Derivatives of “Matrix Derivatives”

Examples

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.

Examples

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.

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

矩阵 向量

Examples

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.

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

By plugging in the definition above:

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

Examples

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.

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

By plugging in the definition above:

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

Examples

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.

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

By plugging in the definition above:

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

$$\frac{d\mathbf{f}}{dx} = A \quad \frac{d(ax)}{dx} = a, \quad \frac{d(ax^2)}{dx} = 2ax$$

Examples

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.

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

Examples

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:

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

Examples

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:

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

Examples

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.

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

By plugging in the definition above:

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

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

Derivatives as Linear Operators

Perhaps we shouldn't stick to the matrix algebra,

Derivatives as Linear Operators

Perhaps we shouldn't stick to the matrix algebra, and simply view derivatives as a **linear operator**:

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

Derivatives as Linear Operators

Perhaps we shouldn't stick to the matrix algebra, and 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:

$$\textcolor{red}{dy} = \mathcal{L}\{dx\} = f'dx$$

Derivatives as Linear Operators

Perhaps we shouldn't stick to the matrix algebra, and 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:

$$\mathcal{L}\{dx\} = f'dx$$

$$\mathcal{L}\{dx\} = dx f'$$

Derivatives as Linear Operators

Perhaps we shouldn't stick to the matrix algebra, and 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:

$$\mathcal{L}\{dx\} = f'dx$$

$$\mathcal{L}\{dx\} = dx f'$$

$$\mathcal{L}\{dx\} = \sqrt{f'}dx\sqrt{f'}.$$

Derivatives as Linear Operators

Perhaps we shouldn't stick to the matrix algebra, and 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:

$$\mathcal{L}\{dx\} = f'dx$$

$$\mathcal{L}\{dx\} = dx f'$$

$$\mathcal{L}\{dx\} = \sqrt{f'}dx\sqrt{f'}.$$

All works as long as it is linear:

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

Derivatives as Linear Operators

But since all linear operators can be represented by matrices, we can, still, represent derivatives as matrices and vectors alike.

Derivatives as Linear Operators

But since all linear operators can be represented by matrices, we can, still, represent derivatives as matrices and vectors alike.

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

Derivatives as Linear Operators

But since all linear operators can be represented by matrices, we can, still, represent derivatives as matrices and vectors alike.

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

Let us see some examples of functions where the representation of

$$\mathcal{L}\{\cdot\} = A(\cdot)$$

holds, where A is matrix.

Differentiate by Vectors

(1) Differentiate scalar function f by vector x :

Differentiate by Vectors

(1) Differentiate scalar function f by vector x :

Given a function $f(x) = f(x_1, \dots, x_n)$ 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:

$$df = \frac{df}{dx_i} dx_i \quad df = \frac{\partial f}{\partial x_1} dx_1 + \cdots + \frac{\partial f}{\partial x_n} dx_n. \quad \text{全微分} \quad (4)$$

But this is just the dot product of ∇f and dx !

$$\text{內積 } \vec{a} \cdot \vec{b} = a_1 b_1 + \cdots + a_n b_n$$

Differentiate by Vectors

Notice that

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

Differentiate by 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 ∇f and $d\mathbf{x}$:

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

Where the **gradient** $\overset{\text{nabla}}{\nabla f}$ is defined as a column vector.

梯度

Differentiate by 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 ∇f and $d\mathbf{x}$:

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

Where the **gradient** ∇f is defined as a column vector. Compare with

$$df = \boxed{\frac{df}{d\mathbf{x}}} d\mathbf{x}. \\ = (\nabla f)^T$$

Differentiate by Vectors

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

$$df = \frac{df}{d\mathbf{x}} d\mathbf{x},$$

then

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

Differentiate by Vectors

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

$$df = \frac{df}{d\mathbf{x}} d\mathbf{x}, \quad = [-\nabla f^T] \begin{bmatrix} d\mathbf{x} \\ 1 \end{bmatrix}$$

then

$$\frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_n} \end{bmatrix} = (\nabla f)^T.$$

But if we define

$$df = d\mathbf{x}^T \frac{df}{d\mathbf{x}^T}, \quad = [d\mathbf{x}^T] \begin{bmatrix} \nabla f \\ 1 \end{bmatrix}$$

Differentiate by Vectors

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

$$df = \frac{df}{d\mathbf{x}} d\mathbf{x},$$

then

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

But if we define

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

then

$$\frac{df}{d\mathbf{x}^\top} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \dots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \nabla f. = \left(\frac{df}{d\mathbf{x}} \right)^\top$$

Differentiate by Vectors

Notice that:

$$\begin{aligned} 1 \times n &\rightarrow \frac{df}{dx} \leftarrow 1 \times 1 \\ &\quad \text{---} \\ n \times 1 &\rightarrow \frac{df}{dx^T} \leftarrow 1 \times n \end{aligned}$$

Differentiate by Vectors

Notice that:

$$\begin{aligned} 1 \times n &\rightarrow \frac{df}{dx} \leftarrow 1 \times 1 \\ &\quad \frac{df}{dx} \leftarrow n \times 1 \\ n \times 1 &\rightarrow \frac{df}{dx^T} \leftarrow 1 \times 1 \\ &\quad \frac{df}{dx^T} \leftarrow 1 \times n \end{aligned}$$

Perhaps we could define the derivative as:

$$\underbrace{\frac{df}{dA}}_{m \times n} = \underbrace{\frac{df}{dA}}_{m \times s} \underbrace{\frac{dA}{s \times n}}_{s \times n} \text{ or } \underbrace{\frac{dB}{dB}}_{m \times s} \underbrace{\frac{df}{dB}}_{s \times n}.$$

$$\underbrace{\frac{df}{dA}}_{m \times n} = \underbrace{M_1}_{m \times p} \underbrace{\frac{dA}{p \times q}}_{p \times q} \underbrace{M_2}_{q \times n}$$

Differentiate by Vectors

(2) Differentiate vector function f by another vector x :

Differentiate by Vectors

(2) Differentiate vector function f by another vector x :

If x is a column vector and f is of the form:

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

$$m \times n \rightarrow \frac{df}{dx} \leftarrow m \times 1 \quad , \quad \underbrace{df}_{m \times 1} = \underbrace{\left(\frac{df}{dx} \right)}_{m \times n} \underbrace{dx}_{n \times 1}$$

Differentiate by Vectors

(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:

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial f_m}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix},$$

Differentiate by Vectors

(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:

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial f_m}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix},$$

Remark: the notations $\frac{df}{dx}$ and $\frac{\partial f}{\partial x}$ are used interchangeably.

Differentiate by Vectors

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

Differentiate by Vectors

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

- ① Gradient (transpose):

$$\nabla f := \frac{\partial f}{\partial \mathbf{x}^T} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}, \quad \nabla^T f := \frac{\partial f}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix} \quad (5)$$

$\left(\frac{\partial f}{\partial \mathbf{x}} \right)^T$

Differentiate by Vectors

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

- ① Gradient (transpose):

$$\nabla f := \frac{\partial f}{\partial \mathbf{x}^T} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}, \quad \nabla^T f := \frac{\partial f}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix} \quad (5)$$

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

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1 x_2 \\ x_1^2 \end{bmatrix}$$

Differentiate by Vectors

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

③ Hessian:

海森矩阵

$$\mathsf{H}(f) := \nabla(\nabla^\top f) = \frac{\partial^2 f}{\partial \mathbf{x}^\top \partial \mathbf{x}} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{bmatrix} \quad (6)$$

[...]

Differentiate by Vectors

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

③ Hessian:

$$\mathsf{H}(f) := \nabla \nabla^T f = \frac{\partial^2 f}{\partial \mathbf{x}^T \partial \mathbf{x}} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{bmatrix} \quad (6)$$

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

$$\nabla f = \begin{bmatrix} 2x_1 x_2 \\ x_1^2 \end{bmatrix}, \quad \frac{\partial^2 f}{\partial \mathbf{x}^T \partial \mathbf{x}} = \begin{bmatrix} 2x_2 & 2x_1 \\ 2x_1 & 0 \end{bmatrix}$$

對稱

Differentiate by Vectors

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

Differentiate by Vectors

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

③ Jacobian:

雅可比

$$J := \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial y_m}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}, \quad (7)$$

Differentiate by Vectors

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

③ Jacobian:

$$J := \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial y_n}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_n} \end{bmatrix}, \quad (7)$$

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

雅可比行列式

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

$$\int f(y) dy = \int f(y(x)) \left| \frac{dy}{dx} \right| dx$$

Differentiate by Matrices

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

$$df = f(A + dA) - f(A)$$

Some examples are as follows:

① $f(A) = Ax$

$$df = (dA)x$$

② $f(A) = \underline{\underline{A^T A}}^{n \times m \quad m \times n}$

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

Differentiate by Matrices

For:

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \textcolor{red}{a_{31}} & a_{32} & \cdots & a_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}_{m \times n}, \quad dA = \begin{bmatrix} da_{11} & da_{12} & \cdots \\ da_{21} & da_{22} & \cdots \\ \vdots & \vdots & \ddots \\ da_{m1} & da_{m2} & \cdots \end{bmatrix}_{m \times n}$$

we have that if f is a scalar function:

$$\cancel{df} \neq \frac{\partial f}{\partial A} dA \quad \frac{\partial f}{\partial A} = \begin{bmatrix} \frac{\partial f}{\partial a_{11}} & \frac{\partial f}{\partial a_{21}} & \frac{\partial f}{\partial a_{31}} & \cdots & \frac{\partial f}{\partial a_{m1}} \\ \frac{\partial f}{\partial a_{12}} & \frac{\partial f}{\partial a_{22}} & \frac{\partial f}{\partial a_{32}} & \cdots & \frac{\partial f}{\partial a_{m2}} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \frac{\partial f}{\partial a_{1n}} & \frac{\partial f}{\partial a_{2n}} & \frac{\partial f}{\partial a_{3n}} & \cdots & \frac{\partial f}{\partial a_{mn}} \end{bmatrix}_{n \times m} \quad (8)$$

$$\underline{\Delta f} = \frac{\partial f}{\partial a_{11}} da_{11} + \frac{\partial f}{\partial a_{12}} da_{12} + \dots \\ \text{+ } \dots + \frac{\partial f}{\partial a_{mn}} da_{mn} = \mathcal{L} \left\{ \begin{bmatrix} da_{11} & da_{12} & \dots \\ da_{21} & da_{22} & \dots \\ \vdots & \vdots & \ddots \\ da_{m1} & da_{m2} & \dots & da_{mn} \end{bmatrix} \right\} = \mathcal{L} \left\{ \underline{dA} \right\}$$

$$\boxed{\Delta f = \text{tr} \left(\frac{\partial f}{\partial A} \cdot dA \right)}$$

That's all folks.

$$= \text{tr} \left(\begin{bmatrix} \frac{\partial f}{\partial a_{11}} & \frac{\partial f}{\partial a_{21}} & \frac{\partial f}{\partial a_{31}} & \dots & \frac{\partial f}{\partial a_{m1}} \\ \frac{\partial f}{\partial a_{12}} & \frac{\partial f}{\partial a_{22}} & \frac{\partial f}{\partial a_{32}} & \dots & \frac{\partial f}{\partial a_{m2}} \\ \vdots & \vdots & \ddots & & \vdots \\ \frac{\partial f}{\partial a_{1n}} & \frac{\partial f}{\partial a_{2n}} & \frac{\partial f}{\partial a_{3n}} & \dots & \frac{\partial f}{\partial a_{mn}} \end{bmatrix} \begin{bmatrix} da_{11} & da_{12} & \dots & da_{1n} \\ da_{21} & da_{22} & \dots & da_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ da_{m1} & da_{m2} & \dots & da_{mn} \end{bmatrix} \right)$$

$$= \text{tr} \left[\begin{array}{c|c} \frac{\partial f}{\partial a_{11}} da_{11} + \dots + \frac{\partial f}{\partial a_{m1}} da_{m1} & \frac{\partial f}{\partial a_{12}} da_{12} + \dots + \frac{\partial f}{\partial a_{m2}} da_{m2} \\ \hline \dots & \dots \end{array} \right] \\ = \underline{\Delta f}$$

Time Derivatives

Time derivatives seem like the easiest, just differentiate term by term:

$$\frac{d}{dt} \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} \frac{dv_1}{dt} \\ \vdots \\ \frac{dv_n}{dt} \end{bmatrix}.$$

Time Derivatives

Time derivatives seem like the easiest, just differentiate term by term:

$$\frac{d}{dt} \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} \frac{dv_1}{dt} \\ \vdots \\ \frac{dv_n}{dt} \end{bmatrix}.$$

By the chain rule, we also have that

$$df = \frac{\partial f}{\partial A} dA \quad \frac{d}{dt} f(A(t)) = \frac{\partial f(A)}{\partial A} \frac{dA}{dt} = \mathcal{L} \left\{ \frac{dA}{dt} \right\}$$

And all the other chain rule, multiplication rule and etc. are satisfied.

Note that, again, the **order of multiplication matters**. And if the matrix derivative is of other form, the time derivative follows suit.

Time Derivatives

Time derivatives seem like the easiest, just differentiate term by term:

$$\frac{d}{dt} \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} \frac{dv_1}{dt} \\ \vdots \\ \frac{dv_n}{dt} \end{bmatrix}.$$

By the chain rule, we also have that

$$\frac{d}{dt} f(A(t)) = \frac{\partial f(A)}{\partial A} \frac{dA}{dt}.$$

E.g. $LHS = \frac{df}{dt} = \frac{d}{dt}(t^4) = 4t^3 = RHS = \frac{\partial f}{\partial \vec{x}} \frac{d\vec{x}}{dt} = [2x_1 x_2 x_1^2] \begin{bmatrix} 1 \\ 2t \end{bmatrix}$

$$\begin{aligned} f(\vec{x}) &= x_1^2 x_2 \\ \vec{x}(t) &= [x_1, x_2]^T = [t, t^2]^T \end{aligned}$$

✓

$$= [2t^3, t^2] \begin{bmatrix} 1 \\ 2t \end{bmatrix} = 4t^3$$

(Recap)

(Recap)

- ① We should view derivatives of $f(x)$ as a **linear approximation** of $f(x + dx)$.

(Recap)

- ① We should view derivatives of $f(x)$ as a **linear approximation** of $f(x + dx)$.
- ② Hence, we can represent derivatives as a **linear operator** $\mathcal{L}\{\cdot\}$ satisfying:

$$df(x) = \mathcal{L}\{dx\}.$$

(Recap)

- ① We should view derivatives of $f(x)$ as a **linear approximation** of $f(x + dx)$.
- ② Hence, we can represent derivatives as a **linear operator** $\mathcal{L}\{\cdot\}$ satisfying:

$$df(x) = \mathcal{L}\{dx\}.$$

- ③ In working with scalars, vectors or matrices, we can represent derivatives as

$$df(x) = \frac{df}{dx}dx,$$

as long as their dimension matches.

Before we take a break, a short remark is needed.

Before we take a break, a short remark is needed.

Notation

Two conventions exist in when differentiating with matrix/vector:
for y of size $m \times 1$ and x of size $n \times 1$,

- ① Numerator layout:

$\frac{\partial y}{\partial x}$ follows the size of $y \times x^T$, i.e. $m \times n$.

Before we take a break, a short remark is needed.

Notation

Two conventions exist in when differentiating with matrix/vector:
for y of size $m \times 1$ and x of size $n \times 1$,

- ① Numerator layout:

$\frac{\partial y}{\partial x}$ follows the size of $y \times x^T$, i.e. $m \times n$.

- ② Denominator layout:

$\frac{\partial y}{\partial x}$ follows the size of $x \times y^T$, i.e. $n \times m$.

$$\text{f}(\tilde{x}) = A\tilde{x} \quad , \quad \frac{\partial f}{\partial x} = A^T$$

Table of Contents

- 1 Differentiation Revisited
- 2 Matrix Differentiation
- 3 Examples
- 4 Derivatives of “Matrix Derivatives”

Vector Functions

① $f(\mathbf{x}) = \mathbf{x}$

Vector Functions

$$\textcircled{1} \quad f(\mathbf{x}) = \mathbf{x}$$

$$df = (\mathbf{x} + d\mathbf{x}) - \mathbf{x} = \mathbf{1}d\mathbf{x}$$

$$\frac{\partial \mathbf{x}}{\partial \mathbf{x}} = 1$$

Vector Functions

① $f(\mathbf{x}) = \mathbf{x}$

$$df = (\mathbf{x} + d\mathbf{x}) - \mathbf{x} = \mathbf{1}d\mathbf{x}$$

$$\frac{\partial \mathbf{x}}{\partial \mathbf{x}} = 1$$

② $f(\mathbf{x}) = \mathbf{x}^\top$

Vector Functions

① $f(\mathbf{x}) = \mathbf{x}$

$$df = (\mathbf{x} + d\mathbf{x}) - \mathbf{x} = \mathbf{1}d\mathbf{x}$$

$$\frac{\partial \mathbf{x}}{\partial \mathbf{x}} = 1$$

② $f(\mathbf{x}) = \mathbf{x}^\top$

$$df = (\mathbf{x} + d\mathbf{x})^\top - (\mathbf{x})^\top = (d\mathbf{x})^\top$$

轉置 \rightarrow linear

Vector Functions

Find the derivative with respect to x of the function

$$f(x) = x \cdot x = x^T x.$$

[Sol.]

$$\begin{bmatrix}] \\] \end{bmatrix} \cdot \begin{bmatrix}] \\] \end{bmatrix} = \text{[} \cancel{\text{——}} \text{]} \begin{bmatrix} | \\ | \end{bmatrix}$$

Vector Functions

Find the derivative with respect to \boldsymbol{x} of the function

$$f(\boldsymbol{x}) = \boldsymbol{x} \cdot \boldsymbol{x} = \boldsymbol{x}^\top \boldsymbol{x}.$$

[Sol.]

$$\mathrm{d}f(\boldsymbol{x}) = \mathrm{d}(\boldsymbol{x}^\top \boldsymbol{x}) =$$

Vector Functions

Find the derivative with respect to x of the function

$$f(x) = x \cdot x = x^T x.$$

[Sol.]

$$\begin{aligned} df(x) &= d(x^T x) = (x + dx)^T (x + dx) - x^T x \\ &= \cancel{x^T x} + x^T dx + (dx)^T x + (dx)^T dx - \cancel{x^T x} \end{aligned}$$

Vector Functions

Find the derivative with respect to \mathbf{x} of the function

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

[Sol.]

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

Vector Functions

Find the derivative with respect to x of the function

$$f(x) = x \cdot x = x^T x.$$

[Sol.]

$$\begin{aligned} df(x) &= d(x^T x) = (x + dx)^T (x + dx) - x^T x \\ &= x^T x + (dx)^T x + x^T dx + (dx)^T dx - x^T x \\ &= \cancel{(dx)^T x} + x^T \cancel{dx} + \boxed{\cancel{(dx)^T dx}} \rightarrow 0 \end{aligned}$$

$$f(x) = x^2$$

$$f'(x) = \lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h} = \lim_{h \rightarrow 0} \frac{2xh + h^2}{h} = \lim_{h \rightarrow 0} 2x + h = 2x$$

Vector Functions

Find the derivative with respect to x of the function

$$f(x) = x \cdot x = x^T x.$$

[Sol.]

$$\begin{aligned} df(x) &= d(x^T x) = (x + dx)^T (x + dx) - x^T x \\ &= x^T x + (dx)^T x + x^T dx + (dx)^T dx - x^T x \\ &= (dx)^T x + x^T dx + \cancel{(dx)^T dx} \end{aligned}$$

$dx^T (\dots)$

Ignoring the terms of $O(dx^2)$ and notice that the transpose of a scalar is still itself ($(dx)^T x = x^T dx$), we hence have:

$$df = 2x^T dx = \mathcal{L}\{dx\}. \quad \blacksquare$$

Quadratic Form

Find the derivative with respect to x of the function

$$f(x) = x^T A x. \quad \text{二次型 Quadratic Form}$$

[Sol.]

Quadratic Form

Find the derivative with respect to \mathbf{x} of the function

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

[Sol.]

$$df = (\mathbf{x} + d\mathbf{x})^T A (\mathbf{x} + d\mathbf{x}) - \mathbf{x}^T \cancel{Ax}$$

Quadratic Form

Find the derivative with respect to \mathbf{x} of the function

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

[Sol.]

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

Quadratic Form

Find the derivative with respect to \mathbf{x} of the function

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

[Sol.]

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

Quadratic Form

Find the derivative with respect to \mathbf{x} of the function

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

[Sol.]

$$\begin{aligned} df &= (\mathbf{x} + d\mathbf{x})^T A (\mathbf{x} + d\mathbf{x}) - \mathbf{x}^T \mathbf{x} \\ &= \mathbf{x}^T A \mathbf{x} + (d\mathbf{x})^T A \mathbf{x} + \mathbf{x}^T A d\mathbf{x} + (d\mathbf{x})^T A d\mathbf{x} - \mathbf{x}^T \mathbf{x} \\ &= \underline{(d\mathbf{x})^T A \mathbf{x}} + \mathbf{x}^T A d\mathbf{x} + \underline{(d\mathbf{x})^T A d\mathbf{x}} \\ &\quad \text{`}\mathbf{x}^T A^T d\mathbf{x}\text{' } \end{aligned}$$

Ignoring the terms of $O(d\mathbf{x}^2)$, we hence have:

$$df = \mathbf{x}^T (A + A^T) d\mathbf{x} = \mathcal{L}\{d\mathbf{x}\}. \quad \blacksquare$$

Inverse

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

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

[Sol.]

Inverse

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

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

[Sol.] Since we have

$$A(t)A^{-1}(t) = 1,$$

(单位矩阵)

Inverse

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

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

[Sol.] Since we have

$$A(t)A^{-1}(t) = 1,$$

by chain rule:

$$\frac{d}{dt} \left(AA^{-1} \right) = \frac{dA}{dt} A^{-1} + A \frac{dA^{-1}}{dt} = 0$$

Inverse

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

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

[Sol.] Since we have

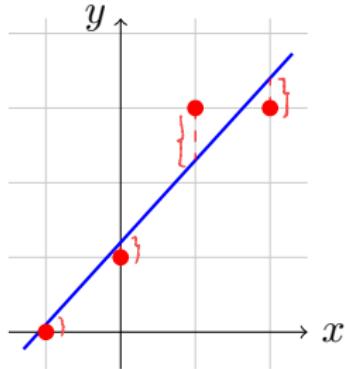
$$A(t)A^{-1}(t) = 1,$$

by chain rule:

$$\begin{aligned}\frac{d}{dt} \left(AA^{-1} \right) &= \frac{dA}{dt} A^{-1} + A \frac{dA^{-1}}{dt} = 0 \\ \frac{dA^{-1}}{dt} &= -A^{-1} \frac{dA}{dt} A^{-1}. \quad \blacksquare\end{aligned}$$

Least Mean Square

Suppose the regression line for the data points is



估測值 $\rightarrow \hat{y} = a_0 + a_1 x$,

we can record the relationship as follows

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \hat{y}_4 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ 1 & x_4 \end{bmatrix}}_X \underbrace{\begin{bmatrix} a_0 \\ a_1 \end{bmatrix}}_w \approx \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

$$\mathcal{E} = \|y - \hat{y}\|^2 = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots$$

Least Mean Square

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:

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

Least Mean Square

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

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

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

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

Least Mean Square

[Sol.]

Least Mean Square

[Sol.]

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

$$1. \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{x} = 2\mathbf{x}^\top = 2(\mathbf{X}\mathbf{w} - \mathbf{y})^\top \mathbf{X}$$

$$2. \frac{\partial}{\partial \mathbf{x}} (\mathbf{A}\mathbf{x} + \mathbf{b}) = \mathbf{A}$$

Least Mean Square

[Sol.]

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

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

Least Mean Square

[Sol.]

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

$$= 2(\mathbf{X}\mathbf{w} - \mathbf{y})^\top \mathbf{X}$$

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

Least Mean Square

[Sol.]

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

$$= 2(\mathbf{X}\mathbf{w} - \mathbf{y})^\top \mathbf{X}$$

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

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

Least Mean Square

[Sol.]

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

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

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

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

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \end{bmatrix}, \quad X^\top X = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$$

$$|X^\top X| = n \sum x_i^2 - (\sum x_i)^2 = n^2 \text{Var}(X) \geq 0, \quad "=\text{" iff } x_i \text{ 同}$$

Least Mean Square

[Sol.]

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

$$\begin{aligned}(\mathbf{X}\mathbf{w}^* - \mathbf{y})^\top \mathbf{X} &= 0 \\ \mathbf{X}^\top \mathbf{X}\mathbf{w}^* &= \mathbf{X}^\top \mathbf{y} \\ \rightarrow \mathbf{w}^* &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad \blacksquare\end{aligned}$$

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

$$\mathbf{X}\mathbf{w} \approx \mathbf{y} \rightarrow \mathbf{w} = \mathbf{X}^{-1} \mathbf{y}$$

Least Mean Square

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:

Least Mean Square

We've only checked that $\mathbf{w}^* = (X^\top X)^{-1} X^\top \mathbf{y}$ is an extrema, but we have yet to check whether its a maxima or a minima. A second derivative test is needed:

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

Least Mean Square

We've only checked that $\mathbf{w}^* = (X^\top X)^{-1} X^\top \mathbf{y}$ is an extrema, but we have yet to check whether its a maxima or a minima. A second derivative test is needed:

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

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

Least Mean Square

We've only checked that $\mathbf{w}^* = (X^\top X)^{-1} X^\top \mathbf{y}$ is an extrema, but we have yet to check whether its a maxima or a minima. A second derivative test is needed:

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

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

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

$$\begin{aligned} \|X(\mathbf{w}^* + d\mathbf{w}) - \mathbf{y}\|^2 - \|X\mathbf{w}^* - \mathbf{y}\|^2 &= (d\mathbf{w})^\top (2X^\top X)d\mathbf{w} > 0, \\ &= 2(Xd\mathbf{w})^\top (Xd\mathbf{w}) \end{aligned}$$

it is therefore a minima.

Definiteness

Def. (Definiteness)

A matrix A is called **positive definite** if for any non-zero vector x it satisfies: $A^T = A$ 正定

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

Moreover, we have:

(positive semi-definite) $x^T Ax \geq 0 \iff A \succeq 0$

(negative definite) $x^T Ax < 0 \iff A \prec 0$

(negative semi-definite) $x^T Ax \leq 0 \iff A \preceq 0$

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

Taylor Expansion

From the derivation of LLMSE solution above, we can find a second order approximation of a scalar-valued function $f(\mathbf{x})$ of vector \mathbf{x} by:

$$f(\mathbf{x}) = f(\mathbf{a}) + \underbrace{\frac{\partial f(\mathbf{a})}{\partial \mathbf{x}}(\mathbf{x} - \mathbf{a})}_{\text{linear approx. term}} + \underbrace{\frac{1}{2!}(\mathbf{x} - \mathbf{a})^T \frac{\partial^2 f(\mathbf{a})}{\partial \mathbf{x} \partial \mathbf{x}^T} (\mathbf{x} - \mathbf{a})}_{\text{quadratic approx. term}} + \dots, \quad (9)$$

or as

$$f(\mathbf{x}) = f(\mathbf{a}) + f'(\mathbf{a})(\mathbf{x} - \mathbf{a}) + \frac{1}{2} f''(\mathbf{a})(\mathbf{x} - \mathbf{a})^2 + \dots$$

$$f(\mathbf{a} + d\mathbf{x}) = f(\mathbf{a}) + \underbrace{\nabla^T f(\mathbf{a}) d\mathbf{x}}_{\text{gradient}} + \frac{1}{2} (d\mathbf{x})^T \underbrace{\mathbf{H}(f(\mathbf{a}))}_{\text{Hessian}} d\mathbf{x}. \quad (10)$$



Regression

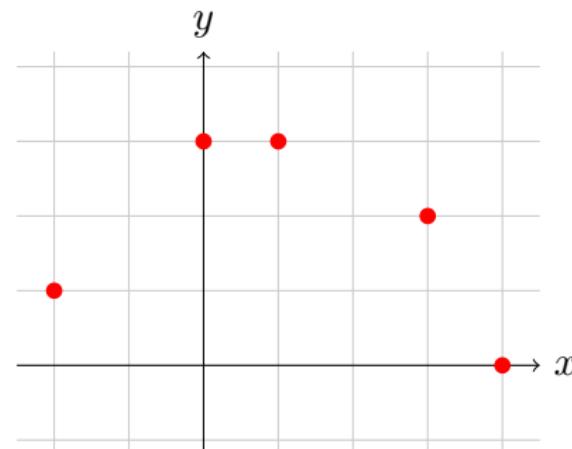
Given a set of data points (x, y) , find a function that interpolates them with the least mean square error.

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.



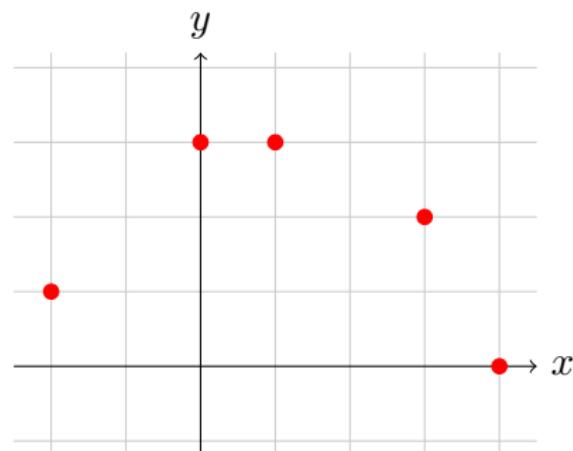
Regression

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:

$$\hat{y} = \begin{bmatrix} 1 & x & x^2 \\ 1 & -2 & 4 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = Xw$$

$$y = [1, 3, 3, 2, 0]^T$$



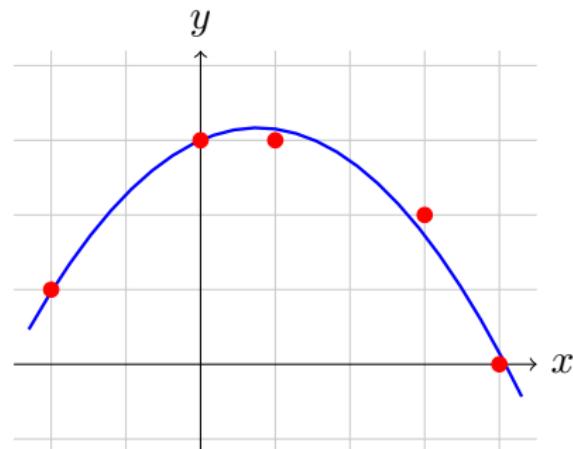
Regression

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

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

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

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



HTML

HTML, HW3

Q6. Let the cross-entropy error function for $E_{\text{in}}(\mathbf{w})$ be:

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n) \right).$$

Find the Hessian of the function. Express it in diagonalized form of
 $\nabla^\top \nabla E_{\text{in}} = XDX^\top$.

HTML

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n) \right)$$

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{\partial E_{\text{in}}}{\partial \mathbf{w}^\top} = \frac{\partial}{\partial \mathbf{w}^\top} \left(\frac{1}{N} \sum_{n=1}^N \ln \left(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n) \right) \right)$$

HTML

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n) \right)$$

$$\begin{aligned}\nabla E_{\text{in}}(\mathbf{w}) &= \frac{\partial E_{\text{in}}}{\partial \mathbf{w}^\top} = \frac{\partial}{\partial \mathbf{w}^\top} \left(\frac{1}{N} \sum_{n=1}^N \ln \left(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n) \right) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)} \frac{\partial}{\partial \mathbf{w}^\top} \left(\exp(-y_n \mathbf{w}^\top \mathbf{x}_n) \right)\end{aligned}$$

HTML

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n) \right)$$

$$\begin{aligned}\nabla E_{\text{in}}(\mathbf{w}) &= \frac{\partial E_{\text{in}}}{\partial \mathbf{w}^\top} = \frac{\partial}{\partial \mathbf{w}^\top} \left(\frac{1}{N} \sum_{n=1}^N \ln \left(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n) \right) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \frac{1}{1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)} \frac{\partial}{\partial \mathbf{w}^\top} \left(\exp(-y_n \mathbf{w}^\top \mathbf{x}_n) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \frac{-y_n \mathbf{x}_n \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}\end{aligned}$$

HTML

$$\frac{\partial^2 E_{\text{in}}}{\partial \mathbf{w} \partial \mathbf{w}^\top} = \frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{N} \sum_{n=1}^N \frac{-y_n \mathbf{x}_n \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)} \right)$$

HTML

$$\begin{aligned}\frac{\partial^2 E_{\text{in}}}{\partial \mathbf{w} \partial \mathbf{w}^T} &= \frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{N} \sum_{n=1}^N \frac{-y_n \mathbf{x}_n \exp(-y_n \mathbf{w}^T \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)} \right) \\ &= \sum_{n=1}^N \frac{-y_n \mathbf{x}_n - y_n \mathbf{x}_n^T \exp(-y_n \mathbf{w}^T \mathbf{x}_n) (1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)) + y_n \mathbf{x}_n^T (\exp(-y_n \mathbf{w}^T \mathbf{x}_n))^2}{N (1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n))^2}\end{aligned}$$

HTML

$$\begin{aligned}\frac{\partial^2 E_{\text{in}}}{\partial \mathbf{w} \partial \mathbf{w}^\top} &= \frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{N} \sum_{n=1}^N \frac{-y_n \mathbf{x}_n \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)} \right) \\ &= \sum_{n=1}^N \frac{-y_n \mathbf{x}_n - y_n \mathbf{x}_n^\top \exp(-y_n \mathbf{w}^\top \mathbf{x}_n) (1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)) + y_n \mathbf{x}_n^\top (\exp(-y_n \mathbf{w}^\top \mathbf{x}_n))^2}{N (1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n))^2} \\ &= \sum_{i=1}^N \left(\frac{y_n^2 \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}{N (1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n))^2} \right) \mathbf{x}_n \mathbf{x}_n^\top.\end{aligned}$$

HTML

$$\frac{\partial^2 E_{\text{in}}}{\partial \mathbf{w} \partial \mathbf{w}^\top} = \sum_{i=1}^N \left(\frac{y_n^2 \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}{N(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n))^2} \right) \mathbf{x}_n \mathbf{x}_n^\top$$

HTML

$$\frac{\partial^2 E_{\text{in}}}{\partial \mathbf{w} \partial \mathbf{w}^\top} = \sum_{i=1}^N \left(\frac{y_n^2 \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}{N(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n))^2} \right) \mathbf{x}_n \mathbf{x}_n^\top$$

Diagonalized (spectral decomposition):

$$E_{\text{in}} = XDX^\top = \sum_{n=1}^N \lambda_n \mathbf{x}_n \mathbf{x}_n^\top. \quad (11)$$

$$\begin{matrix} \mathbf{x}_1 \dots \mathbf{x}_N \\ \parallel \\ \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_N \end{bmatrix} \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} \end{matrix}$$

HTML

$$\frac{\partial^2 E_{\text{in}}}{\partial \mathbf{w} \partial \mathbf{w}^\top} = \sum_{i=1}^N \left(\frac{y_n^2 \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}{N(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n))^2} \right) \mathbf{x}_n \mathbf{x}_n^\top$$

Diagonalized (spectral decomposition):

$$E_{\text{in}} = XDX^\top = \sum_{n=1}^N \lambda_n \mathbf{x}_n \mathbf{x}_n^\top. \quad (11)$$

This is the Hessian in its diagonalized form, with \mathbf{x}_n being the eigenvectors associated with the eigenvalues of

$$\frac{y_n^2 \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}{N(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n))^2}.$$

Table of Contents

- 1 Differentiation Revisited
- 2 Matrix Differentiation
- 3 Examples
- 4 Derivatives of "Matrix Derivatives"

Matrix Derivatives

We'll go over the derivative of objects including the trace, the determinant, eigenvalues, and singular values.

Trace

Def. (Trace)

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

Trace

Def. (Trace)

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

$$\text{tr} \left(\begin{bmatrix} a_{11} & * & \cdots & * \\ * & a_{22} & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ * & * & \cdots & a_{nn} \end{bmatrix} \right) = \sum_{i=1}^n a_{ii}. \quad (12)$$

Trace

Def. (Trace)

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

$$\text{tr} \left(\begin{bmatrix} a_{11} & * & \cdots & * \\ * & a_{22} & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ * & * & \cdots & a_{nn} \end{bmatrix} \right) = \sum_{i=1}^n a_{ii}. \quad (12)$$

It can easily be checked that the derivative operator and trace operator are commutative:

$$\frac{d}{dt} \text{tr}(A(t)) = \text{tr} \left(\frac{d}{dt} A(t) \right).$$

Determinant

Lemma

The following identity holds for all square matrices A :

$$\det(e^A) = e^{\text{tr}(A)}. \quad (13)$$

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

Determinant

Lemma

The following identity holds for all square matrices A :

$$\det(e^A) = e^{\text{tr}(A)}. \quad (13)$$

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

e.g. $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$, $e^A = \sum_{n=0}^{\infty} \frac{A^n}{n!} = \sum_{n=0}^{\infty} \frac{A^{2n}}{(2n)!} + \sum_{n=0}^{\infty} \frac{A^{2n+1}}{(2n+1)!}$

$$A^2 = -\mathbb{1}$$

$$= \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n)!} \mathbb{1} + \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} A$$

$$= (\cos 1) \mathbb{1} + (\sin 1) A = \begin{bmatrix} \cos 1 & \sin 1 \\ -\sin 1 & \cos 1 \end{bmatrix}$$

Determinant

Lemma

The following identity holds for all square matrices A :

$$\det(e^A) = e^{\text{tr}(A)}. \quad (13)$$

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

e.g. $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$ $e^A = \begin{bmatrix} \cos 1 & \sin 1 \\ -\sin 1 & \cos 1 \end{bmatrix}$, $\det(e^A) = 1$

$\text{tr}(A) = 0$, $e^{\text{tr}(A)} = 1$

Determinant

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

$$A(t) = e^{B(t)},$$

then a (pseudo-)proof is as follows:

$$\frac{d}{dt} \det(A(t)) =$$

Determinant

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

$$A(t) = e^{B(t)},$$

then a (pseudo-)proof is as follows:

$$\frac{d}{dt} \det(A(t)) = \frac{d}{dt} e^{\text{tr}B(t)}$$

Determinant

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

$$A(t) = e^{B(t)},$$

then a (pseudo-)proof is as follows:

$$\begin{aligned}\frac{d}{dt} \det(A(t)) &= \frac{d}{dt} e^{\text{tr}B(t)} \\ &= e^{\text{tr}B(t)} \text{tr} \left(\frac{d}{dt} B(t) \right)\end{aligned}$$

$$\therefore \frac{d}{dt} A(t) = e^{B(t)} \frac{d}{dt} B(t)$$

Determinant

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

$$A(t) = e^{B(t)},$$

then a (pseudo-)proof is as follows:

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

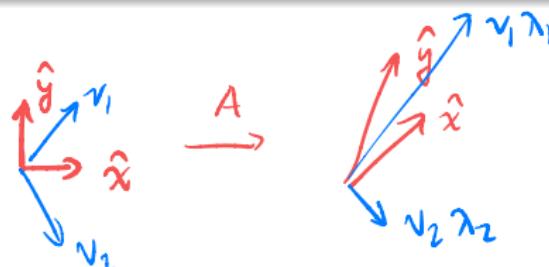
Eigenvalues

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**.



Eigenvalues

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.

Eigenvalues

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.

Written in the language of matrices:

$$AV = A \begin{bmatrix} \mathbf{v}_1 & \cdots & \mathbf{v}_n \end{bmatrix} = \begin{bmatrix} \mathbf{v}_1 & \cdots & \mathbf{v}_n \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} = V\Lambda$$

$$= \begin{bmatrix} \lambda_1 \mathbf{v}_1 & \cdots & \lambda_n \mathbf{v}_n \end{bmatrix}$$

Eigenvalues

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:

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

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

Written in the language of matrices:

$$U^T A = [\mathbf{u}_1 \ \dots \ \mathbf{u}_n]^T A = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} [\mathbf{u}_1 \ \dots \ \mathbf{u}_n]^T = \Lambda U^T$$

Eigenvalues

Some theorems of linear algebra:

Eigenvalues

Some theorems of linear algebra: since we have $AV = V\Lambda$,

Thm. (Diagonalization)

If V is full rank, i.e., V^{-1} exists, then the matrix A can be diagonalized via:

$$\Lambda = V^{-1}AV.$$

$$\begin{cases} AV = V\Lambda \\ V^T A = \Lambda V^T \end{cases} \rightarrow \begin{matrix} \Lambda = V^{-1}AV \\ [\lambda_1 \dots \lambda_n] \quad [**] \end{matrix} \rightarrow \begin{cases} A = V\Lambda V^{-1} \\ A = U^T \Lambda U^T \end{cases}$$

$$\Lambda = U^T A U^{-T} \quad \Lambda^n = \begin{bmatrix} \lambda_1^n & 0 \\ 0 & \lambda_2^n \end{bmatrix} \quad A = V\Lambda U^T$$

$$A^n = V\Lambda^n V^{-1} \quad = [v_1 \dots v_n] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} u_1^T \\ \vdots \\ u_n^T \end{bmatrix}$$

$$= \sum_{i=1}^n \lambda_i v_i u_i^T$$

Eigenvalues

Some theorems of linear algebra: since we have $AV = V\Lambda$,

Thm. (Diagonalization)

If V is full rank, i.e., V^{-1} exists, then the matrix A can be diagonalized via:

$$\Lambda = V^{-1}AV.$$

By expanding out $A = V\Lambda V^{-1}$ and setting $V = [\mathbf{v}_1, \dots, \mathbf{v}_n]$,
 $(V^{-1})^\top = [\mathbf{u}_1, \dots, \mathbf{u}_n]$,

Thm. (Spectral Decomposition)

The $n \times n$ matrix A can be decomposed by

$$A = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{u}_i^\top.$$

Eigenvalues

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

Eigenvalues

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: $v_i(t)$ and $u_i(t)$ that has length satisfying:

$$u_i^T v_i = 1.$$

Eigenvalues

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: $v_i(t)$ and $u_i(t)$ that has length satisfying:

$$u_i^T v_i = 1.$$

Then,

$$u_i^T (t) A(t) \underline{\lambda_i v_i} = \lambda_i(t)$$

Eigenvalues

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,

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

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

Eigenvalues

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

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

Eigenvalues

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

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

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

Singular Values

Def. (Singular Value Decomposition, SVD)

For a real (complex) matrix A of size $m \times n$ (WLOG let $m > n$), then it can be decomposed into a sandwich product of diagonal matrix Σ by two orthogonal (unitary) matrices U and V :

$$A = \underbrace{U}_{m \times m} \underbrace{\Sigma}_{m \times n} \underbrace{V^T}_{n \times n} = \begin{bmatrix} | & & | \\ u_1 & \cdots & u_m \\ | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & \\ \hline & 0 & & 0 \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix} \quad (14)$$

where $U^T U = 1$, $V^T V = 1$, and $\sigma_1 \geq \dots \geq \sigma_r \geq 0$.

Note that for the case of complex A , the transposition are replaced by conjugate-transpose.

Singular Values

Thm. (Spectral Decomposition)

The result of singular value decomposition can also be written as:

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

Singular Values

Thm. (Spectral Decomposition)

The result of singular value decomposition can also be written as:

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

The same derivation can be applied for finding the time derivative of singular values.

Singular Values

Thm. (Spectral Decomposition)

The result of singular value decomposition can also be written as:

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

The same derivation can be applied for finding the time derivative of singular values. By utilizing:

$$AV = U\Sigma V^\top V = U\Sigma, \quad U^\top A = U^\top U\Sigma V^\top = \Sigma V^\top$$

Singular Values

Thm. (Spectral Decomposition)

The result of singular value decomposition can also be written as:

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

The same derivation can be applied for finding the time derivative of singular values. By utilizing:

$$AV = U\Sigma V^\top V = U\Sigma, \quad U^\top A = U^\top U\Sigma V^\top = \Sigma V^\top$$

the solution is given by

$$\frac{d\sigma_i}{dt} = \mathbf{u}_i^\top \frac{dA}{dt} \mathbf{v}_i, \quad (15)$$

where $\mathbf{u}_i(t)$ and $\mathbf{v}_i(t)$ are the left and right singular vector associated with the singular value $\sigma_i(t)$.

That's all folks.