

# VECTOR / MATRIX DERIVATIVES

Everything You Need To Know

KANG GYEONGHUN

ESC, YONSEI UNIVERSITY

March 29, 2020

# Table of Contents

- 1 Vector Derivatives
  - Directional Derivative
  - Jacobian Matrix
  - Vector Derivatives Rules
  - Example: Least Squares Estimators
- 2 Change of Variables
  - Why do we need Jacobian?
  - Example: Geometric Interpretation of MVN
- 3 Matrix Derivatives
  - Directional Derivative in Matrix space
  - Matrix Derivatives for Determinant
  - Matrix Derivatives for Trace
  - Example: MLE for MVN

# I. Vector Derivatives

# DIRECTIONAL DERIVATIVE

## Taylor Expansion for Derivatives

- Suppose  $f(x) \in C^2[a, b]$ . We consider **Taylor expansion** of the function  $f$  at  $w$ ;

$$f(x) \approx f(w) + \frac{(x-w)}{1!} f'(w) + \frac{(x-w)^2}{2!} f''(w)$$

- Then for sufficiently small  $e$ , we have;

$$f(x+e) \approx f(w) + \frac{(x+e-w)}{1!} f'(w) + \frac{(x+e-w)^2}{2!} f''(w)$$

- With respect to  $e$ , we can say

$$f(x+e) \approx O(1) + O(e) + O(e^2)$$

This formulation provides us an intuitive framework in understanding derivatives in various spaces.

# DIRECTIONAL DERIVATIVE

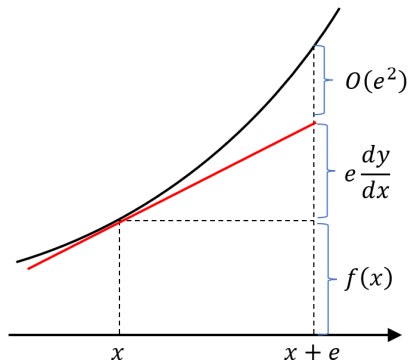
## Recall: 1-dimensional Derivative

- For  $f(x) : \mathbb{R} \mapsto \mathbb{R}$ , an infinitesimal change  $e$  in the input  $x$  leads to change in the output;

$$f(x + e) = f(x) + e \frac{df}{dx} + O(e^2)$$

- Rearrange,  $e \rightarrow 0$ , and we have 1D Derivative;

$$\frac{dy}{dx} = \lim_{e \rightarrow 0} \frac{f(x + e) - f(x)}{e}$$



*Scalar field*

# DIRECTIONAL DERIVATIVE

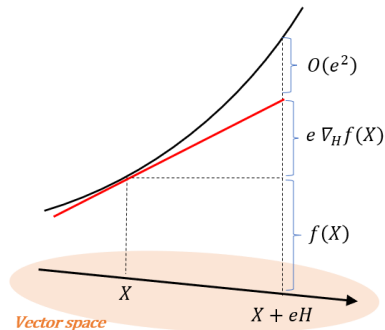
## n-dimensional Derivative

- For  $f(x) : \mathbb{R}^n \mapsto \mathbb{R}$ , an infinitesimal change  $e$  in the input  $\mathbf{x}$  **in the direction of  $\mathbf{H}$**  leads to change in the output;

$$f(\mathbf{X} + e\mathbf{H}) = f(\mathbf{X}) + e\nabla_{\mathbf{H}}f(\mathbf{X}) + O(e^2)$$

- Rearrange,  $e \rightarrow 0$ , and we have nD **Directional Derivative**;

$$\nabla_{\mathbf{H}}f(\mathbf{X}) = \lim_{e \rightarrow 0} \frac{f(\mathbf{X} + e\mathbf{H}) - f(\mathbf{X})}{e}$$



$$\nabla_{\mathbf{H}}f(\mathbf{X}) = \mathbf{H} \cdot \nabla f = \mathbf{H} \cdot \frac{\partial f}{\partial \mathbf{X}}$$

# DIRECTIONAL DERIVATIVE

## Directional Derivative in Vectorspace

- With a gradient of a function  $f; \mathbf{x} \mapsto \mathbb{R}$ , once we specify a direction  $\mathbf{b}$  (unit vector) of variation in the input  $\mathbf{x}$ , we have a **directional derivative of  $f$** ;

$$\begin{aligned}\nabla_{\mathbf{b}} f(\mathbf{x}) &= f'(\mathbf{x}) \cdot \mathbf{b} = \lim_{\epsilon \rightarrow 0} \frac{f(\mathbf{x} + \epsilon \mathbf{b}) - f(\mathbf{x})}{\epsilon} \\ &= \nabla f(\mathbf{x}) \cdot \mathbf{b} = \|\nabla f(\mathbf{x})\| \|\mathbf{b}\| \cos \theta\end{aligned}$$

- Directional derivative is just a  $\mathbb{R}^1$  derivative generalized to  $\mathbb{R}^n$  space. The difference is that, unlike in  $\mathbb{R}^1$  where we had only a number line, in higher dimensions we have to specify which direction  $d\mathbf{x}$  is headed.

$$\text{1D derivatives: } \frac{df(x)}{dx} = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x)}{\epsilon}$$

- $\nabla f(\mathbf{x})$  is called **gradient** and represents a direction in  $\mathbf{x}$  of the maximum change in  $f(\mathbf{x})$  since  $\cos \theta = 1 \iff \theta = 0$ .

# JACOBIAN MATRIX

## Jacobian Matrix for Multidimensional Output

Define  $\psi : \mathbb{R}^n \rightarrow \mathbb{R}^m$  by  $\mathbf{y} = \psi(\mathbf{x})$ . The  $m \times n$  matrix of first-order derivatives of this transformation is called **Jacobian Matrix** and is expressed as;

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \nabla y_1 \\ \nabla y_2 \\ \vdots \\ \nabla y_m \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}_{m \times n}$$

- If  $\psi : \mathbb{R}^n \rightarrow \mathbb{R}^1$ , we have a single row  $\frac{\partial y}{\partial \mathbf{x}} = \nabla \psi(\mathbf{x}) = \left[ \frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$ , which is called **Gradient of  $y$** .
- Jacobian Matrix for  $m \times 1$  vector is essentially a vertical stack of each element's gradient.



# JACOBIAN MATRIX

## Jacobian Matrix for Multidimensional Output

Define  $\psi : \mathbb{R}^n \rightarrow \mathbb{R}^m$  by  $\mathbf{y} = \psi(\mathbf{x})$ . The  $m \times n$  matrix of first-order derivatives of this transformation is called **Jacobian Matrix** and is expressed as;

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}_{m \times n}$$

• If  $\psi : \mathbb{R}^1 \rightarrow \mathbb{R}^m$ , we have a single column  $\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix}$

# VECTOR DERIVATIVES

## Vector Derivatives Rule 1

Let  $\mathbf{y} = \mathbf{A}\mathbf{x}$  where  $\mathbf{y} \in \mathbb{R}^{m \times 1}$ ,  $\mathbf{x} \in \mathbb{R}^{n \times 1}$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Then  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A}$ . ( $\mathbf{A}$  is constant matrix.)

- **pf)** The easiest way is by showing element-wise operations.

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

So we have  $y_i = a_{i1}x_1 + a_{i2}x_2 + \dots + a_{in}x_n$ , and  $\frac{\partial y_i}{\partial x_j} = a_{ij}$ .

$$\text{Therefore, } \mathbf{A} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}}. \quad \square$$

# VECTOR DERIVATIVES

## Vector Derivatives Rule 2

Let  $\mathbf{y} = \mathbf{A}\mathbf{x}$  where  $\mathbf{y} \in \mathbb{R}^{m \times 1}$ ,  $\mathbf{x} \in \mathbb{R}^{n \times 1}$  ( $\mathbf{x}$  is a function of  $\mathbf{z}$ ), and  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Then  $\frac{\partial \mathbf{y}}{\partial \mathbf{z}} = \mathbf{A} \frac{\partial \mathbf{x}}{\partial \mathbf{z}}$ . ( $\mathbf{A}$  is constant matrix.)

- **pf)** Since  $y_i = a_{i1}x_1 + a_{i2}x_2 + \dots + a_{in}x_n$ , we have  $\frac{\partial y_i}{\partial z_j} = a_{i1} \frac{\partial x_1}{\partial z_j} + a_{i2} \frac{\partial x_2}{\partial z_j} + \dots + a_{in} \frac{\partial x_n}{\partial z_j}$ .

Since  $\frac{\partial x_i}{\partial z_j}$  is  $(i, j)$ th element of  $\frac{\partial \mathbf{x}}{\partial \mathbf{z}}$ ,

$$\frac{\partial y}{\partial \mathbf{z}}_{(i,j)} = \frac{\partial y_i}{\partial z_j} = \sum_{k=1}^n a_{ik} \frac{\partial x_k}{\partial z_j} = \sum_{k=1}^n \mathbf{A}_{(i,k)} \frac{\partial \mathbf{x}}{\partial \mathbf{z}}_{(k,j)} = \frac{\partial \mathbf{x}}{\partial \mathbf{z}}_{(i,j)} \quad \square$$

- This serves as a proof of vector chain rule:  $\frac{\partial \mathbf{y}}{\partial \mathbf{z}} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} = \mathbf{A} \frac{\partial \mathbf{x}}{\partial \mathbf{z}}$ .

# VECTOR DERIVATIVES

## Vector Derivatives Rule 3

Let  $\alpha = \mathbf{y}^T \mathbf{A} \mathbf{x}$  where  $\mathbf{y} \in \mathbb{R}^{m \times 1}$ ,  $\mathbf{x} \in \mathbb{R}^{n \times 1}$ , and  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Then  $\frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{y}^T \mathbf{A}$  and  $\frac{\partial \alpha}{\partial \mathbf{y}} = \mathbf{x}^T \mathbf{A}^T$  ( $\mathbf{A}$  is constant matrix.)

- pf) Since  $\alpha = \mathbf{y}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{A}^T \mathbf{y}$ , we only prove  $\frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{y}^T \mathbf{A}$ .

$$\alpha = \mathbf{y}^T \mathbf{A} \mathbf{x} = \mathbf{w}^T \mathbf{x} = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$\frac{\partial \alpha}{\partial x_j} = w_j = [\mathbf{y}^T \mathbf{A}]_j$$

$$\therefore \frac{\partial \alpha}{\partial \mathbf{x}_j} = [\mathbf{y}^T \mathbf{A}]_j \implies \frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{y}^T \mathbf{A}$$



# VECTOR DERIVATIVES

## Vector Derivatives Rule 4 (Quadratic formula)

Let  $\alpha = \mathbf{x}^T \mathbf{A} \mathbf{x}$  where  $\mathbf{x} \in \mathbb{R}^{n \times 1}$ , and  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Then  $\frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T)$ . ( $\mathbf{A}$  is constant matrix.)

- **pf)** It helps to see what  $\alpha$  is made of.

$$\begin{aligned} \alpha &= \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\ &= \begin{bmatrix} \sum_{i=1}^n x_i a_{i1} & \sum_{i=1}^n x_i a_{i2} & \cdots & \sum_{i=1}^n x_i a_{in} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \sum_{j=1}^n x_j \sum_{i=1}^n x_i a_{ij} = \sum_{j=1}^n \sum_{i=1}^n x_i x_j a_{ij} \end{aligned}$$

# VECTOR DERIVATIVES

## Vector Derivatives Rule 4 (Quadratic formula) (continued)

Let  $\alpha = \mathbf{x}^T \mathbf{A} \mathbf{x}$  where  $\mathbf{x} \in \mathbb{R}^{n \times 1}$ , and  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Then  $\frac{\partial \alpha}{\partial \mathbf{x}} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T)$ . ( $\mathbf{A}$  is constant matrix.)

- **pf)** It helps to see what  $\alpha$  is made of.

$$\alpha = \sum_{j=1}^n \sum_{i=1}^n x_i x_j a_{ij}$$

Note that each  $x_j$  has coefficients 1) along the row  $(\sum_{i=1}^n a_{ji} x_i)$  and 2) along the column  $(\sum_{i=1}^n a_{ij} x_i)$ . Therefore,

$$\begin{aligned} \frac{\partial \alpha}{\partial x_j} &= \sum_{i=1}^n a_{ji} x_i + \sum_{i=1}^n a_{ij} x_i \\ \frac{\partial \alpha}{\partial \mathbf{x}} &= \mathbf{x}^T \mathbf{A} + \mathbf{x}^T \mathbf{A}^T = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T) \quad \square \end{aligned}$$

# VECTOR DERIVATIVES

## Vector Derivatives Rule 5 (Quadratic formula)

Let  $\alpha = \mathbf{y}^T \mathbf{x}$  where  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{n \times 1}$ , and each is a function of  $\mathbf{z}$ . Then  $\frac{\partial \alpha}{\partial \mathbf{z}} = \mathbf{x}^T \frac{\partial \mathbf{y}}{\partial \mathbf{z}} + \mathbf{y}^T \frac{\partial \mathbf{x}}{\partial \mathbf{z}}$ .

- pf) In element-wise expansion,

$$\begin{aligned}\alpha &= x_1 y_1 + x_2 y_2 + \dots + x_n y_n \\ \frac{\partial \alpha}{\partial z_k} &= \sum_{i=1}^n y_i \frac{\partial x_i}{\partial z_k} + \sum_{j=1}^n x_j \frac{\partial y_j}{\partial z_k} = \sum_{i=1}^n y_i \frac{\partial \mathbf{x}}{\partial \mathbf{z}}_{(i,k)} + \sum_{j=1}^n x_j \frac{\partial \mathbf{y}}{\partial \mathbf{z}}_{(j,k)} \\ &= [\mathbf{y}^T \frac{\partial \mathbf{x}}{\partial \mathbf{z}}]_k + [\mathbf{x}^T \frac{\partial \mathbf{y}}{\partial \mathbf{z}}]_k \quad \square\end{aligned}$$

Or with the chain rule,  $\frac{\partial \alpha}{\partial \mathbf{z}} = \frac{\partial \alpha}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{z}} + \frac{\partial \alpha}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} = \mathbf{x}^T \frac{\partial \mathbf{y}}{\partial \mathbf{z}} + \mathbf{y}^T \frac{\partial \mathbf{x}}{\partial \mathbf{z}}$ .  $\square$

# VECTOR DERIVATIVES

## Vector Derivatives Rule 6 (Quadratic formula)

Let  $\alpha = \mathbf{y}^T \mathbf{A} \mathbf{x}$  where  $\mathbf{y} \in \mathbb{R}^{m \times 1}$ ,  $\mathbf{x} \in \mathbb{R}^{n \times 1}$ , both a function of  $\mathbf{z}$  and  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Then  $\frac{\partial \alpha}{\partial \mathbf{z}} = \mathbf{y}^T \mathbf{A} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} + \mathbf{x}^T \mathbf{A}^T \frac{\partial \mathbf{y}}{\partial \mathbf{z}}$ . ( $\mathbf{A}$  is constant matrix.)

• pf) Let  $\mathbf{y}^T \mathbf{A} = \mathbf{w}^T$ . By the last rule,

$$\frac{\partial \alpha}{\partial \mathbf{z}} = \mathbf{w}^T \frac{\partial \mathbf{x}}{\partial \mathbf{z}} + \mathbf{x}^T \frac{\partial \mathbf{w}}{\partial \mathbf{z}} = \mathbf{y}^T \mathbf{A} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} + \mathbf{x}^T \frac{\partial \mathbf{A}^T \mathbf{y}}{\partial \mathbf{z}} = \mathbf{y}^T \mathbf{A} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} + \mathbf{x}^T \mathbf{A}^T \frac{\partial \mathbf{y}}{\partial \mathbf{z}} \quad \square$$

Or with the chain rule,

$$\frac{\partial \alpha}{\partial \mathbf{z}} = \frac{\partial \alpha}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} + \frac{\partial \alpha}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{z}} = \mathbf{y}^T \mathbf{A} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} + \mathbf{x}^T \mathbf{A}^T \frac{\partial \mathbf{y}}{\partial \mathbf{z}}. \quad \square$$



# VECTOR DERIVATIVES

## Vector Derivatives Rule 7 (Quadratic formula)

Let  $\alpha = \mathbf{x}^T \mathbf{A} \mathbf{x}$  where  $\mathbf{x} \in \mathbb{R}^{n \times 1}$ , a function of  $\mathbf{z}$ , and  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Then  $\frac{\partial \alpha}{\partial \mathbf{z}} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T) \frac{\partial \mathbf{x}}{\partial \mathbf{z}}$ . ( $\mathbf{A}$  is constant matrix.)

• **pf)** Since if  $\alpha = \mathbf{y}^T \mathbf{A} \mathbf{x}$ , then  $\frac{\partial \alpha}{\partial \mathbf{z}} = \mathbf{y}^T \mathbf{A} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} + \mathbf{x}^T \mathbf{A}^T \frac{\partial \mathbf{y}}{\partial \mathbf{z}}$ , we have

$$\frac{\partial \alpha}{\partial \mathbf{z}} = \mathbf{x}^T \mathbf{A} \frac{\partial \mathbf{x}}{\partial \mathbf{z}} + \mathbf{x}^T \mathbf{A}^T \frac{\partial \mathbf{x}}{\partial \mathbf{z}} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T) \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \quad \square$$

# LEAST SQUARES ESTIMATORS

## Example: Least Squares Estimators

- For a target vector  $\mathbf{t}$  and a design matrix  $\Phi$ , LSE  $\mathbf{w}$  can be obtained by;

$$\mathbf{w} = \arg \min_{\mathbf{w}} \|\mathbf{t} - \Phi \mathbf{w}\|^2 = (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w})$$

Differentiate wrt  $\mathbf{w}$  and we have

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{t}^T \mathbf{t} - 2\mathbf{t}^T \Phi \mathbf{w} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}) \stackrel{set}{=} 0$$

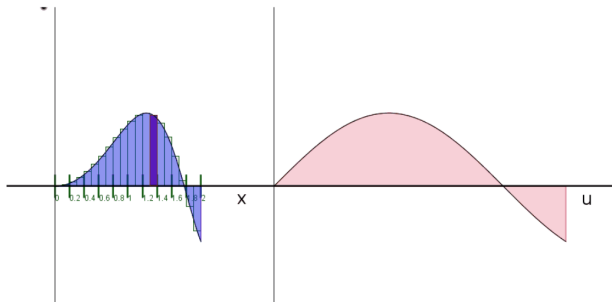
$$\therefore \mathbf{w}_{OLS} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

## II. Change of Variables

# WHY NEED JACOBIAN?

## Recall: Change of Variable Technique for Integration

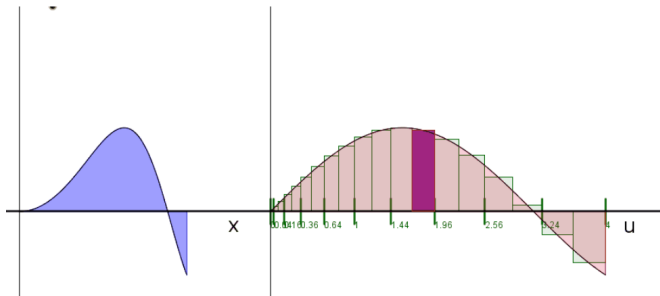
- For  $\int_0^2 \sin(x^2)dx$ , substitute  $u := x^2$  to have  $\int_0^2 \sin(x^2)dx = \int_0^4 \sin(u)dx$ .
- By substituting  $u := x^2$ , the coordinate space has been extended by  $\frac{du}{dx}$ .



# WHY NEED JACOBIAN?

## Recall: Change of Variable Technique for Integration

- We multiply by  $\frac{dx}{du}$  to mitigate this elongation;  $\int_0^4 \sin(u) dx = \int_0^4 \sin(u) \frac{dx}{du} du = \int_0^4 \frac{\sin(u)}{2\sqrt{u}} du$
- Generalization of  $\frac{dx}{du}$  to nd space is det of Jacobian;  $|\frac{\partial \mathbf{x}}{\partial \mathbf{u}}|$ . **It comes from the amount that the area is stretched under the coordinate transformation.** (source)



# WHY NEED JACOBIAN?

## Change of Variable in Probability Distribution

- COV for probability distribution is identical to COV for definite integral. For  $x \sim p(x)$ ,  
 $u := f(x) \sim p(f^{-1}(u)) \left| \frac{dx}{du} \right|$ .
- **Example:** Let  $x \sim 3x^2$  over  $[0, 1]$ . Define  $u := x^2$ . Then  $u \sim 3u \left| \frac{1}{2} y^{-1/2} \right| = \frac{3}{2} u^{1/2}$ .
- For multivariable case where  $\mathbf{y} = \mathbf{f}(\mathbf{x})$ , we have

$$p(\mathbf{y}) = p(\mathbf{x}) \left| \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right|$$

where  $p(\mathbf{x})$  is expressed in terms of  $\mathbf{y}$ .

# MVN: GEOMETRIC INTERPRETATION

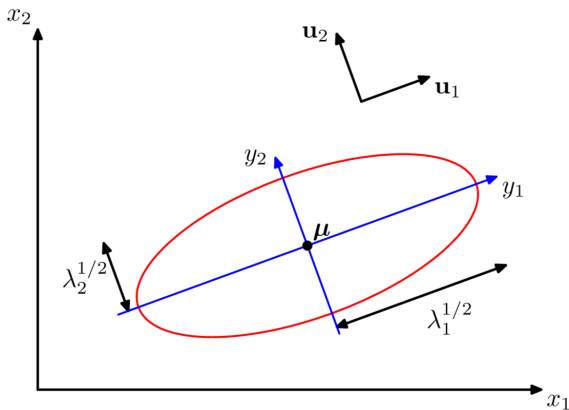
## Change of Basis to Eigenvectors of $\Sigma$

$$\text{MVN: } N(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{(\Sigma)^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\}$$

- Since  $\Sigma$  is symmetric and positive definite, we have eigenvalue decomposition of  $\Sigma$  with strictly positive  $\lambda_i$ ;  $\Sigma = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T$ , and  $\Sigma^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$
- Substitute  $\Sigma^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$  and we have  $(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu) = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$ , where we defined  $y_i := \mathbf{u}_i^T(\mathbf{x} - \mu)$ , which is a coordinate of  $(\mathbf{x} - \mu)$  projected to a unit vector  $\mathbf{u}$ .
- Define  $\mathbf{y} := \mathbf{U}^T(\mathbf{x} - \mu)$  where  $\mathbf{U}$  has  $\mathbf{u}_i$  as column. Then we have a det of jacobian of 1 since  $|\mathbf{J}|^2 = |\mathbf{U}^T|^2 = |\mathbf{U}^T| |\mathbf{U}| = |\mathbf{I}|$ . Thus we have  $p(\mathbf{y}) = \prod_{i=1}^D \frac{1}{(2\pi\lambda_i)^{1/2}} \exp\left\{-\sum_{i=1}^D \frac{y_i^2}{2\lambda_i}\right\}$ .

# MVN: GEOMETRIC INTERPRETATION

## Change of Basis to Eigenvectors of $\Sigma$



PCA derives from the exactly same logic, except that sample cov matrix  $\Phi^T \Phi$  is used in lieu of  $\Sigma$ .



### III. Matrix Derivatives

# MATRIX DERIVATIVES

## Matrix Derivatives

Let  $\mathbf{A}$  be  $m \times n$  matrix. Derivative of  $\mathbf{A}$  with respect to a scalar  $k$  is;

$$\frac{\partial \mathbf{A}}{\partial \alpha} = \begin{bmatrix} \frac{\partial a_{11}}{\partial \alpha} & \frac{\partial a_{12}}{\partial \alpha} & \cdots & \frac{\partial a_{1n}}{\partial \alpha} \\ \frac{\partial a_{21}}{\partial \alpha} & \frac{\partial a_{22}}{\partial \alpha} & \cdots & \frac{\partial a_{2n}}{\partial \alpha} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial a_{m1}}{\partial \alpha} & \frac{\partial a_{m2}}{\partial \alpha} & \cdots & \frac{\partial a_{mn}}{\partial \alpha} \end{bmatrix}_{m \times n}$$

## Matrix Derivatives Rule:

Let  $\mathbf{A}$  be nonsingular. Then  $\frac{\partial \mathbf{A}^{-1}}{\partial \alpha} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \alpha} \mathbf{A}^{-1}$

- **pf)** Since  $\mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$ ,  $\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \alpha} + \frac{\partial \mathbf{A}^{-1}}{\partial \alpha} \mathbf{A} = \mathbf{0}$ . Rearrange and we have  $\frac{\partial \mathbf{A}^{-1}}{\partial \alpha} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \alpha} \mathbf{A}^{-1}$ .



# MATRIX DERIVATIVES: DETERMINANT

## Directional Derivative in Matrix Space

- Since  $\det : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}$ , we first need to define directional derivative in matrix space. For a mapping  $f : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}$  from a matrix to a real number, a differential in  $f$  in response to an infinitesimal change in the input  $\mathbf{X}$  in the direction of  $\mathbf{H}$  is defined as;

$$\text{Directional Derivative in Matrix space: } \nabla_{\mathbf{H}} f(\mathbf{X}) = \lim_{\epsilon \rightarrow 0} \frac{f(\mathbf{X} + \epsilon \mathbf{H}) - f(\mathbf{X})}{\epsilon}$$

There are at least dozens of alternative notations for this;  $Df(\mathbf{X})(\mathbf{H})$ ,  $f'(\mathbf{X})\mathbf{H}$  ( $f'(\mathbf{X})\mathbf{H}$  is NOT a matrix multiplication. In fact, it should be written  $\langle f'(\mathbf{X}), \mathbf{H} \rangle_F$  where  $F$  denotes Frobenius inner product.)

- For comparison, in scalar space and vector space we had;

$$\text{Derivative in scalar space: } f'(x) = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x)}{\epsilon}$$

$$\text{Derivative in Vector space: } \nabla_{\mathbf{h}} f(\mathbf{X}) = \lim_{\epsilon \rightarrow 0} \frac{f(\mathbf{x} + \epsilon \mathbf{h}) - f(\mathbf{x})}{\epsilon}$$

# MATRIX DERIVATIVES: DETERMINANT

## Directional Derivative for $\det$

- Armed with this concept, we can define a directional derivative for  $\det$ ;

$$\nabla_{\mathbf{H}} \det(\mathbf{X}) = \det'(\mathbf{X})\mathbf{H} = \lim_{\epsilon \rightarrow 0} \frac{\det(\mathbf{X} + \epsilon\mathbf{H}) - \det(\mathbf{X})}{\epsilon}$$

and prove some useful facts.

**Fact 1.**  $\det'(\mathbf{I})\mathbf{H} = \text{tr}(\mathbf{H})$

- pf)**  $\det(\mathbf{I} + \epsilon\mathbf{H}) = 1 + \epsilon \text{tr}(\mathbf{H}) + o(\epsilon)$ , where  $o(\epsilon)$  consists of terms with  $2 \leq$  degree in  $\epsilon$ , so when divided by  $\epsilon$ , gets annihilated as  $\epsilon \rightarrow 0$ . To see this, recall that  $\det$  is a sum of signed products of one element of each column that are NOT on the same rows. This makes the diagonal product the only term with  $\epsilon$  of 1st degree.  $\square$
- This tells us that  $\text{tr}$  is more than a boring sum of diagonal element. If  $\mathbf{T}$  moves in the direction of  $\mathbf{H}$ , then its determinant changes accordingly by a factor of  $\text{tr}(\mathbf{H})$ .

# MATRIX DERIVATIVES: DETERMINANT

**Fact 2.**  $\det'(\mathbf{X})\mathbf{H} = \det(\mathbf{X})\text{tr}(\mathbf{X}^{-1}\mathbf{H})$

- **pf)** The proof is similar to that of Fact 1;

$$\begin{aligned}\det(\mathbf{X} + \epsilon\mathbf{H}) &= \det(\mathbf{X}) \det(\mathbf{I} + \epsilon\mathbf{X}^{-1}\mathbf{H}) \\ &= \det(\mathbf{X}) (1 + \epsilon \text{tr}(\mathbf{X}^{-1}\mathbf{H}) + o(\epsilon)) \\ &= \det(\mathbf{X}) + \epsilon \det(\mathbf{X}) \text{tr}(\mathbf{X}^{-1}\mathbf{H}) + o(\epsilon) \quad \square\end{aligned}$$

**Fact 3.**  $\frac{\partial}{\partial \mathbf{X}} \det(\mathbf{X}) = \det(\mathbf{X})\mathbf{X}^{-T}$ . Equivalently,  $\frac{\partial}{\partial \mathbf{X}} \ln \det(\mathbf{X}) = \mathbf{X}^{-T}$ .

- **pf)** We use  $\langle \mathbf{A}, \mathbf{B} \rangle_F = \text{tr}(\mathbf{A}^T \mathbf{B})$ . From the Fact 2,

$$\begin{aligned}\det'(\mathbf{X})\mathbf{H} &= \langle \det'(\mathbf{X}), \mathbf{H} \rangle_F = \det(\mathbf{X})\text{tr}(\mathbf{X}^{-1}\mathbf{H}) \\ &= \langle \det(\mathbf{X})\mathbf{X}^{-T}, \mathbf{H} \rangle_F \quad \square\end{aligned}$$

# MATRIX DERIVATIVES: DETERMINANT

**Fact 4.**  $\frac{\partial}{\partial x} \det(\mathbf{X}) = \det(\mathbf{X}) \text{tr}(\mathbf{X}^{-1} \frac{\partial \mathbf{X}}{\partial x})$ .

• **pf)** This follows from Fact 2. We had;

$$\det'(\mathbf{X})\mathbf{H} = \det(\mathbf{X})\text{tr}(\mathbf{X}^{-1}\mathbf{H})$$

Now put  $\mathbf{H} = \frac{\partial \mathbf{X}}{\partial x}$ , which is a matrix itself, and by the Chain Rule we have;

$$\frac{\partial \det \mathbf{X}}{\partial \mathbf{X}} \frac{\partial \mathbf{X}}{\partial x} = \frac{\partial}{\partial x} \det(\mathbf{X}) = \det(\mathbf{X}) \text{tr}(\mathbf{X}^{-1} \frac{\partial \mathbf{X}}{\partial x}) \quad \square$$

# MATRIX DERIVATIVES: TRACE

These are relatively intuitive, and the proof is literally no more than laying out elementwise multiplications.

$$\frac{\partial}{\partial A_{ij}} \text{tr}(\mathbf{AB}) = B_{ji}$$

$$\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{AB}) = \mathbf{B}^T$$

$$\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{AB}^T) = \frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{A}^T \mathbf{B}) = \mathbf{B}$$

$$\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{A}) = \mathbf{I}$$

$$\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{ABA}^T) = \mathbf{A}(\mathbf{B} + \mathbf{B}^T)$$

# MLE FOR MULTIVARIATE GAUSSIAN

Let's say we have a data  $\mathbf{X}$  consists of  $N$  observations with  $D$  features, with each observation is iid sample of  $MVN(\mu, \Sigma)$ . Then the joint likelihood of  $\mathbf{X}$  is given as;

$$p(\mathbf{X}|\mu, \Sigma) = (2\pi)^{-ND/2} \Sigma^{-N/2} \prod_{n=1}^N \exp\left(-\frac{1}{2}(\mathbf{x}_n - \mu)^T \Sigma^{-1}(\mathbf{x}_n - \mu)\right)$$

Taking a log we have

$$\ln p(\mathbf{X}|\mu, \Sigma) = -\frac{ND}{2} \ln 2\pi - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1}(\mathbf{x}_n - \mu)$$

We can get MLE for  $\mu, \Sigma$  by differentiating the above formula by each and find  $\mu_{ml}, \Sigma_{ml}$  that sets it to zero.



# MLE FOR MULTIVARIATE GAUSSIAN

**Differentiating w.r.t.  $\mu$ :** By the Vector Derivatives Rule 4,

$$\frac{\partial \ln p}{\partial \mu} = \sum_{n=1}^N \Sigma^{-1}(\mathbf{x}_n - \mu) \stackrel{set}{=} 0 \quad \rightarrow \quad \mu_{ml} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

**Differentiating w.r.t.  $\Sigma$ :** This is the tricky part. We have;

$$\frac{\partial \ln p}{\partial \Sigma} = -\frac{N}{2} \frac{\partial \ln |\Sigma|}{\partial \Sigma} - \frac{1}{2} \frac{\partial}{\partial \Sigma} \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu)$$

We know from Fact 3 that  $\frac{\partial \ln |\Sigma|}{\partial \Sigma} = \Sigma^{-T}$ . The second term is a headache. It helps us to know that for any square matrix  $\mathbf{A}$  and vector  $\mathbf{x}$ ,

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \text{tr}(\mathbf{A} \mathbf{x} \mathbf{x}^T)$$

# MLE FOR MULTIVARIATE GAUSSIAN

Differentiating w.r.t.  $\Sigma$ : (continued)

$$\begin{aligned}
 \text{tr}(\mathbf{A}\mathbf{x}\mathbf{x}^T) &= \text{tr}\left( \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1x_1 & x_1x_2 & \dots & x_1x_n \\ x_2x_1 & x_2x_2 & \dots & x_2x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_nx_1 & x_nx_2 & \dots & x_nx_n \end{bmatrix} \right) \\
 &= \sum_{i=1}^n x_1 a_{1i} x_i + \dots + \sum_{i=1}^n x_n a_{ni} x_i \\
 &= \sum_{j=1}^n \sum_{i=1}^n x_j a_{ji} x_i \\
 &= \mathbf{x}^T \mathbf{A} \mathbf{x}
 \end{aligned}$$

# MLE FOR MULTIVARIATE GAUSSIAN

## Differentiating w.r.t. $\Sigma$ : (continued)

Hence we have  $(\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) = \text{tr}(\Sigma^{-1} (\mathbf{x}_n - \mu)(\mathbf{x}_n - \mu)^T)$ , and the second term can be expressed as;  $\sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) = N \text{tr}(\Sigma^{-1} \mathbf{S})$ , where  $\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu)(\mathbf{x}_n - \mu)^T$ .

$$\begin{aligned} \frac{\partial}{\partial \Sigma_{ij}} \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) &= N \frac{\partial}{\partial \Sigma_{ij}} \text{tr}(\Sigma^{-1} \mathbf{S}) = N \text{tr} \left( \frac{\partial}{\partial \Sigma_{ij}} \Sigma^{-1} \mathbf{S} \right) \\ &= -N \text{tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial \Sigma_{ij}} \Sigma^{-1} \mathbf{S} \right) \\ &= -N \text{tr} \left( \frac{\partial \Sigma}{\partial \Sigma_{ij}} \Sigma^{-1} \mathbf{S} \Sigma^{-1} \right) \\ &= -N (\Sigma^{-1} \mathbf{S} \Sigma^{-1})_{ij} \end{aligned}$$

The last line follows since  $\left[ \frac{\partial \Sigma}{\partial \Sigma_{ij}} \right]_{ij} = 0$  for all but the position  $(i, j)$ .

# MLE FOR MULTIVARIATE GAUSSIAN

## Differentiating w.r.t. $\Sigma$ : (continued)

This makes the derivative in the second term

$$-\frac{1}{2} \frac{\partial}{\partial \Sigma} \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) = \frac{1}{2} N (\Sigma^{-1} \mathbf{S} \Sigma^{-1})$$

To sum up, rewriting  $\frac{\partial \ln p}{\partial \Sigma} = -\frac{N}{2} \frac{\partial \ln |\Sigma|}{\partial \Sigma} - \frac{1}{2} \frac{\partial}{\partial \Sigma} \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu)$ , we have

$$-\frac{N}{2} \Sigma^{-T} + \frac{N}{2} \Sigma^{-1} \mathbf{S} \Sigma^{-1} \stackrel{set}{=} \mathbf{0}$$

and this yields

$$\Sigma_{ml} = \mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu)(\mathbf{x}_n - \mu)^T$$

# References

- <https://atmos.washington.edu/~dennis/MatrixCalculus.pdf>
- <https://www.comp.nus.edu.sg/~cs5240/lecture/matrix-differentiation.pdf>
- <https://math.stackexchange.com/questions/1151569/how-to-calculate-the-derivative-of-log-det-matrix>
- <https://math.stackexchange.com/questions/38701/how-to-calculate-the-gradient-of-log-det-matrix-inverse>
- <http://www.math.ucdenver.edu/~esulliva/Calculus3/Taylor.pdf>
- <https://www.quora.com/What-is-the-Jacobian-how-does-it-work-and-what-is-an-intuitive-explanation-of-the-Jacobian>
- Pattern Recognition and machine learning, Christopher M. Bishop
- Advanced Engineering Mathematics, Erwin Kreyszig