

## 2 Convexity, Derivative

### 2.1 Lines, hyperplanes and linear varieties

- The line segment between two points  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  is the set,

$$\{\mathbf{z} \in \mathbb{R}^n : \mathbf{z} = \alpha \mathbf{x} + (1 - \alpha) \mathbf{y}, \alpha \in [0, 1]\}.$$

- A hyperplane of the space  $\mathbb{R}^n$ , is the set of all points  $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$  that satisfy the linear equation

$$u_1 x_1 + u_2 x_2 + \dots + u_n x_n = v,$$

where at least one of the  $u_i$  is nonzero. The hyperplane is defined by

$$\{\mathbf{x} \in \mathbb{R}^n : \mathbf{u}^\top \mathbf{x} = v\},$$

where

$$\mathbf{u} = [u_1, u_2, \dots, u_n]^\top.$$

- Two half-spaces, positive half-space and negative half-space are

$$H_+ = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{u}^\top \mathbf{x} \geq v\},$$

$$H_- = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{u}^\top \mathbf{x} \leq v\}.$$

- A linear variety is a set of form

$$\{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\},$$

for some matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and vector  $\mathbf{b} \in \mathbb{R}^m$ .

### 2.2 Convex sets

- A point  $\mathbf{w} = \alpha \mathbf{u} + (1 - \alpha) \mathbf{v}$  (where  $\alpha \in [0, 1]$ ) is called a convex combination of the points  $\mathbf{u}$  and  $\mathbf{v}$ .

- A set  $\Theta \subset \mathbb{R}^n$  is convex if for all  $\mathbf{u}, \mathbf{v} \in \Theta$ , the *line segment* between  $\mathbf{u}$  and  $\mathbf{v}$  is in  $\Theta$ .

That is,  $\Theta$  is *convex* if and only if  $\alpha \mathbf{u} + (1 - \alpha) \mathbf{v} \in \Theta$  for all  $\mathbf{u}, \mathbf{v} \in \Theta$  and  $\alpha \in (0, 1)$ . Examples of convex sets include the following:

- |                                      |                    |
|--------------------------------------|--------------------|
| – The empty set                      | – A hyperplane     |
| – A set consisting of a single point | – A linear variety |
| – A line or a line segment           | – A half-space     |
| – A subspace                         | – $\mathbb{R}^n$   |

♠ THEOREM 4.3 Convex subsets of  $\mathbb{R}^n$  have the following properties:

- a. If  $\Theta$  is a *convex set* and  $\beta$  is a real number, then the set

$$\beta\Theta = \{\mathbf{x} : \mathbf{x} = \beta\mathbf{v}, \mathbf{v} \in \Theta\}$$

is also convex.

- b. If  $\Theta_1$  and  $\Theta_2$  are *convex sets*, then the set

$$\Theta_1 + \Theta_2 = \{\mathbf{x} : \mathbf{x} = \mathbf{v}_1 + \mathbf{v}_2, \mathbf{v}_1 \in \Theta_1, \mathbf{v}_2 \in \Theta_2\}$$

is also convex.

- c. The intersection of any collection of *convex sets* is convex.

- An extreme point  $\mathbf{x}$  in a *convex set*  $\Theta$ , if there are no two distinct points  $\mathbf{u}$  and  $\mathbf{v}$  in  $\Theta$  such that  $\mathbf{x} = \alpha\mathbf{u} + (1 - \alpha)\mathbf{v}$  for some  $\alpha \in (0, 1)$ .

## 2.3 Differentiation rules

- A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  follows,

$$\begin{aligned} f(\mathbf{x}) &= f\left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}\right) = a_1x_1 + a_2x_2 + \cdots + a_nx_n = \begin{bmatrix} a_1 & \cdots & a_n \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{a}^\top \end{bmatrix} \begin{bmatrix} \vdots \\ \mathbf{x} \\ \vdots \end{bmatrix} = \begin{bmatrix} \mathbf{x}^\top \end{bmatrix} \begin{bmatrix} \vdots \\ \mathbf{a} \\ \vdots \end{bmatrix}. \end{aligned}$$

- A matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{A}$  is  $m \times n$  matrix,

$$\mathbf{A} = \begin{bmatrix} \vdots & \vdots & \cdots & \vdots \\ \mathbf{a}_{*1} & \mathbf{a}_{*2} & \cdots & \mathbf{a}_{*n} \\ \vdots & \vdots & \cdots & \vdots \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^\top \\ \mathbf{a}_2^\top \\ \vdots \\ \mathbf{a}_m^\top \end{bmatrix}.$$

- A function  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{Ax}$  is a column vector whose element is a scalar  $g_*(\mathbf{x})$ .

$$\mathbf{Ax} = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_m^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{a}_1^\top \mathbf{x} \\ \vdots \\ \mathbf{a}_m^\top \mathbf{x} \end{bmatrix} = \begin{bmatrix} \boxed{a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n} \\ \vdots \\ \boxed{a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n} \end{bmatrix} = \begin{bmatrix} g_1(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix} = \mathbf{g}(\mathbf{x}).$$

- To be noted, in this course, we write the **derivative**  $Df(\mathbf{x})$  as a **row vector**, and write the **gradient**  $\nabla f(\mathbf{x})$  as a **column vector**.

# Types of Matrix Derivatives<sup>1</sup>

Types	Scalar		Vector		Matrix	
Scalar	$\frac{dy}{dx}$	$\frac{df(x)}{dx}$ (1)	$\frac{d\mathbf{y}}{dx} = \begin{bmatrix} \frac{\partial y_i}{\partial x} \end{bmatrix}$	$\frac{d\mathbf{g}(t)}{dt}$ (3)	$\frac{d\mathbf{Y}}{dx} = \begin{bmatrix} \frac{\partial y_{ij}}{\partial x} \end{bmatrix}$	$\frac{d\mathbf{A}(t)}{dt}$
Vector	$\frac{d\mathbf{y}}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_j} \end{bmatrix}$	$D_{\mathbf{x}}f(\mathbf{x}) = \begin{bmatrix} \cdot \frac{\partial f(\mathbf{x})}{\partial x_j} \cdot \end{bmatrix}$ $\nabla_{\mathbf{x}}f(\mathbf{x}) = \begin{bmatrix} \cdot \\ \frac{\partial f(\mathbf{x})}{\partial x_j} \\ \cdot \end{bmatrix}$ (2)	$\frac{d\mathbf{y}}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial y_i}{\partial x_j} \end{bmatrix}$	$D_{\mathbf{x}}\mathbf{g}(\mathbf{x}) = \begin{bmatrix} \frac{\partial g_i(\mathbf{x})}{\partial x_j} \end{bmatrix}$ (4)		
Matrix	$\frac{d\mathbf{y}}{d\mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{ji}} \end{bmatrix}$	$D_{\mathbf{X}}f = \begin{bmatrix} \frac{\partial f}{\partial x_{ji}} \end{bmatrix}$				

(1) Given  $f : \mathbb{R} \rightarrow \mathbb{R}$ , if the limit exists, the derivative of  $f$  is a function  $f' : \mathbb{R} \rightarrow \mathbb{R}$  given by

$$D_x(f(x)) = \frac{df}{dx} = f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}.$$

(2) Given  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , consider a scalar  $f(\mathbf{x}) = a_1x_1 + a_2x_2 + \cdots + a_nx_n = \mathbf{a}^\top \mathbf{x}$ .

For **derivative** rule (2),

$$D_{\mathbf{x}}f(\mathbf{x}) = D(\mathbf{a}^\top \mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) & \frac{\partial f}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} a_1 & a_2 & \cdots & a_n \end{bmatrix} = \mathbf{a}^\top.$$

For **gradient** rule (2), if  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable, then the *gradient* of  $f$  is a function  $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  given by

$$\nabla_{\mathbf{x}}f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} = \mathbf{a} = D_{\mathbf{x}}f(\mathbf{x})^\top.$$

(3) Given  $\mathbf{g} : \mathbb{R} \rightarrow \mathbb{R}^m$ , here  $t \in \mathbb{R}$  is a scalar.  $\mathbf{g}(t)$  is a column vector.

$$\mathbf{g}(t) = \begin{bmatrix} g_1(t) \\ \vdots \\ g_m(t) \end{bmatrix}, \quad D_t\mathbf{g}(t) = \begin{bmatrix} \frac{d}{dt}g_1(t) \\ \vdots \\ \frac{d}{dt}g_m(t) \end{bmatrix} = \begin{bmatrix} g'_1(t) \\ \vdots \\ g'_m(t) \end{bmatrix}.$$

(4) Consider  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , here  $\mathbf{x} \in \mathbb{R}^n$  is a vector. Since  $g_i(\mathbf{x})$  is a scalar,  $\mathbf{g} = [g_1, \dots, g_m]^\top$ ,  $\mathbf{g}(\mathbf{x})$  is a column vector.

$$\mathbf{g}(\mathbf{x}) = \begin{bmatrix} g_1(\mathbf{x}) \\ g_2(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix}, \quad D_{\mathbf{x}}\mathbf{g}(\mathbf{x}) = \begin{bmatrix} D_{\mathbf{x}}g_1(x_1, x_2, \dots, x_n) \\ D_{\mathbf{x}}g_2(x_1, x_2, \dots, x_n) \\ \vdots \\ D_{\mathbf{x}}g_m(x_1, x_2, \dots, x_n) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1}g_1 & \frac{\partial}{\partial x_2}g_1 & \cdots & \frac{\partial}{\partial x_n}g_1 \\ \frac{\partial}{\partial x_1}g_2 & \frac{\partial}{\partial x_2}g_2 & \cdots & \frac{\partial}{\partial x_n}g_2 \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial}{\partial x_1}g_m & \frac{\partial}{\partial x_2}g_m & \cdots & \frac{\partial}{\partial x_n}g_m \end{bmatrix} = \mathbf{J}.$$

The matrix  $\mathbf{J}$  is called the Jacobian matrix, or derivative matrix, of function  $\mathbf{g}$ .

<sup>1</sup>Ref: Thomas P. Minka, "Old and New Matrix Algebra Useful for Statistics", 2000

- If all elements in  $\mathbf{g}(\mathbf{x})$  are linear combination of  $\mathbf{x}$ ,

$$\mathbf{g}(\mathbf{x}) = \begin{bmatrix} g_1(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^\top \mathbf{x} \\ \vdots \\ \mathbf{a}_m^\top \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_m^\top \end{bmatrix} \mathbf{x} = \mathbf{A}\mathbf{x}.$$

Then, the derivative of  $\mathbf{A}\mathbf{x}$  is equivalent to  $D_{\mathbf{x}}\mathbf{g}(\mathbf{x})$ ,

$$\mathcal{D}(\mathbf{g}(\mathbf{x})) = \underbrace{\frac{d}{d\mathbf{x}}(\mathbf{A}\mathbf{x})}_{\substack{\text{Notation not used} \\ \text{in this course}}} = D(\mathbf{A}\mathbf{x}) = \begin{bmatrix} D(\mathbf{a}_1^\top \mathbf{x}) \\ D(\mathbf{a}_2^\top \mathbf{x}) \\ \vdots \\ D(\mathbf{a}_m^\top \mathbf{x}) \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^\top \\ \mathbf{a}_2^\top \\ \vdots \\ \mathbf{a}_m^\top \end{bmatrix} = \mathbf{A}.$$

- In summary, the derivative rules are listed as,

$$\mathcal{D}(\mathbf{a}^\top \mathbf{x}) = \mathbf{a}^\top,$$

$$(2) f : \mathbb{R}^n \rightarrow \mathbb{R}, f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x}$$

$$\mathcal{D}(\mathbf{g}(t)) = \begin{bmatrix} \vdots \\ g'_*(t) \\ \vdots \end{bmatrix},$$

$$(3) \mathbf{g} : \mathbb{R} \rightarrow \mathbb{R}^n, \mathbf{g}(t) = \begin{bmatrix} \vdots \\ g_*(t) \\ \vdots \end{bmatrix}$$

$$\mathcal{D}(\mathbf{A}\mathbf{x}) = \mathbf{A},$$

$$(4) \mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m, \mathbf{g}(\mathbf{x}) = \mathbf{A}\mathbf{x}$$

$$\mathcal{D}(\mathbf{A}(\alpha\mathbf{x})) = \alpha\mathbf{A},$$

$$\frac{d}{d\alpha}(\mathbf{A}(\alpha\mathbf{x})) = \mathbf{A}\mathbf{x},$$

$$\nabla \mathbf{a}^\top \mathbf{x} = \mathbf{a},$$

$$(2) f : \mathbb{R}^n \rightarrow \mathbb{R}, f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x}$$

$$\nabla \mathbf{A}\mathbf{x} = \mathbf{A}^\top,$$

$$(4) \mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m, \mathbf{g}(\mathbf{x}) = \mathbf{A}\mathbf{x}$$

$$\nabla \mathbf{A}(\alpha\mathbf{x}) = \alpha\mathbf{A}^\top.$$

- Note that for  $\underline{f : \mathbb{R}^n \rightarrow \mathbb{R}}$ , we have

$$\nabla f(\mathbf{x}) = \mathcal{D}f(\mathbf{x})^\top.$$

## 2.4 Differentiation rules on composite function

- To differentiate the composite function,  $h(t) = f(\mathbf{g}(t))$  is differentiable on  $(a, b)$ , and

$$f(\mathbf{g}(t)) = f\left(\begin{bmatrix} g_1(t) \\ g_2(t) \\ \vdots \\ g_m(t) \end{bmatrix}\right) = a_1g_1(t) + a_2g_2(t) + \cdots + a_mg_m(t).$$

- The differentiated composite function with **derivative** rule is

$$h'(t) = D_{\mathbf{g}}f(\mathbf{g}(t))D_t\mathbf{g}(t) = \nabla f(\mathbf{g}(t))^\top \begin{bmatrix} g'_1(t) \\ \vdots \\ g'_m(t) \end{bmatrix} = \begin{bmatrix} \frac{d}{dg_1}f(\mathbf{g}(t)) & \cdots & \frac{d}{dg_m}f(\mathbf{g}(t)) \end{bmatrix} \begin{bmatrix} g'_1(t) \\ \vdots \\ g'_m(t) \end{bmatrix}.$$

- Consider Hessian matrix, which is second order derivative of scalar. Noted that,  $D(f(\mathbf{x}))$  is spreading the derivative of the polynomials on the horizontal direction. Thus, we would like to ensure each entry is located on a vertical direction, then the entry could be applied to conduct derivative, as

$$D^2(f(\mathbf{x})) = D(Df(\mathbf{x})^\top) = D(\nabla f(\mathbf{x}))$$

$$= \begin{bmatrix} D\left(\frac{\partial f}{\partial x_1}\right) \\ D\left(\frac{\partial f}{\partial x_2}\right) \\ \vdots \\ D\left(\frac{\partial f}{\partial x_n}\right) \end{bmatrix} = \begin{bmatrix} \boxed{\frac{\partial^2 f}{\partial x_1 \partial x_1} \quad \frac{\partial^2 f}{\partial x_2 \partial x_1} \quad \cdots \quad \frac{\partial^2 f}{\partial x_n \partial x_1}} \\ \boxed{\frac{\partial^2 f}{\partial x_1 \partial x_2} \quad \frac{\partial^2 f}{\partial x_2 \partial x_2} \quad \cdots \quad \frac{\partial^2 f}{\partial x_n \partial x_2}} \\ \vdots \\ \boxed{\frac{\partial^2 f}{\partial x_1 \partial x_n} \quad \frac{\partial^2 f}{\partial x_2 \partial x_n} \quad \cdots \quad \frac{\partial^2 f}{\partial x_n \partial x_n}} \end{bmatrix}.$$

- Similarly, the second order gradient of scalar is

$$\nabla^2(f(\mathbf{x})) = \nabla(\nabla f(\mathbf{x})^\top) = \nabla(Df(\mathbf{x})) = \nabla\left(\begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix}\right)$$

$$= \begin{bmatrix} \nabla\left(\frac{\partial f}{\partial x_1}\right) & \nabla\left(\frac{\partial f}{\partial x_2}\right) & \cdots & \nabla\left(\frac{\partial f}{\partial x_n}\right) \end{bmatrix} = \begin{bmatrix} \boxed{\frac{\partial^2 f}{\partial x_1 \partial x_1}} & \boxed{\frac{\partial^2 f}{\partial x_1 \partial x_2}} & \cdots & \boxed{\frac{\partial^2 f}{\partial x_1 \partial x_n}} \\ \boxed{\frac{\partial^2 f}{\partial x_2 \partial x_1}} & \boxed{\frac{\partial^2 f}{\partial x_2 \partial x_2}} & \cdots & \boxed{\frac{\partial^2 f}{\partial x_2 \partial x_n}} \\ \vdots & \vdots & \ddots & \vdots \\ \boxed{\frac{\partial^2 f}{\partial x_n \partial x_1}} & \boxed{\frac{\partial^2 f}{\partial x_n \partial x_2}} & \cdots & \boxed{\frac{\partial^2 f}{\partial x_n \partial x_n}} \end{bmatrix}.$$

## 2.5 Differentiation product rules

- i) Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  and  $g : \mathbb{R} \rightarrow \mathbb{R}$  be two differentiable functions,  $x \in \mathbb{R}$ ,

$$D(f(x)g(x)) = f(x)Dg(x) + g(x)Df(x),$$

$$\nabla(f(x)g(x)) = f(x)\nabla g(x) + g(x)\nabla f(x).$$

- ii) Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  be two differentiable functions,  $\mathbf{x} \in \mathbb{R}^n$ ,

$$D(f(\mathbf{x})g(\mathbf{x})) = f(\mathbf{x}) \begin{bmatrix} Dg(\mathbf{x}) \end{bmatrix} + g(\mathbf{x}) \begin{bmatrix} Df(\mathbf{x}) \end{bmatrix},$$

$$\nabla(f(\mathbf{x})g(\mathbf{x})) = f(\mathbf{x}) \begin{bmatrix} \nabla g(\mathbf{x}) \end{bmatrix} + g(\mathbf{x}) \begin{bmatrix} \nabla f(\mathbf{x}) \end{bmatrix}.$$

- iii) Let  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be two differentiable functions,  $\mathbf{x} \in \mathbb{R}^n$ ,

$$D(\mathbf{f}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})) = \begin{bmatrix} \mathbf{f}(\mathbf{x})^\top \end{bmatrix} \begin{bmatrix} D\mathbf{g}(\mathbf{x}) \end{bmatrix} + \begin{bmatrix} \mathbf{g}(\mathbf{x})^\top \end{bmatrix} \begin{bmatrix} D\mathbf{f}(\mathbf{x}) \end{bmatrix},$$

$$\nabla(\mathbf{f}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})) = \begin{bmatrix} \nabla \mathbf{f}(\mathbf{x}) \end{bmatrix} \begin{bmatrix} \mathbf{g}(\mathbf{x}) \end{bmatrix} + \begin{bmatrix} \nabla \mathbf{g}(\mathbf{x}) \end{bmatrix} \begin{bmatrix} \mathbf{f}(\mathbf{x}) \end{bmatrix}.$$

## 2.6 Differentiation rules

- If  $f = \mathbf{x}^\top \mathbf{Q} \mathbf{x}$  is a scalar,  $f^\top = f$ .
- If  $\mathbf{Q}$  is not symmetric, we can always replace it with a symmetric matrix,

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

Continue with manipulations,

$$\mathbf{x}^\top \mathbf{Q} \mathbf{x} = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} = \mathbf{x}^\top \left( \frac{\mathbf{Q} + \mathbf{Q}^\top}{2} \right) \mathbf{x},$$

where  $\frac{1}{2} (\mathbf{Q} + \mathbf{Q}^\top) = \frac{1}{2} (\mathbf{Q} + \mathbf{Q}^\top)^\top$  is a symmetric matrix.

- Based on the above **derivative** rule, we have
  1. Consider  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a given matrix and  $\mathbf{y} \in \mathbb{R}^m$  a given vector. Then,

$$\begin{aligned} D(\mathbf{y}^\top \mathbf{A} \mathbf{x}) &= \mathbf{y}^\top \mathbf{A}, \\ D(\mathbf{x}^\top \mathbf{A} \mathbf{x}) &= \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top). \end{aligned} \quad \boxed{\text{if } m = n}$$

2. Consider  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a given matrix and  $\mathbf{y} \in \mathbb{R}^n$  a given vector. Then,

$$D(\mathbf{y}^\top \mathbf{x}) = \mathbf{y}^\top.$$

3. Consider if  $\mathbf{Q}$  is a symmetric matrix, then

$$D(\mathbf{x}^\top \mathbf{Q} \mathbf{x}) = 2\mathbf{x}^\top \mathbf{Q}.$$

In particular,

$$D(\mathbf{x}^\top \mathbf{x}) = 2\mathbf{x}^\top.$$

- Based on the above **gradient** rule, we have
  1. Consider  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a given matrix and  $\mathbf{y} \in \mathbb{R}^m$  a given vector. Then,

$$\begin{aligned} \nabla(\mathbf{y}^\top \mathbf{A} \mathbf{x}) &= \mathbf{A}^\top \mathbf{y}, \\ \nabla(\mathbf{x}^\top \mathbf{A} \mathbf{x}) &= (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}. \end{aligned} \quad \boxed{\text{if } m = n}$$

2. Consider  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a given matrix and  $\mathbf{y} \in \mathbb{R}^n$  a given vector. Then,

$$\nabla(\mathbf{y}^\top \mathbf{x}) = \mathbf{y}.$$

3. Consider if  $\mathbf{Q}$  is a symmetric matrix, then

$$\nabla(\mathbf{x}^\top \mathbf{Q} \mathbf{x}) = 2\mathbf{Q} \mathbf{x}.$$

In particular,

$$\nabla(\mathbf{x}^\top \mathbf{x}) = 2\mathbf{x}.$$

## 2.7 Derivative details

- Let  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be two differentiable functions,  $\mathbf{x} \in \mathbb{R}^n$ ,

$$D\left(\mathbf{f}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})\right) = \mathbf{f}(\mathbf{x})^\top D\mathbf{g}(\mathbf{x}) + \mathbf{g}(\mathbf{x})^\top D\mathbf{f}(\mathbf{x}).$$

We can write it into compact matrix form as,

$$D\left(\begin{matrix} 1 \setminus m \\ \mathbf{f}(\mathbf{x})^\top \end{matrix} \begin{matrix} m \setminus 1 \\ \mathbf{g}(\mathbf{x}) \end{matrix}\right) = \begin{matrix} 1 \setminus m \\ \mathbf{f}(\mathbf{x})^\top \end{matrix} \begin{matrix} m \setminus n \\ D\mathbf{g}(\mathbf{x}) \end{matrix} + \begin{matrix} 1 \setminus m \\ \mathbf{g}(\mathbf{x})^\top \end{matrix} \begin{matrix} m \setminus n \\ D\mathbf{f}(\mathbf{x}) \end{matrix}.$$

**Short proof,**

$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\mathbf{x} \in \mathbb{R}^n$ , write the derivative as

$$\begin{aligned} D\left(\mathbf{f}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})\right) &= D\left(\begin{bmatrix} f_1(\mathbf{x}) & \cdots & f_m(\mathbf{x}) \end{bmatrix} \begin{bmatrix} g_1(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix}\right) \\ &= D\left(\begin{bmatrix} f_1(x_1, x_2, \dots, x_n) & \cdots & f_m(x_1, x_2, \dots, x_n) \end{bmatrix} \begin{bmatrix} g_1(x_1, x_2, \dots, x_n) \\ \vdots \\ g_m(x_1, x_2, \dots, x_n) \end{bmatrix}\right) \\ &= D\left(f_1(\mathbf{x})g_1(\mathbf{x}) + \cdots + f_m(\mathbf{x})g_m(\mathbf{x})\right) \\ &= f_1(\mathbf{x})Dg_1(\mathbf{x}) + \cdots + f_m(\mathbf{x})Dg_m(\mathbf{x}) \\ &\quad + g_1(\mathbf{x})Df_1(\mathbf{x}) + \cdots + g_m(\mathbf{x})Df_m(\mathbf{x}) \\ &= \begin{bmatrix} f_1(\mathbf{x}) & \cdots & f_m(\mathbf{x}) \end{bmatrix} \begin{bmatrix} Dg_1(\mathbf{x}) \\ \vdots \\ Dg_m(\mathbf{x}) \end{bmatrix} + \begin{bmatrix} g_1(\mathbf{x}) & \cdots & g_m(\mathbf{x}) \end{bmatrix} \begin{bmatrix} Df_1(\mathbf{x}) \\ \vdots \\ Df_m(\mathbf{x}) \end{bmatrix} \\ &= \begin{matrix} 1 \setminus m \\ \mathbf{f}(\mathbf{x})^\top \end{matrix} \begin{matrix} m \setminus n \\ D\mathbf{g}(\mathbf{x}) \end{matrix} + \begin{matrix} 1 \setminus m \\ \mathbf{g}(\mathbf{x})^\top \end{matrix} \begin{matrix} m \setminus n \\ D\mathbf{f}(\mathbf{x}) \end{matrix}.$$

- Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a given matrix and  $\mathbf{y} \in \mathbb{R}^m$  a given vector. Then,

$$\begin{aligned} D(\mathbf{y}^\top \mathbf{A} \mathbf{x}) &= \mathbf{y}^\top \mathbf{A}, \\ D(\mathbf{x}^\top \mathbf{A} \mathbf{x}) &= \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top). \quad \boxed{\text{if } m = n} \end{aligned}$$

Short proof,

$$\begin{aligned} D(\mathbf{y}^\top \mathbf{A} \mathbf{x}) &= D\left(\mathbf{y}^\top (\mathbf{A} \mathbf{x})\right) = D\left(\mathbf{f}(\mathbf{y})^\top (\mathbf{g}(\mathbf{x}))\right) \\ &= \mathbf{f}(\mathbf{y})^\top D\mathbf{g}(\mathbf{x}) + \mathbf{g}(\mathbf{x})^\top D\mathbf{f}(\mathbf{y}) \\ &= \mathbf{y}^\top D(\mathbf{A} \mathbf{x}) + (\mathbf{A} \mathbf{x})^\top [0] \\ &= \mathbf{y}^\top \mathbf{A}. \end{aligned}$$

If  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,

$$\begin{aligned} D(\mathbf{x}^\top \mathbf{A} \mathbf{x}) &= D\left(\mathbf{x}^\top (\mathbf{A} \mathbf{x})\right) = D\left(\mathbf{f}(\mathbf{x})^\top (\mathbf{g}(\mathbf{x}))\right) \\ &= \mathbf{f}(\mathbf{x})^\top D\mathbf{g}(\mathbf{x}) + \mathbf{g}(\mathbf{x})^\top D\mathbf{f}(\mathbf{x}) \\ &= \mathbf{x}^\top D(\mathbf{A} \mathbf{x}) + (\mathbf{A} \mathbf{x})^\top D(\mathbf{I}_n \mathbf{x}) \\ &= \mathbf{x}^\top \mathbf{A} + \mathbf{x}^\top \mathbf{A}^\top \mathbf{I} \\ &= \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top). \end{aligned}$$

- It follows that if  $\mathbf{Q}$  is a symmetric matrix, then

$$\begin{aligned} D(\mathbf{x}^\top \mathbf{Q} \mathbf{x}) &= 2\mathbf{x}^\top \mathbf{Q}, \\ D(\mathbf{x}^\top \mathbf{x}) &= 2\mathbf{x}^\top. \end{aligned}$$

Short proof,

$$\begin{aligned} D(\mathbf{x}^\top \mathbf{Q} \mathbf{x}) &= D\left(\mathbf{x}^\top (\mathbf{Q} \mathbf{x})\right) \\ &= \mathbf{x}^\top (\mathbf{Q} + \mathbf{Q}^\top) \\ &= 2\mathbf{x}^\top \mathbf{Q}, \\ D(\mathbf{x}^\top \mathbf{x}) &= D(\mathbf{x}^\top \mathbf{I} \mathbf{x}) = D\left(\mathbf{x}^\top (\mathbf{I} \mathbf{x})\right) \\ &= \mathbf{x}^\top (\mathbf{I} + \mathbf{I}^\top) \\ &= 2\mathbf{x}^\top. \end{aligned}$$

- Derivative of scalar by scalar,

$$\begin{aligned} \frac{d}{d\alpha} ((\alpha \mathbf{x})^\top \mathbf{Q} (\alpha \mathbf{x})) &= (\alpha \mathbf{x})^\top \frac{d}{d\alpha} (\mathbf{Q} (\alpha \mathbf{x})) + (\mathbf{Q} (\alpha \mathbf{x}))^\top \frac{d}{d\alpha} (\mathbf{I} (\alpha \mathbf{x})) \\ &= (\alpha \mathbf{x})^\top \mathbf{Q} \mathbf{x} + (\alpha \mathbf{x})^\top \mathbf{Q}^\top \mathbf{x} \\ &= 2(\alpha \mathbf{x})^\top \mathbf{Q} \mathbf{x} \\ &= 2\alpha \mathbf{x}^\top \mathbf{Q} \mathbf{x} \end{aligned}$$



## 2.8 Gradient details

- Consider  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , here  $\mathbf{x} \in \mathbb{R}^n$  is a vector. Since  $g_i(\mathbf{x})$  is a scalar,  $\mathbf{g} = [g_1, \dots, g_m]^\top$ ,  $\mathbf{g}(\mathbf{x})$  is a column vector.

$$\mathbf{g}(\mathbf{x}) = \begin{bmatrix} g_1(\mathbf{x}) \\ g_2(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix}, \nabla_{\mathbf{x}} \mathbf{g}(\mathbf{x}) = \begin{bmatrix} \vdots & \vdots & \vdots \\ \nabla_{\mathbf{x}} g_1 & \cdots & \nabla_{\mathbf{x}} g_m \\ \vdots & \vdots & \vdots \end{bmatrix} = \begin{bmatrix} \boxed{\frac{\partial}{\partial x_1} g_1} & \boxed{\frac{\partial}{\partial x_1} g_2} & \cdots & \boxed{\frac{\partial}{\partial x_1} g_m} \\ \boxed{\frac{\partial}{\partial x_2} g_1} & \boxed{\frac{\partial}{\partial x_2} g_2} & \cdots & \boxed{\frac{\partial}{\partial x_2} g_m} \\ \vdots & \vdots & \ddots & \vdots \\ \boxed{\frac{\partial}{\partial x_n} g_1} & \boxed{\frac{\partial}{\partial x_n} g_2} & \cdots & \boxed{\frac{\partial}{\partial x_n} g_m} \end{bmatrix}.$$

- Not standard derivations in this course, just offer some intuitions on how to obtained the gradients on production rules.
- Let  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be two differentiable functions,  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\nabla \left( \mathbf{f}(\mathbf{x})^\top \mathbf{g}(\mathbf{x}) \right) = \nabla \mathbf{f}(\mathbf{x}) \begin{bmatrix} \mathbf{g}(\mathbf{x}) \end{bmatrix} + \nabla \mathbf{g}(\mathbf{x}) \begin{bmatrix} \mathbf{f}(\mathbf{x}) \end{bmatrix}.$$

Proof.  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\mathbf{x} \in \mathbb{R}^n$ , write the derivative as

$$\begin{aligned} \nabla \left( \mathbf{f}(\mathbf{x})^\top \mathbf{g}(\mathbf{x}) \right) &= \nabla \left( \begin{bmatrix} f_1(\mathbf{x}) & \cdots & f_m(\mathbf{x}) \end{bmatrix} \begin{bmatrix} g_1(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix} \right) \\ &= \nabla \left( \begin{bmatrix} f_1(x_1, x_2, \dots, x_n) & \cdots & f_m(x_1, x_2, \dots, x_n) \end{bmatrix} \begin{bmatrix} g_1(x_1, x_2, \dots, x_n) \\ \vdots \\ g_m(x_1, x_2, \dots, x_n) \end{bmatrix} \right) \\ &= \nabla \left( f_1(\mathbf{x})g_1(\mathbf{x}) + \cdots + f_m(\mathbf{x})g_m(\mathbf{x}) \right) \\ &= f_1(\mathbf{x})\nabla g_1(\mathbf{x}) + \cdots + f_m(\mathbf{x})\nabla g_m(\mathbf{x}) \\ &\quad + g_1(\mathbf{x})\nabla f_1(\mathbf{x}) + \cdots + g_m(\mathbf{x})\nabla f_m(\mathbf{x}) \\ &= \begin{bmatrix} \nabla g_1(\mathbf{x}) & \cdots & \nabla g_m(\mathbf{x}) \end{bmatrix} \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} + \begin{bmatrix} \nabla f_1(\mathbf{x}) & \cdots & \nabla f_m(\mathbf{x}) \end{bmatrix} \begin{bmatrix} g_1(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix} \\ &= \begin{bmatrix} \nabla f_1(\mathbf{x}) & \cdots & \nabla f_m(\mathbf{x}) \end{bmatrix} \begin{bmatrix} g_1(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{bmatrix} + \begin{bmatrix} \nabla g_1(\mathbf{x}) & \cdots & \nabla g_m(\mathbf{x}) \end{bmatrix} \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} \\ &= \begin{matrix} n \setminus m \end{matrix} \begin{bmatrix} \nabla \mathbf{f}(\mathbf{x}) \end{bmatrix} \begin{matrix} m \setminus 1 \end{matrix} \begin{bmatrix} \mathbf{g}(\mathbf{x}) \end{bmatrix} + \begin{matrix} n \setminus m \end{matrix} \begin{bmatrix} \nabla \mathbf{g}(\mathbf{x}) \end{bmatrix} \begin{matrix} m \setminus 1 \end{matrix} \begin{bmatrix} \mathbf{f}(\mathbf{x}) \end{bmatrix} \end{aligned}$$

- For example, if  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,

$$\begin{aligned}
\nabla(\mathbf{x}^\top \mathbf{A} \mathbf{x}) &= \nabla(\mathbf{x}^\top (\mathbf{A} \mathbf{x})) \\
&= \nabla \left( \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix} \begin{bmatrix} \mathbf{a}_1^\top \mathbf{x} \\ \vdots \\ \mathbf{a}_n^\top \mathbf{x} \end{bmatrix} \right) \\
&= \nabla \left( x_1 \mathbf{a}_1^\top \mathbf{x} + \cdots + x_n \mathbf{a}_n^\top \mathbf{x} \right) \\
&= x_1 \nabla(\mathbf{a}_1^\top \mathbf{x}) + \cdots + x_n \nabla(\mathbf{a}_n^\top \mathbf{x}) \\
&\quad + \mathbf{a}_1^\top \mathbf{x} \nabla(1x_1 + 0x_2 + \cdots + 0x_n) + \cdots + \mathbf{a}_n^\top \mathbf{x} \nabla(0x_1 + 0x_2 + \cdots + 1x_n) \\
&= \mathbf{a}_1^\top \mathbf{x} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \cdots + \mathbf{a}_n^\top \mathbf{x} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} + \mathbf{a}_1 x_1 + \cdots + \mathbf{a}_n x_n \\
&= \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} \mathbf{a}_1^\top \mathbf{x} \\ \mathbf{a}_2^\top \mathbf{x} \\ \vdots \\ \mathbf{a}_n^\top \mathbf{x} \end{bmatrix} + \begin{bmatrix} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} \mathbf{x} \end{bmatrix} + \begin{bmatrix} \mathbf{A}^\top \end{bmatrix} \begin{bmatrix} \mathbf{x} \end{bmatrix} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x} \\
&= \begin{bmatrix} \nabla \mathbf{f}(\mathbf{x}) \end{bmatrix} \begin{bmatrix} \mathbf{g}(\mathbf{x}) \end{bmatrix} + \begin{bmatrix} \nabla \mathbf{g}(\mathbf{x}) \end{bmatrix} \begin{bmatrix} \mathbf{f}(\mathbf{x}) \end{bmatrix} \\
&= \nabla(\mathbf{f}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})),
\end{aligned}$$

where  $\mathbf{f}(\mathbf{x}) = \mathbf{x}$ ,  $\mathbf{g}(\mathbf{x}) = \mathbf{A} \mathbf{x}$ .

- Hand writing derivation: given  $\underline{y} \in \mathbb{R}^m$ ,  $\underline{A} \in \mathbb{R}^{m \times n}$ ,  $\underline{x} \in \mathbb{R}^n$ , write the derivative as

$$\begin{aligned}
D_{\underline{x}}(\underline{y}^\top \mathbf{A} \underline{x}) &= D_{\underline{x}}(\underline{y}^\top \underline{A} \underline{x}) \\
&= D \left( {}^{1 \setminus m} \begin{bmatrix} \underline{y}^\top & \end{bmatrix} {}^{m \setminus n} \begin{bmatrix} \underline{A} \end{bmatrix} {}^{n \setminus 1} \begin{bmatrix} \underline{x} \end{bmatrix} \right) \\
&= {}^{1 \setminus m} \begin{bmatrix} \underline{y}^\top & \end{bmatrix} D \left( {}^{m \setminus n} \begin{bmatrix} \underline{A} \end{bmatrix} {}^{n \setminus 1} \begin{bmatrix} \underline{x} \end{bmatrix} \right) \\
&\quad + \left( {}^{m \setminus n} \begin{bmatrix} \underline{A} \end{bmatrix} {}^{n \setminus 1} \begin{bmatrix} \underline{x} \end{bmatrix} \right)^\top {}^{m \setminus n} \begin{bmatrix} D\underline{y} \end{bmatrix} \\
&= {}^{1 \setminus m} \begin{bmatrix} \underline{y}^\top & \end{bmatrix} {}^{m \setminus n} \begin{bmatrix} \underline{A} \end{bmatrix} + {}^{1 \setminus n} \begin{bmatrix} \underline{x}^\top & \end{bmatrix} {}^{n \setminus m} \begin{bmatrix} \underline{A}^\top \end{bmatrix} {}^{m \setminus n} \begin{bmatrix} \underline{0} \end{bmatrix} \\
&= {}^{1 \setminus m} \begin{bmatrix} \underline{y}^\top & \end{bmatrix} {}^{m \setminus n} \begin{bmatrix} \underline{A} \end{bmatrix} \\
&= \underline{y}^\top \mathbf{A}.
\end{aligned}$$

[Ref]: Edwin K.P. Chong, Stanislaw H. Żak, “PART I MATHEMATICAL REVIEW” in “An introduction to optimization”, 4th Edition, John Wiley and Sons, Inc. 2013.