

$$f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$Df(x) = \left[ \frac{\partial f_i(x)}{\partial x_j} \right] = \begin{pmatrix} \frac{\partial f_1(x)}{\partial x_1} & \frac{\partial f_1(x)}{\partial x_2} & \dots & \frac{\partial f_1(x)}{\partial x_n} \\ \frac{\partial f_2(x)}{\partial x_1} & \frac{\partial f_2(x)}{\partial x_2} & \dots & \frac{\partial f_2(x)}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m(x)}{\partial x_1} & \frac{\partial f_m(x)}{\partial x_2} & \dots & \frac{\partial f_m(x)}{\partial x_n} \end{pmatrix}_{m \times n}$$

Define an operator  $\bar{D} : \mathbb{R} \rightarrow \mathbb{R}^n$  by

$$(Df(x))^T = \begin{pmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{pmatrix} (f_1(x) \ f_2(x) \ \dots \ f_m(x)) \quad \bar{D} = \begin{pmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{pmatrix}$$

$$\Rightarrow \underbrace{(Df(x))^T}_{\text{matrix}} = \boxed{\bar{D} \cdot f(x)^T} \quad f(x)^T = f(x)$$

$$\text{If } f : \mathbb{R}^n \rightarrow \mathbb{R}, \quad \nabla f(x) = Df(x)^T = \bar{D} \cdot f(x)$$

$$\begin{aligned} \nabla^2 f(x) &= \nabla(\nabla f(x)) \\ &= \bar{D} \cdot (\nabla f(x))^T \\ &= \boxed{\bar{D} \cdot Df(x)}. \end{aligned}$$

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

$$Ax = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & \ddots & & \\ & \ddots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n \end{pmatrix}$$

Let  $A = \begin{bmatrix} -a_1 \\ -a_2 \\ \vdots \\ -a_m \end{bmatrix}$ .

$$Df(x) = \begin{pmatrix} \frac{\partial a_1 x}{\partial x_1} & \frac{\partial a_1 x}{\partial x_2} & \cdots & \frac{\partial a_1 x}{\partial x_n} \\ \vdots & & & \\ \frac{\partial a_m x}{\partial x_1} & \cdots & & \frac{\partial a_m x}{\partial x_n} \end{pmatrix} = \begin{pmatrix} a_1 x \\ a_2 x \\ \vdots \\ a_m x \end{pmatrix} = A.$$

②  $f(x, y) = \underline{\underline{y^T A x}}$ ,  $A \in \mathbb{R}^{m \times n}$ ,  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^m$ .

$$D_x f(x, y) = \underline{\underline{y^T A}}$$

$$D_y f(x, y) = ?$$

$$y^T A x \in \mathbb{R} = (y^T A x)^T = \boxed{x^T A^T} y$$

$$D_y f(x, y) = x^T A^T.$$

$$f(x) = x^T P x.$$

$$Df(x) = (x^T P^T) + x^T P = 2x^T P \text{ because } P^T = P.$$

$$\underbrace{\frac{d f(x) q(x)}{dx}}_{=} = f(x) q(x) + f'(x) q(x).$$

## Chain rule

Suppose  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is differentiable at  $x \in \text{dom } f$  and  $g: \mathbb{R}^m \rightarrow \mathbb{R}^p$  is diff. at  $f(x) \in \text{int dom } g$ .

Define the composition  $h: \mathbb{R}^n \rightarrow \mathbb{R}^p$  by  $h(x) = g(f(x))$ .

Then  $h$  is differentiable at  $x$ , with derivative

$$Dh(x) = \begin{matrix} Dg(y) \\ p \times n \end{matrix} \begin{matrix} Df(x) \\ p \times m \end{matrix} \quad \left| \begin{matrix} y = f(x) \\ g = f \circ g \end{matrix} \right.$$

Ex.  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  defined by  $f(x) = Ax + b$ .  $A \in \mathbb{R}^{m \times n}$   
Suppose that  $g: \mathbb{R}^m \rightarrow \mathbb{R}^p$  is differentiable.  $b \in \mathbb{R}^m$

Then  $h \stackrel{\Delta}{=} g(f(x)) = g(Ax + b)$ .

$$Dh(x) = Dg(y) \cdot Df(x) = Dg(y) A \quad \left| \begin{matrix} y = Ax + b \\ y = Ax + b \end{matrix} \right.$$

$$\begin{aligned} \text{If } p=1, \quad \nabla h(x) &= (Dh(x))^T = A^T (Dg(y))^T \\ &= A^T \nabla g(y) \quad \left| \begin{matrix} y = Ax + b \\ y = Ax + b \end{matrix} \right. \end{aligned}$$

$$\begin{aligned} n=m=p=1 \\ h(x) = g(f(x)) \rightarrow h(x) = g(f(x)) \cdot f(x) \end{aligned}$$

Example.  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   $\text{dom } f = \mathbb{R}^n$

$$f(y) = \log \left( \sum_{i=1}^m \exp(a_i^T x + b_i) \right), \text{ where } a_i \in \mathbb{R}^n, b_i \in \mathbb{R} \\ i=1, 2, \dots, m.$$

$$\text{Let } h(x) = Ax + b, \quad A = \begin{pmatrix} -a_1^T - \\ \vdots \\ -a_m^T - \end{pmatrix}, \quad b = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}.$$

Let  $g: \mathbb{R}^m \rightarrow \mathbb{R}$  by

$$g(y) = \log \left( \sum_{i=1}^m \exp(y_i) \right).$$

$$\frac{d \log x}{dx} = \frac{1}{x}$$

$$f(x) = g(h(x)).$$

$$\nabla f(x) = A^T \nabla g(y) \quad | y = Ax + b. \quad \text{by chain rule.}$$

$$\frac{\partial g(y)}{\partial y_j} = \frac{1}{\sum_{i=1}^m \exp(y_i)} \exp(y_j)$$

$$\nabla g(y) = \begin{pmatrix} \exp y_1 \\ \vdots \\ \exp y_m \end{pmatrix}$$

$$\nabla^2 f(x) = A^T \left( \nabla^2 g(Ax + b) \right) A \quad \#$$

## Second derivative

For  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , the second derivative of  $f$  (Hessian matrix) at  $x \in \text{int dom } f$ , denoted by  $\nabla^2 f(x)$ , is given by

$$(\nabla^2 f(x))_{i,j} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}, \quad i=1, \dots, n, \quad j=1, \dots, n,$$

provided that  $f$  is twice differentiable at  $x$ .

The second order approximation of  $f$  at or near  $x$ , is the quadratic function of  $z$

defined by

$$\hat{f}(z) = f(x) + [\nabla f(x)]^T(z-x) + \frac{1}{2} (z-x)^T [\nabla^2 f(x)] (z-x)$$

$f: \mathbb{R} \rightarrow \mathbb{R}$ . Taylor series exists if  $f$  is  $m$  times continuously differentiable

$$f(z) = f(x) + f'(x)(z-x) + \frac{1}{2} f''(x)(z-x)^2 + \frac{1}{3!} f'''(x)(z-x)^3 + \dots + \frac{1}{(m-1)!} f^{(m-1)}(x)(z-x)^{m-1} + R_m$$

If  $f$  is differentiable, the gradient mapping is the function  $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  with domain of with value  $\nabla f(x)$  at  $x$ .  $= \text{domain } f$

The derivative of this function is

$$D\nabla f(x) = \nabla^2 f(x)$$

quadratic function.

Ex.  $f(x) = \frac{1}{2}x^T Px + q^T x + r$ ,  $P \in \mathbb{S}^n$ ,  $q \in \mathbb{R}^n$ ,  $r \in \mathbb{R}$

$$Df(x) = x^T P + q^T. \quad \& \boxed{\nabla f(x) = Px + q.}$$

$$\boxed{\nabla^2 f(x) = P.}$$

$$\nabla^2 f(x) = D \cdot Df(x) = \begin{pmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{pmatrix} (x^T P + q^T) \\ = P.$$

The second-order approximation of a quadratic function is itself. (Exercise).

## Chain Rule.

1. Composition with scalar function

Suppose  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $g: \mathbb{R} \rightarrow \mathbb{R}$

and  $h(x) = g(f(x))$ .

Then  $\nabla^2 h(x) = g'(f(x)) \nabla^2 f(x) + g''(f(x)) \nabla f(x) \nabla f(x)^T$ .

If  $n=1$ .  $h'(x) = g'(f(x)) \cdot f'(x)$ .

$h''(x) = g''(f(x)) \cdot f'(x) \cdot f'(x) + g'(f(x)) \cdot f''(x)$ .

Proof.  $\nabla h(x) = g'(f(x)) \cdot \nabla f(x)$ .

$$\nabla^2 h(x) = \bar{D} \cdot (\nabla h(x))^T$$

$$= \bar{D} \cdot (g'(f(x)) \cdot \nabla f(x))$$

$$= (\bar{D} g'(f(x))) \cdot \nabla f(x)^T + g'(f(x)) \bar{D} \cdot (\nabla f(x))^T$$

$$= g''(f(x)) \nabla f(x) \cdot (\nabla f(x))^T + g'(f(x)) \nabla^2 f(x)$$

$$\nabla^2 h(x) = \begin{pmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{pmatrix} (g'(f(x))) \left( \frac{\partial f(x)}{\partial x_1} \quad \frac{\partial f(x)}{\partial x_2} \cdots \quad \frac{\partial f(x)}{\partial x_n} \right)^T$$

2. Composition with an affine function  
 Suppose  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ .  $f(x) = Ax + b$ .  $A \in \mathbb{R}^{m \times n}$   
 $b \in \mathbb{R}^m$ .

$g: \mathbb{R}^m \rightarrow \mathbb{R}^p$  is differentiable.

Define  $h: \mathbb{R}^n \rightarrow \mathbb{R}^p$  by  $h(x) = g(f(x)) = g(Ax + b)$

Then  $D h(x) = Dg(y) \cdot A \quad |_{y=Ax+b}$

Thus  $\nabla h(x) = [A^T \nabla g(y)] \quad |_{y=Ax+b}$

$$\nabla^2 h(x) = \bar{D} \cdot (\nabla h(x))^T$$

$$= \bar{D} \cdot (\nabla g(y))^T \cdot A$$

$$= \underbrace{A^T \cdot \nabla^2 g(y) \cdot A}_{|_{y=Ax+b}}$$

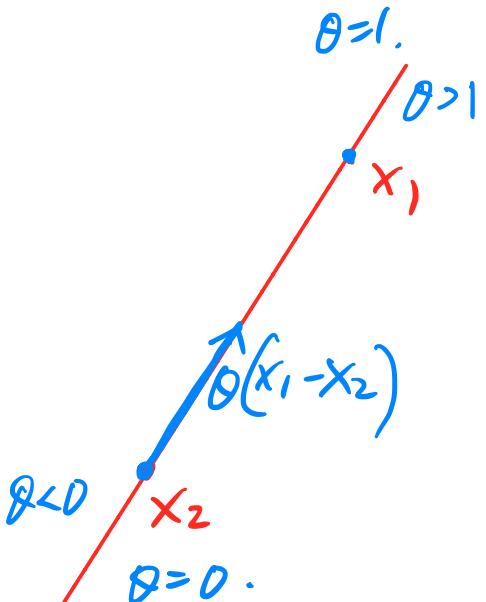
by the chain rule.

## Chapter 2. Convex sets

Line:  $x_1, x_2 \in \mathbb{R}^n, x_1 \neq x_2$

$$y = \theta x_1 + (1-\theta)x_2, \theta \in \mathbb{R}$$

$$= \theta(x_1 - x_2) + x_2$$



line segment

$$y = \theta x_1 + (1-\theta)x_2, \theta \in [0, 1]$$

closed interval  
 $0 \leq \theta \leq 1$ .

## Affine sets

A set  $C \subseteq \mathbb{R}^n$  is affine if for any  $x_1 \neq x_2 \in C$  we have  $\theta x_1 + (1-\theta)x_2 \in C$  for  $\theta \in \mathbb{R}$ .

$\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_k x_k$ , where  $\theta_1 + \theta_2 + \dots + \theta_k = 1$ , is called the affine combination of  $x_1, \dots, x_k$ .

If  $C$  is an affine set,  $x_1, x_2, \dots, x_k \in C$ , then the affine combination of  $x_1, \dots, x_n$  is in  $C$ .

If  $C$  is affine, and  $x_0 \in C$ , then

$V \triangleq C - x_0 = \{x - x_0 : x \in C\}$  is a subspace.

It means for  $v_1, v_2 \in V$ ,  $\alpha, \beta \in \mathbb{R}$ .

$$\underline{\alpha v_1 + \beta v_2 \in V}.$$

proof.  $v_1 = x_1 - x_0 \in V$ .  $v_2 = x_2 - x_0 \in V$ .  
 $x_1, x_2 \in C$ .

$$\begin{aligned}\alpha v_1 + \beta v_2 &= \alpha(x_1 - x_0) + \beta(x_2 - x_0) \\ &= \alpha x_1 + \beta x_2 - (\alpha + \beta)x_0 \\ &= \underbrace{\alpha x_1 + \beta x_2 - (\alpha + \beta - 1)x_0}_{\text{affine combination of } x_1, x_2, x_0 \in C} - x_0 \\ &\in V.\end{aligned}$$

Thus an affine set  $C = V + x_0$  for  
a subspace  $V$  and  $x_0 \in C$ .

$$\dim C \leq \dim(V).$$

We say  $k+1$  points  $v_0, v_1, \dots, v_k$  are  
~~affinely independent~~ if  $v_1 - v_0, v_2 - v_0, \dots, v_k - v_0$  are  
linearly independent.

Affine hull: for some set  $C \subseteq \mathbb{R}^n$

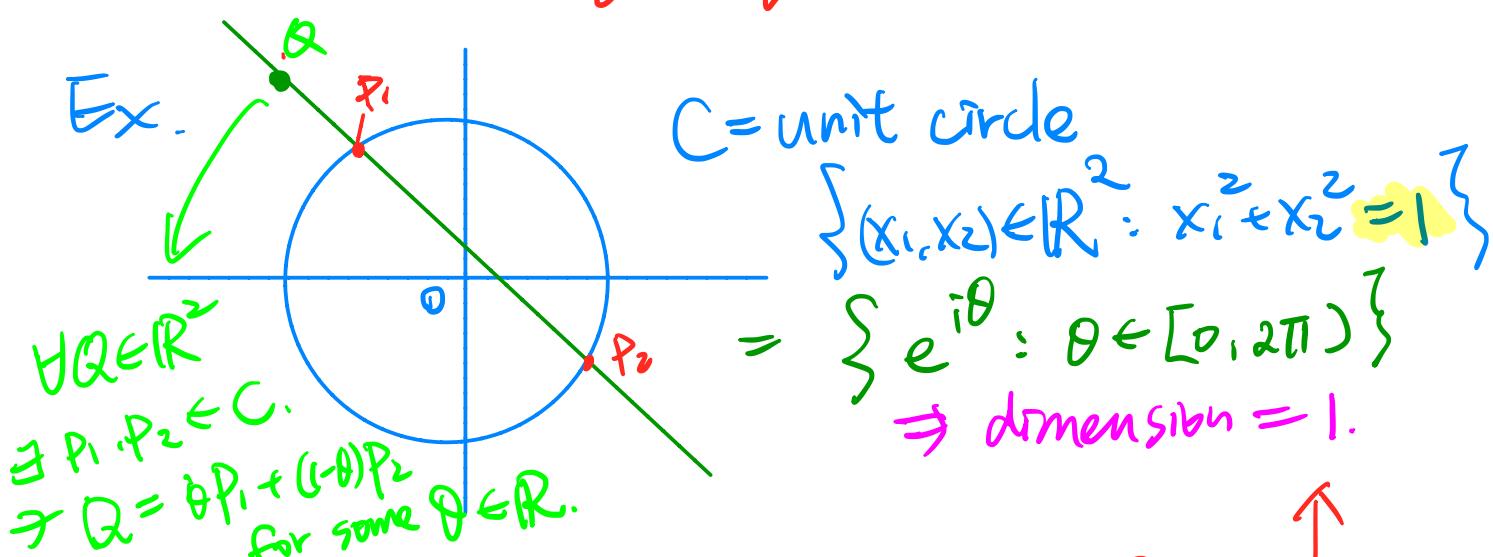
$$\text{aff}(C) \triangleq \left\{ \sum_{i=1}^k \theta_i x_i : x_1, \dots, x_k \in C \right. \\ \left. \quad \sum_i \theta_i = 1, \quad \theta_i \in \mathbb{R} \right\}$$

which is the smallest affine set containing  $C$ .

Affine dimension of  $C$  is the dimension of its affine hull.  $\text{aff}(C)$ .

$\text{aff}(C)$  is affine  $\Rightarrow \exists V \ni \text{aff}(C) = V + x_0$   
for  $x_0 \in C$ .

$$\text{aff dim}(C) \triangleq \dim(\text{aff}(C)) \\ = \dim V$$



$$\text{aff}(C) = \mathbb{R}^2 \text{ has dimension 2.}$$

$$\Rightarrow \text{aff dim } C = 2.$$

thus

$$\mathbb{R}^2 \subseteq \text{aff}(C).$$

Also,  $\text{aff}(C) \subseteq \mathbb{R}^2$  trivially.

## Relative interior -

If  $\text{aff dim}(C) < n$ . for  $C \subset \mathbb{R}^n$ ,

then  $C \subseteq \boxed{\text{aff}(C) \neq \mathbb{R}^n}$

We define the relative interior of  $C$ ,  
denoted by  $\text{relint } C$ , as its interior  
relative to  $\text{aff}(C)$ :

$$\text{relint } C \equiv \left\{ x \in C : \exists B(x, r), r > 0, \text{ such that } B(x, r) \cap \underline{\text{aff}(C)} \subseteq C \right\}$$