

# Introduction

Oscar Dalmau  
dalmau@cimat.mx

Centro de Investigación en Matemáticas  
CIMAT A.C. Mexico

January 2018

# Outline

- ① Introduction
- ② Level sets and gradient
- ③ Taylor Series

# Optimization Problem

- **Optimization Problem:**

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

where  $f(\cdot)$  is class  $\mathcal{C}^1$  or continuously differentiable that consists of all differentiable functions whose derivative is continuous or  $f(\cdot)$  could be class  $\mathcal{C}^2$

# Optimization Problem

- Optimization Problem:**

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

- Examples:**  $\mathbf{x} = [x_1, \dots, x_n]^T \in \mathbb{R}^n$ ,

$$\min_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}, \mathbf{A} \text{ is symmetric}$$

$$\min_{\mathbf{x}} \sum_{i=1}^n (x_i - y_i)^2 + \lambda \sum_{i=1}^{n-1} (x_{i+1} - x_i)^2$$

$$\min_{\mathbf{x}} \sum_{i=1}^{N-1} 100(x_{i+1} - x_i^2)^2 + (1 - x_i)^2$$

# Notation

$$\mathbf{1} = [1, 1, \dots, 1]^T$$

$$\mathbf{0} = [0, 0, \dots, 0]^T$$

$$\mathbf{I} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \ddots & \dots \\ 0 & 0 & \dots & 1 \end{bmatrix}$$

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

$$f(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x})]^T$$

# Definitions

- A real  $n \times n$  matrix  $\mathbf{A}$  is said to be positive definite if the scalar  $\mathbf{x}^\top \mathbf{A} \mathbf{x}$  is positive for every non-zero column vector  $\mathbf{x}$
- A real  $n \times n$  matrix  $\mathbf{A}$  is said to be the **negative definite**, **positive semi-definite**, and **negative semi-definite** matrices are defined in the same way, except that the expression  $\mathbf{x}^\top \mathbf{A} \mathbf{x}$  is required to be always negative, non-negative, and non-positive, respectively.
- We can determine that a **symmetric matrix** is **positive definite** by computing its eigenvalues and verifying that they are all positive, or by performing a Cholesky factorization (Cholesky factorization gives error for non-positive-definite matrices)

## Some theorems

**Lema:** (Sign Preserving Property) Let  $f$  be continuous at  $a$  and  $f(a) \neq 0$ . Then there is an interval  $(a - \delta, a + \delta)$  about  $a$  in which  $f$  has the same sign as  $f(a)$ .

**Proof:** Without loss of generality, assume that  $f(a) > 0$ . Using the continuity of  $f$ . For all  $\epsilon > 0$  there exist  $\delta > 0$  such that

$$|x - a| < \delta \Rightarrow |f(x) - f(a)| < \epsilon$$

Then  $f(x) > f(a) - \epsilon$  and taking  $0 < \epsilon < f(a)$  (for example,  $\epsilon = \frac{f(a)}{2}$ ) one obtains

$$f(x) > f(a) - \epsilon > 0$$

for  $|x - a| < \delta$ , ie,  $x \in (a - \delta, a + \delta)$

## Some theorems

**Theorem:** (Mean value theorem) Let  $f : [a, b] \rightarrow \mathbb{R}$  be a continuous function on the closed interval  $[a, b]$ , and differentiable on the open interval  $(a, b)$ , where  $a < b$ . Then there exists some  $c \in (a, b)$  such that

$$f'(c) = \frac{f(b) - f(a)}{b - a}.$$

The *mean value theorem* is a generalization of *Rolle's theorem*, which assumes  $f(a) = f(b)$ , so that the right-hand side above is zero.



# Norms

For a vector  $\mathbf{x} \in \mathbb{R}^n$ , we have the following norms:

$$\|\mathbf{x}\|_1 \stackrel{def}{=} \sum_{i=1}^n |x_i| \quad (1)$$

$$\|\mathbf{x}\|_2 \stackrel{def}{=} \sqrt{\sum_{i=1}^n x_i^2} \quad (2)$$

$$\|\mathbf{x}\|_\infty \stackrel{def}{=} \max_{i=1,2,\dots,n} |x_i| \quad (3)$$

The norm  $\|\cdot\|_2$  is often called the *Euclidean norm*,  $\|\cdot\|_1$  is called the  $\ell_1$ -norm and  $\|\cdot\|_\infty$  is called the  $\ell_\infty$ -norm.

# Norms

In general, a norm is any mapping  $\|\cdot\|$  from  $\mathbb{R}^n$  to the nonnegative real numbers that satisfies the following properties:

$$\begin{aligned}\|\mathbf{x}\| &= 0 \Leftrightarrow \mathbf{x} = \mathbf{0}; \mathbf{x} \in \mathbb{R}^n \\ \|\alpha\mathbf{x}\| &= |\alpha|\|\mathbf{x}\|; \mathbf{x} \in \mathbb{R}^n; \alpha \in \mathbb{R} \\ \|\mathbf{x} + \mathbf{y}\| &\leq \|\mathbf{x}\| + \|\mathbf{y}\|; \mathbf{x}, \mathbf{y} \in \mathbb{R}^n\end{aligned}$$

# Cauchy-Schwarz inequality

A property that holds for the Euclidean norm is the Cauchy-Schwarz inequality, which states that:

$$|\mathbf{x}^T \mathbf{y}| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

and the equality holds if and only if one of these vectors is a multiple of the other, ie, the vectors are parallels.

# Matrix norms

Let  $\|\cdot\|$  be generic notation for the three norms listed above, we define the corresponding matrix norm as:

$$\|\mathbf{A}\| \stackrel{def}{=} \sup_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}$$

# Matrix norms

## Explicit formulas

$$\|\mathbf{A}\|_1 = \max_{j=1,2,\dots,m} \sum_{i=1}^n |A_{ij}|$$

$$\|\mathbf{A}\|_\infty = \max_{i=1,2,\dots,n} \sum_{j=1}^m |A_{ij}|$$

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(A^T A)} = \sigma_{\max}(A)$$

# Matrix norms

- The Frobenius norm is defined as follows

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j}^{n,m} A_{ij}^2}$$

- Inner product  $\langle A, B \rangle = \text{tr}(A^T B)$

$$\begin{aligned}\|\mathbf{A}\|_F^2 &= \langle A, A \rangle = \text{tr}(A^T A) = \sum_j (A^T A)_{jj} \\ &= \sum_j \sum_i A_{ji}^T A_{ij} = \sum_i \sum_j A_{ij} A_{ij} = \sum_i \sum_j A_{ij}^2\end{aligned}$$

# Condition number of a Matrix

The condition number of a nonsingular matrix is defined as

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|,$$

where any matrix norm can be used in the definition. We can use a subscript for the different norms, ie,  $\kappa_1(\cdot)$ ,  $\kappa_2(\cdot)$ , and  $\kappa_\infty(\cdot)$  respectively. We use  $\kappa$  to denote  $\kappa_2(\cdot)$ .

# Definitions

In many optimization methods the information of the first and/or second derivative of  $f(\cdot)$  is required

- **Gradient:** If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  such that  $f \in \mathcal{C}^1$ , i.e.,  $f$  has continuous partial derivatives of first order, then the *gradient* of  $f(\cdot)$  is defined as

$$\begin{aligned}\nabla f(\mathbf{x}) &= \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]^T \stackrel{\text{def}}{=} \mathbf{g}(\mathbf{x}) \\ Df(\mathbf{x}) &= \nabla f(\mathbf{x})^T\end{aligned}$$



# Definitions

In many optimization methods the information of the first and/or second derivative of  $f(\cdot)$  is required

- **Hessian:** If  $f \in \mathcal{C}^2$ , i.e.,  $f$  has continuous partial derivatives of second order, then the *Hessian* of  $f(\cdot)$  is defined as

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \stackrel{def}{=} \mathbf{H}(\mathbf{x})$$

$\mathbf{H}(\mathbf{x})$  is a symmetric square matrix, ie,  $\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$ . For the *equality of mixed partial derivatives Theorem*, due to the mixed partial derivatives exist and are continuous.

# Definitions

In many optimization methods the information of the first and/or second derivative of  $f(\cdot)$  is required

- The function  $f(\cdot)$ , the gradient  $\mathbf{g}(\cdot)$  and the Hessian  $\mathbf{H}(\cdot)$  at  $\mathbf{x} = \mathbf{x}_k$  are denoted as  $f_k$ ,  $\mathbf{g}_k$  and  $\mathbf{H}_k$  respectively in order to simplify notation, i.e.,

$$\begin{aligned} f_k &\stackrel{\text{def}}{=} f(\mathbf{x}_k) \\ \mathbf{g}_k &\stackrel{\text{def}}{=} \mathbf{g}(\mathbf{x}_k) \\ \mathbf{H}_k &\stackrel{\text{def}}{=} \mathbf{H}(\mathbf{x}_k) \end{aligned}$$

# Definitions

**Directional derivative:** Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is defined on an open ball about a point  $\mathbf{x}_0$ . Given a unit vector  $\mathbf{v}$ , we call

$$D_{\mathbf{v}}f(\mathbf{x}_0) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x}_0 + h\mathbf{v}) - f(\mathbf{x}_0)}{h},$$

provided the limit exists, the **directional derivative** of  $f$  in the direction of  $\mathbf{v}$  at  $\mathbf{x}_0$ . For example

$$\frac{\partial f}{\partial x_i}(\mathbf{x}_0) = D_{\mathbf{e}_i}f(\mathbf{x}_0)$$

where  $\mathbf{e}_i = [0, 0, \dots, 1, \dots, 0]^T$ .

# Differentiation Rules

- **Chain rule:** Let  $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^k$ ,  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be differentiable functions at  $\mathbf{x} \in \mathbb{R}^n$  then

$$D_{\mathbf{x}}(\mathbf{f} \circ \mathbf{g}) = D\mathbf{f}(\mathbf{g}(\mathbf{x}))D\mathbf{g}(\mathbf{x})$$

where  $D\mathbf{f}(\mathbf{g}(\mathbf{x})) \in \mathbb{R}^{k \times m}$ ,  $D\mathbf{g}(\mathbf{x}) \in \mathbb{R}^{m \times n}$  and  $D_{\mathbf{x}}(\mathbf{f} \circ \mathbf{g}) \in \mathbb{R}^{k \times n}$  are matrices and  $(\mathbf{f} \circ \mathbf{g}) : \mathbb{R}^n \rightarrow \mathbb{R}^k$ .

# Differentiation Rules

- Example:**  $\mathbf{f}(x, y) = [x^2 + y, x + y^2]^T$ ,  $\mathbf{g}(t) = [t^2, t]^T$ .  
Compute  $D_t(\mathbf{f} \circ \mathbf{g})$ .

**Solution 1:** Using the chain rule

$$D\mathbf{f}(x, y) = \begin{bmatrix} 2x & 1 \\ 1 & 2y \end{bmatrix}$$

$$D\mathbf{g}(t) = [2t, 1]^T$$

$$D\mathbf{f}(\mathbf{g}(t)) = \begin{bmatrix} 2t^2 & 1 \\ 1 & 2t \end{bmatrix}$$

$$D\mathbf{f}(\mathbf{g}(t))D\mathbf{g}(t) = \begin{bmatrix} 4t^3 + 1 \\ 4t \end{bmatrix}$$

# Differentiation Rules

- **Example:**  $\mathbf{f}(x, y) = [x^2 + y, x + y^2]^T$ ,  $\mathbf{g}(t) = [t^2, t]^T$ .  
Compute  $D_t(\mathbf{f} \circ \mathbf{g})$ .

**Solution 2:** Computing  $(\mathbf{f} \circ \mathbf{g})(t)$

$$\begin{aligned} h(t) = (\mathbf{f} \circ \mathbf{g})(t) &= \begin{bmatrix} t^4 + t \\ 2t^2 \end{bmatrix} \\ h'(t) &= \begin{bmatrix} 4t^3 + 1 \\ 4t \end{bmatrix} \end{aligned}$$

- **Chain rule** (case 1): Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{z} : \mathbb{R} \rightarrow \mathbb{R}^n$  be differentiable functions at  $t \in \mathbb{R}$  then

$$D_t(f \circ \mathbf{z}) = Df(\mathbf{z}(t))D\mathbf{z}(t)$$

Note that:

$$\begin{aligned} f \circ \mathbf{z} : \mathbb{R} &\rightarrow \mathbb{R} \\ Df(\mathbf{z}(t)) &= \nabla f(\mathbf{z}(t))^T \\ D\mathbf{z}(t) &= [z'_1(t), z'_2(t), \dots, z'_n(t)]^T \\ D_t(f \circ \mathbf{z}) &= \nabla f(\mathbf{z}(t))^T D\mathbf{z}(t) = \langle \nabla f(\mathbf{z}(t)), D\mathbf{z}(t) \rangle \in \mathbb{R} \end{aligned}$$

- **Chain rule** (case 1 Example):  $f(\mathbf{x}) = \|\mathbf{x}\|^2$ ,  $\mathbf{z}(t) = \mathbf{x}_0 + t\mathbf{v}$ .

Compute  $D_t(f \circ \mathbf{z})$ .

Solution:

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

$$D\mathbf{z}(t) = [z'_1(t), z'_2(t), \dots, z'_n(t)]^T = \mathbf{v}$$

Then  $D_t(f \circ \mathbf{z}) = 2\mathbf{z}(t)^T \mathbf{v} = 2(\mathbf{x}_0 + t\mathbf{v})^T \mathbf{v}$



- **Product rule:** Let  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be two differentiable functions. Define the function  $h : \mathbb{R}^n \rightarrow \mathbb{R}$  by  $h(\mathbf{x}) = \mathbf{f}(\mathbf{x})^T \mathbf{g}(\mathbf{x})$  then

$$D_{\mathbf{x}}h(\mathbf{x}) = \mathbf{f}(\mathbf{x})^T D\mathbf{g}(\mathbf{x}) + \mathbf{g}(\mathbf{x})^T D\mathbf{f}(\mathbf{x})$$

where  $D\mathbf{g}(\mathbf{x})$  and  $D\mathbf{f}(\mathbf{x})$  are matrices (the Jacobian matrix).

- **Product rule** (Example):  $f(x) = x$ ,  $g(x) = Ax$ . Compute  $D_x h(x)$  with  $h(x) = f(x)^T g(x) = x^T Ax$ .  
Solution:

$$Df(x) = I$$

$$Dg(x) = A$$

Then

$$\begin{aligned} D_x h(x) &= f(x)^T Dg(x) + g(x)^T Df(x) \\ &= x^T A + x^T A^T I = x^T (A + A^T) \end{aligned}$$

and  $\nabla h(x) = (A + A^T)x$ .

**Note:** If  $A^T = A$  then  $\nabla h(x) = 2Ax$  and  $\nabla^2 h(x) = 2A$ .

# Level sets

The level set of a function  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$  at level  $l$  is the set of points

$$S = \{\mathbf{x} : f(\mathbf{x}) = l\}$$

- If  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  then  $S$  is a curve
- If  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  then  $S$  is a surface

# Level sets

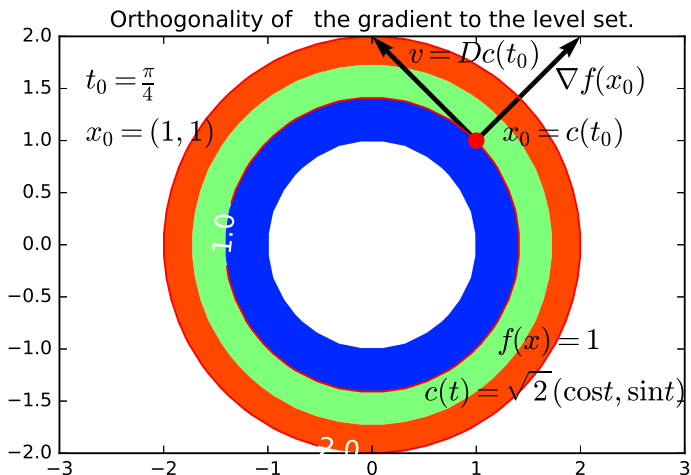
- Let  $\mathbf{x}_0$  a point in the level set  $S$ , i.e.,  $f(\mathbf{x}_0) = l$
- Let  $\mathbf{c}(t)$  be a parametrization of a curve  $\gamma$  that lies on  $S$ , such that

$$\begin{aligned}\mathbf{c}(t) &: \mathbb{R} \rightarrow \mathbb{R}^n \\ \mathbf{c}(t_0) &= \mathbf{x}_0 \\ D\mathbf{c}(t_0) &= \mathbf{v} \neq \mathbf{0}\end{aligned}$$

where  $\mathbf{v}$  is the tangent vector to  $\gamma$  at  $\mathbf{x}_0$

- Example:  $f(\mathbf{x}) = \frac{1}{2}(x_1^2 + x_2^2)$  and the level set  $f(\mathbf{x}) = 1$ , ie,  $x_1^2 + x_2^2 = 2$ . A parametrization could be  $\mathbf{c}(t) = [x_1(t), x_2(t)]^T = \sqrt{2}[\cos t, \sin t]^T$ , see next slide.

# Level sets and gradient



**Theorem:** The vector  $\nabla f(\mathbf{x}_0)$  is orthogonal to the tangent vector to an arbitrary smooth curve passing through  $\mathbf{x}_0$  on the level set determined by  $f(\mathbf{x}) = f(\mathbf{x}_0)$ .

**Proof:** (It is straightforward)

Applying the chain rule to the function  $h(t) = f(\mathbf{c}(t))$

$$h'(t_0) = Df(\mathbf{c}(t_0))D\mathbf{c}(t_0) = Df(\mathbf{x}_0)\mathbf{v} = \nabla f(\mathbf{x}_0)^T \mathbf{v}$$

$$h(t) = f(\mathbf{c}(t)) = \text{constant}$$

$$h'(t) = 0$$

then  $\nabla f(\mathbf{x}_0)^T \mathbf{v} = 0$ , ie,  $\nabla f(\mathbf{x}_0) \perp \mathbf{v}$

# Comments

- It is said that  $\nabla f(\mathbf{x}_0)$  is orthogonal or normal to the level set  $S$  corresponding to  $\mathbf{x}_0$
- The tangent plane to  $S$  at  $\mathbf{x}_0$  the set of all points  $\mathbf{x}$  satisfying

$$\nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) = 0$$

if  $\nabla f(\mathbf{x}_0) \neq \mathbf{0}$ .

# Comments

**Theorem:** Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\mathcal{C}^1$  on an open ball containing the point  $\mathbf{x}_0$ . Then for any unit vector  $\mathbf{v}$ ,  $D_{\mathbf{v}}f(\mathbf{x}_0)$  exists and  $D_{\mathbf{v}}f(\mathbf{x}_0) = Df(\mathbf{x}_0)\mathbf{v} = \nabla f(\mathbf{x}_0)^T \mathbf{v}$ .

**Proof:** Note that

$$D_{\mathbf{v}}f(\mathbf{x}_0) = \left[ \frac{d}{d\alpha} f(z(\alpha)) \right]_{\alpha=0} = \left[ \frac{d}{d\alpha} f(\mathbf{x}_0 + \alpha \mathbf{v}) \right]_{\alpha=0}$$

with  $z(\alpha) = \mathbf{x}_0 + \alpha \mathbf{v}$ . Using the Chain rule

$$D_{\mathbf{v}}f(\mathbf{x}_0) = \left[ Df(z(\alpha))z'(\alpha) \right]_{\alpha=0} = Df(z(0))\mathbf{v} = Df(\mathbf{x}_0)\mathbf{v}$$



# Comments

**Theorem:** Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\mathcal{C}^1$  on an open ball containing the point  $\mathbf{x}_0$ . Then for any unit vector  $\mathbf{v}$ ,  $D_{\mathbf{v}}f(\mathbf{x}_0)$  exists and  $D_{\mathbf{v}}f(\mathbf{x}_0) = Df(\mathbf{x}_0)\mathbf{v} = \nabla f(\mathbf{x}_0)^T \mathbf{v}$ .

**Note:** Using the Cauchy-Schwarz inequality

$$\begin{aligned} |D_{\mathbf{v}}f(\mathbf{x}_0)| &= |\nabla f(\mathbf{x}_0)^T \mathbf{v}| \leq \|\nabla f(\mathbf{x}_0)\| \|\mathbf{v}\| \\ |D_{\mathbf{v}}f(\mathbf{x}_0)| &\leq \|\nabla f(\mathbf{x}_0)\| \end{aligned}$$

if  $\nabla f(\mathbf{x}_0) \neq \mathbf{0}$  and taking  $\mathbf{v} = \frac{\nabla f(\mathbf{x}_0)}{\|\nabla f(\mathbf{x}_0)\|}$

$$\begin{aligned} D_{\mathbf{v}}f(\mathbf{x}_0) &= \|\nabla f(\mathbf{x}_0)\| \\ D_{-\mathbf{v}}f(\mathbf{x}_0) &= -\|\nabla f(\mathbf{x}_0)\| \end{aligned}$$

## Comments

**Proposition:** Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\mathcal{C}^1$  on an open ball containing the point  $\mathbf{x}_0$ . Then  $D_{\mathbf{v}}f(\mathbf{x}_0)$  has a maximum value of  $\|\nabla f(\mathbf{x}_0)\|$  when  $\mathbf{v}$  is the direction of  $\nabla f(\mathbf{x}_0)$  and a minimum value of  $-\|\nabla f(\mathbf{x}_0)\|$  when  $\mathbf{v}$  is the direction of  $-\nabla f(\mathbf{x}_0)$ .

**Proof:**

$$\begin{aligned} D_{\mathbf{v}}f(\mathbf{x}_0) &= \nabla f(\mathbf{x}_0)^T \mathbf{v} \\ &= \|\nabla f(\mathbf{x}_0)\| \|\mathbf{v}\| \cos \angle(\nabla f(\mathbf{x}_0), \mathbf{v}) \\ &= \|\nabla f(\mathbf{x}_0)\| \cos \angle(\nabla f(\mathbf{x}_0), \mathbf{v}) \end{aligned}$$

The maximum and minimum of  $D_{\mathbf{v}}f(\mathbf{x}_0)$  is obtained when  $\cos \angle(\nabla f(\mathbf{x}_0), \mathbf{v}) = 1$  and  $\cos \angle(\nabla f(\mathbf{x}_0), \mathbf{v}) = -1$  respectively, i.e., when  $\angle(\nabla f(\mathbf{x}_0), \mathbf{v}) = 0$  and  $\angle(\nabla f(\mathbf{x}_0), \mathbf{v}) = \pi$

## Comments

- 1 The gradient ,  $\nabla f(\cdot)$ , of a real-valued differentiable function at a point is orthogonal to the level set of the function at that point.
- 2 The gradient vector points in the direction of the maximum rate of increase of the function and the negative of the gradient vector points in the direction of the maximum rate of decrease of the function.
- 3 The length of the gradient vector tells us the rate of increase in the direction of maximum increase and its negative tells us the rate of decrease in the direction of maximum decrease.

**Theorem:** Taylor's Theorem. Assume that a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is  $m + 1$  times continuously differentiable (i.e.,  $f \in \mathcal{C}^{m+1}$ ) at  $a \in \mathbb{R}$ . Denote  $h = x - a$ . Then,

$$f(x) = f(a) + \frac{h}{1!} f^{(1)}(a) + \frac{h^2}{2!} f^{(2)}(a) + \cdots + \frac{h^m}{m!} f^{(m)}(a) + R_{m+1},$$

(called Taylor's formula) where  $f^{(i)}$  is the  $i$ -th derivative of  $f$ , and

$$R_{m+1} = \frac{h^{m+1}}{(m+1)!} f^{(m+1)}(a + \theta h), \quad \theta \in (0, 1)$$

**Theorem:** Taylor's Theorem. Assume that a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is  $m + 1$  times continuously differentiable (i.e.,  $f \in \mathcal{C}^{m+1}$ ) at  $a \in \mathbb{R}$ . Denote  $h = x - a$ . Then,

$$f(x) = f(a) + \frac{h}{1!} f^{(1)}(a) + \frac{h^2}{2!} f^{(2)}(a) + \cdots + \frac{h^m}{m!} f^{(m)}(a) + R_{m+1},$$

(called Taylor's formula) where  $f^{(i)}$  is the  $i$ -th derivative of  $f$ , and

$$R_{m+1} = \frac{h^{m+1}}{(m+1)!} f^{(m+1)}(a + \theta h), \quad \theta \in (0, 1)$$

The remainder  $R_{m+1}$  using **Big O** and **little o** notation?

# Big O and Little o Notation

- **Definition.** We say  $f(x) = O(g(x))$  as  $x \rightarrow a$  if there exists a constant  $C$  such that

$$|f(x)| \leq C|g(x)|$$

in some neighborhood of  $a$ , that is, for  $x \in (a - \delta, a + \delta) \setminus \{a\}$  for some  $\delta > 0$ .

- We say  $f(x) = O(g(x))$  as  $x \rightarrow \infty$  if there exist positive constants  $x_0$  and  $C$  such that  $|f(x)| \leq C|g(x)|$  for all  $x > x_0$ .

# Big O and Little o Notation

**Example:**  $x^3 - 2x + 1$  is  $O(x^3)$  as  $x \rightarrow \infty$

$$|x^3 - 2x + 1| \leq |x^3| + |2x| + |1|$$

for  $x \geq 1$  we have

$$|x^3| \leq 1|x^3|$$

$$|2x| \leq 2|x^3|$$

$$|1| \leq |x^3|$$

$$|x^3 - 2x + 1| \leq 4|x^3|$$

then  $f(x) = O(g(x))$  as  $x \rightarrow \infty$ .

# Big O and Little o Notation

- **Theorem 1:** If  $\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = L$  then  $f(x) = O(g(x))$  as  $x \rightarrow a$
- **Example:**  $x^3 - 2x + 1$  is  $O(x^3)$  as  $x \rightarrow \infty$
- **By the theorem:**

$$\lim_{x \rightarrow \infty} \frac{x^3 - 2x + 1}{x^3} = 1$$

then  $f(x) = O(g(x))$  as  $x \rightarrow \infty$



# Big O and Little o Notation

- **Definition.** We say  $f(x) = o(g(x))$  as  $x \rightarrow a$  if

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = 0$$

- Generally, some results assume that  $a = 0$  since changing  $x - a$  to  $x$  is just a change of coordinates.

# Big O and Little o Notation

- The previous notation is standard, but it is not good. The statement  $f(x) = O(g(x))$  has a meaning, but  $O(g(x)) = f(x)$  is meaningless, i.e., the equal sign ( $=$ ) is not symmetric, then this is an abuse of notation.
- A better option would be  $f(x) \in O(g(x))$  or  $f(x)$  is  $O(g(x))$
- The limit  $a$  is important, for example,  $\frac{1}{x} = o(1)$  as  $x \rightarrow \infty$ , but  $\frac{1}{x} \neq o(1)$  as  $x \rightarrow 0$

# Big O and Little o Notation. Examples

- Big O examples

$$x = O(x), \text{ as } x \rightarrow \infty$$

$$x = O(x^2), \text{ as } x \rightarrow \infty$$

$$ax^n = O(x^m), m \geq n, \text{ as } x \rightarrow \infty$$

$$ax^n \neq O(x^m), m < n, \text{ as } x \rightarrow \infty$$

- Little o examples

$$x^2 = o(x), \text{ as } x \rightarrow 0$$

$$x \neq o(x^2), \text{ as } x \rightarrow 0$$

$$x - \sin x = o(x), \text{ as } x \rightarrow 0$$

$$x - \sin x = o(x^2), \text{ as } x \rightarrow 0$$

# Big O and Little o Notation. Properties

- ❶  $f(x) = O(f(x))$
- ❷ If  $f(x) = O(g(x))$ , then  $cf(x) = O(g(x))$  for any constant  $c$ .
- ❸ If  $f_1(x)$  and  $f_2(x)$  are both  $O(g(x))$ , then so is  $f_1(x) + f_2(x)$ .
- ❹ If  $f(x) = o(g(x))$ , then  $f(x) = O(g(x))$ .
- ❺ If  $f(x) = O(g(x))$ , then  $O(f(x)) + O(g(x)) = O(g(x))$ .
- ❻ If  $f(x) = O(g(x))$ , then  $o(f(x)) + o(g(x)) = o(g(x))$ .
- ❼ If  $f_1(x) = O(g(x))$  but  $f_2(x) = o(g(x))$ , then  
 $f_1(x) + f_2(x) = O(g(x))$
- ❽ If  $f(x) = O(g(x))$ , and  $g(x) = o(h(x))$ , then  $f(x) = o(h(x))$ .
- ❾ Let  $c \neq 0$ , then  $cO(g(x)) = O(g(x))$  and  $co(g(x)) = o(g(x))$ .

# Big O and Little o Notation. Properties

- ①  $O(f(x))O(g(x)) = O(f(x)g(x))$
- ②  $o(f(x))O(g(x)) = o(f(x)g(x))$
- ③ if  $\lim_{x \rightarrow 0} \frac{h(x)}{g(x)} = L$  then  $h(x) = O(g(x))$ . (We can use the definition of limit). then for any  $\epsilon > 0$  there exists  $\delta > 0$ , such that for  $0 < |x| < \delta$  it holds  $|\frac{f(x)}{g(x)} - L| < \epsilon$ . Using the inequality  $|a| \leq |a - b| + |b|$

$$\begin{aligned} \left| \frac{f(x)}{g(x)} \right| &\leq \left| \frac{f(x)}{g(x)} - L \right| + |L| < \epsilon + |L| \stackrel{\text{def}}{=} M \\ |f(x)| &< M|g(x)|, \text{ for } 0 < |x| < \delta \end{aligned}$$

therefore  $f(x) \in O(g(x))$  as  $x \rightarrow 0$ !

# Taylor's formula with Big o and little o notation

$$R_{m+1} = \frac{h^{m+1}}{(m+1)!} f^{(m+1)}(a + \theta h), \quad \theta \in (0, 1)$$

Therefore,

$$R_{m+1} = o(h^m), \text{ as } h \rightarrow 0$$

$$R_{m+1} = O(h^{m+1}), \text{ as } h \rightarrow 0$$

Then, if  $f(x) \in \mathcal{C}^{m+1}$ , we may write Taylor's formula as

$$f(x) = f(a) + \frac{h}{1!} f^{(1)}(a) + \frac{h^2}{2!} f^{(2)}(a) + \cdots + \frac{h^m}{m!} f^{(m)}(a) + o(h^m),$$

$$f(x) = f(a) + \frac{h}{1!} f^{(1)}(a) + \frac{h^2}{2!} f^{(2)}(a) + \cdots + \frac{h^m}{m!} f^{(m)}(a) + O(h^{m+1})$$

If we assume that  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R} \in \mathcal{C}^3$ , we have the formula for the remainder term  $R_3$

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}_0) + \frac{1}{1!} Df(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2!} (\mathbf{x} - \mathbf{x}_0)^T D^2 f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \\ &\quad + o(\|\mathbf{x} - \mathbf{x}_0\|^2) \\ f(\mathbf{x}) &= f(\mathbf{x}_0) + \frac{1}{1!} Df(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2!} (\mathbf{x} - \mathbf{x}_0)^T D^2 f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \\ &\quad + O(\|\mathbf{x} - \mathbf{x}_0\|^3) \end{aligned}$$

# Taylor Theorem

Suppose that  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R} \in \mathcal{C}^2$ . Then we have that

$$f(\mathbf{x}) = f(\mathbf{x}_0) + Df(\mathbf{x}_0 + \theta\mathbf{h})\mathbf{h}$$

$$f(\mathbf{x}) = f(\mathbf{x}_0) + Df(\mathbf{x}_0)\mathbf{h} + \frac{1}{2}\mathbf{h}^T D^2 f(\mathbf{x}_0 + \theta\mathbf{h})\mathbf{h}$$

for some  $\theta \in (0, 1)$  and  $\mathbf{h} = \mathbf{x} - \mathbf{x}_0$ .



# Taylor Theorem

Let  $g : \mathbb{R} \rightarrow \mathbb{R} \in \mathcal{C}^1$ . By Taylor's theorem in 1D (or mean value theorem)

$$g(1) = g(0) + g'(\theta)$$

with  $\theta \in (0, 1)$ . Let define  $g(t) = f(\mathbf{z}(t))$  with  $\mathbf{z}(t) = \mathbf{x}_0 + t\mathbf{h}$  and  $\mathbf{h} = \mathbf{x} - \mathbf{x}_0$ . Therefore,

$$g(1) = f(\mathbf{z}(1)) = f(\mathbf{x})$$

$$g(0) = f(\mathbf{z}(0)) = f(\mathbf{x}_0)$$

$$g'(\theta) = Df(\mathbf{z}(\theta))\mathbf{h} = \nabla f(\mathbf{x}_0 + \theta\mathbf{h})^T \mathbf{h}$$

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0 + \theta\mathbf{h})^T \mathbf{h}$$

# Taylor Theorem

The proof of the second part is similar. Let  $g : \mathbb{R} \rightarrow \mathbb{R} \in \mathcal{C}^1$ . By Taylor's theorem in 1D

$$g(1) = g(0) + g'(0) + \frac{1}{2}g''(\theta)$$

with  $\theta \in (0, 1)$ . Let define  $g(t) = f(\mathbf{z}(t))$  with  $\mathbf{z}(t) = \mathbf{x}_0 + t\mathbf{h}$  and  $\mathbf{h} = \mathbf{x} - \mathbf{x}_0$ . Therefore, for  $t = \|\mathbf{h}\|$

$$\begin{aligned} g(1) &= f(\mathbf{z}(1)) = f(\mathbf{x}) \\ g(0) &= f(\mathbf{z}(0)) = f(\mathbf{x}_0) \\ g'(0) &= Df(\mathbf{z}(0))\mathbf{h} = \nabla f(\mathbf{x}_0)^T \mathbf{h} \end{aligned}$$

# Taylor Theorem

On the other hand

$$\begin{aligned}g'(\theta) &= Df(z(\theta))\mathbf{h} = \nabla f(z(\theta))^T \mathbf{h} \\g''(\theta) &= \mathbf{h}^T D^2 f(z(t)) D\mathbf{z}(t) = \mathbf{h}^T D^2 f(z(\theta)) \mathbf{h} \\&= \mathbf{h}^T D^2 f(\mathbf{x}_0 + \theta \mathbf{h}) \mathbf{h}\end{aligned}$$

now using  $g(1) = g(0) + g'(0) + \frac{1}{2}g''(\theta)$  and substituting one obtains the result.

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T D^2 f(\mathbf{x}_0 + \theta \mathbf{h}) \mathbf{h}$$