

# Sandwich theorem for $\beta$ -smooth convex function

## Andersen Ang

ECS, Uni. Southampton, UK  
andersen.ang@soton.ac.uk

Homepage [angms.science](http://angms.science)

Version: April 6, 2023

First draft: June 6, 2017

## Content

$$\frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

Sandwich theorem for  $\beta$ -smooth convex function

Proving the upper bound

Proving the lower bound

Application - proving gradient is co-coercive

## Convex $\beta$ -smooth function

- ▶ Setup: given a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  that is differentiable
- ▶  $f$  is convex if  $f$  and  $\text{dom}f (= \mathbb{R}^n)$  are convex.
- ▶  $f$  is  $\beta$ -smooth if  $\nabla f$  is  $\beta$ -Lipschitz. I.e., for any  $\mathbf{a}, \mathbf{b} \in \text{dom}f$ , we have a constant  $\beta$  such that:

$$\|\nabla f(\mathbf{a}) - \nabla f(\mathbf{b})\| \leq \beta \|\mathbf{a} - \mathbf{b}\|.$$

Details [here](#).

- ▶ **Theorem.** If  $f$  is convex  $\beta$ -smooth, then  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$\frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \leq \left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \right| \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2. \quad (1)$$

- ▶  $f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) =$  1st-order Taylor series of  $f$  at  $\mathbf{x}$
- ▶  $f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$  is the **Bregman divergence** (error made by the approximation)

It means: for a convex  $\beta$ -smooth  $f$ , its 1st-order Taylor approximation has bounded error.

- ▶ This PDF: prove (1).

## Prerequisite: directional derivative

▶ Given a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  as  $f(\mathbf{x}) = f(x_1, x_2, \dots, x_n)$ .

- ▶ We are in simple Euclidean space  $\mathbb{R}^n$
- ▶ Input of  $f$  is a vector  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  in  $\mathbb{R}^n$
- ▶ Output of  $f$  is a scalar in  $\mathbb{R}$

▶ **Directional derivative** of  $f$  along a vector  $\mathbf{v} = (v_1, v_2, \dots, v_n)$ , denoted as  $D_{\mathbf{v}}f(\mathbf{x})$ , is defined as

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

- ▶  $\mathbf{v}$  specifies a direction where  $\mathbf{x}$  is moving, and  $h$  is the stepsize
- ▶  $\mathbf{v}$  is possibly not a unit vector

▶ **Theorem:** let  $\nabla : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a scalar to vector mapping and let  $\nabla f$  be the gradient of  $f$ , then

$$D_{\mathbf{v}}f(\mathbf{x}) = \nabla f(\mathbf{x})^{\top} \mathbf{v}.$$

(Point,direction) to (point,point) form

$$D_v f(\mathbf{x}) := \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h} = \nabla f(\mathbf{x})^\top \mathbf{v}$$

- ▶ Now instead of  $(\mathbf{x}, \mathbf{v})$ , we are given two points  $(\mathbf{x}, \mathbf{y})$ . Let  $\mathbf{v} = \mathbf{y} - \mathbf{x}$  be the direction pointing from  $\mathbf{x}$  to  $\mathbf{y}$ . The directional derivative is then

$$D_v f(\mathbf{x}) = D_{\mathbf{y}-\mathbf{x}} f(\mathbf{x}) := \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{h} = \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}).$$

- ▶ Now let  $h \in [0, 1]$  and we consider the line segment  $[\mathbf{x}, \mathbf{y}]$

- ▶ At one end we have  $\mathbf{x} = \mathbf{x} + 0 \cdot (\mathbf{y} - \mathbf{x})$

$$h = 0$$

- ▶ At another end we have  $\mathbf{y} = \mathbf{x} + 1 \cdot (\mathbf{y} - \mathbf{x})$

$$h = 1$$

This means

- ▶  $D_v f(\mathbf{x}) = D_v (\mathbf{x} + h(\mathbf{y} - \mathbf{x})) \Big|_{h=0} = \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) \Big|_{h=0}$

- ▶  $D_v f(\mathbf{y}) = D_v (\mathbf{x} + h(\mathbf{y} - \mathbf{x})) \Big|_{h=1} = \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) \Big|_{h=1}$

- ▶ Why do these? Because now we have

$$f(\mathbf{y}) - f(\mathbf{x}) = \int_0^1 \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) dh$$

to prove the upper bound in the sandwich theorem.

Proving the upper bound  $\left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \right| \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \dots 1/3$

► Start with the integral

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &= \int_0^1 \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) dh \\ &= \int_0^1 \left[ \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) + \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}) \right]^\top (\mathbf{y} - \mathbf{x}) dh \\ &= \int_0^1 \left[ \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \right]^\top (\mathbf{y} - \mathbf{x}) dh + \int_0^1 \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) dh \\ &= \int_0^1 \left[ \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \right]^\top (\mathbf{y} - \mathbf{x}) dh + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \underbrace{\int_0^1 dh}_{=1} \end{aligned}$$

► Now we have

$$f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) = \int_0^1 \left[ \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \right]^\top (\mathbf{y} - \mathbf{x}) dh.$$

Compare it with the bound we see that we put absolute value sign next.

Proving the upper bound  $\left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \right| \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \dots 2/3$

► Put absolute value sign on

$$f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) = \int_0^1 \left[ \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \right]^\top (\mathbf{y} - \mathbf{x}) dh$$

gives

$$0 \leq \left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \right| = \left| \int_0^1 \left[ \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \right]^\top (\mathbf{y} - \mathbf{x}) dh \right|$$

► By  $\left| \int_a^b f(x) dx \right| \leq \int_a^b |f(x)| dx$

$$0 \leq \left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \right| \leq \int_0^1 \left| \left[ \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \right]^\top (\mathbf{y} - \mathbf{x}) \right| dh$$

► By  $|\mathbf{a}^\top \mathbf{b}| \leq \|\mathbf{a}\| \cdot \|\mathbf{b}\|$

$$\begin{aligned} 0 \leq \left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \right| &\leq \int_0^1 \left\| \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \right\| \cdot \|\mathbf{y} - \mathbf{x}\| dh \\ &= \|\mathbf{y} - \mathbf{x}\| \int_0^1 \left\| \nabla f(\mathbf{x} + h(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \right\| dh \end{aligned}$$

Proving the upper bound  $\left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \right| \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \dots 3/3$

►  $f$  is  $\beta$ -smooth  $\iff \|\nabla f(\mathbf{a}) - \nabla f(\mathbf{b})\| \leq \beta \|\mathbf{a} - \mathbf{b}\|$ , so

$$0 \leq \left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \right| \leq \|\mathbf{y} - \mathbf{x}\| \int_0^1 \beta \|h(\mathbf{y} - \mathbf{x})\| dh \leq \beta \|\mathbf{y} - \mathbf{x}\|^2 \int_0^1 |h| dh$$

►  $\int_0^1 |t| dt = \frac{1}{2}$

$$0 \leq \left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \right| \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

►  $f$  is convex  $\iff$  the Bregman divergence  $f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$  is nonnegative so

$$0 \leq f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad \square \quad (2)$$

Note: if  $f$  is non-convex,  $0 \leq f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$  does not hold and we need the absolute sign.

► (2) gives

$$f(\mathbf{y}) - f(\mathbf{x}) \leq \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad (\#)$$

Proving the lower bound  $\frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \dots 1/2$

$$f(\mathbf{y}) - f(\mathbf{x}) = f(\mathbf{z}) - f(\mathbf{x}) - [f(\mathbf{z}) - f(\mathbf{y})] \quad \text{not so trivial}$$

$$\geq \nabla f(\mathbf{x})^\top (\mathbf{z} - \mathbf{x}) - [f(\mathbf{z}) - f(\mathbf{y})] \quad f \text{ is cvx}$$

$$\geq \nabla f(\mathbf{x})^\top (\mathbf{z} - \mathbf{x}) - \left[ \nabla f(\mathbf{y})^\top (\mathbf{z} - \mathbf{y}) + \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 \right] \quad \text{by (\#)}$$

$$= \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x} + \mathbf{z} - \mathbf{y}) - \nabla f(\mathbf{y})^\top (\mathbf{z} - \mathbf{y}) - \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|_2^2$$

$$= \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \left( \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \right)^\top (\mathbf{z} - \mathbf{y}) - \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|_2^2$$

We now have

$$\left( \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \right)^\top (\mathbf{z} - \mathbf{y}) - \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$$

Compare with the bound, we need to form  $\frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2$  in the left-hand-side



Proving the lower bound  $\frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \dots 2/2$

$$\left(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\right)^\top (\mathbf{z} - \mathbf{y}) - \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$$

When we introduce  $\mathbf{z}$ ,  $\mathbf{z}$  is free so we can set any value for  $\mathbf{z}$ .

**A very tricky step:** let  $\mathbf{z} = \mathbf{y} - \frac{1}{\beta} (\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}))$

$$\mathbf{z} - \mathbf{y} = -\frac{1}{\beta} (\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}))$$

$$\left(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\right)^\top (\mathbf{z} - \mathbf{y}) = \frac{1}{\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2$$

$$-\frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 = -\frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2$$

$$\left(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\right)^\top (\mathbf{z} - \mathbf{y}) - \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 = \frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \quad \square$$

## Application - proving gradient is coercive

- Consider the lower bound of the theorem on  $\mathbf{x}, \mathbf{y}$  and  $\mathbf{y}, \mathbf{x}$

$$\frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$$

$$\frac{1}{2\beta} \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|^2 \leq f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y})$$

Sum them up

$$\frac{1}{\beta} \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|^2 \leq (\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) \quad (\dagger)$$

this inequality is called co-coercivity of gradient.

- Remark: apply Cauchy-Schwartz inequality  $|ab| \leq |a||b|$  on  $(\dagger)$  gives the Lipschitz continuity of  $\nabla f$ :  
 $\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \leq \beta \|\mathbf{y} - \mathbf{x}\|$ .  
So  $f$  is  $\beta$ -smooth  $\iff \nabla f$  is Lipschitz continuous  $\iff \nabla f$  is coercive.

## Last page - summary

1. For a  $\beta$ -smooth  $f$ ,  $\forall \mathbf{x}, \mathbf{y} \in \text{dom} f$ ,

$$0 \leq |f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})| \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

2. Sandwich theorem for  $\beta$ -smooth convex function: if  $f$  is convex and  $\beta$ -smooth,  $\forall \mathbf{x}, \mathbf{y} \in \text{dom} f$ ,

$$\frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

Implication: 1st-order Taylor series has bounded error.

3. Application: proving gradient is co-coercive

$$\frac{1}{\beta} \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|^2 \leq \left( \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) \right)^\top (\mathbf{y} - \mathbf{x}).$$

End of document