

# Special classes of function in optimization in machine learning

## Andersen Ang

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

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

Version: April 28, 2023

First draft: June 6, 2017

## Content

Convex

$\alpha$ -strongly convex

Lipschitz

Smooth / Lipschitz gradient

Relatively-smooth

Lipschitz continuous Hessian

Strongly convex & smooth

Other properties

Lower semicontinuous

Closed and proper

argmin

Polyak-Łojasiewicz and Kurdyka-Łojasiewicz

## Some old-school terminology

Notation used by Nesterov, Mordukhovich, or any classical real analysis textbooks:

- ▶  $f \in C^0$  :  $f(\mathbf{x})$  is continuous
- ▶  $f \in C^1$  :  $f(\mathbf{x})$  and  $\nabla f(\mathbf{x})$  are continuous
- ▶  $f \in C^2$  :  $f(\mathbf{x})$ ,  $\nabla f(\mathbf{x})$  and  $\nabla^2 f(\mathbf{x})$  are continuous
- ▶  $f \in C^{1,1}$  :  $f(\mathbf{x})$  and  $\nabla f(\mathbf{x})$  are continuous,  $\nabla f(\mathbf{x})$  is  $L$ -Lipschitz with  $L < +\infty$
- ▶  $f \in C_L^{k,p}$  :  $f$  is  $k$  times continuously differentiable and  $p$ th derivative is  $L$ -Lipschitz
- ▶  $f \in \mathcal{F}_L^k$  :  $f$  is  $\mathcal{C}_L^k$  and convex
- ▶  $f \in \mathcal{S}_{M,L}^k$  :  $f$  is  $\mathcal{F}_L^k$  and  $M$ -strongly convex

# Table of Contents

## Convex

- $\alpha$ -strongly convex

## Lipschitz

- Smooth / Lipschitz gradient

- Relatively-smooth

- Lipschitz continuous Hessian

## Strongly convex & smooth

## Other properties

- Lower semicontinuous

- Closed and proper

- argmin

- Polyak-Łojasiewicz and Kurdyka-Łojasiewicz

## Convex function

A function  $f(\mathbf{x}) : \text{dom } f \rightarrow \mathbb{R}$  is **convex** if :

- ▶  $\text{dom } f$  is a convex set<sup>1</sup>
- ▶  $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f$ , we have any one of the following
  1. Jensen's inequality:  $f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y})$ .
  2. Epigraph of  $f$  is a convex set.
  3. 1st-order Taylor series at  $\mathbf{x}$  is a global under-estimator:  $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ .
  4. Gradient is monotonic:  $\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq 0$ .(For 3,4, if  $f$  is not differentiable, we replace gradient by subgradient.)

- ▶ The 4 definitions are equivalent ("if and only if"). See optimization books for proof. See [here](#) for proof of  $1 \iff 3$ .
- ▶ If  $f$  is twice differentiable, it is convex iff  $\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$ .
- ▶  $f$  is **strictly convex** if  $\leq, \geq$  became  $<, >$  (i.e. strict inequality).

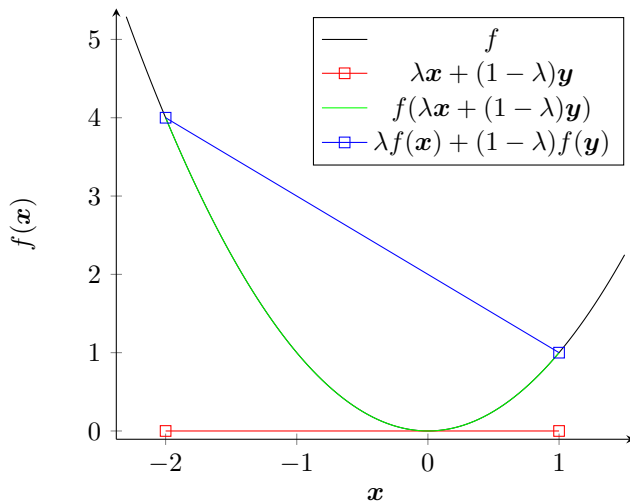
---

<sup>1</sup> $\text{dom } f$  can be open set. However, in optimization usually  $\text{dom } f$  is closed because optimization over an open set has no solution.

## Convexity: the geometry of Jensen's inequality

$f : \text{dom } f \rightarrow \mathbb{R}$  is **convex** IF

- (1)  $\text{dom } f$  is a convex set and
- (2)  $\forall x, y \in \text{dom } f, f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$

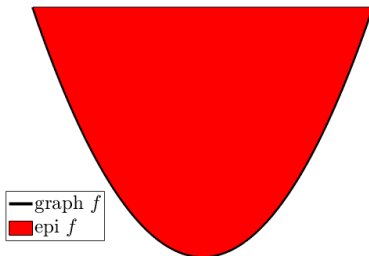


# Convexity: the convex geometry of epigraph

$f : \text{dom } f \rightarrow \mathbb{R}$  is **convex** IF **epigraph of  $f$**  is a convex set

Visualization of  $\text{graph } f$  and  $\text{epi } f$

- ▶  $\text{epi } f =$  **all** the points of  $\mathbb{R}^{n+1}$  lying **on or above**  $\text{graph } f$ .
- ▶ Example:  $f(x) = x^2$ 
  - ▶  $n = 1$  (1-dimensional)
  - ▶  $\text{graph } f := \{(x, y) \in \mathbb{R} \times \mathbb{R} : y = f(x)\}$  is a 1d curve in a 2d space.
  - ▶  $\text{epi } f := \{(x, \alpha) \in \mathbb{R} \times \mathbb{R} : \alpha \geq f(x)\}$  is a 2d set in a 2d space.



## Convexity: the geometry of 1st-order Taylor series

- $f : \text{dom } f \rightarrow \mathbb{R}$  is **convex** if :

1.  $\text{dom } f$  is a convex set
2.  $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f$ , we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle. \quad (*)$$

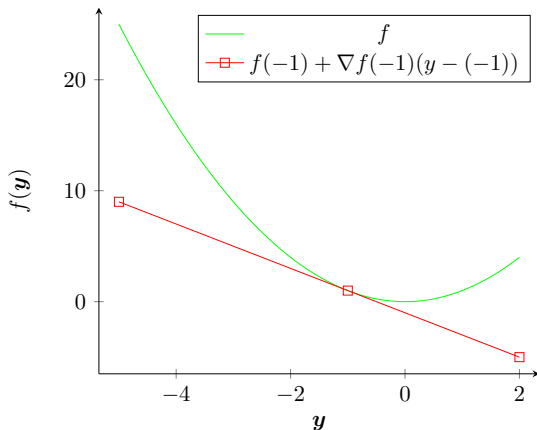
- (\*) assumes  $f$  is differentiable at  $\mathbf{x}$ . If  $f$  is not differentiable at  $\mathbf{x}$ , we generalize gradient with **subgradient**:

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{q}, \mathbf{y} - \mathbf{x} \rangle. \quad (\#)$$

i.e., we replace  $\nabla f(\mathbf{x})$  by any vector  $\mathbf{q}$  that (#) holds.

- In fact, subgradient is defined using (#)

- The gap between  $f$  and the 1st-order Taylor series is known as the **Bregman Divergence**.



# Why convex and differentiable $f$ is lower-bounded by their own 1st-order Taylor series?

- Consider a pedagogical case:  $f$  is (twice) differentiable of single variable, then

$$\begin{aligned} f(y) &= f(x) + f'(x)(y - x) + o(y - x) && \text{Taylor series} \\ &= f(x) + f'(x)(y - x) + \frac{f''(\xi)}{2}(y - x)^2 && \text{see 1} \\ &\geq f(x) + f'(x)(y - x) && \text{see 2} \end{aligned}$$

1. Lagrange remainder theorem: using mean-value theorem, the remainder term  $o(y - x) = \frac{f''(\xi)}{2}(y - x)^2$  for some  $\xi$  in the interval  $[x, y]$ .
2. As  $f$  is convex, which means  $f'' \geq 0$  so the last term is nonnegative.

- The arguments above generalize to multi-variable  $f$ .
- Note: **this is not a prove** but an illustration, because
  - apart from assuming  $f$  is differentiable, we assumed  $f$  is twice differentiable,
  - we didn't show that  $f$  is convex  $\iff$  its Hessian is positive semi-definite.



## $\alpha$ -strongly convex function

A function  $f : \text{dom } f \rightarrow \mathbb{R}$  is  $\alpha$ -strongly convex if:

- ▶  $\text{dom } f$  is a convex set.
- ▶  $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f$ , we have any one of the following

1. Jensen's inequality with an additional quadratic term with  $\alpha > 0$

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) - \frac{\alpha}{2} \lambda (1 - \lambda) \|\mathbf{x} - \mathbf{y}\|_2^2.$$

2.  $\text{grad } f$  is monotonic with an additional quadratic term with  $\alpha > 0$

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \alpha \|\mathbf{x} - \mathbf{y}\|_2^2 \geq 0.$$

3. 1st-order Taylor series at  $\mathbf{x}$  is global under-estimator with an additional quadratic term with  $\alpha > 0$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2,$$

or we say  $f$  is lower bounded by a quadratic function.

4. With  $\alpha > 0$ , the function  $f(\mathbf{x}) - \frac{\alpha}{2} \|\mathbf{x}\|_2^2$  is convex.

- ▶ These definitions are equivalent.
- ▶ If  $f$  is twice differentiable, it is  $\alpha$ -strongly convex iff  $\nabla^2 f(\mathbf{x}) \succeq \alpha \mathbf{I}$ .

## Illustrating equivalence between definitions of strong convexity

For  $\alpha > 0$  and  $f$  twice differentiable,  $\nabla^2 f(\mathbf{x}) \succeq \alpha \mathbf{I} \implies \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \alpha \|\mathbf{x} - \mathbf{y}\|_2^2$ .

- **Proof.** Recall from calculus  $G(b) - G(a) = \int_a^b g(\theta) d\theta$ . Next, a smart step, let  $\theta = \mathbf{y} + \tau(\mathbf{x} - \mathbf{y})$ , then  $d\theta = (\mathbf{x} - \mathbf{y}) d\tau$ . Consider integral range from 0 to 1 for  $\tau$  we let  $G$  be  $\nabla f$  and  $g$  be  $\nabla^2 f$ , this gives

$$\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) = \int_0^1 \nabla^2 f(\mathbf{y} + \tau(\mathbf{x} - \mathbf{y})) (\mathbf{x} - \mathbf{y}) d\tau.$$

(left hand side is a vector, right hand side is matrix-vector product, also a vector)

- Take dot product with  $\mathbf{x} - \mathbf{y}$  on the whole equation on both sides

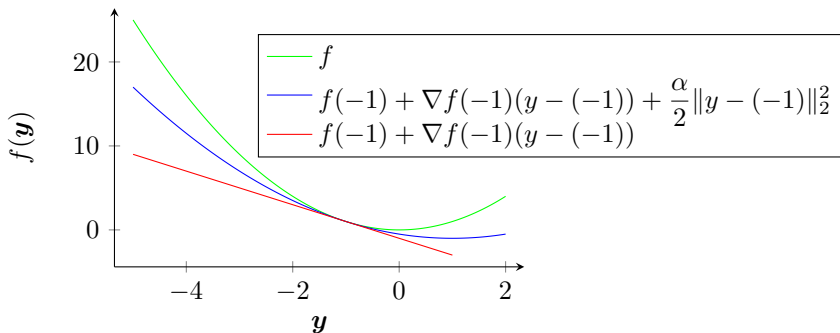
$$\begin{aligned} \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle &= \left\langle \mathbf{x} - \mathbf{y}, \int_0^1 \nabla^2 f(\mathbf{y} + \tau(\mathbf{x} - \mathbf{y})) (\mathbf{x} - \mathbf{y}) d\tau \right\rangle \\ &\geq \left\langle \mathbf{x} - \mathbf{y}, \int_0^1 \alpha (\mathbf{x} - \mathbf{y}) d\tau \right\rangle \\ &= \alpha \|\mathbf{x} - \mathbf{y}\|_2^2, \end{aligned}$$

where the inequality is due to  $\nabla^2 f(\mathbf{x}) \succeq \alpha \mathbf{I}$  for all  $\mathbf{x}$ : we have  $\nabla^2 f(\mathbf{y} + \tau(\mathbf{x} - \mathbf{y})) \succeq \alpha \mathbf{I}$ . ■

$\alpha$ -strongly convex: the geometry of the lower bounded

$f(x) : \text{dom } f \rightarrow \mathbb{R}$  is  $\alpha$ -**strongly convex** if

(1)  $\text{dom } f$  is a convex and (2)  $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f$ :  $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$



**Meaning:**  $f$  is lower bounded by a **quadratic curve with some curvature**, which is also lower bounded by the **1st order Taylor series (zero curvature)**

$\implies f$  is not “too flat” / at least “as curved as” **the lower bound**  
In other words:  $f$  is at least  $\alpha$ -amount of “bumpy”.

## Remarks on convexity

- ▶ Strongly convex  $\implies$  strictly convex  $\implies$  convex.

The opposite is false.

- ▶ e.g.,  $x^4$  is strictly convex but not strongly convex.  
Why:  $x^4$  is not globally lower-bounded by  $x^2$ .

- ▶ Convexity function needs not to be differentiable.
  - ▶ That's why we have Jansen's definition

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}),$$

which does not involve  $\nabla f$ .

- ▶ Strongly convex functions are **coercive**.

# Table of Contents

Convex

$\alpha$ -strongly convex

Lipschitz

Smooth / Lipschitz gradient

Relatively-smooth

Lipschitz continuous Hessian

Strongly convex & smooth

Other properties

Lower semicontinuous

Closed and proper

argmin

Polyak-Łojasiewicz and Kurdyka-Łojasiewicz

## Lipschitz continuity

A function  $f(\mathbf{x}) : \text{dom } f \rightarrow \mathbb{R}$  is *Lipschitz* if for any  $\mathbf{x}, \mathbf{y} \in \text{dom } f$ , there exists a constant  $L \geq 0$  (the Lipschitz constant) such that

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|.$$

- Re-arrange gives

$$\frac{|f(\mathbf{x}) - f(\mathbf{y})|}{\|\mathbf{x} - \mathbf{y}\|} \leq L \quad \mathbf{y} \xrightarrow{\approx} \mathbf{x} \quad \text{size of } \nabla f(\mathbf{x}) \leq L$$

$\implies f$  is Lipschitz means the “slope” (rate of change) of  $f$  is bounded above globally by  $L$ .

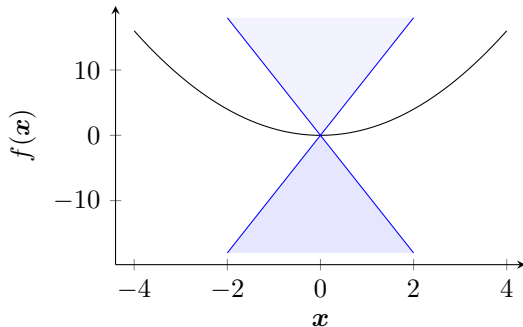
- Removing the absolute value sign:

$$\begin{cases} f(\mathbf{x}) \leq f(\mathbf{y}) + L\|\mathbf{x} - \mathbf{y}\| \\ f(\mathbf{x}) \geq f(\mathbf{y}) - L\|\mathbf{x} - \mathbf{y}\| \end{cases}$$

means that  $f$  for all  $\mathbf{x}$  is bounded above and below by a linear function constructed at  $\mathbf{y}$ .

# The geometry of Lipschitz continuity

$f$  is Lipschitz  $\iff f$  does not have sharp change everywhere:  $\forall x$  the curve  $f$  is entirely outside a cone which is modeled by the linear functions in the last page.



Important note: such property is **global**, such cone exists for all points on  $f$ . i.e. the cone can “slide” along the curve and the argument still holds.

## Lipschitz continuity and differentiability

- ▶ **Q:** If  $f$  is Lipschitz continuous, is  $f$  differentiable?

**A:** No.

- ▶ **Rademacher's theorem:** Lipschitz function is *almost everywhere* differentiable.

Almost everywhere  $\neq$  everywhere.

- ▶ Example.  $|x|$

- ▶  $|x|$  is 1-Lipschitz but not differentiable at  $x = 0$ .

- ▶ However, the single point  $x = 0$  has a measure zero<sup>2</sup> on  $\mathbb{R}$ , this is what “almost everywhere” means in Rademacher's theorem.

---

<sup>2</sup>The probability of getting this number in a random guess on the real line is zero, because there are infinitely many real numbers.



## Composition of Lipschitz functions

- ▶ Suppose  $f_1$  is  $L_1$ -Lipschitz and  $f_2$  is  $L_2$ -Lipschitz. Then  $f_1 \circ f_2$  is  $L_1 L_2$ -Lipschitz.
- ▶  $f_1 \circ f_2$  means the composition of  $f_1$  and  $f_2$ , i.e.,  $f_1(f_2)$
- ▶ The proof: direct proof

$$\begin{aligned}\|(f_1 \circ f_2)(\mathbf{x}) - (f_1 \circ f_2)(\mathbf{y})\| &\leq \|f_1(f_2(\mathbf{x})) - f_1(f_2(\mathbf{y}))\| \\ &\leq L_1 \|f_2(\mathbf{x}) - f_2(\mathbf{y})\| && f_1 \text{ is } L_1\text{-Lipschitz} \\ &\leq L_1 L_2 \|\mathbf{x} - \mathbf{y}\| && f_2 \text{ is } L_2\text{-Lipschitz}\end{aligned}$$

(The proof holds for any norm, not only for  $\ell_2$  norm)

- ▶ This result is commutative:  $f_1 \circ f_2$  and  $f_2 \circ f_1$  are both  $L_1 L_2$ -Lipschitz
- ▶ A small detail: in Euclidean space  $f_1 \circ f_2$  assumes the output dimension of  $f_2$  match the input dimension of  $f_1$
- ▶ Corollary:  $f_1 \circ f_2 \circ \dots \circ f_n$  is  $L_1 L_2 \dots L_n$ -Lipschitz

## Sum of Lipschitz functions

► Suppose  $f_1$  is  $L_1$ -Lipschitz and  $f_2$  is  $L_2$ -Lipschitz. Then  $\alpha_1 f_1 + \alpha_2 f_2$  is  $|\alpha_1|L_1 + |\alpha_2|L_2$ -Lipschitz.

► **Proof** First we group the terms

$$\left\| \alpha_1 f_1(\mathbf{x}) + \alpha_2 f_2(\mathbf{x}) - \alpha_1 f_1(\mathbf{y}) + \alpha_2 f_2(\mathbf{y}) \right\| \leq \left\| \alpha_1 (f_1(\mathbf{x}) - f_1(\mathbf{y})) + \alpha_2 (f_1(\mathbf{y}) - f_2(\mathbf{y})) \right\|$$

Use triangle inequality<sup>3</sup>

$$\begin{aligned} \left\| \alpha_1 f_1(\mathbf{x}) + \alpha_2 f_2(\mathbf{x}) - \alpha_1 f_1(\mathbf{y}) + \alpha_2 f_2(\mathbf{y}) \right\| &\leq \left\| \alpha_1 (f_1(\mathbf{x}) - f_1(\mathbf{y})) \right\| + \left\| \alpha_2 (f_1(\mathbf{y}) - f_2(\mathbf{y})) \right\| \\ &\leq |\alpha_1| \|f_1(\mathbf{x}) - f_1(\mathbf{y})\| + |\alpha_2| \|f_1(\mathbf{y}) - f_2(\mathbf{y})\| \\ &\leq |\alpha_1| L_1 \|\mathbf{x} - \mathbf{y}\| + |\alpha_2| L_2 \|\mathbf{x} - \mathbf{y}\| \\ &= (|\alpha_1| L_1 + |\alpha_2| L_2) \|\mathbf{x} - \mathbf{y}\| \end{aligned}$$

---

<sup>3</sup>First for the squared term  $\|\mathbf{a} + \mathbf{b}\|^2 \leq \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 + 2\langle \mathbf{a}, \mathbf{b} \rangle \leq \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 + 2\|\mathbf{a}\|\|\mathbf{b}\| = (\|\mathbf{a}\| + \|\mathbf{b}\|)^2$ .  
Remove the square we have  $\|\mathbf{a} + \mathbf{b}\| \leq \|\mathbf{a}\| + \|\mathbf{b}\|$

# Max of Lipschitz functions

- Suppose  $f_1$  is  $L_1$ -Lipschitz and  $f_2$  is  $L_2$ -Lipschitz. Then  $\max\{f_1, f_2\}$  is  $\max\{L_1, L_2\}$ -Lipschitz.
- Tools we need

$$a \leq |a|$$

$$a \leq \max\{a, b\}$$

$$\begin{cases} a \leq M \\ b \leq M \end{cases} \iff \max\{a, b\} \leq M$$

$$a \leq M \text{ and } -a \leq M \text{ imply } |a| \leq M$$

- **Proof**  $f_1$  is Lipschitz so  $|f_1(\mathbf{x}) - f_1(\mathbf{y})| \leq L_1 \|\mathbf{x} - \mathbf{y}\|$ . By  $\blacksquare$   $f_1(\mathbf{x}) - f_1(\mathbf{y}) \leq L_1 \|\mathbf{x} - \mathbf{y}\|$ , which gives

$$f_1(\mathbf{x}) \leq f_1(\mathbf{y}) + L_1 \|\mathbf{x} - \mathbf{y}\| \iff f_1(\mathbf{x}) \leq \max\{f_1(\mathbf{y}), f_2(\mathbf{y})\} + \max\{L_1, L_2\} \|\mathbf{x} - \mathbf{y}\| \quad (1)$$

Similarly,

$$f_2(\mathbf{x}) \leq \max\{f_1(\mathbf{y}), f_2(\mathbf{y})\} + \max\{L_1, L_2\} \|\mathbf{x} - \mathbf{y}\| \quad (2)$$

By  $\blacksquare$ , (1) and (2) gives

$$\max\{f_1(\mathbf{x}), f_2(\mathbf{x})\} \leq \max\{f_1(\mathbf{y}), f_2(\mathbf{y})\} + \max\{L_1, L_2\} \|\mathbf{x} - \mathbf{y}\| \quad (3)$$

(3) holds by swapping  $(\mathbf{x}, \mathbf{y})$  as  $(\mathbf{y}, \mathbf{x})$ , we have

$$\max\{f_1(\mathbf{y}), f_2(\mathbf{y})\} \leq \max\{f_1(\mathbf{x}), f_2(\mathbf{x})\} + \max\{L_1, L_2\} \|\mathbf{x} - \mathbf{y}\| \quad (4)$$

$$(3) \iff \underbrace{\max\{f_1(\mathbf{x}), f_2(\mathbf{x})\} - \max\{f_1(\mathbf{y}), f_2(\mathbf{y})\}}_a \leq \max\{L_1, L_2\} \|\mathbf{x} - \mathbf{y}\|$$

$$(4) \iff \underbrace{\max\{f_1(\mathbf{y}), f_2(\mathbf{y})\} - \max\{f_1(\mathbf{x}), f_2(\mathbf{x})\}}_{-a} \leq \max\{L_1, L_2\} \|\mathbf{x} - \mathbf{y}\|$$

By  $\blacksquare$ ,

$$\left| \max\{f_1(\mathbf{x}), f_2(\mathbf{x})\} - \max\{f_1(\mathbf{y}), f_2(\mathbf{y})\} \right| \leq \max\{L_1, L_2\} \|\mathbf{x} - \mathbf{y}\|. \quad \blacksquare$$

## $L$ -smooth function / Lipschitz continuous gradient

A function  $f : \text{dom } f \rightarrow \mathbb{R}$  is  $L$ -smooth if for any two points  $\mathbf{x}, \mathbf{y} \in \text{dom } f$ , there exists a constant  $L < +\infty$  such that

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

- ▶ This assume  $f$  is differentiable.
- ▶ “ $f$  is  $L$ -smooth” is also called  $L$ -Lipschitz gradient, or  $\mathcal{C}^{1,1}$ .
- ▶ “ $f$  is  $L$ -smooth” is equivalent to

$$|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

Removing the absolute value sign gives

$$\begin{cases} f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \\ f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \end{cases}$$

meaning that  $f$  is bounded above and below by a quadratic function.

## Equivalent definitions of $L$ -smooth function

A function  $f(x)$  is  $L$ -smooth if

- $\text{grad} f$  is  $L$ -Lipschitz with  $L \geq 0$ . I.e.  $\forall \mathbf{x}, \mathbf{y} \in \text{dom} f$  we have  $L \geq 0$

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

- $f$  is bounded by a quadratic function with  $L > 0$ :

$$|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

- the gradient of  $f$  is monotonic with additional term with  $L > 0$ :

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \frac{1}{L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2.$$

- the norm of the slope of  $\nabla f$  (which is  $\nabla^2 f$ ) is bounded above.
- If  $f$  is twice differentiable,  $\nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}$ , or all the eigenvalue of  $\nabla^2 f(\mathbf{x})$  is below  $L$ .

These definitions are equivalent. See [here](#) for more about the 2nd definition.

## Proof of equivalence

We show for  $L > 0$ ,  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$  implies  $|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$ .

Recall calculus  $G(b) - G(a) = \int_a^b g(\theta) d\theta$ . Next, a smart step, let  $g(\tau) = \langle \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle$  be a function in  $\tau$  and  $d\theta = d\tau$ . Consider the definite integral of  $g(\tau)$  from 0 to 1, let  $G(b) = f(\mathbf{y})$  and  $G(a) = f(\mathbf{x})$ , hence

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &= \int_0^1 \langle \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle d\tau \\ &= \int_0^1 \langle \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) + \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle d\tau. \end{aligned}$$

As  $\nabla f(\mathbf{x})$  is independent of  $\tau$ , can take out from the integral

$$f(\mathbf{y}) - f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 \langle \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle d\tau.$$

The idea is to create the term  $\langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$  so that we can move it to the left and get  $|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle|$

## Proof of equivalence - continue

$$\begin{aligned} |f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| &= \left| \int_0^1 \langle \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle d\tau \right| \\ &\leq \int_0^1 | \langle \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle | d\tau \\ &\stackrel{\text{Cauchy - Schwarz}}{\leq} \int_0^1 \| \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \| \cdot \| \mathbf{y} - \mathbf{x} \| d\tau. \end{aligned}$$

Look at  $\| \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \|$ , this is exactly where we can apply the Lipschitz gradient inequality

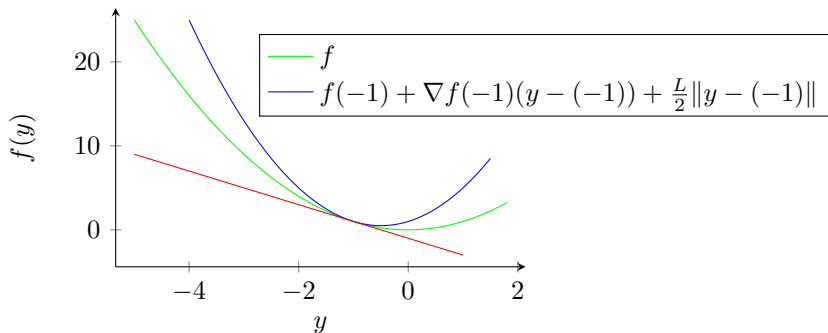
$$\| \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \| \leq L \| \tau(\mathbf{y} - \mathbf{x}) \| \leq L |\tau| \| \mathbf{y} - \mathbf{x} \| = L \tau \| \mathbf{y} - \mathbf{x} \|^2$$

where  $\| \tau(\mathbf{y} - \mathbf{x}) \| = |\tau| \| \mathbf{y} - \mathbf{x} \|$  as norm is non-negative. Note that the integral range is from 0 to 1 so the absolute sign in  $\tau$  can be removed. Lastly

$$|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \int_0^1 L \tau d\tau \cdot \| \mathbf{y} - \mathbf{x} \|^2 = \frac{L}{2} \| \mathbf{y} - \mathbf{x} \|^2. \quad \square$$

## $L$ -smoothness: the geometry of the upper bound

$f$  is  $L$ -**smooth** if  $\forall x, y \in \text{dom } f$ ,  $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|_2^2$



**Meaning:**  $f$  is globally bounded above by a quadratic function.

i.e.  $f$  cannot be “too sharp” ( $f$  is flatter than the upper bound), or  $f$  cannot grow “too fast”.



## Relatively-smooth function

- ▶  $f$  is  $L$ -smooth

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + L \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

- ▶  $f$  is  $L$ -smooth relative to the distance kernel  $h$

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + L D_h(\mathbf{x}, \mathbf{y}),$$

where  $D_h$  is the Bregman divergence on the distance kernel  $h$ .

- ▶ Why relative smoothness

- ▶ for proving convergence of gradient descent on non-Euclidean geometry
- ▶ for function that is not uniformly smooth,  
e.g. the slope of  $x^2 - \log(x)$  approaches to  $\infty$  as  $x \rightarrow 0$ , the value  $L$  change dramatically as  $x$  moves.
- ▶ application in minimizing  $\frac{1}{4} \|\mathbf{Ax} - \mathbf{b}\|_4^4$ .
- ▶ mirror descent

## Lipschitz continuous Hessian

A function  $f(\mathbf{x}) : \text{dom } f \rightarrow \mathbb{R}$  has  $L$ -Lipschitz Hessian, if  $\forall \mathbf{x}, \mathbf{y} \in \text{dom } f, \exists L < \infty$  such that

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

- ▶ This assumes  $f$  is twice differentiable.
- ▶ This means the norm of  $\nabla^3 f(\mathbf{x})$  is bounded above by  $L$ .
- ▶  $f$  has  $L$ -Lipschitz Hessian is equivalent to

$$\left| f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \langle \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \right| \leq \frac{L}{6} \|\mathbf{y} - \mathbf{x}\|_2^3$$

see [here](#) for the proof.

Removing the absolute value sign, and make  $\mathbf{y}$  the subject:

$$\begin{cases} f(\mathbf{y}) \geq f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \langle \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \frac{L}{6} \|\mathbf{y} - \mathbf{x}\|_2^3 \\ f(\mathbf{y}) \leq f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \langle \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{6} \|\mathbf{y} - \mathbf{x}\|_2^3 \end{cases}$$

which means  $f(\mathbf{y})$  is bounded above and below by two cubic functions parameterized at the point  $\mathbf{x}$  for all  $\mathbf{y}$ .

# Table of Contents

Convex

$\alpha$ -strongly convex

Lipschitz

Smooth / Lipschitz gradient

Relatively-smooth

Lipschitz continuous Hessian

**Strongly convex & smooth**

Other properties

Lower semicontinuous

Closed and proper

argmin

Polyak-Łojasiewicz and Kurdyka-Łojasiewicz

## Strongly convex smooth function

- ▶ A function  $f : \text{dom} \rightarrow \mathbb{R}$  is  $\alpha$ -strongly convex and  $\beta$ -smooth if
  - ▶  $f$  is  $\beta$ -smooth, which means  $f$  is differentiable and  $\nabla f$  is monotone

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \frac{1}{\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2.$$

- ▶  $f$  is  $\alpha$ -strongly convex, which means gradient is strongly monotone

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \alpha \|\mathbf{x} - \mathbf{y}\|_2^2.$$

- ▶ As  $f$  satisfies both monotone inequalities, so we have

$$\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \frac{\alpha\beta}{\alpha + \beta} \|\mathbf{x} - \mathbf{y}\|_2^2 + \frac{1}{\alpha + \beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2.$$

Details [here](#).

# Table of Contents

## Convex

$\alpha$ -strongly convex

## Lipschitz

Smooth / Lipschitz gradient

Relatively-smooth

Lipschitz continuous Hessian

## Strongly convex & smooth

## Other properties

Lower semicontinuous

Closed and proper

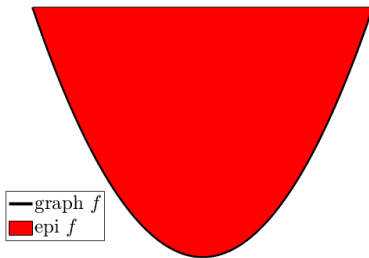
argmin

Polyak-Łojasiewicz and Kurdyka-Łojasiewicz

# Epigraph

## Visualization of $\text{graph } f$ and $\text{epi } f$

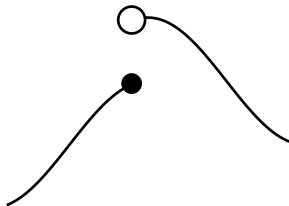
- ▶  $\text{epi } f =$  **all** the points of  $\mathbb{R}^{n+1}$  lying **on or above**  $\text{graph } f$ .
- ▶ Example:  $f(x) = x^2$ 
  - ▶  $n = 1$  (1-dimensional)
  - ▶  $\text{graph } f := \{(x, y) \in \mathbb{R} \times \mathbb{R} : y = f(x)\}$  is a 1d curve in a 2d space.
  - ▶  $\text{epi } f := \{(x, \alpha) \in \mathbb{R} \times \mathbb{R} : \alpha \geq f(x)\}$  is a 2d set in a 2d space.



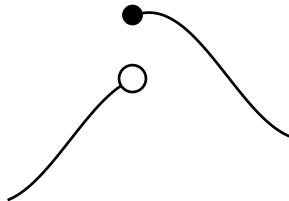
## Lower semicontinuous function (l.s.c.)

- ▶  $\bar{\mathbb{R}} := \mathbb{R} \cup \{+\infty\}$  is the extended real line.
- ▶ A function is continuous means it has no “jump”.
- ▶ A function is l.s.c. means it allow jump but still continuous if viewed from below.
- ▶ A function  $f$  is l.s.c. if its epigraph is closed.

L.S.C. (epi  $f$  is closed)



Not L.S.C. (epi  $f$  is open)



- ▶ Why care about l.s.c. function: indicator function of a closed convex set are all l.s.c..

## Closed and proper function

- ▶ A function  $f$  is proper if it never takes the value  $-\infty$  and  $\text{dom } f \neq \emptyset$

OR euivalently,

$\text{epi } f \neq \emptyset$  without svertical line<sup>4</sup>.

- ▶ A proper function  $f$  is closed if  $\text{dom } f$  is closed and  $f$  is lower semicontinuous at each  $x \in \text{dom } f$

OR equivalently,

$\text{epi } f$  is closed.

---

<sup>4</sup>which can move downward and touch  $-\infty$



argmin (argument of minimum = set of minimizer)

- argmin is a set defined as

$$\operatorname{argmin} f := \left\{ \boldsymbol{x} \in \operatorname{dom} f \mid f(\boldsymbol{x}) = \inf_{\boldsymbol{z} \in \operatorname{dom} f} f(\boldsymbol{z}) \right\}.$$

- If  $f$  is closed convex proper, then  $\operatorname{argmin} f$  is closed convex and possibly empty<sup>5</sup>

---

<sup>5</sup>If  $\operatorname{argmin} f$  is an empty set that means there is no minimizer for  $f$

# Polyak-Łojasiewicz and Kurdyka-Łojasiewicz

- ▶  $f$  is Polyak-Łojasiewicz (PŁ) if  $\exists \mu > 0$  such that  $\|\nabla f(\mathbf{x})\|_2^2 \geq \mu(f(\mathbf{x}) - f^*)$  for all  $\mathbf{x} \in \text{dom } f$ .
  - ▶ PŁ is weaker than strong convexity.
  - ▶ If  $f$  is  $\mu$ -strongly convex, then  $f$  is  $\mu$ -PŁ.
  - ▶ PŁ can be used as a tool to prove convergence of gradient descent, see [here](#) for more.
- ▶ Kurdyka-Łojasiewicz
  - ▶ Generalized PŁ: it can handle nonsmooth functions
  - ▶ KŁ is a tool for proving convergence of gradient method on nonsmooth optimization.
  - ▶ Very technical. The original full definition is long, so we give a simplified one here.  
 $f$  is KŁ at a point  $\bar{\mathbf{x}}$  if there exists  $c > 0$  and  $\mu \in [0, 1)$  such that  $\|\partial f(\mathbf{x})\|_2 \geq \frac{1}{c(1-\mu)} (f(\mathbf{x}) - f(\bar{\mathbf{x}}))^\mu$  holds for all  $\mathbf{x}$  within a neighbourhood of  $\bar{\mathbf{x}}$ . For  $\partial f(\mathbf{x})$ , we use the norm of the subgradient with smallest  $\ell_2$  norm to define  $\|\partial f(\mathbf{x})\|_2$ .
  - ▶ If  $f$  is a semi-algebraic function, then  $f$  is KŁ.
- ▶ Semi-algebraic function
  - ▶ A function is semi-algebraic if  $\text{epi } f$  is a semialgebraic set.
  - ▶ A set is semialgebraic if it is defined by polynomial equations and polynomial inequalities

$f$  is proper if  $\text{epi } f$  is non-empty and has no vertical line  
 proper  $f$  is closed if  $\text{epi } f$  is closed  
 $f$  is l.s.c. if  $\text{epi } f$  is closed.  
 $\text{argmin } f$  is closed convex if  $f$  is closed convex proper

proper  
 closed of proper  $f$   
 Lower semicontinuous  
 $\text{argmin } f$  closed convex

$f$  is convex if  $\text{dom } f$  is convex and

1.  $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$
2.  $\langle x - y, \nabla f(x) - \nabla f(y) \rangle \geq 0$
3.  $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$
4.  $\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$ , if  $f$  is twice differentiable
5.  $\text{epi } f$  is convex

Jensen  
 Gradient is monotone  
 1st-order Taylor series is global support  
 Hessian argument  
 epigraph is convex set

$f$  is  $\alpha$ -strongly convex if  $\text{dom } f$  is convex and

1.  $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) - \frac{\alpha}{2} \lambda(1 - \lambda) \|x - y\|_2^2$
2.  $\langle x - y, \nabla f(x) - \nabla f(y) \rangle \geq \alpha \|x - y\|_2^2$
3.  $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\alpha}{2} \|x - y\|_2^2$
4.  $f(x) - \frac{\alpha}{2} \|x\|_2^2$  is convex
5.  $\nabla^2 f(\mathbf{x}) \succeq \alpha \mathbf{I}$ , if  $f$  is twice differentiable

Jensen  
 Strongly monotone  
 Global quadratic lower bound  
 Convexity  
 Hessian argument

$f$  is  $L$ -Lipschitz gradient ( $L$ -smooth) if  $f$  is differentiable and

1.  $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$
2.  $|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq \frac{L}{2} \|y - x\|_2^2$
3.  $\langle x - y, \nabla f(x) - \nabla f(y) \rangle \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|_2^2$
4.  $\nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}$ , if  $f$  is twice differentiable

Definition of Lipschitz  
 Quadratic inequality  
 monotone  
 Hessian argument

$f$  is  $L$ -Lipschitz Hessian if  $f$  is twice differentiable and

1.  $\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L \|x - y\|$
2.  $|f(x) - f(y) - \langle \nabla f(x), y - x \rangle - \langle \nabla^2 f(x)(y - x), y - x \rangle| \leq \frac{L}{6} \|y - x\|_2^3$

Definition of Lipschitz  
 Cubic inequality

$f$  is  $\alpha$ -strongly convex and  $\beta$ -smooth

1.  $\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \frac{\alpha\beta}{\alpha + \beta} \|\mathbf{x} - \mathbf{y}\|_2^2 + \frac{1}{\alpha + \beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2$

End of document