

# Introductory Lectures on Optimization

## Beyond The Black-box Model (1)

Hui Qian  
qianhui@zju.edu.cn

College of Computer Science, Zhejiang University

December 4, 2024

# Outline

## 1 Proximal Gradient Method

- Proximal Operator
- Properties of Proximal Operator
- Analysis for Proximal Gradient Method
- Accelerated Proximal Gradient Method
- Special case: Proximal Point Method

## 2 Douglas-Rachford Splitting

- Different Settings for Convex Problem
- Fixed Point for Nonsmooth Composition
- Splitting Algorithm

## 3 Reference

# Part I

## Proximal Gradient Method

# Proximal Operator

The **proximal operator** ( or **proximal mapping**) of a convex function  $h$  is defined as

$$\text{prox}_h(\mathbf{x}) = \underset{\mathbf{u}}{\operatorname{argmin}} \left( h(\mathbf{u}) + \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|_2^2 \right).$$

## Example 1

- 1 If  $h(\mathbf{x}) = 0$ , we have  $\text{prox}_h(\mathbf{x}) = \mathbf{x}$ .
- 2 If  $h(\mathbf{x})$  is an **indicator function** of a closed convex set  $\mathcal{C}$ ,  $\text{prox}_h$  is projection on  $\mathcal{C}$ , that is

$$\text{prox}_h(\mathbf{x}) = \underset{\mathbf{u} \in \mathcal{C}}{\operatorname{argmin}} \|\mathbf{u} - \mathbf{x}\|_2^2 = \pi_{\mathcal{C}}(\mathbf{x}).$$

## Proximal Operator

3 If  $h(\mathbf{x}) = \|\mathbf{x}\|_1$ ,  $\text{prox}_h$  is the "soft-threshold" (shrinkage) operation:

$$\text{prox}_h(\mathbf{x})_i = \begin{cases} \mathbf{x}_i - 1, & \mathbf{x}_i \geq 1, \\ 0, & |\mathbf{x}_i| \leq 1, \\ \mathbf{x}_i + 1, & \mathbf{x}_i \leq -1. \end{cases}$$

Remark. The problem

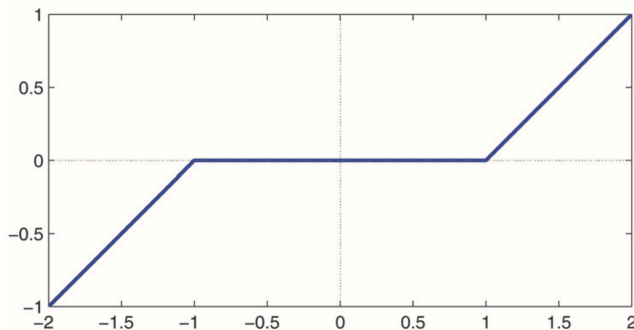
$$\underset{\mathbf{u}}{\text{argmin}} \left\{ \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|_2^2 + \|\mathbf{u}\|_1 \right\}$$

is separable with respect to both  $\mathbf{u}$  and  $\mathbf{x}$  hence one could solve the following problem:

$$\underset{u_i}{\text{argmin}} \left\{ \frac{1}{2} (u_i - x_i)^2 + |u_i| \right\}.$$

# Proximal Operator

- 1 If  $u_i > 0$ , we have  $u_i^* - x_i + 1 = 0$ , then  $u_i^* = x_i - 1$  with  $x_i > 1$ .
- 2 If  $u_i < 0$ , we have  $u_i^* - x_i - 1 = 0$ , then  $u_i^* = x_i + 1$  with  $x_i < -1$ .
- 3 If  $u_i = 0$ , we have  $0 \in 0 - x_i + \partial|0|$ , then  $x_i \in \partial|0| = [-1, 1]$ .



# Proximal Gradient Method

Unconstrained optimization with objective split into two components:

$$\min \left\{ f(\mathbf{x}) \triangleq g(\mathbf{x}) + h(\mathbf{x}) \right\}.$$

where  $g$  is convex and differentiable, and  $\text{dom } g = \mathbb{R}^n$ ;  $h$  is convex with *inexpensive* proximal operator.

Proximal Gradient algorithm:

$$\mathbf{x}_{k+1} = \text{prox}_{t_k h}(\mathbf{x}_k - t_k \nabla g(\mathbf{x}_k)).$$

- 1  $t_k > 0$  is step size, constant or determined by line search.
- 2 can start at infeasible  $\mathbf{x}_0$  (However  $\mathbf{x}_k \in \text{dom } f = \text{dom } h$  for  $k \geq 1$ .)

# Proximal Gradient Method

How to interpret proximal algorithm? From the definition

$$\begin{aligned}\mathbf{x}_{k+1} &= \text{prox}_{th}(\mathbf{x}_k - t\nabla g(\mathbf{x}_k)) \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \left( h(\mathbf{u}) + \frac{1}{2t} \|\mathbf{u} - \mathbf{x}_k + t\nabla g(\mathbf{x}_k)\|_2^2 \right) \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \left( h(\mathbf{u}) + g(\mathbf{x}_k) + \nabla g(\mathbf{x}_k)^\top (\mathbf{u} - \mathbf{x}_k) + \frac{1}{2t} \|\mathbf{u} - \mathbf{x}_k\|_2^2 \right).\end{aligned}$$

$\mathbf{x}_{k+1}$  minimize  $h(\mathbf{u})$  plus a simple quadratic local model of  $g(\mathbf{u})$  around  $\mathbf{x}_k$ .



# Examples

Example 2 (Gradient Method:)

Special case with  $h(\mathbf{x}) = 0$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t \nabla g(\mathbf{x}).$$

If  $h(\mathbf{x}) = 0$ , then  $\text{prox}_h(\mathbf{x}) = \mathbf{x}$ .

# Examples

## Example 3 (Projected Gradient Method)

Special case with  $h(\mathbf{x}) = I_C(\mathbf{x})$

$$\mathbf{x}_{k+1} = \pi_C(\mathbf{x}_k - t\nabla g(\mathbf{x}_k)).$$

# Examples

## Example 4 (Soft-thresholding Method)

Special case with  $h(\mathbf{x}) = \|\mathbf{x}\|_1$

$$\mathbf{x}_{k+1} = \text{prox}_{th}(\mathbf{x} - t\nabla g(\mathbf{x})),$$

where

$$\text{prox}_{th}(\mathbf{u})_i = \begin{cases} \mathbf{u}_i - t, & \mathbf{u}_i \geq t, \\ 0, & |\mathbf{u}_i| \leq t, \\ \mathbf{u}_i + t, & \mathbf{u}_i \leq -t. \end{cases}$$

# Properties

## Proposition 5

If  $h$  is convex and closed (has a closed epigraph), then

$$\text{prox}_h(\mathbf{x}) = \underset{\mathbf{u}}{\operatorname{argmin}} \left( h(\mathbf{u}) + \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|_2^2 \right)$$

exists and is unique for all  $\mathbf{x}$ .

**Proof.**

$h(\mathbf{u}) + \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|_2^2$  is strongly convex.



# Properties

## Proposition 6

$\mathbf{u} = \text{prox}_h(\mathbf{x})$  is equivalent to the following

- 1  $\mathbf{x} - \mathbf{u} \in \partial h(\mathbf{u})$ ,
- 2  $h(\mathbf{z}) \geq h(\mathbf{u}) + (\mathbf{x} - \mathbf{u})^\top (\mathbf{z} - \mathbf{u})$  for all  $\mathbf{z}$ .

**Proof.**  $0 \in \partial\{h(\mathbf{u}) + \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|_2^2\}$ . Thus

$$0 \in \partial h(\mathbf{u}) + \mathbf{u} - \mathbf{x}.$$

Also we have for  $g \in \partial h(\mathbf{u})$ ,

$$h(\mathbf{z}) \geq h(\mathbf{u}) + g^\top (\mathbf{z} - \mathbf{u}).$$



# Properties

## Proposition 7

Proximal mapping of indicator function  $I_{\mathcal{C}}$  is Euclidean projection on  $\mathcal{C}$

$$\mathbf{u}^* = \text{prox}_{I_{\mathcal{C}}}(\mathbf{x}) = \underset{\mathbf{u} \in \mathcal{C}}{\text{argmin}} \|\mathbf{u} - \mathbf{x}\|_2^2 = \pi_{\mathcal{C}}(\mathbf{x}),$$

and

$$(\mathbf{x} - \mathbf{u}^*)^\top (\mathbf{z} - \mathbf{u}^*) \leq 0, \forall \mathbf{z} \in \mathcal{C},$$

for all  $\mathbf{z} \in \mathcal{C}$ .

**Proof.** First of all,  $h(\mathbf{z}) \geq h(\mathbf{u}^*) + (\mathbf{x} - \mathbf{u}^*)^\top (\mathbf{z} - \mathbf{u}^*)$  for all  $\mathbf{z}$ . And from the definition of indicator function, we have  $h(\mathbf{z}) = h(\mathbf{u}^*) = 0$ . □

# Properties

## Proposition 8 (Fixed Point)

Let  $f$  be a convex function, we have that a point  $\mathbf{x}_*$  minimizes  $f(\mathbf{x})$  if and only if  $\mathbf{x}_* = \text{prox}_f(\mathbf{x}_*)$ .

**proof.** First, if  $\mathbf{x}_*$  minimize  $f(\mathbf{x})$ , we have  $f(\mathbf{x}) \geq f(\mathbf{x}_*)$ . Hence,

$$f(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{x}_*\|_2^2 \geq f(\mathbf{x}_*) + \frac{1}{2} \|\mathbf{x}_* - \mathbf{x}_*\|_2^2.$$

This implies that

$$\mathbf{x}_* = \underset{\mathbf{x}}{\text{argmin}} \left\{ f(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{x}_*\|_2^2 \right\} = \text{prox}_f(\mathbf{x}_*).$$

# Properties

## Proof. (Continued.)

To prove the converse, consider if

$$\mathbf{x}_* = \text{prox}_f(\mathbf{x}_*) = \underset{\mathbf{x}}{\text{argmin}} \left\{ f(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{x}_*\|_2^2 \right\}.$$

By the optimality condition, this implies that

$$0 \in \partial f(\mathbf{x}_*) + (\mathbf{x}_* - \mathbf{x}_*) \Rightarrow 0 \in \partial f(\mathbf{x}_*).$$

Therefore,  $\mathbf{x}_*$  minimizes  $f$ .





# Properties

## Proposition 9 (Non-expansive)

$$\|\text{prox}_f(\mathbf{x}) - \text{prox}_f(\mathbf{y})\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2.$$

**Proof.** Let us denote  $\mathbf{u} = \text{prox}_f(\mathbf{x})$  and  $\mathbf{v} = \text{prox}_f(\mathbf{y})$ , then in view of Prop. 6,

$$\mathbf{x} - \mathbf{u} \in \partial f(\mathbf{u}), \quad \mathbf{y} - \mathbf{v} \in \partial f(\mathbf{v}).$$

Combining this with **monotonicity of subdifferential** gives  $\langle \mathbf{x} - \mathbf{u} - \mathbf{y} + \mathbf{v}, \mathbf{u} - \mathbf{v} \rangle \geq 0$ .  
Hence, we have

$$\langle \mathbf{x} - \mathbf{y}, \mathbf{u} - \mathbf{v} \rangle \geq \|\mathbf{u} - \mathbf{v}\|_2^2.$$

By Cauchy-Schwartz inequality, this also leads to  $\|\text{prox}_f(\mathbf{x}) - \text{prox}_f(\mathbf{y})\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2$ .



# Properties

Remark.

The subdifferential of a convex function is a monotone operator:

$$\langle g_u - g_v, u - v \rangle \geq 0$$

for all  $u, v, g_u \in \partial f(u)$ , and  $g_v \in \partial f(v)$ .

The proof is very straightforward. Combining the following two inequalities shows it:

$$f(v) \geq f(u) + g_u^\top (v - u), \text{ and } f(u) \geq f(v) + g_v^\top (u - v).$$

## Problem

For problem

$$\min \left\{ f(\mathbf{x}) \triangleq g(\mathbf{x}) + h(\mathbf{x}) \right\},$$

we assume that

- 1  $h$  is closed and convex (so that  $\text{prox}_{th}$  is well-defined).
- 2  $g$  is differentiable with  $\text{dom } g = \mathbb{R}^n$ , and belong to  $\mathcal{S}_{1,1}^{L,m}(\mathbb{R}^n)$ .
- 3 The optimal value  $f^*$  is finite and attained at  $\mathbf{x}^*$  ( not necessarily unique).

From Lemma 1.2.3 and definition 2.1.2, for all  $\mathbf{x}, \mathbf{y}$ ,

$$g(\mathbf{y}) \geq g(\mathbf{x}) + \nabla g(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{m}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \quad (1)$$

and

$$g(\mathbf{y}) \leq g(\mathbf{x}) + \nabla g(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \quad (2)$$

## Update rule for proximal gradient

The **update rule** for proximal gradient:

$$\begin{aligned}\mathbf{x}' &= \text{prox}_{th}(\mathbf{x} - t\nabla g(\mathbf{x})) \\ &= \mathbf{x} - tG_t(\mathbf{x}),\end{aligned}$$

where  $G_t(\mathbf{x})$  is the **gradient map** in this settings, defined as

$$G_t(\mathbf{x}) = \frac{1}{t} (\mathbf{x} - \text{prox}_{th}(\mathbf{x} - t\nabla g(\mathbf{x}))).$$

(The following properties of  $G_t(\mathbf{x})$  will be left as excercises)

- 1  $G_t(\mathbf{x}) \in \nabla g(\mathbf{x}) + \partial h(\mathbf{x} - tG_t(\mathbf{x})).$
- 2  $G_t(\mathbf{x}_*) = 0$  if and only if  $\mathbf{x}_*$  minimizes  $f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x}).$

## Analysis for Fixed Step Size

For  $0 < t \leq 1/L$ ,

**1** The upper bound (2) implies (substitute  $\mathbf{y} = \mathbf{x} - tG_t(\mathbf{x})$  into bound (2))

$$\begin{aligned} g(\mathbf{x} - tG_t(\mathbf{x})) &\leq g(\mathbf{x}) + \nabla g(\mathbf{x})^\top (-tG_t(\mathbf{x})) + \frac{L}{2} \|tG_t(\mathbf{x})\|_2^2; \\ &= g(\mathbf{x}) - t\nabla g(\mathbf{x})^\top (G_t(\mathbf{x})) + \frac{t \cdot t \cdot L}{2} \|G_t(\mathbf{x})\|_2^2; \\ \Rightarrow \quad g(\mathbf{x} - tG_t(\mathbf{x})) &\leq g(\mathbf{x}) - t\nabla g(\mathbf{x})^\top (G_t(\mathbf{x})) + \frac{t}{2} \|G_t(\mathbf{x})\|_2^2. \end{aligned} \tag{3}$$

The last inequality comes from  $t \cdot L \leq 1$  (since  $0 < t \leq 1/L$ ).

## Analysis for Fixed Step Size

For  $0 < t \leq 1/L$ ,

**2** If the inequality (1) is also satisfied, we have  $\boxed{mt \leq 1}$ , since

$$\frac{mt^2}{2} \|G_t(\mathbf{x})\|_2^2 \leq \frac{t}{2} \|G_t(\mathbf{x})\|_2^2 \text{ and } t \neq 0.$$

**3** We also establish an upper bound for  $f(\mathbf{x} - tG_t(\mathbf{x}))$ , that is

$$f(\mathbf{x} - tG_t(\mathbf{x})) \leq f(\mathbf{z}) + G_t(\mathbf{x})^\top (\mathbf{x} - \mathbf{z}) - \frac{t}{2} \|G_t(\mathbf{x})\|_2^2 - \frac{m}{2} \|\mathbf{x} - \mathbf{z}\|_2^2. \quad (4)$$

(Please refer to the remark below for the establishing process.)

# Analysis for Fixed Step Size

From inequality(3), we have

$$\begin{aligned}
 f(\mathbf{x} - tG_t(\mathbf{x})) &= g(\mathbf{x} - tG_t(\mathbf{x})) + h(\mathbf{x} - tG_t(\mathbf{x})) \\
 &\leq \underbrace{g(\mathbf{x}) - t\nabla g(\mathbf{x})^\top G_t(\mathbf{x}) + \frac{t}{2} \|G_t(\mathbf{x})\|_2^2}_{\text{using bound (3) for } g(\mathbf{x} - tG_t(\mathbf{x}))} + h(\mathbf{x} - tG_t(\mathbf{x})).
 \end{aligned}$$

Using upper bound of  $g$ , we arrive at

$$\begin{aligned}
 f(\mathbf{x} - tG_t(\mathbf{x})) &= g(\mathbf{x} - tG_t(\mathbf{x})) + h(\mathbf{x} - tG_t(\mathbf{x})) \\
 &\leq \underbrace{g(\mathbf{z}) - \nabla g(\mathbf{x})^\top (\mathbf{z} - \mathbf{x}) - \frac{m}{2} \|\mathbf{z} - \mathbf{x}\|_2^2}_{\text{using bound (1) for } g(\mathbf{z})} - t\nabla g(\mathbf{x})^\top G_t(\mathbf{x}) + \frac{t}{2} \|G_t(\mathbf{x})\|_2^2 \\
 &\quad + h(\mathbf{x} - tG_t(\mathbf{x})).
 \end{aligned}$$

## Analysis for Fixed Step Size

Recall that

$$G_t(\mathbf{x}) = \frac{1}{t} (\mathbf{x} - \text{prox}_{th}(\mathbf{x} - t\nabla g(\mathbf{x}))).$$

In view of properties of gradient mapping, we have

$$\mathbf{x} - tG_t(\mathbf{x}) = \text{prox}_{th}(\mathbf{x} - t\nabla g(\mathbf{x})), \quad (5)$$

$$\Rightarrow (\mathbf{x} - t\nabla g(\mathbf{x})) - (\mathbf{x} - tG_t(\mathbf{x})) \in \partial h(\mathbf{x} - tG_t(\mathbf{x})), \quad (6)$$

$$\Rightarrow t(G_t(\mathbf{x}) - \nabla g(\mathbf{x})) \in \partial th(\mathbf{x} - tG_t(\mathbf{x})), \quad (7)$$

$$\Rightarrow G_t(\mathbf{x}) - \nabla g(\mathbf{x}) \in \partial h(\mathbf{x} - tG_t(\mathbf{x})). \quad (8)$$



## Analysis for Fixed Step Size

Thus,

$$\begin{aligned} f(\mathbf{x} - tG_t(\mathbf{x})) &= g(\mathbf{x} - tG_t(\mathbf{x})) + \boxed{h(\mathbf{x} - tG_t(\mathbf{x}))} \\ &\leq g(\mathbf{z}) - \nabla g(\mathbf{x})^\top (\mathbf{z} - \mathbf{x}) - \frac{m}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 \\ &\quad - t \nabla g(\mathbf{x})^\top G_t(\mathbf{x}) + \frac{t}{2} \|G_t(\mathbf{x})\|_2^2 \\ &\quad + \boxed{h(\mathbf{z}) - (G_t(\mathbf{x}) - \nabla g(\mathbf{x}))^\top (\mathbf{z} - \mathbf{x} + tG_t(\mathbf{x}))} \\ &= \underbrace{g(\mathbf{z}) + h(\mathbf{z})}_{f(\mathbf{z})} + G_t(\mathbf{x})^\top (\mathbf{x} - \mathbf{z}) - \frac{t}{2} \|G_t(\mathbf{x})\|_2^2 - \frac{m}{2} \|\mathbf{x} - \mathbf{z}\|_2^2. \end{aligned}$$

## Analysis for Fixed Step Size

$$f(\mathbf{x} - tG_t(\mathbf{x})) \leq f(\mathbf{z}) + G_t(\mathbf{x})^\top (\mathbf{x} - \mathbf{z}) - \frac{t}{2} \|G_t(\mathbf{x})\|_2^2 - \frac{m}{2} \|\mathbf{x} - \mathbf{z}\|_2^2. \quad (4)$$

- inequality (4) with  $\mathbf{z} = \mathbf{x}$  shows that the algorithm is a **descent method**:

$$\left\{ f(\mathbf{x}^+) \triangleq f(\mathbf{x} - tG_t(\mathbf{x})) \right\} \leq f(\mathbf{x}) - \frac{t}{2} \|G_t(\mathbf{x})\|_2^2.$$

# Analysis for Fixed Step Size

$$f(\mathbf{x} - tG_t(\mathbf{x})) \leq f(\mathbf{z}) + G_t(\mathbf{x})^\top (\mathbf{x} - \mathbf{z}) - \frac{t}{2} \|G_t(\mathbf{x})\|_2^2 - \frac{m}{2} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

- inequality (4) with  $\mathbf{z} = \mathbf{x}_*$  shows that

$$\begin{aligned} f(\mathbf{x}^+) - f(\mathbf{x}_*) &\leq G_t(\mathbf{x})^\top (\mathbf{x} - \mathbf{x}_*) - \frac{t}{2} \|G_t(\mathbf{x})\|_2^2 - \frac{m}{2} \|\mathbf{x} - \mathbf{x}_*\|_2^2 \\ &= \frac{1}{2t} \left( \|\mathbf{x} - \mathbf{x}_*\|_2^2 - \|\mathbf{x} - \mathbf{x}_* - tG_t(\mathbf{x})\|_2^2 \right) - \frac{m}{2} \|\mathbf{x} - \mathbf{x}_*\|_2^2 \\ &= \frac{1}{2t} \left( (1 - mt) \|\mathbf{x} - \mathbf{x}_*\|_2^2 - \|\mathbf{x}^+ - \mathbf{x}_*\|_2^2 \right) \end{aligned} \quad (9)$$

$$\leq \frac{1}{2t} \left( \|\mathbf{x} - \mathbf{x}_*\|_2^2 - \|\mathbf{x}^+ - \mathbf{x}_*\|_2^2 \right) \quad (\text{since } mt \leq 1) \quad (10)$$

## Analysis for Fixed Step Size

Add inequalities (10) with  $\mathbf{x} = \mathbf{x}_i$ ,  $\mathbf{x}^+ = \mathbf{x}_{i+1}$ ,  $t = t_i = 1/L$ , from  $i = 0, \dots, k-1$ ,

$$\begin{aligned} \sum_{i=1}^k (f(\mathbf{x}_i) - f(\mathbf{x}_*)) &\leq \frac{1}{2t} \sum_{i=1}^k \left( \|\mathbf{x}_i - \mathbf{x}_*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}_*\|_2^2 \right) \\ &= \frac{1}{2t} \left( \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 - \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 \right) \\ &\leq \frac{1}{2t} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2. \end{aligned}$$

Since  $f(\mathbf{x}_i)$  is nonincreasing,

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \frac{1}{k} \sum_{i=1}^k (f(\mathbf{x}_i) - f(\mathbf{x}_*)) \leq \frac{1}{2kt} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 \quad (11)$$

## Analysis for Fixed Step Size

**Distance to Optimal Set:** from (9) and  $f(\mathbf{x}^+) \geq f(\mathbf{x}_*)$ , we have

$$\begin{aligned} 0 \leq f(\mathbf{x}^+) - f(\mathbf{x}_*) &\leq \frac{1}{2t} \left( (1 - mt) \|\mathbf{x} - \mathbf{x}_*\|_2^2 - \|\mathbf{x}^+ - \mathbf{x}_*\|_2^2 \right) \\ \Rightarrow \|\mathbf{x}^+ - \mathbf{x}_*\|_2^2 &\leq (1 - mt) \|\mathbf{x} - \mathbf{x}_*\|_2^2 \\ &\leq \|\mathbf{x} - \mathbf{x}_*\|_2^2 \quad (\text{since } 0 < mt \leq 1). \end{aligned}$$

For fixed step size  $t_k = 1/L$ ,

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2^2 \leq \left(1 - \frac{m}{L}\right)^k \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2.$$

This implies the linear convergence if  $g$  is strongly convex ( $m > 0$ ).

## Analysis with line search

From inequality (10), if (3) holds in iteration  $i$ , then  $f(\mathbf{x}_{i+1}) < f(\mathbf{x}_i)$  and

$$t_i(f(\mathbf{x}_{i+1}) - f(\mathbf{x}_*)) \leq \frac{1}{2} \left( \|\mathbf{x}_i - \mathbf{x}_*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}_*\|_2^2 \right)$$

Adding inequalities for  $i = 0$  to  $k - 1$  gives

$$\left( \sum_{i=0}^{k-1} t_i \right) (f(\mathbf{x}_k) - f(\mathbf{x}_*)) \leq \sum_{i=0}^{k-1} t_i (f(\mathbf{x}_{i+1}) - f(\mathbf{x}_*)) \leq \frac{1}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2.$$

First inequality holds because  $f(\mathbf{x}_i)$  is nonincreasing. Since  $t_i \geq t_{\min}$ , we obtain a similar  $1/k$  bound as for fixed step size

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \frac{1}{2kt_{\min}} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2.$$

## Analysis with line search

**Distance to Optimal Set:** from inequality (10), if (3) holds in iteration  $i$ , then

$$\begin{aligned} 0 \leq f(\mathbf{x}_{i+1}) - f(\mathbf{x}_*) &\leq \frac{1}{2t} \left( (1 - mt_i) \|\mathbf{x}_i - \mathbf{x}_*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}_*\|_2^2 \right) \\ \Rightarrow \|\mathbf{x}_{i+1} - \mathbf{x}_*\|_2^2 &\leq (1 - mt_i) \|\mathbf{x}_i - \mathbf{x}_*\|_2^2 \\ &\leq (1 - mt_{\min}) \|\mathbf{x}_i - \mathbf{x}_*\|_2^2 \quad (\text{if } 0 < mt_{\min} \leq 1). \end{aligned}$$

Thus, for  $c = 1 - mt_{\min}$ ,

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2^2 \leq c^k \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2.$$

This implies the linear convergence if  $g$  is strongly convex ( $m > 0$ ).

# Analysis with line search

Backtracking line-search for Lipschitz constant:

(1) We initialize  $L_0 = 1$  and some  $\alpha > 1$ .

(2) At each iteration  $i$ , we find the smallest integer  $t$  such that  $L = \alpha^t L_{i-1}$ , specifically:

$$f(\mathbf{x}^+) \leq f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)(\mathbf{x}^+ - \mathbf{x}_i) + \frac{L}{2} \|\mathbf{x}^+ - \mathbf{x}_i\|_2^2,$$

where  $\mathbf{x}^+ = \text{prox}_h(\mathbf{x}_i - \frac{1}{L} \nabla f(\mathbf{x}_i))$ .

(3) Update  $L_i = L$  and  $\mathbf{x}_{i+1} = \mathbf{x}^+$



## Accelerated Proximal Gradient Method

Originally developed by [Nesterov, 2013a] and [Beck and Teboulle, 2009], we can accelerate the proximal gradient method simply as follows

$$\begin{aligned}\mathbf{x}_{i+1} &= \text{prox}_{th}(\mathbf{y}_i - t\nabla g(\mathbf{y}_i)), \\ \mathbf{y}_{i+1} &= \mathbf{x}_{i+1} + \beta_i(\mathbf{x}_{i+1} - \mathbf{x}_i).\end{aligned}$$

Some simple choice for  $\beta$ :

1 NESTEROV07[Nesterov, 2013a]:  $\beta_i = \frac{i}{i+3}$ .

2 FISTA[Beck and Teboulle, 2009]:  $\beta_i = \frac{\lambda_i - 1}{\lambda_{i+1}}$ , where  $\lambda_0 = 0$ ,  $\lambda_{t+1} = \frac{1 + \sqrt{1 + 4\lambda_i^2}}{2}$ .

# Accelerated Proximal Gradient Method

It was shown in [Beck and Teboulle, 2009] that

## Theorem 10

The sequences  $\mathbf{x}_k, f(\mathbf{x}_k)$  generated via FISTA with either a constant or backtracking (with a ratio  $\alpha \geq 1$ ) stepsize rule satisfy

$$f(\mathbf{x}_k) - f^* \leq \frac{2\alpha L}{k^2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2.$$

# Proximal Point Method

Proximal Point Method is an algorithm for minimizing a closed convex function  $f$ :

$$\begin{aligned}\mathbf{x}_{k+1} &= \text{prox}_{tf}(\mathbf{x}_k) \\ &= \underset{\mathbf{u}}{\text{argmin}} \left( f(\mathbf{u}) + \frac{1}{2t} \|\mathbf{u} - \mathbf{x}_k\|_2 \right).\end{aligned}$$

- 1 PPM can be viewed as [proximal gradient method](#) with  $g(\mathbf{x}) = 0$ .
- 2 PPM is basis of the [augmented Lagrangian method](#) (coming sections).
- 3 PPM is related to [Moreau-Yosida smoothing](#) (coming sections).

# Proximal Point Method

## Theorem 11

For PPM, if  $f$  is closed and convex and optimal value  $f^*$  is finite and attained at  $\mathbf{x}_*$ . We have

$$f(\mathbf{x}_k) - f^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{2 \sum_{i=0}^{k-1} t_i}, \quad \text{for } k \geq 1.$$

## Remark.

- (1) Implies convergence if  $\sum_i t_i \rightarrow \infty$
- (2) Convergence rate is  $1/k$  if  $t_i$  is fixed, or variable but bounded, away from zero.
- (3)  $t_i$  is arbitrary; however cost of PPM evaluations will depend on  $t_i$

## Part II

### Splitting Method

## Different Settings for Convex Problem

### 1 Nonsmooth Minimization:

$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ , where  $f(\mathbf{x})$  is convex and nonsmooth. We show the optimal solution has the following fixed point property:

$$\mathbf{x}_* \text{ is optimal.} \Leftrightarrow 0 \in \partial f(\mathbf{x}_*) \Leftrightarrow \forall \lambda > 0, \mathbf{x}_* = \text{prox}_{\lambda f}(\mathbf{x}_*)$$

The fixed point iteration gives rise to the **Proximal Point Algorithm**:

$$\mathbf{x}_{t+1} = \text{prox}_{\lambda_t f}(\mathbf{x}_t).$$

## Different Settings for Convex Problem

### 2 Smooth + Nonsmooth Minimization:

$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \triangleq g(\mathbf{x}) + h(\mathbf{x})$ , where  $g(\mathbf{x})$  is convex and smooth and  $h(\mathbf{x})$  is convex and nonsmooth. We show the optimal solution has the following fixed point property:

$$\mathbf{x}_* \text{ is optimal.} \Leftrightarrow 0 \in \nabla g(\mathbf{x}_*) + \partial h(\mathbf{x}_*) \Leftrightarrow \forall \lambda > 0, \mathbf{x}_* = \text{prox}_{\lambda h}(\mathbf{x}_* - \lambda \nabla g(\mathbf{x}_*))$$

The fixed point iteration gives rise to the **Proximal Gradient Algorithm**:

$$\mathbf{x}_{t+1} = \text{prox}_{\lambda_t h}(\mathbf{x}_t - \lambda_t \nabla g(\mathbf{x}_t)).$$

## Different Settings for Convex Problem

### 3 Nonsmooth + Nonsmooth Minimization:

$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \triangleq g(\mathbf{x}) + h(\mathbf{x})$ , where both  $g(\mathbf{x})$  and  $h(\mathbf{x})$  are convex and nonsmooth.

We show the optimal solution has the following fixed point property:

$$\mathbf{x}_* \text{ is optimal.} \Leftrightarrow 0 \in \partial g(\mathbf{x}_*) + \partial h(\mathbf{x}_*) \Leftrightarrow \forall \lambda > 0, \mathbf{x}_* = \text{prox}_{\lambda(g+h)}(\mathbf{x}_*)$$

The fixed point iteration gives rise to the **Splitting Algorithms**:

$$\mathbf{x}_{t+1} = \text{prox}_{\lambda_t(g+h)}(\mathbf{x}_t).$$

However, this would require the proximal operator of the sum of two convex function, which is not always easy to compute, even if the proximal operators of both functions **separately** may be easy to compute.



## Different Settings for Convex Problem

### Example 12

Let  $g(\mathbf{x}) = \|\mathbf{x}\|_1$ , and  $h(\mathbf{x}) = \|A\mathbf{x}\|_2^2$ . The proximal operator of  $g$  is given by

$$\text{prox}_g(\mathbf{y}) = \underset{\mathbf{x}}{\text{argmin}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \|\mathbf{x}\|_1 \right\}.$$

The proximal operator of  $h$  ([see 6.2.3 Convex Quadratic of \[Beck, 2017\]](#)) is given by

$$\text{prox}_h(\mathbf{y}) = \underset{\mathbf{x}}{\text{argmin}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \|A\mathbf{x}\|_2^2 \right\}.$$

Both are easy to compute. However, the proximal operator of the sum of  $g$  and  $h$  is given by

$$\text{prox}_{g+h}(\mathbf{y}) = \underset{\mathbf{x}}{\text{argmin}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \|\mathbf{x}\|_1 + \|A\mathbf{x}\|_2^2 \right\}.$$

Therefore, the above fixed point property is not really useful.

## Fixed Point for Nonsmooth Composition

### Theorem 13 (Fixed Point)

Consider  $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \triangleq g(\mathbf{x}) + h(\mathbf{x})$ , where both  $g(\mathbf{x})$  and  $h(\mathbf{x})$  are convex and nonsmooth. If  $\mathbf{x}_*$  is optimal if and only if for any  $\lambda > 0$  and  $\rho \in \mathbb{R}$ ,

$$\mathbf{x}_* = \text{prox}_{\lambda g}(\mathbf{y}_*) \quad \text{and} \quad \mathbf{y}_* = \mathbf{y}_* + \rho [\text{prox}_{\lambda h}(\boxed{2\text{prox}_{\lambda g}(\mathbf{y}_*) - \mathbf{y}_*}) - \text{prox}_{\lambda h}(\mathbf{y}_*)]$$

**Proof.** Suppose that  $\mathbf{x}_*$  is optimal.

- 1  $\Leftrightarrow 0 \in \partial g(\mathbf{x}_*) + \partial h(\mathbf{x}_*)$ .
- 2  $\Leftrightarrow \forall \lambda > 0$ , there exists  $\mathbf{z}$  such that  $\mathbf{z} \in \partial(\lambda g)(\mathbf{x}_*)$  and  $-\mathbf{z} \in \partial(\lambda h)(\mathbf{x}_*)$ .
- 3  $\Leftrightarrow \forall \lambda > 0$ , there exists  $\mathbf{y}$  such that  $\mathbf{y} - \mathbf{x}_* \in \partial(\lambda g)(\mathbf{x}_*)$  and  $\mathbf{x}_* - \mathbf{y} \in \partial(\lambda h)(\mathbf{x}_*)$ .

## Fixed Point for Nonsmooth Composition

4  $\Leftrightarrow \forall \lambda > 0$ , there exists  $y$  such that  $y - x_* \in \partial(\lambda g)(x_*)$  and  $2x_* - y \in x_* + \partial(\lambda h)(x_*)$ .

5  $\Leftrightarrow x_* = \text{prox}_{\lambda g}(y)$ , thus we have  $\boxed{2\text{prox}_{\lambda g}(y) - y} \in x_* + \partial(\lambda h)(x_*)$ .

6  $\Leftrightarrow x_* = \text{prox}_{\lambda g}(y)$ , and  $x_* = \text{prox}_{\lambda h}(\boxed{2\text{prox}_{\lambda g}(y) - y})$ .

7  $\Leftrightarrow x_* = \text{prox}_{\lambda g}(y)$ , and  $\forall \rho$ ,

$$y = y + \rho \left[ \underbrace{\text{prox}_{\lambda h}(\boxed{2\text{prox}_{\lambda g}(y) - y})}_{x_*} - \underbrace{\text{prox}_{\lambda g}(y)}_{x_*} \right].$$

All the statements are equivalent. □

## Fixed Point for Nonsmooth Composition

Define  $\text{refl}_f(x) \triangleq 2\text{prox}_f(x) - x$ .

1 When  $\rho = 1$ , we have

$$\begin{aligned} \mathbf{y}_* &= \mathbf{y}_* + [\text{prox}_{\lambda h}(2\text{prox}_{\lambda g}(\mathbf{y}_*) - \mathbf{y}_*) - \text{prox}_{\lambda g}(\mathbf{y}_*)] \\ &= \frac{1}{2}[\mathbf{y}_* + 2\text{prox}_{\lambda h}(\boxed{2\text{prox}_{\lambda g}(\mathbf{y}_*) - \mathbf{y}_*}) - (\boxed{2\text{prox}_{\lambda g}(\mathbf{y}_*) - \mathbf{y}_*})] \\ &= \frac{1}{2}[\mathbf{y}_* + \text{refl}_{\lambda h} \circ \text{refl}_{\lambda g}(\mathbf{y}_*)]. \end{aligned}$$

Hence,  $\mathbf{y}_* = \mathcal{T}_{\lambda,g,h}(\mathbf{y}_*)$  with operator

$$\mathcal{T}_{\lambda,g,h} = \frac{1}{2}[I + \text{refl}_{\lambda h} \circ \text{refl}_{\lambda g}].$$

This is known as the **Douglas-Rachford operator** [Lions and Mercier, 1979].

## Fixed Point for Nonsmooth Composition

2 When  $\rho = 2$ , we have

$$\begin{aligned} \mathbf{y}_* &= \mathbf{y}_* + 2 [\text{prox}_{\lambda h}(2\text{prox}_{\lambda g}(\mathbf{y}_*) - \mathbf{y}_*) - \text{prox}_{\lambda g}(\mathbf{y}_*)] \\ &= 2\text{prox}_{\lambda h}(\boxed{2\text{prox}_{\lambda g}(\mathbf{y}_*) - \mathbf{y}_*}) - (\boxed{2\text{prox}_{\lambda g}(\mathbf{y}_*) - \mathbf{y}_*}) \\ &= \text{refl}_{\lambda h} \circ \text{refl}_{\lambda g}(\mathbf{y}_*). \end{aligned}$$

Hence,  $\mathbf{y}_* = \mathcal{T}_{\lambda,g,h}(\mathbf{y}_*)$  with operator

$$\mathcal{T}_{\lambda,g,h} = \text{refl}_{\lambda h} \circ \text{refl}_{\lambda g}.$$

This is known as the **Peaceman-Rachford operator**[Lions and Mercier, 1979].

# Non-expansive

## Lemma 14

The reflection operator  $\text{refl}_{\lambda f}(\cdot)$  is non-expansive for any  $\lambda > 0$ .

**Proof.** This is because

$$\begin{aligned}\|\text{refl}_{\lambda f}(\mathbf{x}) - \text{refl}_{\lambda f}(\mathbf{y})\|_2^2 &= \|2\text{prox}_{\lambda f}(\mathbf{x}) - 2\text{prox}_{\lambda f}(\mathbf{y}) - (\mathbf{x} - \mathbf{y})\|_2^2 \\ &= 4 \|\text{prox}_{\lambda f}(\mathbf{x}) - \text{prox}_{\lambda f}(\mathbf{y})\|_2^2 \\ &\quad - 4\langle \text{prox}_{\lambda f}(\mathbf{x}) - \text{prox}_{\lambda f}(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \|\mathbf{x} - \mathbf{y}\|_2^2 \\ &\leq 4 \|\text{prox}_{\lambda f}(\mathbf{x}) - \text{prox}_{\lambda f}(\mathbf{y})\|_2^2 - 4 \|\text{prox}_{\lambda f}(\mathbf{x}) - \text{prox}_{\lambda f}(\mathbf{y})\|_2^2 \\ &\quad + \|\mathbf{x} - \mathbf{y}\|_2^2 \\ &= \|\mathbf{x} - \mathbf{y}\|_2^2.\end{aligned}$$

# Non-expansive

The proximal operator is firmly nonexpansiv, i.e.,

$$\|\text{prox}_{\lambda f}(\boldsymbol{x}) - \text{prox}_{\lambda f}(\boldsymbol{y})\|_2^2 \leq \langle \text{prox}_{\lambda f}(\boldsymbol{x}) - \text{prox}_{\lambda f}(\boldsymbol{y}), \boldsymbol{x} - \boldsymbol{y} \rangle.$$

# Non-expansive

## Lemma 15

Both the Peaceman-Rachford operator  $\mathcal{T}_1$  and the Douglas-Rachford operator  $\mathcal{T}_2$  are non-expansive.

**Proof.** For  $\mathcal{T}_1$ ,

$$\begin{aligned}\|\mathcal{T}_1 \mathbf{x} - \mathcal{T}_1 \mathbf{y}\|_2^2 &= \|\text{refl}_{\lambda h} \circ \text{refl}_{\lambda g}(\mathbf{x}) - \text{refl}_{\lambda h} \circ \text{refl}_{\lambda g}(\mathbf{y})\|_2^2 \\ &\leq \|\text{refl}_{\lambda g}(\mathbf{x}) - \text{refl}_{\lambda g}(\mathbf{y})\|_2^2 \\ &\leq \|\mathbf{x} - \mathbf{y}\|_2^2.\end{aligned}$$

For  $\mathcal{T}_2$ ,

$$\|\mathcal{T}_2 \mathbf{x} - \mathcal{T}_2 \mathbf{y}\|_2^2 \leq \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \frac{1}{2} \|\mathcal{T}_1 \mathbf{x} - \mathcal{T}_1 \mathbf{y}\|_2^2 \leq \|\mathbf{x} - \mathbf{y}\|_2^2.$$





# Splitting Algorithm

The fixed point iteration corresponding to these non-expansive operators lead to the following algorithm:

## 1 Douglas-Rachford Splitting Algorithm:

$$\mathbf{y}_{t+1} = \frac{1}{2}[\mathbf{y}_t + \text{refl}_{\lambda h} \circ \text{refl}_{\lambda g}(\mathbf{y}_t)].$$

## 2 Peaceman-Rachford Splitting Algorithm:

$$\mathbf{y}_{t+1} = \text{refl}_{\lambda h} \circ \text{refl}_{\lambda g}(\mathbf{y}_t).$$

## 3 relaxed Peaceman-Rachford Splitting Algorithm: let $\gamma_t \in [0, 1]$ ,

$$\mathbf{y}_{t+1} = (1 - \gamma_t)\mathbf{y}_t + \gamma_t \text{refl}_{\lambda h} \circ \text{refl}_{\lambda g}(\mathbf{y}_t).$$

When  $\gamma_t = 1$ , Relaxed PR is PR. And if  $\gamma_t = 1/2$ , Relaxed PR is DR.

# Splitting Algorithm

Let us initialize  $\mathbf{y}_1$  and  $\mathbf{x}_1 = \text{prox}_{\lambda g}(\mathbf{y}_1)$ , the DR algorithm can be rewritten as

$$\begin{aligned}\mathbf{y}_{t+1} &= \mathbf{y}_t + \text{prox}_{\lambda h}(2\mathbf{x}_t - \mathbf{y}_t) - \mathbf{x}_t, \\ \mathbf{x}_{t+1} &= \text{prox}_{\lambda g}(\mathbf{y}_{t+1}).\end{aligned}$$

Let  $\mathbf{z}_{t+1} = \text{prox}_{\lambda h}(2\mathbf{x}_t - \mathbf{y}_t)$  and switch  $\mathbf{x}$  and  $\mathbf{y}$  updates, this can be further formulated as

$$\begin{aligned}\mathbf{z}_{t+1} &= \text{prox}_{\lambda h}(2\mathbf{x}_t - \mathbf{y}_t), \\ \mathbf{x}_{t+1} &= \text{prox}_{\lambda g}(\mathbf{y}_t + \mathbf{z}_{t+1} - \mathbf{x}_t), \\ \mathbf{y}_{t+1} &= \mathbf{y}_t + \mathbf{z}_{t+1} - \mathbf{x}_t.\end{aligned}$$

# Splitting Algorithm

Let  $\mathbf{u}_t = \mathbf{x}_t - \mathbf{y}_t$ . We have

$$\mathbf{z}_{t+1} = \text{prox}_{\lambda h}(\mathbf{x}_t + \mathbf{u}_t),$$

$$\mathbf{x}_{t+1} = \text{prox}_{\lambda g}(\mathbf{z}_{t+1} - \mathbf{u}_t),$$

$$\mathbf{u}_{t+1} = \mathbf{u}_t + (\mathbf{x}_{t+1} - \mathbf{z}_{t+1}),$$

The above is a special case of the Alternating Direction Methods of Multipliers (ADMM).

## Related to ADMM

We consider the following optimization problem

$$\begin{aligned} & \min g(\mathbf{x}) + h(\mathbf{z}) \\ \text{s.t. } & A\mathbf{x} + B\mathbf{z} = c. \end{aligned}$$

The iteration of ADMM is

$$\begin{aligned} \mathbf{z}_{t+1} &= \operatorname{argmin}_z \{h(\mathbf{z}) + \frac{\rho}{2} \|A\mathbf{x}_t + B\mathbf{z}_t - c + \mathbf{u}_t\|_2^2\}, \\ \mathbf{x}_{t+1} &= \operatorname{argmin}_x \{h(\mathbf{z}) + \frac{\rho}{2} \|A\mathbf{x}_t + B\mathbf{z}_{t+1} - c + \mathbf{u}_t\|_2^2\}, \\ \mathbf{u}_{t+1} &= \mathbf{u}_t + (A\mathbf{x}_{t+1} + B\mathbf{z}_{t+1} - c). \end{aligned}$$

Let  $A = I, B = -I, c = 0$ . One can see that the above problem is exactly the same as the Douglas-Rachford splitting algorithm in this case.

## Convergence Analysis

Let  $\mathcal{T} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a nonexpansive operator. We consider the relaxed fixed point algorithm

$$\mathbf{x}_{t+1} = (1 - \gamma_t)\mathbf{x}_t + \gamma_t \cdot \mathcal{T}\mathbf{x}_t, \quad \text{for all } t \geq 0$$

where  $\gamma_t \in (0, 1]$ . This is known as the **Krasnosel'skii-Mann (KM)** algorithm. Note that when  $\gamma_t = 1$ , this reduce to the usual fixed point algorithm.

### Theorem 16

Let  $\mathcal{T}$  be a nonexpansive operator and not a self-map. The KM algorithm satisfies that

$$\|\mathcal{T}\mathbf{x}_t - \mathbf{x}_t\|_2^2 \leq \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{\sum_{\tau=0}^t \gamma_\tau (1 - \gamma_\tau)}.$$

## Convergence Analysis

**Proof.** First,  $\|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2$  is non-increasing, since

$$\begin{aligned}\|\mathbf{x}_{t+1} - \mathcal{T}\mathbf{x}_{t+1}\|_2^2 &= \|(1 - \gamma_t)\mathbf{x}_t + \gamma_t\mathcal{T}\mathbf{x}_t - \mathcal{T}\mathbf{x}_{t+1}\|_2^2 \\ &= \|(1 - \gamma_t)(\mathbf{x}_t - \mathcal{T}\mathbf{x}_t) + \mathcal{T}\mathbf{x}_t - \mathcal{T}\mathbf{x}_{t+1}\|_2^2 \\ &\leq (1 - \gamma_t)\|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2 + \|\mathcal{T}\mathbf{x}_t - \mathcal{T}\mathbf{x}_{t+1}\|_2^2 \quad (\text{Triangular inequality}) \\ &\leq (1 - \gamma_t)\|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2 + \|\mathbf{x}_t - \mathbf{x}_{t+1}\|_2^2 \quad (\mathcal{T} \text{ is a nonexpansive operator}) \\ &= (1 - \gamma_t)\|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2 + \gamma_t\|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2 \\ &= \|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2.\end{aligned}$$

To be continued ...

## Convergence Analysis

**Proof. (continued)** We now show that

$$\begin{aligned}\|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2 &= \|(1 - \gamma_t)\mathbf{x}_t + \gamma_t\mathcal{T}\mathbf{x}_t - (1 - \gamma_t)\mathbf{x}_* - \gamma_t\mathcal{T}\mathbf{x}_*\|_2^2 \\ &= \|(1 - \gamma_t)(\mathbf{x}_t - \mathbf{x}_*) + \gamma_t(\mathcal{T}\mathbf{x}_t - \mathcal{T}\mathbf{x}_*)\|_2^2 \\ &= (1 - \gamma_t)\|\mathbf{x}_t - \mathbf{x}_*\|_2^2 + \gamma_t\|\mathcal{T}\mathbf{x}_t - \mathcal{T}\mathbf{x}_*\|_2^2 - \gamma_t(1 - \gamma_t)\|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2 \\ &= \|\mathbf{x}_t - \mathbf{x}_*\|_2^2 - \underbrace{\gamma_t\left(\|\mathbf{x}_t - \mathbf{x}_*\|_2^2 - \|\mathcal{T}\mathbf{x}_t - \mathcal{T}\mathbf{x}_*\|_2^2\right)}_{\geq 0 \text{ since } \mathcal{T} \text{ is non-expansive.}} - \gamma_t(1 - \gamma_t)\|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2 \\ &\leq \|\mathbf{x}_t - \mathbf{x}_*\|_2^2 - \gamma_t(1 - \gamma_t)\|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2.\end{aligned}$$

The third equality is due to the fact that

$$\|(1 - \gamma_t)\mathbf{u} + \gamma_t\mathbf{v}\|_2^2 = (1 - \gamma_t)\|\mathbf{u}\|_2^2 + \gamma_t\|\mathbf{v}\|_2^2 - \gamma_t(1 - \gamma_t)\|\mathbf{u} - \mathbf{v}\|_2^2.$$

# Convergence Analysis

**Proof. (continued)** We now have that

$$\|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2 \leq \|\mathbf{x}_t - \mathbf{x}_*\|_2^2 - \gamma_t(1 - \gamma_t) \|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2.$$

Taking summation over  $t$  leads to

$$\begin{aligned} \sum_{\tau=0}^t \gamma_{\tau}(1 - \gamma_{\tau}) \|\mathbf{x}_{\tau} - \mathcal{T}\mathbf{x}_{\tau}\|_2^2 &\leq \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2 \\ &\leq \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2. \end{aligned}$$

Thus, we obtain

$$\left( \sum_{\tau=0}^t \gamma_{\tau}(1 - \gamma_{\tau}) \right) \|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2 \leq \sum_{\tau=0}^t \gamma_{\tau}(1 - \gamma_{\tau}) \|\mathbf{x}_{\tau} - \mathcal{T}\mathbf{x}_{\tau}\|_2^2 \leq \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2.$$



# Convergence Analysis

If we set  $\gamma_t = \gamma \in (0, 1)$ , we have

$$\|\mathbf{x}_t - \mathcal{T}\mathbf{x}_t\|_2^2 \leq \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{\sum_{\tau=0}^t \gamma_\tau (1 - \gamma_\tau)} = \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{\gamma(1 - \gamma)(t + 1)}.$$

## Appendix: Contraction Mapping

Let  $(\mathcal{X}, D)$  be a metric space, and let  $\mathcal{T} : \mathcal{X} \rightarrow \mathcal{X}$  be an operator.  $\mathcal{T}$  is considered a **contraction mapping** if there exists a constant  $0 \leq k < 1$  such that, for any two points  $\mathbf{x}$  and  $\mathbf{y}$  in  $\mathcal{X}$ , the distance between  $\mathcal{T}(\mathbf{x})$  and  $\mathcal{T}(\mathbf{y})$  is less than or equal to the contraction factor times the distance between  $\mathbf{x}$  and  $\mathbf{y}$ :

$$D(\mathcal{T}(\mathbf{x}), \mathcal{T}(\mathbf{y})) \leq k \cdot D(\mathbf{x}, \mathbf{y}).$$

let  $(\mathcal{X}, D)$  be a metric space, and let  $\mathcal{T} : \mathcal{X} \rightarrow \mathcal{X}$  be an operator.  $\mathcal{T}$  is considered **non-expansive** if, for any two points  $\mathbf{x}$  and  $\mathbf{y}$  in  $\mathcal{X}$ , the distance between  $\mathcal{T}(\mathbf{x})$  and  $\mathcal{T}(\mathbf{y})$  is less than or equal to the distance between  $\mathbf{x}$  and  $\mathbf{y}$ :

$$D(\mathcal{T}(\mathbf{x}), \mathcal{T}(\mathbf{y})) \leq D(\mathbf{x}, \mathbf{y}).$$

## Appendix: Fixed-Point Theorem

The Banach Fixed-Point Theorem states that if a mapping  $\mathcal{T}$  is a **contraction mapping** on a complete metric space, then it has a **unique** fixed point and any sequence generated by iteratively applying the mapping will converge to that fixed point.

In the case of a **non-expansive** operator, which is a **weaker** condition than contraction, the guarantee for convergence depends on the specific properties of the operator and the space it operates on. If the non-expansive operator is defined on a **compact set** or operates in a **finite-dimensional space**, convergence is guaranteed. This is because compact sets and finite-dimensional spaces have certain properties that ensure convergence for non-expansive mappings.

## Appendix: Fixed-Point Theorem

However, in an infinite-dimensional space, the convergence of a non-expansive operator is **NOT** always guaranteed. Additional conditions, such as convexity or strong monotonicity, may be required to ensure convergence in such cases.

Therefore, while a non-expansive operator with a fixed point indicates the potential for convergence, the actual convergence depends on the specific properties of the operator and the space in which it operates.

## References I

- Yu Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1):125–161, 2013a.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- Amir Beck. *First-order methods in optimization*, volume 25. SIAM, 2017.
- Pierre-Louis Lions and Bertrand Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM Journal on Numerical Analysis*, 16(6):964–979, 1979.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Yurii Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018.
- Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013b.

## References II

L. Vandenberghe. Proximal gradient method, course ece236c. URL <http://www.seas.ucla.edu/~vandenbe/236C/lectures/proxgrad.pdf>.

Niao He. Big data optimizatou course ece236c. URL [http://niaohe.ise.illinois.edu/IE598\\_2016/index.html](http://niaohe.ise.illinois.edu/IE598_2016/index.html).

# Thank You!

Email: [qianhui@zju.edu.cn](mailto:qianhui@zju.edu.cn)