

# Modern Optimization

Jean-Luc Bouchot

School of Mathematics and Statistics  
Beijing Institute of Technology  
jlouchot@bit.edu.cn

Spring 2021

# 1 Proximal methods

# Outline

## 1 Proximal methods

### Definition 2.1 (Composite model)

Let  $f(x) = g(x) + h(x)$  where

- $g$  is nice (i.e. for which the analysis from the previous sections carry over)
- $h$  is simple – which we will describe later on

This is called a **composite** model.

### Example 2.1

Assume we are trying to solve the following constrained optimization problem

$$\begin{aligned} \min f_0(x) \\ \text{s.t. } x \in \Omega \end{aligned}$$

where  $\Omega$  is a convex body.

This can be rewritten in the form of a composite function with

- $g = f_0$
- $h = \chi_\Omega$  (which is 0 for points in  $\Omega$  and  $\infty$  elsewhere)

### Example 2.2

Assume we are trying to solve the following constrained optimization problem

$$\begin{aligned} \min f_0(x) \\ \text{s.t. } Ax = 0 \end{aligned}$$

where  $A \in \mathbb{R}^{m \times n}$ .

This can be approximated via a composite function with

- $g = f_0$
- $h = \|Ax\|$

### Remark 2.1

Note that if both functions  $g$  and  $h$  are differentiable, we're good to go!  
The interesting part is if  $h$  is not differentiable (e.g. indicator function)

## Remark 2.2

At each iterations, we will (try to) solve:

$$x^{k+1} := \operatorname{argmin} \left\{ \frac{1}{2\gamma} \|y - (x^k - \gamma \nabla g(x^k))\|^2 + h(y) \right\}$$



### Definition 2.2

Let  $f$  be a function and  $\gamma > 0$  a given parameter. We define the **proximal operator** as

$$\text{prox}_{f,\gamma}(x) := \operatorname{argmin}\left\{f(y) + \frac{1}{2\gamma}\|y - x\|^2\right\}.$$

### Example 2.3

Let  $C$  be a nonempty closed convex body and define

$$\chi_C(x) := \begin{cases} 0 & \text{if } x \in C, \\ \infty & \text{elsewhere.} \end{cases}$$

Its proximal operator is precisely the projection.

### Definition 2.3

We define the **proximal gradient descent** as the sequence of iterates

$$x^{k+1} := \text{prox}_{h,\gamma}(x^k - \gamma \nabla g(x^k))$$

with a certain starting point  $x^0$ .

### Proposition 2.1 (Admitted)

*Under some very general assumptions (e.g.  $f$  is proper closed and convex or  $f$  is proper closed and coercive) the proximal operator admits a unique valued and is defined.*

### Example 2.4

Let  $A \in \mathbb{R}^{d \times d}$  be symmetric positive definite,  $b \in \mathbb{R}^d$  be a constant vector and  $c \in \mathbb{R}$  a scalar and define  $f(x) = \frac{1}{2}x^T Ax + b^T x + c$ . Then, for  $\gamma > 0$ ,

$$\text{prox}_{f,\gamma}(x) = \left( A + \frac{1}{\gamma} I \right)^{-1} \left( \frac{1}{\gamma} x - b \right).$$

### Remark 2.3

The proximal gradient descent algorithm is a generalization of both the gradient descent and the projected gradient descent.

### Definition 2.4

Let  $f = g + h$  with  $g$  convex differentiable (and smooth) and  $h$  simple. We define its **generalized gradient** as the operator

$$G_{h,\gamma}(x) := \frac{1}{\gamma} (x - \text{prox}_{h,\gamma}(x - \gamma \nabla g(x))).$$

## Proposition 2.2

*The proximal gradient descent can also be written as a generalized gradient descent*

$$x^{k+1} = x^k - \gamma G_{h,\gamma}(x^k).$$