

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 Ω 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 Ω and ∞ 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

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

Theorem 2.1

Let $f = g + h$ be a composite function such that g is convex (proper closed) and L -smooth and h is convex (and proper closed). Let $\{x^k\}$ be the sequence of iterates generated by the proximal gradient descent algorithm with stepsize $\gamma = 1/L$ and starting at $x^0 \in \mathbb{R}^d$. Assume moreover that the function f admits a minimum point x^ . Then for any $K \geq 1$ it holds*

$$f(x^K) - f(x^*) \leq \frac{L}{2K} \|x^0 - x^*\|^2$$

Lemma 1

Let f be a (proper closed) convex function and $\gamma > 0$. For any x in the domain

$$u = \text{prox}_{f,\gamma}(x) \Rightarrow \frac{1}{\gamma} \langle x - u, y - u \rangle \leq f(y) - f(u), \quad \forall y.$$

Theorem 2.2 (Prox of separable functions)

Suppose that $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is **separable**:

$$f(x) = f(x_1, \dots, x_d) = \sum_{i=1}^d f_i(x_i),$$

where all the f_i 's a proper closed and convex univariate functions. Then

$$\text{prox}_{f,\gamma}(x) = (\text{prox}_{f_i,\gamma}(x_i))_{i=1}^d.$$

Remark 2.4

This previous result is a consequence of the more general result: if

$$f : \begin{cases} \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \times \cdots \times \mathbb{R}^{d_m} & \rightarrow (-\infty, \infty] \\ (x_1, \cdots, x_m) & \mapsto \sum_{i=1}^m f_i(x_i) \end{cases}$$

then

$$\text{prox}_{f,\gamma}(x) = \text{prox}_{f_1,\gamma}(x_1) \times \text{prox}_{f_2,\gamma}(x_2) \times \cdots \times \text{prox}_{f_m,\gamma}(x_m)$$

Example 2.5

Let $f(x) = \|x\|_1$. For any $\gamma > 0$, its proximal operator is given by

$$\text{prox}_{f,\gamma}(x) = S_\gamma(x)$$

where S_γ denotes the soft-thresholding operator applied component-wise and defined as

$$S_\gamma(x)_i = \max\{|x_i| - \gamma, 0\} \text{sign}(x_i), \quad \forall 1 \leq i \leq d.$$

Theorem 2.3 (Prox with scaling and translation)

Let g be a proper function. For any scaling parameter $\lambda \neq 0$ and translation $a \in \mathbb{R}^d$, define $f(x) = g(\lambda x + a)$. It follows that

$$\text{prox}_{f,\gamma}(x) = \frac{1}{\lambda} \left(\text{prox}_{\lambda^2 g, \gamma}(\lambda x + a) - a \right).$$

Theorem 2.4 (Proved as part of the convergence of the proximal gradient descent)

Let f be a (proper closed) convex function. Then for any $x, y \in \mathbb{R}^d$ the following statements are equivalent

- ① $y = \text{prox}_{f,\gamma}(x)$
- ② $x - y \in \partial f(y)$ (the subgradient – see exercises)
- ③ $\langle x - y, z - y \rangle \leq f(z) - f(y)$ for all $z \in \mathbb{R}^d$.

Proposition 2.3

Let f be a (proper closed) convex function. Then x is a minimizer of f if and only if it is a fixed point of its proximal operator:

$$x = \operatorname{prox}_{f,1}(x).$$