

Proximal Gradient Descent: convergence [2025 update]

Friday, March 7, 2025

10:44 AM

$$\min_x F(x) := g(x) + h(x)$$

smooth easy proximity operator

Includes gradient descent
as special case (but doesn't
analyze strongly convex case)

Assume ∇g is Lipschitz continuous

WLOG let Lipschitz constant be $L=1$ for simplicity

(i.e., redefine $\tilde{F}(x) = \frac{1}{L} \cdot F(x)$)

Assume $g, h \in \Gamma_0(\mathbb{R}^n)$

Algorithm: $x_{k+1} = \text{prox}_h(\underbrace{x_k - \nabla g(x_k)}_{\tilde{x}})$... or if $L \neq 1$, use stepsize $t = \frac{1}{L}$
 $x_{k+1} = \text{prox}_{th}(x_k - t \nabla g(x_k))$

Analysis:

introduce the gradient map $G(x) = x - \text{prox}_h(x - \nabla g(x))$

ex: $h(x)=0 \Rightarrow \text{prox}_h(y)=y$ so $G(x) = \nabla g(x)$

thus the algo. can be written as

$$x_{k+1} = x_k - G(x_k) \quad \dots \text{ looks like gradient descent.}$$

Property of prox

$$\text{let } y = \text{prox}_h(\tilde{x}) = \arg\min_y \frac{1}{2} \|y - \tilde{x}\|^2 + h(y)$$

i.e. if $y = x_{k+1}$,

i.e. (Fermat's rule) $0 \in y - \tilde{x} + \partial h(y)$

$$y = \text{prox}_h(\tilde{x} := x_k - \nabla g(x_k))$$

$$\text{so } \boxed{\tilde{x} - y \in \partial h(y) \text{ if } y = \text{prox}_h(\tilde{x})} \quad (*)$$

Key inequality (via descent lemma)

Since g is 1-Lipschitz, the descent lemma says

$$g(y) \leq g(x) + \langle \nabla g(x), y - x \rangle + \frac{1}{2} \|y - x\|^2$$

hence

$$F(y) = g(y) + h(y) \leq g(x) + h(y) + \langle \nabla g(x), y - x \rangle + \frac{1}{2} \|y - x\|^2$$

So, thinking of x as x_k , and $y = x - G(x)$, this means
 $= \text{prox}_h(\tilde{x})$

Proximal Gradient Descent: convergence (p. 2)

Friday, March 7, 2025 10:44 AM

use $\forall z, g(z) \geq g(x) + \langle \nabla g(x), z - x \rangle$

$$F(y) \leq g(x) + \underbrace{h(x - G(x))}_{\substack{\text{via convexity and} \\ \text{definition of subgradients}}} + \langle \nabla g(x), -G(x) \rangle + \frac{1}{2} \|G(x)\|^2$$

$g(z) + \langle \nabla g(x), x - z \rangle$ $h(z) + \langle v, y - z \rangle$ where $v \in \partial h(y)$

combine to form $f(z)$

Well, note that $v = G(x) - \nabla g(x) \in \partial h(y)$

Since $G(x) = x - \text{prox}_h(x - \nabla g(x))$, i.e.,

$y = x - G(x) = \text{prox}_h(x - \nabla g(x))$ so via (*)

$$(x - \nabla g(x)) - (x - G(x)) \in \partial h(x - G(x)) = \partial h(y)$$

So...

$$\begin{aligned} F(y) &\leq \underbrace{g(z) + \langle \nabla g(x), x - z \rangle}_{\substack{\text{cancel} \\ \text{cancel}}} + \underbrace{h(z) + \langle G(x) - \nabla g(x), y - z \rangle}_{\substack{\text{cancel} \\ \text{cancel}}} - \langle \nabla g(x), G(x) \rangle + \frac{1}{2} \|G(x)\|^2 \\ &= F(z) + \langle \nabla g(x), x - z \rangle + \langle G(x) - \nabla g(x), x - z \rangle - \langle G(x) - \nabla g(x), G(x) \rangle - \langle \nabla g(x), G(x) \rangle + \frac{1}{2} \|G(x)\|^2 \end{aligned}$$

$$= F(z) + \langle G(x), x - z \rangle - \langle G(x), G(x) \rangle + \frac{1}{2} \|G(x)\|^2$$

$$= F(z) + \langle G(x), x - z \rangle - \frac{1}{2} \|G(x)\|^2$$

$\therefore \forall z$

$$F(y) \leq F(z) + \langle G(x), x - z \rangle - \frac{1}{2} \|G(x)\|^2 \quad \star$$

Proximal Gradient Descent: convergence (p. 3)

Friday, March 7, 2025 10:44 AM

Recall $F = g + h$

$$\forall z \quad F(y) \leq F(z) + \langle G(x), x - z \rangle - \frac{1}{2} \|G(x)\|^2 \quad \star$$

recall
 $y = x - G(x)$
 $x = x_k$
 x_{k+1}

if $z = x \Rightarrow F(y) \leq F(x) - \frac{1}{2} \|G(x)\|^2$, i.e., a descent method.

if $z = x^* \in \arg\min F(x)$

$$\begin{aligned} \Rightarrow F(y) - F^* &\leq \langle G(x), x - x^* \rangle - \frac{1}{2} \|G(x)\|^2 \quad \text{Complete-the-square} \\ &= \frac{1}{2} \left(\|x - x^*\|^2 - \|x - x^* - G(x)\|^2 \right) \\ &= \frac{1}{2} \left(\|x - x^*\|^2 - \|y - x^*\|^2 \right) \end{aligned}$$

For $y = x_{k+1}$, $x = x_k$, now sum for $k = 1, \dots, K$

$$\sum_{k=1}^K F(x_k) - F^* \leq \frac{1}{2} \sum_{k=1}^K \left(\|x_{k-1} - x^*\|^2 - \|x_k - x^*\|^2 \right)$$

telescopes!

since a descent method,

$$F(x_k) - F^* \leq \frac{1}{K} \sum_{k=1}^K F(x_k) - F^*$$

$$\leq \frac{1}{2} \left(\|x_0 - x^*\|^2 - \|x_K - x^*\|^2 \right)$$

unknown, but ≥ 0

$$\leq \frac{1}{2} \|x_0 - x^*\|^2$$

$$\Rightarrow F(x_K) - F^* \leq \frac{1}{2K} \|x_0 - x^*\|^2$$

or if we "undo" our scaling by $1/L$,

$$F(x_k) - F^* \leq \frac{L}{2K} \|x_0 - x^*\|^2$$

i.e. $O(1/k)$
convergence rate

