

CSED490Y: Optimization for Machine Learning

Week 07-2: Proximal gradient descent

Namhoon Lee

POSTECH

Spring 2022

Proximal gradient method

Minimizing composite function:

$$\min_x f(x) = \overbrace{g(x)}^{\text{differentiable}} + \overbrace{h(x)}^{\text{non-differentiable}}$$

Proximal gradient:

$$x_{t+1} = \text{prox}_{\eta h}(\underline{x_t} - \eta \nabla g(x_t))$$

How fast is it? \rightarrow rate of convergence

Gradient mapping

$$x_{t+1} = \text{prox}_{\eta h} (x_t - \eta \nabla g(x_t))$$

Define gradient mapping:

- $G_\eta(x) \stackrel{\text{def.}}{=} \frac{1}{\eta} (x - \text{prox}_{\eta h}(x - \eta \nabla g(x)))$

Then we can rewrite the proximal gradient method into something that looks more like a gradient descent update step:

- $x_{t+1} = x_t - \eta G_\eta(x_t)$

- ▶ G_η is called the gradient map of proximal gradient method, and we treat this as if it's a gradient, but G_η is not a (sub)gradient of f in general.
- ▶ We do this to make analyzing convergence behavior easier.

Gradient mapping and optimality condition $x = \text{prox}_{\eta h}(y) = \underset{x}{\text{argmin}} \frac{1}{2\eta} \|x-y\|_2^2 + h(x)$

Optimal solutions are the only fixed points of the prox grad update $\Leftrightarrow 0 \in \frac{1}{\eta}(x-y) + \partial h(x)$

$$G_{\eta}(x_t) \notin \partial f(x_t)$$

$$x_{t+1} = x_t - \eta \underline{G_{\eta}(x_t)}$$

$$\Leftrightarrow \boxed{y-x \in \eta \partial h(x)}$$

i.e., $G_{\eta}(x) = 0$ when x is a minimizer of f .

$$G_{\eta}(\hat{x}) = 0 \Leftrightarrow \frac{1}{\eta}(\hat{x} - \text{prox}_{\eta h}(\hat{x} - \eta \nabla g(\hat{x}))) = 0$$

$$\Leftrightarrow \hat{x} = \text{prox}_{\eta h}(\hat{x} - \eta \nabla g(\hat{x}))$$

$$\Leftrightarrow \hat{x} - \eta \nabla g(\hat{x}) - \hat{x} \in \eta \partial h(\hat{x})$$

$$\Leftrightarrow 0 \in \nabla g(\hat{x}) + \partial h(\hat{x})$$

$$\min_x f(x) = g(x) + h(x)$$

Key lemma

$$\min_x f(x) = g(x) + h(x) \quad \text{convex}$$

Lemma

β -smooth
 α -strongly convex.

$$f(x_t - \eta G_\eta(x_t)) \leq f(z) + \langle G_\eta(x_t), x_t - z \rangle - \frac{\eta}{2} \|G_\eta(x_t)\|_2^2 - \frac{\alpha}{2} \|x_t - z\|_2^2$$

Setting $z = x_t$, $x = x_t$ gives

$$\underline{f(x_{t+1})} = \underline{f(x_t - \eta G_\eta(x_t))} \leq \underline{f(x_t) - \frac{\eta}{2} \|G_\eta(x_t)\|_2^2}$$

which looks the progress bound for smooth functions.

$$\circ \quad f(x_{t+1}) = f(x_t) - \frac{1}{2\beta} \|\nabla f(x_t)\|_2^2$$

Key lemma – proof (1/3)

$$f(x - \eta G_{\eta}(x)) = \underbrace{g(x - \eta G_{\eta}(x))}_{(A)} + \underbrace{h(x - \eta G_{\eta}(x))}_{(B)} \leq ? \quad \eta = \frac{1}{L}$$

(1) β -smooth: $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|y - x\|_2^2$

$$g(x - \eta G_{\eta}(x)) \leq \underbrace{g(x)} + \langle \nabla g(x), -\eta G_{\eta}(x) \rangle + \frac{\beta}{2} \|G_{\eta}(x)\|_2^2$$

(2) α -strongly convex: $f(y) \geq \underline{f(x) + \langle \nabla f(x), y - x \rangle + \frac{\alpha}{2} \|y - x\|_2^2}$

$$\underbrace{g(x - \eta G_{\eta}(x))}_{(A)} \leq g(z) - \langle \nabla g(x), z - x \rangle - \frac{\alpha}{2} \|z - x\|_2^2$$

$$+ \langle \nabla g(x), -\eta G_{\eta}(x) \rangle + \frac{\beta}{2} \|G_{\eta}(x)\|_2^2$$

Key lemma – proof (2/3)

$$\textcircled{B}: h(x - \eta G_\eta(x))$$

$$(1) \quad \underline{x - \eta G_\eta(x)} = \cancel{x - \eta \frac{1}{\eta} (x - \text{prox}_{\eta h}(x - \eta \nabla g(x)))} = \underline{\text{prox}_{\eta h}(x - \eta \nabla g(x))}$$

$$x - \eta G_\eta(x) = \text{prox}_{\eta h}(x - \eta \nabla g(x))$$

$$\Leftrightarrow \cancel{x - \eta \nabla g(x)} - \cancel{(x - \eta G_\eta(x))} \in \cancel{\eta} \partial h(x - \eta G_\eta(x))$$

$$\Leftrightarrow \underline{G_\eta(x) - \nabla g(x)} \in \partial h(x - \eta G_\eta(x))$$

$$(2) \quad h(x - \eta G_\eta(x)) \leq h(z) - \langle G_\eta(x) - \nabla g(x), z - (x - \eta G_\eta(x)) \rangle$$

Key lemma – proof (3/3)

$$f(x) = \underbrace{g(x)}_{\textcircled{A}} + \underbrace{h(x)}_{\textcircled{B}}$$

$$\textcircled{A} = g(x - \eta \nabla g(x)) \leq \underbrace{g(z)}_{\textcircled{1}} - \underbrace{\langle \nabla g(x), z - x \rangle}_{\textcircled{4}} - \underbrace{\frac{d}{2} \|z - x\|_2^2}_{\textcircled{6}} + \underbrace{\langle \nabla g(x), -\eta \nabla g(x) \rangle}_{\textcircled{5}} + \underbrace{\frac{\eta}{2} \|\nabla g(x)\|_2^2}_{\textcircled{3}}$$

$$\textcircled{B} = h(x - \eta \nabla g(x)) \leq \underbrace{h(z)}_{\textcircled{1}} - \underbrace{\langle \nabla g(x) - \nabla h(x), z - (x - \eta \nabla g(x)) \rangle}_{\substack{\textcircled{2} \textcircled{3} \textcircled{4} \textcircled{5} \\ \textcircled{2} \textcircled{4} \textcircled{3} \textcircled{5}}}$$

$$\textcircled{A} + \textcircled{B} \leq f(z) + \langle \nabla g(x), x - z \rangle - \frac{\eta}{2} \|\nabla g(x)\|_2^2 - \frac{d}{2} \|x - z\|_2^2 \quad \square$$

Convergence for smooth functions

$g : \beta$ -smooth. $d=0$

$$f(x - \eta G_{\eta}(x)) \leq f(z) + \langle G_{\eta}(x), x - z \rangle - \frac{\eta}{2} \|G_{\eta}(x)\|_2^2 - \underline{\frac{\beta}{2} \|x - z\|_2^2}$$

substitute $x \rightarrow x_t$, $z \rightarrow x^*$

$$f(x_{t+1}) \leq f(x^*) + \langle G_{\eta}(x_t), x_t - x^* \rangle - \frac{\eta}{2} \|G_{\eta}(x_t)\|_2^2$$

$$\Leftrightarrow \underbrace{f(x_{t+1}) - f(x^*)}_{\text{sum over } T \text{ iterations}} \leq \langle G_{\eta}(x_t), x_t - x^* \rangle - \frac{\eta}{2} \|G_{\eta}(x_t)\|_2^2$$

$$= \frac{1}{2\eta} \left(2\eta \langle G_{\eta}(x_t), x_t - x^* \rangle - \eta^2 \|G_{\eta}(x_t)\|_2^2 - \|x_t - x^*\|_2^2 + \underbrace{\|x_t - x^*\|_2^2} \right)$$

sum over T iterations

both sides & divide by T

$$= \frac{1}{2\eta} \left(\|x_t - x^*\|_2^2 - \|x_t - x^* - \eta G_{\eta}(x_t)\|_2^2 \right)$$

$$\sim o\left(\frac{1}{T}\right)$$

$$= \frac{1}{2\eta} \left(\underbrace{\|x_t - x^*\|_2^2 - \|x_{t+1} - x^*\|_2^2}_{\text{telescoping sum}} \right)$$

Convergence for smooth and strongly convex functions

g - β -smooth & α -strongly convex

$$\min_x f(x) = g(x) + h(x)$$

$$\begin{aligned} 0 \leq f(x_{t+1}) - f(x^*) &\leq \langle G_\eta(x_t), x_t - x^* \rangle - \frac{\eta}{2} \|G_\eta(x_t)\|_v^2 - \frac{\alpha}{2} \|x_t - x^*\|_v^2 \\ &= \frac{1}{2\eta} \left(\|x_t - x^*\|_v^2 - \underline{\|x_{t+1} - x^*\|_v^2} - \eta\alpha \|x_t - x^*\|_v^2 \right) \end{aligned}$$

$$\Leftrightarrow \|x_{t+1} - x^*\|_v^2 \leq (1 - \eta\alpha) \|x_t - x^*\|_v^2$$

$$\Leftrightarrow \|x_{t+1} - x^*\|_v^2 \leq \underbrace{(1 - \eta\alpha)^t}_{R^t} \|x_1 - x^*\|_v^2 \quad : \quad \sim O(\rho^t)$$

Example: ISTA

L1-regularized least squares or Lasso:

$$\min_x \underbrace{\frac{1}{2} \|Ax - y\|_2^2} + \lambda \underbrace{\|x\|_1}$$

- ▶ Recall that this can achieve sparse solutions.
- ▶ Apply subgradient method? It will converge with $\mathcal{O}(1/\sqrt{t})$.
- ▶ f is a composite function of g smooth and h non-smooth functions, so consider using proximal gradient method.

Example: ISTA

(cont'd) Apply proximal gradient:

$$\begin{aligned}x_{t+1} &= \text{prox}_{\eta h}(y_{t+1}) \\ &= \arg \min_x \frac{1}{2\eta} \|x - y_{t+1}\|_2^2 + \lambda \|x\|_1\end{aligned}$$

where $y_{t+1} = x_t - \eta \nabla g(x_t) = x_t - \eta(A^\top(Ax_t - y))$.

The solution to the prox will be

$$x^+ : x_i = \begin{cases} x_i + \eta\lambda & \text{for } y_i < -\eta\lambda \\ 0 & \text{for } |y_i| \leq \eta\lambda \\ x_i - \eta\lambda & \text{for } y_i > \eta\lambda . \end{cases}$$


“A fast iterative shrinkage-thresholding algorithm for linear inverse problems” (Beck and Teboulle 2009)

Thank you

Any questions?

A lot of material in this course is borrowed or derived from the following:

- ▶ Numerical Optimization, Jorge Nocedal and Stephen J. Wright.
- ▶ Convex Optimization, Stephen Boyd and Lieven Vandenberghe.
- ▶ Convex Optimization, Ryan Tibshirani.
- ▶ Optimization for Machine Learning, Martin Jaggi and Nicolas Flammarion.
- ▶ Optimization Algorithms, Constantine Caramanis.
- ▶ Advanced Machine Learning, Mark Schmidt.

-  Beck, Amir and Marc Teboulle (2009). “A fast iterative shrinkage-thresholding algorithm for linear inverse problems”. In: *SIAM journal on imaging sciences*.