

CSIE5410 Optimization algorithms

Lecture 7: Convergence of the proximal gradient method

Yen-Huan Li (yenhuan.li@csie.ntu.edu.tw)

08.11.2018

Department of Computer Science and Information Engineering
National Taiwan University

We have seen the proximal gradient method. This lecture addresses the following questions.

- Under what conditions does the algorithm provably converge?
- How fast does it converge?

Recommended reading

- A. Beck and M. Teboulle. 2009. A fast iterative shrinkage-thresholding algorithm for linear inverse problems.
- Yu. Nesterov and A. Nemirovski. 2013. On first-order algorithms for ℓ_1 /nuclear norm minimization.
- H. Lu *et al.* 2018. Relatively smooth convex optimization by first-order methods, and applications.
- M. Teboulle. 2018. A simplified view of first order methods for optimization.

Table of contents

1. Smooth case
2. Acceleration for the smoothness case
3. Relatively smooth case
4. Conclusions

Smooth case

Recap: Proximal gradient method

Consider the optimization problem

$$\varphi^* = \min_x \{ f(x) + g(x) \mid x \in \mathbb{R}^p \},$$

for some L -smooth convex function f and proper closed convex function g .

Proximal gradient method

$$\begin{aligned} x_{t+1} &\leftarrow (I + \eta_t \partial g)^{-1} (I - \eta_t \nabla f) x_t \\ &= \text{prox}_{\eta_t g} (x_t - \eta_t \nabla f(x_t)) \end{aligned}$$

Equivalent formulation

$$x_{t+1} \leftarrow \text{prox}_{\eta_t g}(x_t - \eta_t \nabla f(x_t)).$$

Proposition. Equivalently, we write

$$x_{t+1} = \arg \min_x \left\{ \langle \nabla f(x_t), x - x_t \rangle + g(x) + \frac{1}{2\eta_t} \|x - x_t\|_2^2 \mid x \in \mathbb{R}^p \right\}.$$

Proof. Plug in the definition of a proximal mapping.

Remark. Recall that when g is the indicator function of a closed convex set \mathcal{X} , the iteration rule of the projected gradient method becomes a special case.

Proximal set-up (1/3)

Consider the problem

$$\varphi^{\star} = \min_x \{ f(x) + g(x) \mid x \in \mathcal{X} \},$$

for some convex function f that is L -smooth on \mathcal{X} , proper closed convex function g , and closed convex set $\mathcal{X} \subseteq \mathbb{R}^p$.

Consider the Banach space $(\mathbb{R}^p, \|\cdot\|)$.

Proximal set-up (2/3)

Definition. We say that a function $h : \mathcal{X} \rightarrow \mathbb{R}$ is a *distance generating function (DGF)*, if and only if the following are satisfied.

- The function h is continuous and convex on \mathcal{X} .
- There exists a function h' that is continuous on $\mathcal{X}^\circ := \mathcal{X} \cap \text{dom}(\partial h)$ and satisfies

$$h'(x) \in \partial h(x), \quad \forall x \in \mathcal{X}^\circ.$$

- The function h is 1-strongly convex on \mathcal{X}° , i.e.,

$$\langle h'(x) - h'(y), x - y \rangle \geq \|x - y\|^2, \quad \forall x, y \in \mathcal{X}^\circ.$$

Proximal set-up (3/3)

Definition. The *Bregman divergence* is given by

$$D_h(y, x) := h(y) - h(x) - \langle h'(x), y - x \rangle, \quad \forall x \in \mathcal{X}^\circ, y \in \mathcal{X}.$$

Definition. The *(composite) prox-mapping* is given by

$$T_L(x) := \arg \min_{y \in \mathcal{X}} \{ \langle \nabla f(x), y - x \rangle + g(y) + LD_h(y, x) \mid y \in \mathcal{X} \}.$$

Bregman proximal gradient method.

$$x_t \leftarrow T_L(x_{t-1}).$$

Interpretation of the Bregman proximal gradient method (1/3)

Recall the derivation of the proximal point method:

$$\begin{aligned}0 \in \partial f(x^*) &\Leftrightarrow x^* \in (I + \partial f)x^* \\ &\Leftrightarrow x^* = (I + \partial f)^{-1}x^*\end{aligned}$$

How about the following derivation? Let h be a differentiable μ -strongly convex function (w.r.t. the 2-norm).

$$\begin{aligned}0 \in \partial f(x^*) &\Leftrightarrow \nabla h(x^*) \in (\nabla h + \partial f)x^* \\ &\Leftrightarrow x^* = (\nabla h + \partial f)^{-1}(\nabla h(x^*))\end{aligned}$$

J. Eckstein. 1993. Nonlinear proximal point algorithms using Bregman functions, with applications to convex programming.

H. H. Bauschke *et al.* 2003. Bregman monotone optimization algorithms.

Interpretation of the Bregman proximal gradient method (2/3)

Consider the algorithm:

$$x_{t+1} \leftarrow (\nabla h + \partial f)^{-1}(\nabla h(x_t)).$$

Then, we write

$$\nabla h(x_t) \in (\nabla h + \partial f)x_{t+1} \Leftrightarrow 0 \in \partial f(x_{t+1}) + \nabla h(x_{t+1}) - \nabla h(x_t),$$

which is equivalent to

$$x_{t+1} \in \arg \min_x \{ f(x) + D_h(x, x_t) \mid x \in \mathbb{R}^p \}.$$

J. Eckstein. 1993. Nonlinear proximal point algorithms using Bregman functions, with applications to convex programming.

H. H. Bauschke *et al.* 2003. Bregman monotone optimization algorithms.

Interpretation of the Bregman proximal gradient method (3/3)

Similarly, consider the monotone inclusion problem

$$0 \in \nabla f(x^\star) + \partial g(x^\star).$$

We may write

$$\begin{aligned}(\nabla h - \nabla f)(x^\star) &\in (\nabla h + \partial g)(x^\star) \\ \Leftrightarrow x^\star &= (\nabla h + \partial g)^{-1}(\nabla h - \nabla f)(x^\star),\end{aligned}$$

which motivates the algorithm:

$$\begin{aligned}x_{t+1} &\leftarrow (\nabla h + \partial g)^{-1}(\nabla h - \nabla f)(x_t) \\ &= \arg \min_x \{ \langle \nabla f(x_t), x - x_t \rangle + g(x) + D_h(x, x_t) \mid x \in \mathbb{R}^p \}.\end{aligned}$$

Convergence of the Bregman proximal gradient method

Algorithm Bregman proximal gradient method

- 1: Set $x_0 \in \mathcal{X}^\circ$.
 - 2: **for** $t = 0, 1, \dots, T$ **do**
 - 3: $x_{t+1} \leftarrow \arg \min_{x \in \mathcal{X}} \{ \langle \nabla f(x_t), x - x_t \rangle + LD(x, x_t) + g(x) \}$
 - 4: **end for**
-

Theorem. Suppose that g is finite on $\text{ri } \mathcal{X}$. Then, for all $t \geq 1$, it holds that $x_t \in \mathcal{X}^\circ$, and

$$\varphi(x_t) - \varphi(x) \leq \frac{LD_h(x, x_0)}{t}, \quad \forall x \in \mathcal{X}.$$

Yu. Nesterov and A. Nemirovski. 2013. On first-order algorithms for ℓ_1 /nuclear norm minimization.

Proof of the theorem (1/5)

Lemma. Let ψ be a proper closed convex function that is finite on $\text{ri } \mathcal{X}$. Define

$$x_\psi^\star \in \arg \min_x \{ \psi(x) + h(x) \mid x \in \mathcal{X} \}.$$

Then, $x_\psi^\star \in \mathcal{X}^\circ$, and there exists some $\xi \in \partial\psi(x_\psi^\star)$, such that

$$\langle \xi + h'(x_\psi^\star), x - x_\psi^\star \rangle \geq 0, \quad \forall x \in \mathcal{X}.$$

Proof. See the reference.

Yu. Nesterov and A. Nemirovski. 2013. On first-order algorithms for ℓ_1 /nuclear norm minimization.

Proof of the theorem (2/5)

Lemma. Suppose that $\tilde{\psi} := \psi - Lh$ is convex. Let $x_\psi^\star \in \arg \min_{x \in \mathcal{X}} \psi(x)$. Then $x_\psi^\star \in \mathcal{X}^\circ$, and

$$\psi(x) \geq \psi(x_\psi^\star) + LD_h(x, x_\psi^\star), \quad \forall x \in \mathcal{X}.$$

Proof. By the previous lemma, $x_\psi^\star \in \mathcal{X}^\circ$, and there exists some $\eta \in \partial \tilde{\psi}(x_\psi^\star)$ such that

$$\langle \eta + Lh'(x_\psi^\star), x - x_\psi^\star \rangle \geq 0, \quad \forall x \in \mathcal{X}.$$

Proof of the theorem (3/5)

Proof (of the lemma) continued. We write

$$\begin{aligned}\psi(x) &= \tilde{\psi}(x) + Lh(x) \\ &\geq \tilde{\psi}(x_\psi^\star) + \langle \eta, x - x_\psi^\star \rangle + Lh(x) \\ &\geq \tilde{\psi}(x_\psi^\star) - \langle Lh'(x_\psi^\star), x - x_\psi^\star \rangle + Lh(x) \\ &= \tilde{\psi}(x_\psi^\star) + Lh(x_\psi^\star) + Lh(x) - (Lh(x_\psi^\star) + \langle Lh'(x_\psi^\star), x - x_\psi^\star \rangle) \\ &= \psi(x_\psi^\star) + LD_h(x, x_\psi^\star).\end{aligned}$$

Proof of the theorem (4/5)

Proof of the theorem. Define

$$m_L(x) := f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + g(x) + LD_h(x, x_t).$$

By smoothness of f and strong convexity of h ,

$$m_L(x) \geq \varphi(x) := f(x) + g(x).$$

By the previous lemma, $x_{t+1} \in \mathcal{X}^\circ$, and

$$m_L(x) \geq m_L(x_{t+1}) + LD_h(x, x_{t+1}), \quad \forall x \in \mathcal{X}.$$

Then, by smoothness of f ,

$$\begin{aligned} \varphi(x_t) = m_L(x_t) &\geq m_L(x_{t+1}) + LD_h(x_t, x_{t+1}) \\ &\geq \varphi(x_{t+1}), \end{aligned}$$

showing that the sequence $(\varphi(x_t))_{t \in \mathbb{N}}$ is non-increasing.

Proof of the theorem (5/5)

Proof of the theorem continued.

Moreover, we have

$$\begin{aligned}\varphi(x) + LD_h(x, x_t) &\geq m_L(x) \geq m_L(x_{t+1}) + LD_h(x, x_{t+1}) \\ &\geq \varphi(x_{t+1}) + LD_h(x, x_{t+1}).\end{aligned}$$

Summing up the inequalities over t , we obtain

$$\begin{aligned}t(\varphi(x_t) - \varphi(x)) &\leq \sum_{\tau=0}^{t-1} (\varphi(x_\tau) - \varphi(x)) \\ &\leq L(D_h(x, x_0) - D_h(x, x_t)).\end{aligned}$$

Acceleration for the smoothness case

Estimate sequence approach (1/2)

Estimate sequence approach. Consider an iterative algorithm that maintains three sequences:

1. A sequence of iterates $(y_t)_{t \geq 0}$.
2. A sequence of increasing numbers $(A_t)_{t \geq 0}$ such that

$$A_0 = 0, \quad A_t := A_{t-1} + a_t, \quad \forall t \in \mathbb{N}.$$

3. A sequence of *estimate functions*

$$\psi_t := \langle \xi_k, \cdot \rangle + A_k g(x) + \frac{1}{2} \|\cdot - x_0\|^2, \quad \forall t \geq 0.$$

Estimate sequence approach (2/2)

Conditions. The algorithm has to satisfy for all t :

1. $A_t \varphi(y_t) \leq \psi_t^* := \min_x \psi_t(x)$.
2. $\psi_t(x) \leq A_t \varphi(x) + LD_h(x, y_0)$ for all $x \in \mathbb{R}^p$.

Theorem. If the conditions are satisfied, then it holds that

$$\varphi(y_t) - \varphi(x) \leq \frac{LD_h(x, y_0)}{A_t}, \quad \forall x \in \mathcal{X}.$$

Yu. Nesterov. 2008. Accelerating the cubic regularization of Newton's method on convex problems.

Accelerated Bregman proximal gradient method

Algorithm Accelerated Bregman proximal gradient method

- 1: Set $y_0 = \arg \min_{x \in \mathcal{X}} h(x)$ and $\psi_0(x) = LD_h(x, y_0)$.
 - 2: **for** $t = 0, 1, \dots, T$ **do**
 - 3: $z_t \leftarrow \arg \min_{x \in \mathcal{X}} \psi_t(x)$.
 - 4: $\gamma_t \leftarrow \frac{2(t+2)}{(t+1)(t+4)}$.
 - 5: $x_{t+1} \leftarrow \gamma_t z_t + (1 - \gamma_t) y_t$.
 - 6: $\hat{x}_{t+1} \leftarrow T_{\frac{2L}{t+2}}(x_{t+1})$.
 - 7: $y_{t+1} \leftarrow \gamma_t \hat{x}_{t+1} + (1 - \gamma_t) y_t$.
 - 8: $\psi_{t+1} \leftarrow \psi_t + \frac{t+2}{2} [f(x_{t+1}) + \langle \nabla f(x_{t+1}), x - x_{t+1} \rangle + g(x)]$
 - 9: **end for**
-

Yu. Nesterov. 2013. Gradient methods for minimizing composite functions.

Yu. Nesterov. 2013. On first-order algorithms for ℓ_1 /nuclear norm minimization.

Theorem. It holds that for all $t, y_t \in \mathcal{X}$, and

$$\varphi(y_t) - \varphi^* \leq \frac{4LD_h(x^*, y_0)}{t(t+3)}.$$

Remark. The algorithm requires evaluating two prox-mappings per iteration.

Remark. A similar algorithm for the case where $g \equiv 0$ can be found in the first reference below.

Yu. Nesterov. 2005. Smooth minimization of non-smooth functions.

Yu. Nesterov. 2013. Gradient methods for minimizing composite functions.

Yu. Nesterov. 2013. On first-order algorithms for ℓ_1 /nuclear norm minimization.

Proof of the theorem (1/5)

Lemma. It holds that $y_t \in \mathcal{X}$ for all t .

Proof. Notice that $\hat{x}_{t+1} \in \mathcal{X}$ as the output of a prox-mapping. The lemma follows by induction.

Define $a_{t+1} := \frac{t+2}{2}$. Then, $A_t = \frac{t(t+3)}{4}$.

Lemma. It holds that

$$\psi_t(x) \leq A_t \varphi(x) + LD(x, y_0), \quad \forall x \in \mathcal{X}.$$

Proof. The lemma follows by induction with the observation that

$$\begin{aligned} \psi_{t+1} &= \psi_t + a_{t+1} [f(x_{t+1}) + \langle \nabla f(x_{t+1}), x - x_{t+1} \rangle + g(x)] \\ &\leq \psi_t + a_{t+1} [f(x) + g(x)]. \end{aligned}$$

Proof of the theorem (2/5)

Lemma. It holds that

$$A_t \varphi(y_t) \leq \psi_t^\star := \min_x \psi_t(x).$$

Proof. We prove by induction. The inequality obvious holds for $t = 0$. Suppose now that the inequality holds for some $t \geq 0$.

Since $\psi_t - Lh$ is obviously convex, we have

$$\psi_t(x) \geq \psi_t^\star + LD_h(x, z_t).$$

By the convexity of f , we write

$$\begin{aligned} \psi_t(x) &\geq A_t \varphi(y_t) + LD_h(x, z_t) \\ &\geq A_t [f(x_{t+1}) + \langle \nabla f(x_{t+1}), y_t - x_{t+1} \rangle + g(y_t)] + LD_h(x, z_t). \end{aligned}$$

Proof of the theorem (3/5)

Proof continued. Then, by the algorithm, we obtain for every $x \in \mathcal{X}$,

$$\begin{aligned}\psi_{t+1}(x) \geq A_t [f(x_{t+1}) + \langle \nabla f(x_{t+1}), y_t - x_{t+1} \rangle + g(y_t)] + LD_h(x, z_t) + \\ a_{t+1} [f(x_{t+1}) + \langle \nabla f(x_{t+1}), x - x_{t+1} \rangle + g(x)].\end{aligned}$$

Also by the algorithm, we have

$$A_t(y_t - x_{t+1}) - a_{t+1}x_{t+1} = -a_{t+1}z_t.$$

Therefore, we obtain

$$\begin{aligned}\psi_{t+1}(x) \geq a_{t+1} [\langle \nabla f(x_{t+1}), x - z_t \rangle + g(x)] + \\ A_{t+1}f(x_{t+1}) + A_tg(y_t) + LD_h(x, z_t).\end{aligned}$$

Proof of the theorem (4/5)

Proof continued. Then, we write

$$\begin{aligned}\psi_{t+1}^* &\geq \min_{x \in \mathcal{X}} \{a_{t+1} [\langle \nabla f(x_{t+1}), x - z_t \rangle + g(x)] + \\ &\quad A_{t+1}f(x_{t+1}) + A_t g(y_t) + LD_h(x, z_t)\} \\ &\geq a_{t+1} [\langle \nabla f(x_{t+1}), \hat{x}_{t+1} - z_t \rangle + g(\hat{x}_{t+1})] + \\ &\quad A_{t+1}f(x_{t+1}) + A_t g(y_t) + LD_h(\hat{x}_{t+1}, z_t) \\ &\geq a_{t+1} \langle \nabla f(x_{t+1}), \hat{x}_{t+1} - z_t \rangle + A_{t+1}g(y_{t+1}) + A_{t+1}f(x_{t+1}) + \\ &\quad \frac{L}{2} \|\hat{x}_{t+1} - z_t\|^2,\end{aligned}$$

where we have exploited the definition of \hat{x}_{t+1} , strong convexity of h , convexity of g , and the definition of y_{t+1} .

Proof of the theorem (5/5)

Proof continued. Notice that

$$\hat{x}_{t+1} - z_t = \frac{y_{t+1} - x_{t+1}}{\gamma_t} = \frac{A_{t+1}}{a_{t+1}} (y_{t+1} - x_{t+1}).$$

Then, we obtain

$$\begin{aligned}\psi_{t+1}^* &\geq A_{t+1} \langle \nabla f(x_{t+1}), y_{t+1} - x_{t+1} \rangle + A_{t+1} g(y_{t+1}) + \\ &\quad A_{t+1} f(x_{t+1}) + A_{t+1} \frac{L}{2} \|y_{t+1} - x_{t+1}\|^2 \\ &\geq A_{t+1} [f(y_{t+1}) + g(y_{t+1})].\end{aligned}$$

The lemma follows.

Proof of the theorem. Directly apply the estimation sequence result.

Remark. Arguably, FISTA looks more natural, but its proof is even less principled.

Assumption. Set $\mathcal{X} = \mathbb{R}^p$. The function f is L -smooth (w.r.t the 2-norm) on \mathbb{R}^p .

Algorithm Fast iterative shrinkage thresholding algorithm (FISTA)

- 1: Set $y_1 = x_0 \in \mathbb{R}^p$, $\gamma_1 = 1$, and $\eta = 1/L$.
 - 2: **for** $t = 1, \dots, T$ **do**
 - 3: $x_t \leftarrow (I + \eta \partial g)^{-1}(I - \eta \nabla f)y_t$.
 - 4: $\gamma_{t+1} \leftarrow \frac{1 + \sqrt{1 + 4\gamma_t^2}}{2}$.
 - 5: $y_{t+1} = x_t + \frac{\gamma_t - 1}{\gamma_{t+1}}(x_t - x_{t-1})$.
 - 6: **end for**
-

A. Beck and M. Teboulle. 2009. A fast iterative shrinkage thresholding algorithm for linear inverse problems.

Relatively smooth case

Problem formulation

Consider the optimization problem

$$\varphi^{\star} = \min_x \{ f(x) + g(x) \mid x \in \mathcal{X} \},$$

where f is convex and L -smooth relative to a proper closed convex function h on $\text{int}(\text{dom } h)$, g is a proper closed convex function, and \mathcal{X} is a closed convex set.

Definition. We say that f is L -smooth relative to a convex function h on a set \mathcal{X} , if and only if

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + LD_h(y, x), \quad \forall x, y \in \mathcal{X}.$$

Another prox set-up (1/3)

Definition. We say that h is *essentially smooth*, if and only if $\text{int}(\text{dom } h) \neq \emptyset$, h is differentiable on $\text{int}(\text{dom } h)$, and $\|\nabla h(x_k)\| \rightarrow +\infty$ for any sequence $(x_k)_{k \in \mathbb{N}}$ in $\text{int}(\text{dom } h)$ converging to a point on the boundary of $\text{dom } h$.

Theorem. The function h is essentially smooth, if and only if $\text{dom } \partial h = \text{int}(\text{dom } h) \neq \emptyset$.

Definition. We say that h is *strictly convex*, if and only if

$$h(\alpha x + (1 - \alpha)y) < \alpha h(x) + (1 - \alpha)h(y), \quad \forall \alpha \in]0, 1[, x, y.$$

Definition. We say that h is *Legendre*, if and only if it is essentially smooth and strictly convex.

Another prox set-up (2/3)

$$\varphi^* = \min_x \{ f(x) + g(x) \mid x \in \mathcal{X} \},$$

Assumptions.

1. h is Legendre with $\text{cl}(\text{dom } h) = \mathcal{X}$.
2. $\text{dom } f \supset \text{dom } h$ and f is differentiable on $\text{int}(\text{dom } h)$.
3. $\text{dom } g \cap \text{int}(\text{dom } h) \neq \emptyset$.

Question. How do we specialize the set-up for entropic mirror descent?

Another prox set-up (3/3)

Theorem. Let $x \in \text{int}(\text{dom } h)$. Let ψ be a proper closed convex function. Consider the optimization problem

$$x_+ = \arg \min_y \{ \psi(y) + D_h(y, x) \mid y \in \mathcal{X} \}.$$

1. If $\psi(x) > -\infty$ on \mathcal{X} , then x_+ exists and is unique.
2. If in addition, $\text{ri}(\text{dom } \psi \cap \mathcal{X}) \subset \text{int}(\text{dom } h)$, then $x^+ \in \text{dom } \psi \cap \text{int}(\text{dom } h)$, and

$$0 \in \partial\psi(x_+) + \nabla h(x_+) - \nabla h(x).$$

Bregman proximal inequality

This is indeed the “key lemma” in Lecture 4.

Theorem. Let $x \in \text{int}(\text{dom } h)$. Let ψ be proper closed convex.
Define

$$x_+ = \arg \min_y \{ \psi(y) + D_h(y, x) \mid y \in \mathcal{X} \}.$$

Then $x_+ \in \text{dom } \psi \cap \text{int}(\text{dom } h)$, and

$$\psi(x_+) - \psi(u) \leq D_h(u, x) - D_h(u, x_+) - D_h(x_+, x), \quad \forall u \in \text{dom } h.$$

Algorithm Bregman proximal gradient method

- 1: Set $x_0 \in \mathcal{X}^\circ$.
 - 2: **for** $t = 0, 1, \dots, T$ **do**
 - 3: $x_{t+1} \leftarrow \arg \min_{x \in \mathcal{X}} \{ \langle \nabla f(x_t), x - x_t \rangle + LD(x, x_t) + g(x) \}$
 - 4: **end for**
-

Theorem. Define $\varphi := f + g$. For all $t \geq 1$, it holds that $x_t \in \text{dom } g \cap \text{int}(\text{dom } h)$, and

$$\varphi(x_t) - \varphi(u) \leq \frac{LD_h(u, x_0)}{t}, \quad \forall u \in \text{dom } h.$$

Proof of the theorem (1/2)

Proof. We write, by L -relative smoothness,

$$f(x_t) \leq f(x_{t-1}) + \langle \nabla f(x_{t-1}), x_t - x_{t-1} \rangle + LD_h(x_t, x_{t-1}).$$

By the Bregman proximal inequality, we write, for any $x \in \text{dom } h$,

$$\begin{aligned} LD_h(x_t, x_{t-1}) &\leq \langle \nabla f(x_{t-1}), x - x_t \rangle + g(x) - g(x_t) + \\ &\quad L(D_h(x, x_{t-1}) - D_h(x, x_t)). \end{aligned}$$

Then, we obtain

$$\begin{aligned} \varphi(x_t) &\leq f(x_{t-1}) + \langle \nabla f(x_{t-1}), x - x_{t-1} \rangle + g(x) + \\ &\quad L(D_h(x, x_{t-1}) - D_h(x, x_t)) \\ &\leq \varphi(x) + L(D_h(x, x_{t-1}) - D_h(x, x_t)). \end{aligned}$$

Proof of the theorem (2/2)

Proof continued. Notice that then the sequence $(\varphi(x_t))_{t \in \mathbb{N}}$ is non-increasing.

Summing over all t , we obtain

$$\sum_{\tau=1}^t \varphi(x_\tau) \leq t\varphi(x) + LD_h(x, x_0) - LD_h(x, x_t).$$

Since $(\varphi(x_t))_{t \in \mathbb{N}}$ is non-increasing, we have

$$t(\varphi(x_t) - \varphi(x)) \leq LD_h(x, x_0).$$

Remark. The proof strategy is exact the same as that for the “mirror descent”.

Is the proof easy?

Below is a slightly earlier convergence guarantee due to Nesterov.

Theorem. Suppose that f is L -smooth on \mathcal{X} . Suppose that

$$\|x - x^*\| \leq R, \quad \forall x \in \mathcal{X} \text{ such that } \varphi(x) \leq \varphi(x_0).$$

If $\varphi(x_0) - \varphi(x^*) \geq LR^2$, then

$$\varphi(x_1) - \varphi^* \leq \frac{LR^2}{2}.$$

Otherwise, for every $t \in \mathbb{N}$,

$$\varphi(x_t) - \varphi^* \leq \frac{2LR^2}{t+2}.$$

Yu. Nesterov. 2013. Gradient methods for minimizing composite functions.

Conclusions

Summary

Two new notions:

1. Prox-mapping.
2. Estimate sequence approach.

Two algorithms:

1. Bregman proximal gradient algorithm.
2. Accelerated Bregman proximal gradient algorithm.

Two things to notice:

1. Subtle difference between the two proximal set-ups.
2. Similarity of the proof strategies.

Next lecture

- Frank-Wolfe algorithm.