

Lecture 7. Proximal Gradient Methods

Bao Wang

Department of Mathematics

Scientific Computing and Imaging Institute

University of Utah

Math 5750/6880, Fall 2021

Loss function

So far, we have formulated training machine learning models as

$$\min f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_i(\mathbf{x}) + R(\mathbf{x})$$

where \mathbf{x} is the parameter of the machine learning model, $\mathcal{L}_i(\mathbf{x})$ is the loss of the i th training instance, and $R(\mathbf{x})$ is the regularization term.

How to find the optimal \mathbf{x}^* if $R(\mathbf{x})$ is not differentiable everywhere, e.g. ℓ_1 -regularization?

Subgradient methods or **proximal gradient methods**.

Recap: Subgradient method

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{g}^t, \quad (1)$$

where \mathbf{g}^t is any subgradient of f at \mathbf{x}^t .

We say \mathbf{g} is a subgradient of f at the point \mathbf{x} if

$$f(\mathbf{z}) \geq \underbrace{f(\mathbf{x}) + \mathbf{g}^\top(\mathbf{z} - \mathbf{x})}_{\text{a linear under-estimate of } f}, \quad \forall \mathbf{z}. \quad (2)$$

The set of all subgradients of f at \mathbf{x} is called the subdifferential of f at \mathbf{x} , denoted by $\partial f(\mathbf{x})$.

Recap: Convergence of subgradient methods

	stepsize rule	convergence rate	iteration complexity
convex & Lipschitz problems	$\eta_t \asymp \frac{1}{\sqrt{t}}$	$O\left(\frac{1}{\sqrt{t}}\right)$	$O\left(\frac{1}{\varepsilon^2}\right)$
strongly convex & Lipschitz problems	$\eta_t \asymp \frac{1}{t}$	$O\left(\frac{1}{t}\right)$	$O\left(\frac{1}{\varepsilon}\right)$

Subgradient method is very slow!

Proximal gradient descent for composite functions

Consider the composite model

$$\min_{\mathbf{x}} F(\mathbf{x}) := f(\mathbf{x}) + h(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^n,$$

let $F^{opt} := \min_{\mathbf{x}} F(\mathbf{x})$ be the optimal cost.

1. ℓ_1 regularized minimization for promoting sparsity (e.g., lasso)




$$\min_{\mathbf{x}} f(\mathbf{x}) + \|\mathbf{x}\|_1$$

2. nuclear norm (sum of the singular values) regularized minimization for promoting low-rank structure (Netflix competition)

$$\min_{\mathbf{X}} f(\mathbf{X}) + \|\mathbf{X}\|_*$$

Matrix completion

Movies

				
Bob	4	?	?	4
Alice	?	5	4	?
Joe	?	5	?	?
Sam	5	?	?	?

Users

Recommender system through matrix completion! (User privacy?)

Nuclear norm?

The rank is the L0 norm of the spectrum (the vector of singular values) and so the best approximation is the L1 norm of the spectrum. However, for a positive semidefinite matrix, the singular values are the eigenvalues and the sum of eigenvalues is the trace.

The connection between the rank and the nuclear norm can be understood as follows. The convex hull of all rank-1 matrices with norm equal to 1 is the nuclear ball. This is why the tightest convex approximation to the rank functional is the nuclear norm.

A proximal view of gradient descent

We note that the gradient descent iteration

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t)$$

can be written as

$$\mathbf{x}^{t+1} = \arg \min_{\mathbf{x}} \left\{ \underbrace{f(\mathbf{x}^t) + \langle \nabla f(\mathbf{x}^t), \mathbf{x} - \mathbf{x}^t \rangle}_{\text{first-order approximation}} + \underbrace{\frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}^t\|_2^2}_{\text{proximal term}} \right\}.$$

Motivation. GD can be considered as find the optimal solution of the linear approximation of $f(\mathbf{x}^t)$, and the linear approximation is accurate when \mathbf{x} and \mathbf{x}^t is close to each other.

Proximal gradient algorithm

We note that the gradient descent iteration

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t)$$

can be written as

$$\mathbf{x}^{t+1} = \arg \min_{\mathbf{x}} \left\{ \underbrace{f(\mathbf{x}^t) + \langle \nabla f(\mathbf{x}^t), \mathbf{x} - \mathbf{x}^t \rangle}_{\text{first-order approximation}} + \underbrace{\frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}^t\|_2^2}_{\text{proximal term}} \right\}.$$

\Leftrightarrow

$$\mathbf{x}^{t+1} = \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x} - (\mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t))\|_2^2 \right\}.$$

Proximal gradient algorithm

Define the proximal operator

$$\text{prox}_h(\mathbf{x}) := \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 + h(\mathbf{z}) \right\}$$

for any convex function h .

This allows one to express GD update as (set $h(\mathbf{z}) = 0$),

$$\mathbf{x}^{t+1} = \text{prox}_0(\mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t)). \quad (3)$$

One can generalize (3) to accommodate more general h ,

$$\mathbf{x}^{t+1} = \text{prox}_{\eta_t h}(\mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t)).$$

The proximal gradient algorithm alternates between gradient updates on f and proximal minimization on h , and it will be useful if prox_h is inexpensive.

Proximal gradient descent

Consider the composite model

$$\min_{\mathbf{x}} F(\mathbf{x}) := f(\mathbf{x}) + h(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^n,$$

Proximal gradient descent

for $k = 0, 1, \dots$

$$\mathbf{x}^{t+1} = \text{prox}_{\eta_t h}(\mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t))$$

Proximal mapping/operator

The proximal operator is define by

$$\text{prox}_h(\mathbf{x}) := \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 + h(\mathbf{z}) \right\}.$$

- > well-defined under very general conditions (including nonsmooth convex functions)
- > can be evaluated efficiently for many widely used functions (in particular, regularizers)
- > this abstraction is conceptually and mathematically simple, and covers many well-known optimization algorithms.

Example (ℓ_1 norm)

If $h(\mathbf{x}) = \|\mathbf{x}\|_1$, then $(\text{prox}_{\lambda h}(\mathbf{x}))_i = \psi_{st}(x_i; \lambda)$ (soft-thresholding) where

$$\psi_{st}(x; \lambda) = \begin{cases} x - \lambda & \text{if } x \geq \lambda \\ x + \lambda & \text{if } x \leq -\lambda \\ 0 & \text{else} \end{cases}$$

Why?

$$\text{prox}_{\lambda\|\cdot\|_1}(\mathbf{x}) = \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 + \lambda \|\mathbf{z}\|_1 \right\} = \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{z}\|_2^2 - \langle \mathbf{z}, \mathbf{x} \rangle + \lambda \|\mathbf{z}\|_1 \right\}$$

Note that

$$\arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{z}\|_2^2 - \langle \mathbf{z}, \mathbf{x} \rangle + \lambda \|\mathbf{z}\|_1 \right\} = \sum_i \mathcal{L}_i,$$

where

$$\mathcal{L}_i := \frac{1}{2} z_i^2 - z_i x_i + \lambda |z_i|.$$

If $x_i > 0$, then we must have $z_i \geq 0$, otherwise, let $z_i^* < 0$ minimizes \mathcal{L}_i , then $-z_i^*$ enables even smaller \mathcal{L}_i .

If $x_i < 0$, then we must have $z_i \leq 0$.

If $x_i > 0$, since $z_i \geq 0$, then we have

$$\mathcal{L}_i = -x_i z_i + \frac{1}{2} z_i^2 + \lambda z_i,$$

$$\frac{\partial \mathcal{L}}{\partial z_i} = 0 \Rightarrow -x_i + z_i + \lambda = 0 \Rightarrow z_i = x_i - \lambda.$$

Here, we require the RHS is positive (we require $z_i \geq 0$), i.e., $x_i \geq \lambda$.

If $x_i < 0$, since $z_i \leq 0$, then we have

$$\mathcal{L}_i = -x_i z_i + \frac{1}{2} z_i^2 - \lambda z_i,$$

$$\frac{\partial \mathcal{L}}{\partial z_i} = 0 \Rightarrow -x_i + z_i - \lambda = 0 \Rightarrow z_i = x_i + \lambda.$$

Here, we require the RHS is negative (we require $z_i \leq 0$), i.e., $x_i \leq -\lambda$.

Finally, let us consider the case when $-\lambda < x_i < \lambda$, our goal is

$$\arg \min \mathcal{L}_i := -x_i z_i + \frac{1}{2} z_i^2 + \lambda |z_i|.$$

1. $z_i = 0 \Rightarrow \mathcal{L}_i = 0$
2. $z_i > 0 \Rightarrow \mathcal{L}_i = -x_i z_i + \frac{1}{2} z_i^2 + \lambda z_i$ and the minimum is obtained when $z_i = 1 - \lambda$, in this case we have

$$\mathcal{L}_i = -x_i(1 - \lambda) + \frac{1}{2}(1 - \lambda)^2 + \lambda(1 - \lambda) > 0$$

3. $z_i < 0 \Rightarrow \mathcal{L}_i = -x_i z_i + \frac{1}{2} z_i^2 - \lambda z_i$ and the minimum is obtained when $z_i = 1 + \lambda$, in this case we have

$$\mathcal{L}_i = -x_i(1 + \lambda) + \frac{1}{2}(1 + \lambda)^2 + \lambda(1 + \lambda) > 0.$$

Basic rules

If $f(\mathbf{x}) = ag(\mathbf{x}) + b$ with $a > 0$, then

$$\text{prox}_f(\mathbf{x}) = \text{prox}_{ag}(\mathbf{x}).$$

Basic rules, Affine addition

If $f(\mathbf{x}) = g(\mathbf{x}) + \mathbf{a}^\top \mathbf{x} + b$, then

$$\text{prox}_f(\mathbf{x}) = \text{prox}_g(\mathbf{x} - \mathbf{a})$$

Basic rules, Quadratic addition

If $f(\mathbf{x}) = g(\mathbf{x}) + \frac{\rho}{2}\|\mathbf{x} - \mathbf{a}\|_2^2$, then

$$\text{prox}_f(\mathbf{x}) = \text{prox}_{\frac{1}{1+\rho}g}\left(\frac{1}{1+\rho}\mathbf{x} + \frac{\rho}{1+\rho}\mathbf{a}\right)$$

Proof.

$$\begin{aligned}\text{prox}_f(\mathbf{x}) &= \arg \min_{\mathbf{z}} \left\{ \frac{1}{2}\|\mathbf{z} - \mathbf{x}\|_2^2 + g(\mathbf{z}) + \frac{\rho}{2}\|\mathbf{z} - \mathbf{a}\|_2^2 \right\} \\ &= \arg \min_{\mathbf{z}} \left\{ \frac{1+\rho}{2}\|\mathbf{z}\|_2^2 - \langle \mathbf{z}, \mathbf{x} + \rho\mathbf{a} \rangle + g(\mathbf{z}) \right\} \\ &= \arg \min_{\mathbf{z}} \left\{ \frac{1}{2}\|\mathbf{z}\|_2^2 - \frac{1}{1+\rho}\langle \mathbf{z}, \mathbf{x} + \rho\mathbf{a} \rangle + \frac{1}{1+\rho}g(\mathbf{z}) \right\} \\ &= \arg \min_{\mathbf{z}} \left\{ \frac{1}{2}\|\mathbf{z} - \left(\frac{1}{1+\rho}\mathbf{x} + \frac{\rho}{1+\rho}\mathbf{a}\right)\|_2^2 + \frac{1}{1+\rho}g(\mathbf{z}) \right\} \\ &= \text{prox}_{\frac{1}{1+\rho}g}\left(\frac{1}{1+\rho}\mathbf{x} + \frac{\rho}{1+\rho}\mathbf{a}\right)\end{aligned}$$

Basic rules, Scaling and translation

If $f(\mathbf{x}) = g(a\mathbf{x} + \mathbf{b})$ with $a \neq 0$, then

$$\text{prox}_f(\mathbf{x}) = \frac{1}{a} \left(\text{prox}_{a^2 g}(a\mathbf{x} + \mathbf{b}) - \mathbf{b} \right)$$

Why?

$$\begin{aligned}
\text{prox}_f(\mathbf{x}) &= \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 + g(a\mathbf{z} + \mathbf{b}) \right\} \\
&= \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \left\| \frac{\mathbf{z}' - \mathbf{b}}{a} - \mathbf{x} \right\|_2^2 + g(\mathbf{z}') \right\} \quad (\text{Let } \mathbf{z}' = a\mathbf{z} + \mathbf{b}) \\
&= \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{z}' - (a\mathbf{x} + \mathbf{b})\|_2^2 + a^2 g(\mathbf{z}') \right\}
\end{aligned}$$

Next, consider

$$\mathbf{z}'^* = \arg \min_{\mathbf{z}'} \left\{ \frac{1}{2} \|\mathbf{z}' - (a\mathbf{x} + \mathbf{b})\|_2^2 + a^2 g(\mathbf{z}') \right\} = \text{prox}_{a^2 g}(a\mathbf{x} + \mathbf{b}).$$

Moreover, we have $\mathbf{z}^* = \frac{\mathbf{z}'^* - \mathbf{b}}{a}$, thus

$$\text{prox}_f(\mathbf{x}) = \frac{1}{a} \left(\text{prox}_{a^2 g}(a\mathbf{x} + \mathbf{b}) - \mathbf{b} \right).$$

Basic rules, Orthogonal mapping

If $f(\mathbf{x}) = g(\mathbf{Q}\mathbf{x})$ with \mathbf{Q} orthogonal ($\mathbf{Q}\mathbf{Q}^\top = \mathbf{Q}^\top\mathbf{Q} = \mathbf{I}$), then

$$\text{prox}_f(\mathbf{x}) = \mathbf{Q}^\top \text{prox}_g(\mathbf{Q}^\top \mathbf{x})$$

$$\begin{aligned} \text{prox}_f(\mathbf{x}) &= \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|^2 + f(\mathbf{z}) \right\} \\ &= \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|^2 + g(\mathbf{Q}\mathbf{z}) \right\} \\ &= \arg \min_{\mathbf{z}'} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{Q}^\top \mathbf{z}'\|_2^2 + g(\mathbf{z}') \right\} \end{aligned}$$

Let $\mathbf{z}'^* = \arg \min_{\mathbf{z}'} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{Q}^\top \mathbf{z}'\|_2^2 + g(\mathbf{z}') \right\} = \text{prox}_g(\mathbf{Q}^\top \mathbf{x})$ and we have $\mathbf{z}^* = \mathbf{Q}^\top \mathbf{z}'^*$, therefore

$$\text{prox}_f(\mathbf{x}) = \mathbf{Q}^\top \text{prox}_g(\mathbf{Q}^\top \mathbf{x})$$

Basic rules, Orthogonal affine mapping

If $f(\mathbf{x}) = g(\mathbf{Q}\mathbf{x} + \mathbf{b})$ with $\underbrace{\mathbf{Q}\mathbf{Q}^\top = \alpha^{-1}\mathbf{I}}$, then
does not require $\mathbf{Q}^\top\mathbf{Q} = \alpha^{-1}\mathbf{I}$

$$\text{prox}_f(\mathbf{x}) = (\mathbf{I} - \alpha\mathbf{Q}^\top\mathbf{Q})\mathbf{x} + \alpha\mathbf{Q}^\top(\text{prox}_{\alpha^{-1}g}(\mathbf{Q}\mathbf{x} + \mathbf{b}) - \mathbf{b})$$

Basic rules, Norm composition

If $f(\mathbf{x}) = g(\|\mathbf{x}\|_2)$ with $\text{domain}(g) = [0, \infty)$, then

$$\text{prox}_f(\mathbf{x}) = \text{prox}_g(\|\mathbf{x}\|_2) \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \quad \forall \mathbf{x} \neq 0$$

Basic rules, Norm composition – cont'd

Proof. Observe that

$$\begin{aligned}\min_{\mathbf{z}} \left\{ f(\mathbf{z}) + \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 \right\} &= \min_{\mathbf{z}} \left\{ g(\|\mathbf{z}\|_2) + \frac{1}{2} \|\mathbf{z}\|_2^2 - \mathbf{z}^\top \mathbf{x} + \frac{1}{2} \|\mathbf{x}\|_2^2 \right\} \\ &= \min_{\alpha \geq 0} \min_{\|\mathbf{z}\|_2 = \alpha} \left\{ g(\alpha) + \frac{1}{2} \alpha^2 - \mathbf{z}^\top \mathbf{x} + \frac{1}{2} \|\mathbf{x}\|_2^2 \right\} \\ &\stackrel{\text{Cauchy-Schwarz}}{=} \min_{\alpha \geq 0} \left\{ g(\alpha) + \frac{1}{2} \alpha^2 - \alpha \|\mathbf{x}\|_2 + \frac{1}{2} \|\mathbf{x}\|_2^2 \right\} \\ &= \min_{\alpha \geq 0} \left\{ g(\alpha) + \frac{1}{2} (\alpha - \|\mathbf{x}\|_2)^2 \right\}\end{aligned}$$

From the above calculation, we know the optimal point is

$$\alpha^* = \text{prox}_g(\|\mathbf{x}\|_2) \quad \text{and} \quad \mathbf{z}^* = \alpha^* \frac{\mathbf{x}}{\|\mathbf{x}\|_2} = \text{prox}_g(\|\mathbf{x}\|_2) \frac{\mathbf{x}}{\|\mathbf{x}\|_2},$$

thus concluding proof.

Convergence analysis

Lemma 5. [Cost monotonicity] Suppose f is convex and L -smooth. If $\eta_t \equiv 1/L$, then

$$F(\mathbf{x}^{t+1}) \leq F(\mathbf{x}^t).$$

Remark. This is different from subgradient methods (for which objective value might be non-monotonic in t). Also, constant stepsize rule is recommended when f is convex and smooth.

Fundamental Inequality

Lemma 6. (key lemma) Let $\mathbf{y}^+ = \text{prox}_{\frac{1}{L}h}\left(\mathbf{y} - \frac{1}{L}\nabla f(\mathbf{y})\right)$, then

$$F(\mathbf{y}^+) - F(\mathbf{x}) \leq \frac{L}{2}\|\mathbf{x} - \mathbf{y}\|_2^2 - \frac{L}{2}\|\mathbf{x} - \mathbf{y}^+\|_2^2 - \underbrace{g(\mathbf{x}, \mathbf{y})}_{\geq 0 \text{ by convexity}}$$

where $g(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$.

Take $\mathbf{x} = \mathbf{y} = \mathbf{x}^t$ and hence $\mathbf{y}^+ = \mathbf{x}^{t+1}$ to complete the proof of Lemma 5.

Proof of Lemma 6. Define $\phi(\mathbf{z}) = f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 + h(\mathbf{z})$. It is easily seen that $\mathbf{y}^+ = \arg \min_{\mathbf{z}} \phi(\mathbf{z})$. Two important properties:

1. Since $\phi(\mathbf{z})$ is L -strongly convex, one has

$$\phi(\mathbf{x}) \geq \phi(\mathbf{y}^+) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}^+\|_2^2.$$

2. From smoothness,

$$\phi(\mathbf{y}^+) = \underbrace{f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{y}^+ - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{y}^+ - \mathbf{y}\|_2^2}_{\text{upper bound on } f(\mathbf{y}^+) \text{ (L-smoothness)}} + h(\mathbf{y}^+) \geq f(\mathbf{y}^+) + h(\mathbf{y}^+) = F(\mathbf{y}^+).$$

Proof of Lemma 6 (cont'd). Taken collectively, these yield

$$\phi(\mathbf{x}) \geq F(\mathbf{y}^+) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}^+\|_2^2,$$

which together with the definition of $\phi(\mathbf{x})$ gives

$$\underbrace{f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + h(\mathbf{x})}_{=f(\mathbf{x})+h(\mathbf{x})-g(\mathbf{x},\mathbf{y})=F(\mathbf{x})-g(\mathbf{x},\mathbf{y})} + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \geq F(\mathbf{y}^+) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}^+\|_2^2$$

which finishes the proof.

Monotonicity in estimation error

Lemma 7. Suppose f is convex and L -smooth. If $\eta_t \equiv 1/L$, then

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 \leq \|\mathbf{x}^t - \mathbf{x}^*\|_2.$$

Proof. From Lemma 6, taking $\mathbf{x} = \mathbf{x}^*$, $\mathbf{y} = \mathbf{x}^t$ (and hence $\mathbf{y}^+ = \mathbf{x}^{t+1}$) yields

$$\underbrace{F(\mathbf{x}^{t+1}) - F(\mathbf{x}^*)}_{\geq 0} + \underbrace{g(\mathbf{x}, \mathbf{y})}_{\geq 0} \leq \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^t\|_2^2 - \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^{t+1}\|_2^2$$

which immediately concludes the proof.

Remark. Proximal gradient iterates are not only monotonic w.r.t. cost, but also monotonic in estimation error.

Convergence for convex problems

Theorem. [Convergence of proximal gradient methods for convex problems] Suppose f is convex and L -smooth. If $\eta_t \equiv 1/L$, then

$$F(\mathbf{x}^t) - F^{opt} \leq \frac{L \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{2t}.$$

Convergence for convex problems

Proof. With Lemma 6 in mind, set $\mathbf{x} = \mathbf{x}^*$, $\mathbf{y} = \mathbf{x}^t$ to obtain

$$\begin{aligned} F(\mathbf{x}^{t+1}) - F(\mathbf{x}^*) &\leq \frac{L}{2} \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \frac{L}{2} \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 - \underbrace{g(\mathbf{x}^*, \mathbf{x}^t)}_{\geq 0 \text{ by convexity}} \\ &\leq \frac{L}{2} \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \frac{L}{2} \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \end{aligned}$$

Apply it recursively and add up all inequalities to get

$$\sum_{k=0}^{t-1} \left(F(\mathbf{x}^{k+1}) - F(\mathbf{x}^*) \right) \leq \frac{L}{2} \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 - \frac{L}{2} \|\mathbf{x}^t - \mathbf{x}^*\|_2^2.$$

This combines with monotonicity of $F(\mathbf{x}^t)$ (cf. Lemma 6) yields

$$F(\mathbf{x}^t) - F(\mathbf{x}^*) \leq \frac{\frac{L}{2} \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{t}.$$

Convergence for convex problems

Theorem. [Convergence of proximal gradient methods for strongly convex problems]
Suppose f is μ -strongly convex and L -smooth. If $\eta_t \equiv 1/L$, then

$$\|\mathbf{x}^t - \mathbf{x}^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^t \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2.$$

Convergence for convex problems

Proof. Taking $\mathbf{x} = \mathbf{x}^*$, $\mathbf{y} = \mathbf{x}^t$ (and hence $\mathbf{y}^+ = \mathbf{x}^{t+1}$) in Lemma 6 gives

$$\begin{aligned} F(\mathbf{x}^{t+1}) - F(\mathbf{x}^*) &\leq \frac{1}{L} \|\mathbf{x}^* - \mathbf{x}^t\|_2^2 - \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^{t+1}\|_2^2 - \underbrace{g(\mathbf{x}^*, \mathbf{x}^t)}_{\geq \frac{\mu}{2} \|\mathbf{x}^* - \mathbf{x}^t\|_2^2} \\ &\leq \frac{L - \mu}{2} \|\mathbf{x}^t - \mathbf{x}^*\|_2^2 - \frac{L}{2} \|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2. \end{aligned}$$

This taken collectively with $F(\mathbf{x}^{t+1}) - F(\mathbf{x}^*) \geq 0$ yields

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}^t - \mathbf{x}^*\|_2^2.$$

Applying it recursively concludes the proof.