



EECS 559 Optimization Methods for SIPML

Lecture 7 – Smoothing & Proximal Methods

Instructor: Prof. Qing Qu (qingqu@umich.edu)

Lecture Agenda

- Smoothing for Nonsmooth Problems
- Moreau Envelope and Proximal Operator
- Proximal Gradient Methods

Lecture Agenda

- Smoothing for Nonsmooth Problems
- Moreau Envelope and Proximal Operator
- Proximal Gradient Methods

Nonsmooth Optimization

$$\min_{x \in \mathbb{R}^n} f(x)$$

Can the convergence gap be closed?

- **For convex and smooth problems,** (accelerated) gradient descent (GD) converges to ε -accuracy in

$$O\left(\frac{1}{\sqrt{\varepsilon}}\right) \text{iterations, or } O\left(\frac{1}{k^2}\right) \text{convergence rate.}$$

- **For convex but non-smooth problems,** subgradient method converges to ε -accuracy in

$$O\left(\frac{1}{\varepsilon^2}\right) \text{iterations, or } O\left(\frac{1}{\sqrt{k}}\right) \text{convergence rate.}$$

Smooth Approximation

Definition. A convex function f is called (α, β) -smoothable, if for any $\mu > 0$, there exists a convex function f_μ , such that

- $f_\mu(\mathbf{x}) \leq f(\mathbf{x}) \leq f_\mu(\mathbf{x}) + \beta\mu, \quad \forall \mathbf{x}$
- f_μ is $\frac{\alpha}{\mu}$ -smooth, with gradient:
$$\|\nabla f_\mu(\mathbf{x}) - \nabla f_\mu(\mathbf{x}')\|_2 \leq \frac{\alpha}{\mu} \|\mathbf{x} - \mathbf{x}'\|_2, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$$

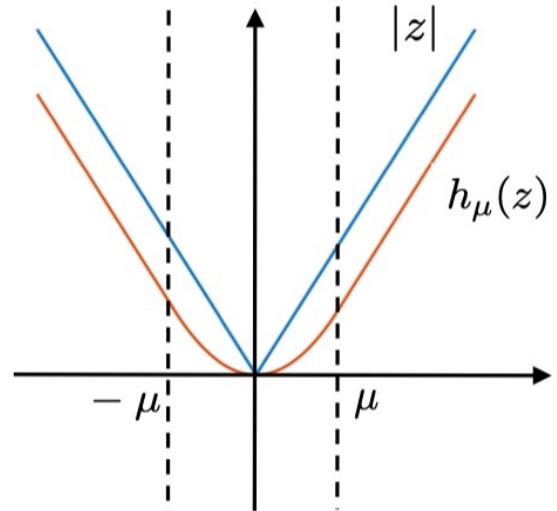
f_μ is called a $\frac{1}{\mu}$ -smooth approximation of f with parameters (α, β) .
 μ : tradeoff between approximation accuracy ϵ and smoothness

Example: ℓ_1 -Norm

Consider the Huber function

$$f_\mu(\mathbf{x}) := \sum_{i=1}^n h_\mu(x_i),$$

$$h_\mu(x) = \begin{cases} |x| - \frac{\mu}{2} & |x| \geq \mu \\ \frac{x^2}{2\mu} & |x| < \mu \end{cases}$$



which satisfies

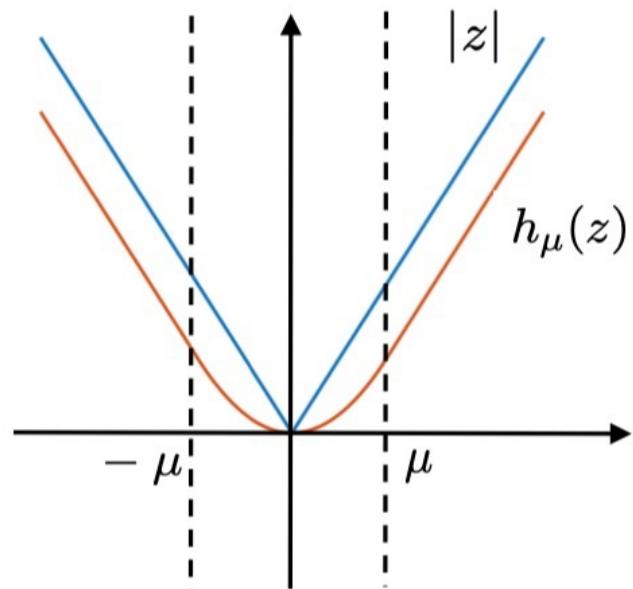
$$h_\mu(x) \leq |x| \leq h_\mu(x) + \mu/2, \text{ and } h_\mu(x) \text{ is } \frac{1}{\mu}\text{-smooth.}$$

Example: ℓ_1 -Norm

Therefore, $f_\mu(\mathbf{x}) = \sum_{i=1}^n h_\mu(x_i)$ is $\frac{1}{\mu}$ -smooth,
and obeys

$$f_\mu(\mathbf{x}) \leq \|\mathbf{x}\|_1 \leq f_\mu(\mathbf{x}) + \frac{n\mu}{2}$$

so that $\|\mathbf{x}\|_1$ is $(1, n/2)$ -smoothable.

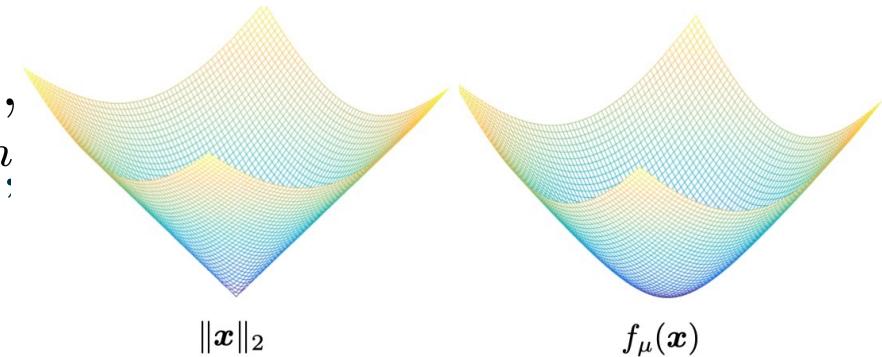


Example: ℓ_2 -Norm

Consider $f_\mu(\mathbf{x}) := \sqrt{\|\mathbf{x}\|_2^2 + \mu^2} - \mu$,
then for any $\mu > 0$ and any $\mathbf{x} \in \mathbb{R}^n$,

$$f_\mu(\mathbf{x}) \leq (\|\mathbf{x}\|_2 + \mu) - \mu = \|\mathbf{x}\|_2$$

$$\|\mathbf{x}\|_2 \leq \sqrt{\|\mathbf{x}\|_2^2 + \mu^2} = f_\mu(\mathbf{x}) + \mu$$

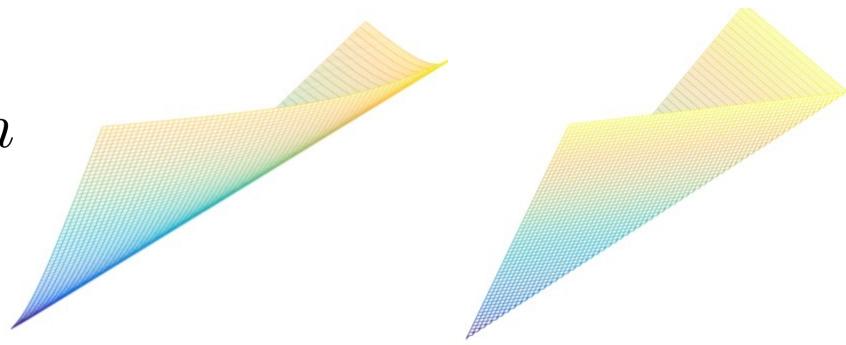


Additionally, we can show that $f_\mu(\mathbf{x})$ is $\frac{1}{\mu}$ -smooth.
Therefore, $\|\mathbf{x}\|_2$ is $(1,1)$ -smoothable.

Example: Max Function

$$f_\mu(\mathbf{x}) := \mu \log \left(\sum_{i=1}^n e^{x_i/\mu} \right) - \mu \log n$$

For any $\mu > 0$ and any $\mathbf{x} \in \mathbb{R}^n$,



$$f_\mu(\mathbf{x}) \leq f(\mathbf{x}) = \max_i x_i \leq f_\mu(\mathbf{x}) + \mu \log n.$$

Additionally, we can show that $f_\mu(\mathbf{x})$ is $\frac{1}{\mu}$ -smooth.
Therefore, $f(\mathbf{x}) = \max_{1 \leq i \leq n} x_i$ is $(1, \log n)$ -smoothable.

Basic Composition Rules

- **Addition:**

Let $f_{\mu,1}$ be a $\frac{1}{\mu}$ -smooth approximation of f_1 with (α_1, β_1) , and let $f_{\mu,2}$ be a $\frac{1}{\mu}$ -smooth approximation of f_2 with (α_2, β_2) . Then

$\lambda_1 f_{\mu,1} + \lambda_2 f_{\mu,2}$ ($\lambda_1, \lambda_2 > 0$) is a $\frac{1}{\mu}$ -smooth approximation of $\lambda_1 f_1 + \lambda_2 f_2$ with parameters $(\lambda_1 \alpha_1 + \lambda_2 \alpha_2, \lambda_1 \beta_1 + \lambda_2 \beta_2)$.

Basic Composition Rules

- **Affine Transformation:**

h_μ is a $\frac{1}{\mu}$ -smooth approximation of h with parameters (α, β) and let $f(\mathbf{x}) := h(\mathbf{A}\mathbf{x} + \mathbf{b})$. Then $h_\mu(\mathbf{A}\mathbf{x} + \mathbf{b})$ is a $\frac{1}{\mu}$ -smooth approximation of f with parameters $(\alpha \|\mathbf{A}\|^2, \beta)$.

Example: $f(x) = |x|$

Claim. $f_\mu(x) := \mu \log(e^{x/\mu} + e^{-x/\mu}) - \mu \log 2$ is a $\frac{1}{\mu}$ -smooth approximation of $f(x) = |x|$ with $(2, \log 2)$.

Example: $f(\mathbf{x}) = \|\mathbf{Ax} + \mathbf{b}\|_2$

Claim. $f_\mu(\mathbf{x}) = \sqrt{\|\mathbf{Ax} + \mathbf{b}\|_2^2 + \mu^2} - \mu$ is a $\frac{1}{\mu}$ -smooth approximation of $\|\mathbf{Ax} + \mathbf{b}\|_2$ with $(\|\mathbf{A}\|_2^2, 1)$.

Lecture Agenda

- Smoothing for Nonsmooth Problems
- Moreau Envelope and Proximal Operator
- Proximal Gradient Methods

Minimizing a Nonsmooth Function?

Suppose we want to minimize

$$\min_{\substack{x \\ \text{nonsmooth}}} f(x) \quad , \quad \text{s.t.} \quad x \in \mathbb{R}^n.$$

- We know that subgradient method converges slowly;
- Can we improve the convergence by minimizing a smooth surrogate of f ?

Smoothing via the Moreau Envelope

Definition. The Moreau envelope of a convex function f with parameter $\mu > 0$ is defined as

$$M_{\mu f}(\mathbf{x}) := \inf_{\mathbf{z}} \left\{ f(\mathbf{z}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{z}\|_2^2 \right\}$$

- $M_{\mu f}$ is a smoothed approximation of f ;
- Minimizers of $f =$ minimizers of $M_{\mu f}$.

Smoothing via the Moreau Envelope

$$M_{\mu f}(\mathbf{x}) := \inf_{\mathbf{z}} \left\{ f(\mathbf{z}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{z}\|_2^2 \right\}$$

Show the following (homework):

- $M_{\mu f}$ is convex;
- $M_{\mu f}$ is $\frac{1}{\mu}$ -smooth;
- If f is L -Lipschitz, then $M_{\mu f}$ is a $\frac{1}{\mu}$ -smooth approximation of f with parameters $(1, L^2/2)$.

Smoothing via the Moreau Envelope

$$M_{\mu f}(\mathbf{x}) := \inf_{\mathbf{z}} \left\{ f(\mathbf{z}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{z}\|_2^2 \right\}$$

Therefore, we can minimize f by minimizing $M_{\mu f}(\mathbf{x})$

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \tau_k \nabla M_{\mu f}(\mathbf{x}_k) \\ &= \text{prox}_{\mu f}(\mathbf{x}_k), \quad \text{if } \tau_k = \mu \end{aligned}$$

Proximal Operator & Connections

Definition. The proximal operator for a convex function f is defined as

$$\text{prox}_{\mu f}(\mathbf{x}) = \arg \inf_{\mathbf{z}} \left(f(\mathbf{z}) + \frac{1}{2\mu} \|\mathbf{z} - \mathbf{x}\|_2^2 \right).$$

Therefore, we have

$$M_{\mu f}(\mathbf{x}) = f(\text{prox}_{\mu f}(\mathbf{x})) + \frac{1}{2\mu} \|\mathbf{x} - \text{prox}_{\mu f}(\mathbf{x})\|_2^2$$

Proximal Operator

Definition. The proximal operator for a convex function f at the point \mathbf{x} is defined as

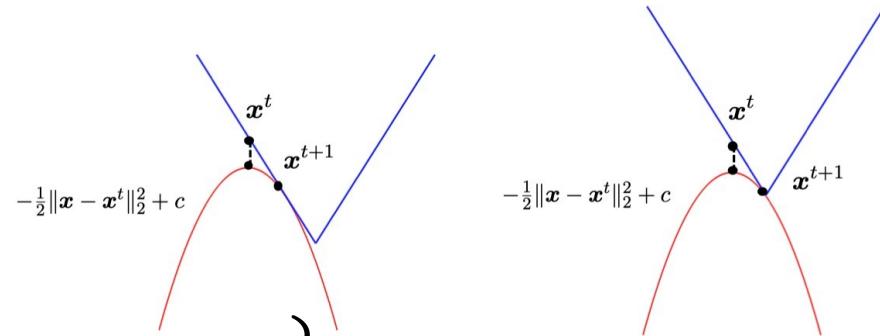
$$\text{prox}_{\mu f}(\mathbf{x}) = \arg \inf_{\mathbf{z}} \left(f(\mathbf{z}) + \frac{1}{2\mu} \|\mathbf{z} - \mathbf{x}\|_2^2 \right).$$

- well-defined under very general conditions;
- can be evaluated efficiently for many widely used functions;
- the abstraction is conceptually & mathematically simple, covers many well-known optimization methods.

Proximal Operator: ℓ_1 -Norm

If $f(\mathbf{x}) = \|\mathbf{x}\|_1$ with

$$\text{prox}_{\mu f}(\mathbf{x}) = \arg \min_{\mathbf{z}} \left\{ \frac{1}{2\mu} \|\mathbf{z} - \mathbf{x}\|_2^2 + \|\mathbf{z}\|_1 \right\}.$$



Then the proximal operator is the *soft-thresholding operator*

$$\text{prox}_{\mu f}(\mathbf{x}) = S_\mu(\mathbf{x}), \quad [S_\mu(\mathbf{x})]_i = \begin{cases} x_i - \mu & \text{if } x_i > \mu, \\ 0 & \text{if } |x_i| \leq \mu, \\ x_i + \mu & \text{if } x_i < -\mu. \end{cases}$$

Proximal Operator: Nuclear Norm

If $f(\mathbf{X}) = \|\mathbf{X}\|_*$ with $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$, and

$$\text{prox}_{\mu f}(\mathbf{X}) = \arg \min_{\mathbf{Z}} \left\{ \frac{1}{2\mu} \|\mathbf{Z} - \mathbf{X}\|_F^2 + \|\mathbf{Z}\|_* \right\}.$$

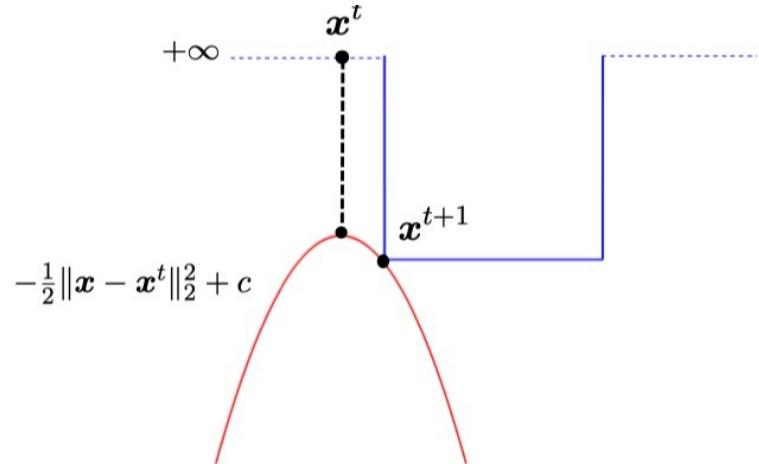
Then the proximal operator is the *singular value thresholding operator*

$$\text{prox}_{\mu f}(\mathbf{X}) = \mathbf{U}\mathbf{S}_\mu(\Sigma)\mathbf{V}^\top.$$

Proximal Operator: Indicator Function

If $f(\mathbf{x}) = \mathbb{1}_{\mathbf{x} \in \mathcal{C}}$ with

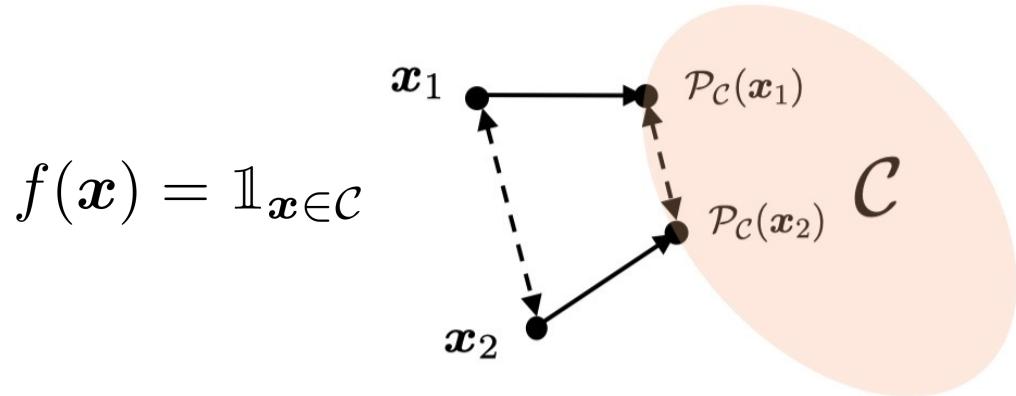
$$f(\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{x} \in \mathcal{C} \\ \infty & \text{else.} \end{cases}$$



The proximal operator is the *Euclidean projection*:

$$\text{prox}_{\mu f}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2$$

Nonexpansiveness of Proximal Operators

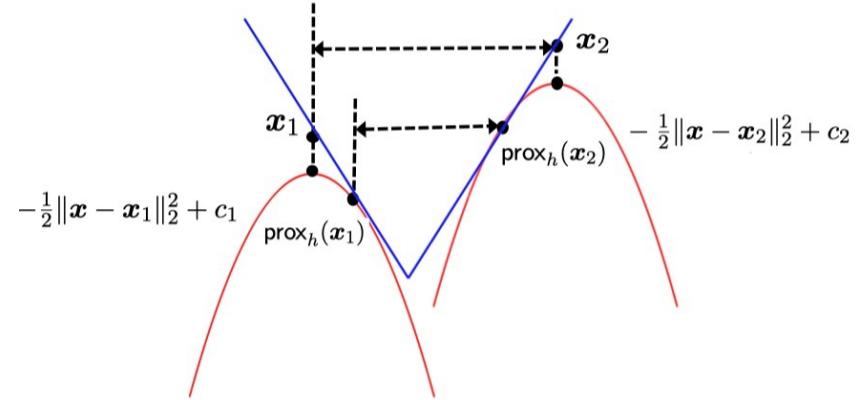


When $\text{prox}_{\mu f}(\mathbf{x})$ is the Euclidean projection onto \mathcal{C} , it is *nonexpansive* for a *convex* \mathcal{C} :

$$\|\mathcal{P}_{\mathcal{C}}(\mathbf{x}_1) - \mathcal{P}_{\mathcal{C}}(\mathbf{x}_2)\|_2 \leq \|\mathbf{x}_1 - \mathbf{x}_2\|_2.$$

Nonexpansiveness of Proximal Operators

Proximal operator behaves like the projection operator.



- **Nonexpansiveness** (homework):

$$\|\text{prox}_{\mu f}(\mathbf{x}_1) - \text{prox}_{\mu f}(\mathbf{x}_2)\|_2 \leq \|\mathbf{x}_1 - \mathbf{x}_2\|_2.$$

- **Firm nonexpansiveness** (homework):

$$\langle \text{prox}_{\mu f}(\mathbf{x}_1) - \text{prox}_{\mu f}(\mathbf{x}_2), \mathbf{x}_1 - \mathbf{x}_2 \rangle \geq \|\text{prox}_{\mu f}(\mathbf{x}_1) - \text{prox}_{\mu f}(\mathbf{x}_2)\|_2^2$$

Proximal Operator & Moreau Envelope

Theorem. (Gradient of Moreau envelope)

The Moreau envelope

$$M_{\mu f}(\mathbf{x}) = f(\text{prox}_{\mu f}(\mathbf{x})) + \frac{1}{2\mu} \|\mathbf{x} - \text{prox}_{\mu f}(\mathbf{x})\|_2^2$$

is continuously differentiable with gradient

$$\nabla M_{\mu f}(\mathbf{x}) = \frac{1}{\mu} (\mathbf{x} - \text{prox}_{\mu f}(\mathbf{x})).$$

- Proximal operator is a gradient step for minimizing $M_{\mu f}$:

$$\text{prox}_{\mu f}(\mathbf{x}) = \mathbf{x} - \mu \cdot \nabla M_{\mu f}(\mathbf{x}).$$

Proximal Operator & Connections

Proof. First of all, let us introduce the conjugate of a given function f , which is defined by

$$f^*(\mathbf{x}) := \sup_{\mathbf{z} \in \text{dom } f} \mathbf{z}^\top \mathbf{x} - f(\mathbf{z})$$

and we can show that

$$\nabla f^*(\mathbf{x}) = \arg \max_{\mathbf{z}} \{ \mathbf{z}^\top \mathbf{x} - f(\mathbf{z}) \}$$

Proximal Operator & Connections

For our case, we can rearrange terms in $M_{\mu f}(\mathbf{x})$ as

$$\begin{aligned} M_{\mu f}(\mathbf{x}) &= \frac{1}{2\mu} \|\mathbf{x}\|_2^2 - \frac{1}{\mu} \sup_{\mathbf{z}} \left\{ \mathbf{z}^\top \mathbf{x} - \mu f(\mathbf{z}) - \frac{1}{2} \|\mathbf{z}\|_2^2 \right\} \\ &= \frac{1}{2\mu} \|\mathbf{x}\|_2^2 - \frac{1}{\mu} \left(\mu f + \frac{1}{2} \|\cdot\|_2^2 \right)^*(\mathbf{x}) \end{aligned}$$

Therefore, we have

$$\nabla M_{\mu f}(\mathbf{x}) = \frac{\mathbf{x}}{\mu} - \frac{1}{\mu} \arg \max_{\mathbf{z}} \left\{ \mathbf{z}^\top \mathbf{x} - \mu f(\mathbf{z}) - \frac{1}{2} \|\mathbf{z}\|_2^2 \right\} = \frac{1}{\mu} (\mathbf{x} - \text{prox}_{\mu f}(\mathbf{x}))$$

which implies that $\text{prox}_{\mu f}(\mathbf{x}) = \mathbf{x} - \mu \cdot \nabla M_{\mu f}(\mathbf{x})$.

Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma. (**Chapter 8.2**)
- *Proximal Algorithms, Foundations & Trends in Optimization.* Neal Parikh, Stephen Boyd, 2014.
https://web.stanford.edu/~boyd/papers/prox_algs.html

Lecture Agenda

- Smoothing for Nonsmooth Problems
- Moreau Envelope and Proximal Operator
- **Proximal Gradient Methods**

Nonsmooth Composite Problems

$$\min_{\boldsymbol{x}} F(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})$$

- f : convex, continuously differentiable, and L -smooth,

$$\|\nabla f(\boldsymbol{x}') - \nabla f(\boldsymbol{x})\|_2 \leq L \|\boldsymbol{x}' - \boldsymbol{x}\|_2, \quad \forall \boldsymbol{x}, \boldsymbol{x}' \in \mathbb{R}^n.$$

- g : convex but may *not* be differentiable.

Example I: ℓ_1 -Regularized Problems

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}) + \underbrace{\lambda \cdot \|\boldsymbol{x}\|_1}_{g(\boldsymbol{x})}$$

- **Lasso problem:** $f(\boldsymbol{x}) = \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2$
- **Sparse logistic regression:**

$$f(\boldsymbol{x}) = \sum_{i=1}^m \left(-y_i \boldsymbol{x}^\top \boldsymbol{z}_i + \log \left(1 + e^{y_i \boldsymbol{x}^\top \boldsymbol{z}_i} \right) \right)$$

Example II: Nuclear Norm Regularization

- Stable low-rank matrix recovery:

$$\min_{\mathbf{X}} F(\mathbf{X}) = \underbrace{\frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_F^2}_{f(\mathbf{X})} + \underbrace{\lambda \|\mathbf{X}\|_*}_{g(\mathbf{X})}.$$

- Stable principal component pursuit (PCP):

$$\min_{\mathbf{L}, \mathbf{S}} F(\mathbf{L}, \mathbf{S}) = \underbrace{\frac{1}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2}_{f(\mathbf{L}, \mathbf{S})} + \underbrace{\lambda \|\mathbf{L}\|_* + \mu \|\mathbf{S}\|_1}_{g(\mathbf{L}, \mathbf{S})}.$$

Closing the Gap?

$$\min_{\mathbf{x}} F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$$

- **For convex and smooth problems,** (accelerated) gradient descent (GD) converges to ε -accuracy in

$$O\left(\frac{1}{\sqrt{\varepsilon}}\right) \text{iterations, or } O\left(\frac{1}{k^2}\right) \text{convergence rate.}$$

- **For convex but non-smooth problems,** subgradient method converges to ε -accuracy in

$$O\left(\frac{1}{\varepsilon^2}\right) \text{iterations, or } O\left(\frac{1}{\sqrt{k}}\right) \text{convergence rate.}$$

Taylor Approximation Theory I

Theorem (first-order approximation).

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be continuously differentiable, and $\nabla f(\mathbf{x})$ is $\gamma_L(\mathbf{x}_0)$ -Lipschitz continuous at \mathbf{x}_0 , then

$$\left| f(\mathbf{x}) - \hat{f}_L(\mathbf{x}; \mathbf{x}_0) \right| \leq \frac{\gamma_L(\mathbf{x}_0)}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2,$$

where we define

$$\hat{f}_L(\mathbf{x}; \mathbf{x}_0) := f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle.$$

Proximal Gradient Method

Consider optimizing the following *quadratic upper bound*

$$\begin{aligned}\widehat{F}_\mu(\boldsymbol{x}, \boldsymbol{x}_k) &= \underbrace{f(\boldsymbol{x}_k) + \langle \nabla f(\boldsymbol{x}_k), \boldsymbol{x} - \boldsymbol{x}_k \rangle + g(\boldsymbol{x})}_{=: \ell_F(\boldsymbol{x}, \boldsymbol{x}_k)} + \frac{1}{2\mu} \|\boldsymbol{x} - \boldsymbol{x}_k\|_2^2\end{aligned}$$

$$\boldsymbol{x}_{k+1} = \arg \min_{\boldsymbol{x}} \widehat{F}_\mu(\boldsymbol{x}, \boldsymbol{x}_k).$$

Proximal Gradient Method

$$\begin{aligned}\widehat{F}_\mu(\mathbf{x}, \mathbf{x}_k) &= g(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{x} - (\mathbf{x}_k - \mu \nabla f(\mathbf{x}_k))\|_2^2 \\ &\quad + f(\mathbf{x}_k) - \frac{\mu}{2} \|\nabla f(\mathbf{x}_k)\|_2^2\end{aligned}$$

Let $\mathbf{w}_k = \mathbf{x}_k - \mu \nabla f(\mathbf{x}_k)$, then

$$\begin{aligned}\mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \left\{ g(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{w}_k\|_2^2 \right\} \\ &= \text{prox}_{\mu g}(\mathbf{w}_k).\end{aligned}$$

Proximal Gradient Method

Proximal Gradient (PG)

Problem Class:

$$\min_{\mathbf{x}} F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$$

$f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ are convex, ∇f L -Lipschitz and g (maybe) non-smooth.

Basic Iteration: set $\mathbf{x}_0 \in \mathbb{R}^n$.

Repeat:

$$\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k),$$

$$\mathbf{x}_{k+1} \leftarrow \text{prox}_{g/L}[\mathbf{w}_k].$$

Convergence Guarantee:

$F(\mathbf{x}_k) - F(\mathbf{x}_*)$ converges at a rate of $O(1/k)$.

Proximal Gradient Method

In summary, one step of PG is

$$\boldsymbol{x}_{k+1} = \text{prox}_{\mu g}(\boldsymbol{x}_k - \mu \nabla f(\boldsymbol{x}_k))$$

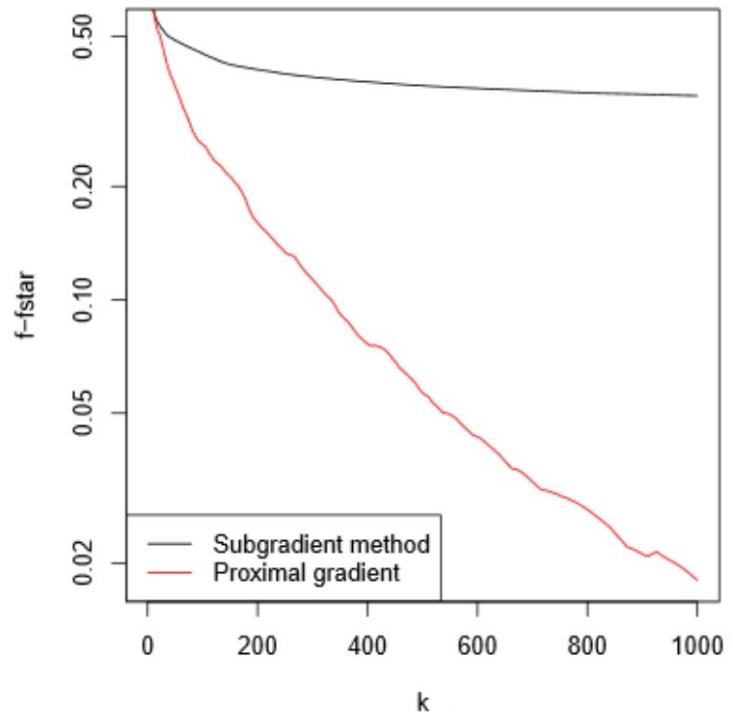
where $\mu \in (0, 1/L]$.

Let $G_\mu(\boldsymbol{x}) := \frac{1}{\mu} (\boldsymbol{x} - \text{prox}_{\mu g}(\boldsymbol{x} - \mu \nabla f(\boldsymbol{x})))$, then

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \mu \cdot \underbrace{G_\mu(\boldsymbol{x}_k)}_{\text{proximal gradient}}.$$

Comparison with Subgradient

$$\min_{\mathbf{x}} F(\mathbf{x}) = \underbrace{\frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2}_{f(\mathbf{x})} + \underbrace{\lambda \|\mathbf{x}\|_1}_{g(\mathbf{x})}$$



Convergence of Proximal Gradient

Theorem. Let $F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$, with

- f is convex, differentiable, and L -smooth;
- g is convex, but possibly nonsmooth.

Consider the sequence of proximal gradient iterates

$$\mathbf{x}_{k+1} = \text{prox}_{\mu_k g} (\mathbf{x}_k - \mu_k \nabla f(\mathbf{x}_k)), \quad \mu_k \equiv 1/L.$$

Let \mathbf{x}_* be a minimizer of F . Then for any $k \geq 1$,

$$F(\mathbf{x}_k) - F(\mathbf{x}_*) \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{2k}.$$

Convergence of Proximal Gradient

Proof. By convexity of f and g , we have

$$f(\mathbf{x}) \geq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle$$

$$g(\mathbf{x}) \geq g(\mathbf{x}_{k+1}) + \langle \gamma, \mathbf{x} - \mathbf{x}_{k+1} \rangle \quad \gamma \in \partial g$$

Thus, we obtain

$$\begin{aligned} F(\mathbf{x}) &= f(\mathbf{x}) + g(\mathbf{x}) \\ &\geq f(\mathbf{x}_k) + g(\mathbf{x}_{k+1}) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \langle \gamma, \mathbf{x} - \mathbf{x}_{k+1} \rangle \end{aligned}$$

On the other hand, for $\mu \in (0, 1/L]$ we have

$$\widehat{F}_\mu(\mathbf{x}, \mathbf{x}_k) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + g(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \geq F(\mathbf{x})$$

→ $\widehat{F}_\mu(\mathbf{x}_{k+1}, \mathbf{x}_k) \geq F(\mathbf{x}_{k+1})$

Convergence of Proximal Gradient

$$\begin{aligned} F(\mathbf{x}) - F(\mathbf{x}_{k+1}) &\geq F(\mathbf{x}) - \widehat{F}_\mu(\mathbf{x}_{k+1}, \mathbf{x}_k) \\ &\geq f(\mathbf{x}_k) + g(\mathbf{x}_{k+1}) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \langle \boldsymbol{\gamma}, \mathbf{x} - \mathbf{x}_{k+1} \rangle \\ &\quad - \left(f(\mathbf{x}_k) + g(\mathbf{x}_{k+1}) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \right) \\ &= \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \langle \boldsymbol{\gamma}, \mathbf{x} - \mathbf{x}_{k+1} \rangle - \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle \\ &\quad - \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \\ &= \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_{k+1} \rangle + \langle \boldsymbol{\gamma}, \mathbf{x} - \mathbf{x}_{k+1} \rangle - \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \\ &= \langle \nabla f(\mathbf{x}_k) + \boldsymbol{\gamma}, \mathbf{x} - \mathbf{x}_{k+1} \rangle - \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \end{aligned}$$

Convergence of Proximal Gradient

On the other hand, we have the following

$$\boldsymbol{x}_{k+1} = \min_{\boldsymbol{x}} \left\{ g(\boldsymbol{x}) + \frac{L}{2} \left\| \boldsymbol{x} - \left(\boldsymbol{x}_k - \frac{1}{L} \nabla f(\boldsymbol{x}_k) \right) \right\|_2^2 \right\} = \min_{\boldsymbol{x}} h(\boldsymbol{x})$$

Thus, we have

$$\mathbf{0} \in \partial h(\boldsymbol{x}_{k+1}) \implies \boldsymbol{\gamma} + \nabla f(\boldsymbol{x}_k) + L(\boldsymbol{x}_{k+1} - \boldsymbol{x}_k) = \mathbf{0}, \quad \boldsymbol{\gamma} \in \partial g$$

Therefore, we have

$$\begin{aligned} F(\boldsymbol{x}) - F(\boldsymbol{x}_{k+1}) &\geq \langle \nabla f(\boldsymbol{x}_k) + \boldsymbol{\gamma}, \boldsymbol{x} - \boldsymbol{x}_{k+1} \rangle - \frac{L}{2} \|\boldsymbol{x}_{k+1} - \boldsymbol{x}_k\|_2^2 \\ &\geq -\frac{L}{2} \|\boldsymbol{x}_{k+1} - \boldsymbol{x}_k\|_2^2 + \langle \nabla f(\boldsymbol{x}_k) + \boldsymbol{\gamma}, \boldsymbol{x} - \boldsymbol{x}_{k+1} \rangle \end{aligned}$$

Convergence of Proximal Gradient

Therefore, we have

$$\begin{aligned} F(\mathbf{x}) - F(\mathbf{x}_{k+1}) &\geq \langle \nabla f(\mathbf{x}_k) + \gamma, \mathbf{x} - \mathbf{x}_{k+1} \rangle - \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 \\ &= -\frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 + L \langle \mathbf{x}_k - \mathbf{x}_{k+1}, \mathbf{x} - \mathbf{x}_{k+1} \rangle - \frac{L}{2} \|\mathbf{x} - \mathbf{x}_{k+1}\|_2^2 \\ &\quad + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_{k+1}\|_2^2 \\ &= \frac{L}{2} \|\mathbf{x} - \mathbf{x}_{k+1}\|_2^2 - \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 \end{aligned}$$

If $\mathbf{x} = \mathbf{x}_\star$, we have

$$0 \geq F(\mathbf{x}_\star) - F(\mathbf{x}_{k+1}) \geq \frac{L}{2} \|\mathbf{x}_\star - \mathbf{x}_{k+1}\|_2^2 - \frac{L}{2} \|\mathbf{x}_\star - \mathbf{x}_k\|_2^2$$

Convergence of Proximal Gradient

If $\mathbf{x} = \mathbf{x}_\star$, we have

$$0 \geq F(\mathbf{x}_\star) - F(\mathbf{x}_{k+1}) \geq \frac{L}{2} \|\mathbf{x}_\star - \mathbf{x}_{k+1}\|_2^2 - \frac{L}{2} \|\mathbf{x}_\star - \mathbf{x}_k\|_2^2$$

→ $\|\mathbf{x}_\star - \mathbf{x}_{k+1}\|_2^2 \leq \|\mathbf{x}_\star - \mathbf{x}_k\|_2^2$

Finally, we obtain

$$\begin{aligned} k(F(\mathbf{x}_k) - F(\mathbf{x}_\star)) &\leq \sum_{i=0}^{k-1} (F(\mathbf{x}_i) - F(\mathbf{x}_\star)) \leq \frac{L}{2} \sum_{i=0}^k \left(\|\mathbf{x}_\star - \mathbf{x}_{i+1}\|_2^2 - \|\mathbf{x}_\star - \mathbf{x}_i\|_2^2 \right) \\ &\leq \frac{L}{2} \|\mathbf{x}_\star - \mathbf{x}_0\|_2^2 \end{aligned}$$

So that $F(\mathbf{x}_k) - F(\mathbf{x}_\star) \leq \frac{L}{2k} \|\mathbf{x}_0 - \mathbf{x}_k\|_2^2$

PG Optimizes Smooth FBE Envelope

Definition. (Forward-backward envelope (FBE) $F(\mathbf{x})$
of)

$$Q_\mu(\mathbf{x}) := \min_{\mathbf{z}} \hat{F}_\mu(\mathbf{z}, \mathbf{x})$$

$$\begin{aligned} Q_\mu(\mathbf{x}) &= \underbrace{\min_{\mathbf{z}} \left\{ g(\mathbf{z}) + \frac{1}{2\mu} \|\mathbf{z} - (\mathbf{x} - \mu \nabla f(\mathbf{x}))\|_2^2 \right\}}_{\text{Moreau envelope } M_{\mu g}(\mathbf{x} - \mu \nabla f(\mathbf{x})) \text{ of } g(\mathbf{x})} \\ &\quad + f(\mathbf{x}) - \frac{\mu}{2} \|\nabla f(\mathbf{x})\|_2^2 \end{aligned}$$

PG Optimizes Smooth FBE Envelope

Theorem. (Differentiability of FBE $Q_\mu(\mathbf{x})$)

Suppose f is twice continuously differentiable and L -smooth.
Then $Q_\mu(\mathbf{x})$ is *continuously differentiable* with

$$\nabla Q_\mu(\mathbf{x}) = (\mathbf{I} - \mu \nabla^2 f(\mathbf{x})) \cdot G_\mu(\mathbf{x}).$$

If $\mu \in (0, 1/L)$, $G_\mu(\mathbf{x})$ is a descent direction on $Q_\mu(\mathbf{x})$, and
the set of stationary points of $Q_\mu(\mathbf{x})$ equals

$$\text{zer } \partial F = \{\mathbf{x} \in \mathbb{R}^n \mid 0 \in \partial F(\mathbf{x})\}.$$

PG Optimizes Smooth FBE Envelope

Whenever $\mu \in (0, 1/L)$, $Q_\mu(\mathbf{x})$ is continuously differentiable with

$$\langle G_\mu(\mathbf{x}), \nabla Q_\mu(\mathbf{x}) \rangle > 0$$

The proximal gradient

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mu \cdot G_\mu(\mathbf{x}_k)$$

can be viewed as a descent method on the smooth envelope $Q_\mu(\mathbf{x})$ of $F(\mathbf{x})$.

Convergence for Strongly Convex Functions

Theorem. Let $F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$, with

- f is γ -strongly convex, differentiable, and L -smooth;
- g is convex, but possibly nonsmooth.

Consider the sequence of proximal gradient iterates

$$\mathbf{x}_{k+1} = \text{prox}_{\mu_k g} (\mathbf{x}_k - \mu_k \nabla f(\mathbf{x}_k)), \quad \mu_k \equiv 1/L.$$

Let \mathbf{x}_* be a minimizer of F . Then for any $k \geq 1$,

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2^2 \leq \left(1 - \frac{\gamma}{L}\right)^k \cdot \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2.$$

Rate of Convergence

Methods	Convex Nonsmooth	Convex smooth	Strongly Convex
Subgradient Method	$O(1/\sqrt{k})$		
Gradient Descent		$O(1/k)$	Linear
Proximal Method	$O(1/k)$		Linear

Backtracking Linesearch for L

- Proximal gradient:

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \frac{1}{L} \mathcal{G}_\mu(\boldsymbol{x}_k)$$

In many cases, the Lipschitz constant L is hard to compute

- Example:

$$F(\boldsymbol{x}) = \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2 + \lambda \|\boldsymbol{x}\|_1$$

Backtracking Linesearch for L

Smooth case: checking Armijo condition for gradient descent, to guarantee sufficient decrease

$$f(\mathbf{x}_k - \tau_k \nabla f(\mathbf{x}_k)) \leq f(\mathbf{x}_k) - \frac{\tau_k}{2} \|\nabla f(\mathbf{x}_k)\|_2^2.$$

which is equivalent to updating $\tau_k = 1/L_k$ until

$$\begin{aligned} \underbrace{f(\mathbf{x}_k - \tau_k \nabla f(\mathbf{x}_k))}_{\mathbf{x}_{k+1}} &\leq f(\mathbf{x}_k) - \frac{1}{L_k} \langle \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_k) \rangle + \frac{1}{2L_k} \|\nabla f(\mathbf{x}_k)\|_2^2 \\ &\leq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{L_k}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2. \end{aligned}$$

Backtracking Stepsize Rule

- Similarly, for PG we have

$$\boldsymbol{x}_{k+1} = T_{\mu_k}(\boldsymbol{x}_k) := \boldsymbol{x}_k - \mu_k G_{\mu_k}(\boldsymbol{x}_k).$$

- Thus, we would like to select μ_k such that

$$f(T_{\mu_k}(\boldsymbol{x}_k)) \leq f(\boldsymbol{x}_k) - \langle \nabla f(\boldsymbol{x}_k), \mu_k G_{\mu_k}(\boldsymbol{x}_k) \rangle + \frac{\mu_k}{2} \|G_{\mu_k}(\boldsymbol{x}_k)\|_2^2.$$

Backtracking Stepsize Rule

Algorithm 1 Backtracking Stepsize Rule for PG

Initialize $\mu_k = 1, \beta \in (0, 1)$

while $f(T_{\mu_k}(\mathbf{x}_k)) \geq f(\mathbf{x}_k) - \mu_k \langle \nabla f(\mathbf{x}_k), G_{\mu_k}(\mathbf{x}_k) \rangle + \frac{\mu_k}{2} \|G_{\mu_k}(\mathbf{x}_k)\|_2^2$ **do**
 $\mu_k \leftarrow \beta \cdot \mu_k$

end while

In practice, we often do *warm restart* for μ_k to accelerate the speed

$$\mu_k = \alpha \cdot \mu_{k-1}, \quad \alpha > 1.$$

Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma.
(Chapter 8.2)
- *Proximal Algorithms, Foundations & Trends in Optimization.*
Neal Parikh, Stephen Boyd, 2014.
https://web.stanford.edu/~boyd/papers/prox_algs.html
- *A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems.* Amir Beck and Marc Teboulle, 2009.
<https://pubs.siam.org/doi/abs/10.1137/080716542?mobileUi=0>