

An Introduction to Non-smooth Optimization

Lecture 03 - Proximal Gradient Descent

Jingwei LIANG

Institute of Natural Sciences, Shanghai Jiao Tong University

Email: optimization.sjtu@gmail.com

Office: Room 355, No. 6 Science Building



饮水思源 · 爱国荣校

We have seen two problems

- Non-negative least square (NLS)

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{y}\|^2 \quad \text{such that} \quad x_i \geq 0, i = 1, 2, \dots, n.$$

- Sparse logistic regression

$$\min_{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R}} \mu \|\mathbf{x}\|_1 + \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-b_i(\mathbf{x}^T \mathbf{a}_i + y)}).$$

We have seen two problems

- Non-negative least square (NLS)

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{y}\|^2 \quad \text{such that} \quad x_i \geq 0, i = 1, 2, \dots, n.$$

- Sparse logistic regression

$$\min_{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R}} \mu \|\mathbf{x}\|_1 + \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-b_i(\mathbf{x}^T \mathbf{a}_i + y)}).$$

Let $S = \{\mathbf{x} \in \mathbb{R}^n : x_i \geq 0, i = 1, 2, \dots, n\}$ and define

$$\iota_S(\mathbf{x}) = \begin{cases} 0 & : \mathbf{x} \in S, \\ +\infty & : \mathbf{x} \notin S. \end{cases}$$

The NLS problem can be equivalently written as

$$\min_{\mathbf{x} \in \mathbb{R}^n} \iota_S(\mathbf{x}) + \|\mathbf{Ax} - \mathbf{y}\|^2.$$

We have seen two problems

- Non-negative least square (NLS)

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{y}\|^2 \quad \text{such that} \quad x_i \geq 0, i = 1, 2, \dots, n.$$

- Sparse logistic regression

$$\min_{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R}} \mu \|\mathbf{x}\|_1 + \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-b_i(\mathbf{x}^T \mathbf{a}_i + y)}).$$

Problem - Non-smooth optimization problem

Let $F, R \in \Gamma_0(\mathbb{R}^n)$, consider

$$\min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \Phi(\mathbf{x}) \stackrel{\text{def}}{=} F(\mathbf{x}) + R(\mathbf{x}) \right\},$$

with

- F : smooth differentiable with ∇F being L -Lipschitz continuous.
- R : non-smooth with proximal mapping easy to compute.

We have seen two problems

- Non-negative least square (NLS)

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{y}\|^2 \quad \text{such that} \quad x_i \geq 0, i = 1, 2, \dots, n.$$

- Sparse logistic regression

$$\min_{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R}} \mu \|\mathbf{x}\|_1 + \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-b_i(\mathbf{x}^T \mathbf{a}_i + y)}).$$

Proposition - Optimality condition

Suppose $\text{zer}(\nabla F + \partial R)$ is non-empty, and let $\mathbf{x}^* \in \text{zer}(\nabla F + \partial R)$. Then

$$\mathbf{0} \in \partial F(\mathbf{x}^*) + \partial R(\mathbf{x}^*).$$

Outline

- 1 Gradient descent
- 2 Projected gradient descent
- 3 Proximal gradient descent
- 4 Convergence analysis



Projected gradient descent

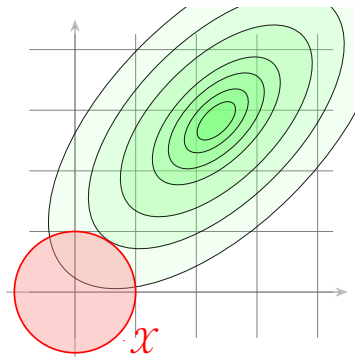
Constrained smooth optimization



Problem - Constrained smooth optimization

Let $S \subset \mathbb{R}^n$ be closed and convex and $F \in C_L^1(\mathbb{R}^n)$,

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \quad \text{such that} \quad \mathbf{x} \in S.$$



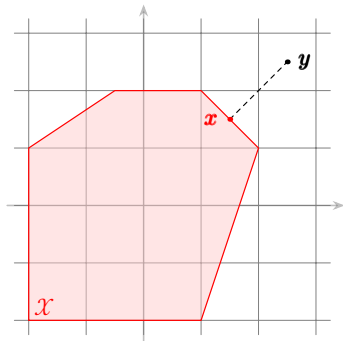
Projection of \mathbf{y} onto S :

$$\min_{\mathbf{x} \in S} \|\mathbf{x} - \mathbf{y}\|.$$

Definition - Projection

Projection mapping onto a set is defined by

$$\mathcal{P}_S(\mathbf{y}) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{x} \in S} \|\mathbf{x} - \mathbf{y}\|.$$



The projection is unique for closed and convex, and S

$$\mathcal{P}_S(\mathbf{y}) \in S \quad \text{and} \quad \forall \mathbf{x} \in S \quad \langle \mathbf{x} - \mathcal{P}_S(\mathbf{y}) \mid \mathbf{y} - \mathcal{P}_S(\mathbf{y}) \rangle \leq 0.$$

Algorithm - Projected Gradient descent

initial: $\mathbf{x}^{(0)} \in \text{dom}(F)$;

repeat:

1. Choose step-size $\gamma_k > 0$
2. GD: $\mathbf{x}^{(k+1/2)} = \mathbf{x}^{(k)} - \gamma_k \nabla F(\mathbf{x}^{(k)})$
3. Projection: $\mathbf{x}^{(k+1)} = \mathcal{P}_S(\mathbf{x}^{(k+1/2)})$

until: stopping criterion is satisfied.

In a compact form

$$\mathbf{x}^{(k+1)} = \mathcal{P}_S(\mathbf{x}^{(k)} - \gamma_k \nabla F(\mathbf{x}^{(k)})).$$

- The same as gradient descent, only one parameter which is γ_k .

Proximal gradient descent

Proximal mapping and algorithm



From projection to proximal mapping



Previously

$$\mathcal{P}_S(\mathbf{y}) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{x} \in S} \|\mathbf{x} - \mathbf{y}\|.$$

Previously

$$\mathcal{P}_S(\mathbf{y}) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{x} \in S} \|\mathbf{x} - \mathbf{y}\|.$$

The following are equivalent: $\iota_S(\mathbf{x}) \in \Gamma_0(\mathbb{R}^n)$ for closed convex S

$$\begin{aligned} \min_{\mathbf{x} \in S} \|\mathbf{x} - \mathbf{y}\| &\iff \min_{\mathbf{x} \in S} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2 \\ &\iff \min_{\mathbf{x} \in \mathbb{R}^n} \iota_S(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

From projection to proximal mapping



Previously

$$\mathcal{P}_S(\mathbf{y}) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{x} \in S} \|\mathbf{x} - \mathbf{y}\|.$$

The following are equivalent: $\iota_S(\mathbf{x}) \in \Gamma_0(\mathbb{R}^n)$ for closed convex S

$$\begin{aligned} \min_{\mathbf{x} \in S} \|\mathbf{x} - \mathbf{y}\| &\iff \min_{\mathbf{x} \in S} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2 \\ &\iff \min_{\mathbf{x} \in \mathbb{R}^n} \iota_S(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

Given any $R(\mathbf{x}) \in \Gamma_0(\mathbb{R}^n)$

$$\boxed{\min_{\mathbf{x} \in \mathbb{R}^n} R(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2.}$$

Definition - Proximal mapping

The proximal mapping (or proximity operator) of a function $R \in \Gamma_0(\mathbb{R}^n)$ is defined by

$$\text{prox}_{\gamma R}(\mathbf{y}) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \gamma R(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

Definition - Proximal mapping

The proximal mapping (or proximity operator) of a function $R \in \Gamma_0(\mathbb{R}^n)$ is defined by

$$\text{prox}_{\gamma R}(\mathbf{y}) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \gamma R(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

- $\text{prox}_{\gamma R}(\mathbf{y})$ is unique for $R \in \Gamma_0(\mathbb{R}^n)$.

Definition - Proximal mapping

The proximal mapping (or proximity operator) of a function $R \in \Gamma_0(\mathbb{R}^n)$ is defined by

$$\text{prox}_{\gamma R}(\mathbf{y}) \stackrel{\text{def}}{=} \underset{\mathbf{x} \in \mathbb{R}^n}{\text{argmin}} \gamma R(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

- $\text{prox}_{\gamma R}(\mathbf{y})$ is unique for $R \in \Gamma_0(\mathbb{R}^n)$.
- Alternative characterization let $\mathbf{x} \stackrel{\text{def}}{=} \text{prox}_{\gamma R}(\mathbf{y})$,

$$\begin{aligned} \mathbf{0} \in \gamma \partial R(\mathbf{x}) + \mathbf{x} - \mathbf{y} &\iff \mathbf{y} - \mathbf{x} \in \gamma \partial R(\mathbf{x}) \\ &\iff \mathbf{y} \in (\mathbf{Id} + \gamma \partial R)(\mathbf{x}) \\ &\iff \mathbf{x} = (\mathbf{Id} + \gamma \partial R)^{-1}(\mathbf{y}). \end{aligned}$$

- $(\mathbf{Id} + \gamma \partial R)^{-1}$ is called the **resolvent** of $\gamma \partial R$.

Projection $R(\mathbf{x}) = \iota_S(\mathbf{x})$, then

$$\partial \iota_S(\mathbf{x}) = \mathcal{N}_S(\mathbf{x}) = \{\mathbf{g} : \langle \mathbf{g} | \mathbf{u} - \mathbf{x} \rangle \leq 0, \forall \mathbf{u} \in S\}$$

and

$$\mathcal{P}_S(\mathbf{y}) = (\mathbf{Id} + \mathcal{N}_S)^{-1}(\mathbf{y}).$$

Projection $R(\mathbf{x}) = \iota_S(\mathbf{x})$, then

$$\partial \iota_S(\mathbf{x}) = \mathcal{N}_S(\mathbf{x}) = \{\mathbf{g} : \langle \mathbf{g} | \mathbf{u} - \mathbf{x} \rangle \leq 0, \forall \mathbf{u} \in S\}$$

and

$$\mathcal{P}_S(\mathbf{y}) = (\mathbf{Id} + \mathcal{N}_S)^{-1}(\mathbf{y}).$$

Examples

- Hyperplane: $S = \{\mathbf{x} : \mathbf{a}^T \mathbf{x} = b\}$, $\mathbf{a} \neq 0$

$$\mathcal{P}_S(\mathbf{y}) = \mathbf{y} + \frac{b - \mathbf{a}^T \mathbf{y}}{\|\mathbf{a}\|^2} \mathbf{a}.$$

- Affine subspace: $S = \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ with $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\text{rank}(\mathbf{A}) = m < n$

$$\mathcal{P}_S(\mathbf{y}) = \mathbf{y} + \mathbf{A}^T (\mathbf{A}\mathbf{A}^T)^{-1} (\mathbf{b} - \mathbf{A}\mathbf{y}).$$

- Nonnegative orthant: $S = \mathbb{R}_+^n$

$$\mathcal{P}_S(\mathbf{y}) = (\max\{0, y_i\})_i.$$

Quadratic function $R(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T\mathbf{A}\mathbf{x} + \mathbf{b}^T\mathbf{x} + c$ with $\mathbf{A} \in \mathbb{R}^{n \times n}$ being symmetric and positive semi-definite

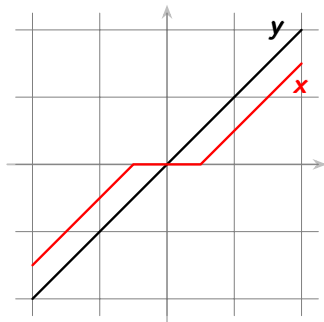
$$\text{prox}_{\gamma R}(\mathbf{y}) = (\mathbf{Id} + \gamma\mathbf{A})^{-1}(\mathbf{y} - \gamma\mathbf{b}).$$

Quadratic function $R(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T\mathbf{A}\mathbf{x} + \mathbf{b}^T\mathbf{x} + c$ with $\mathbf{A} \in \mathbb{R}^{n \times n}$ being symmetric and positive semi-definite

$$\text{prox}_{\gamma R}(\mathbf{y}) = (\mathbf{Id} + \gamma\mathbf{A})^{-1}(\mathbf{y} - \gamma\mathbf{b}).$$

Soft-threshold: $R(x) = |x|,$

$$\text{prox}_{\gamma R}(y) = \mathcal{T}_{\gamma}(y) = \begin{cases} y - \gamma : y > \gamma, \\ 0 : y \in [-\gamma, \gamma], \\ y + \gamma : y < -\gamma. \end{cases}$$



Quadratic function $R(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T\mathbf{A}\mathbf{x} + \mathbf{b}^T\mathbf{x} + c$ with $\mathbf{A} \in \mathbb{R}^{n \times n}$ being symmetric and positive semi-definite

$$\text{prox}_{\gamma R}(\mathbf{y}) = (\mathbf{Id} + \gamma\mathbf{A})^{-1}(\mathbf{y} - \gamma\mathbf{b}).$$

Soft-threshold: $R(x) = |x|$,

$$\text{prox}_{\gamma R}(\mathbf{y}) = \mathcal{T}_{\gamma}(\mathbf{y}) = \begin{cases} \mathbf{y} - \gamma : \mathbf{y} > \gamma, \\ 0 : \mathbf{y} \in [-\gamma, \gamma], \\ \mathbf{y} + \gamma : \mathbf{y} < -\gamma. \end{cases}$$

Euclidean norm $R(\mathbf{x}) = \|\mathbf{x}\|_2$

$$\text{prox}_{\gamma R}(\mathbf{y}) = \begin{cases} (1 - \frac{\gamma}{\|\mathbf{y}\|})\mathbf{y} : \|\mathbf{y}\| > \gamma, \\ \mathbf{0} : \text{o.w.} \end{cases}$$

Quadratic function $R(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T\mathbf{A}\mathbf{x} + \mathbf{b}^T\mathbf{x} + c$ with $\mathbf{A} \in \mathbb{R}^{n \times n}$ being symmetric and positive semi-definite

$$\text{prox}_{\gamma R}(\mathbf{y}) = (\mathbf{Id} + \gamma\mathbf{A})^{-1}(\mathbf{y} - \gamma\mathbf{b}).$$

Soft-threshold: $R(x) = |x|$,

$$\text{prox}_{\gamma R}(\mathbf{y}) = \mathcal{T}_{\gamma}(\mathbf{y}) = \begin{cases} \mathbf{y} - \gamma : \mathbf{y} > \gamma, \\ 0 : \mathbf{y} \in [-\gamma, \gamma], \\ \mathbf{y} + \gamma : \mathbf{y} < -\gamma. \end{cases}$$

Euclidean norm $R(\mathbf{x}) = \|\mathbf{x}\|_2$

$$\text{prox}_{\gamma R}(\mathbf{y}) = \begin{cases} (1 - \frac{\gamma}{\|\mathbf{y}\|})\mathbf{y} : \|\mathbf{y}\| > \gamma, \\ \mathbf{0} : \text{o.w.} \end{cases}$$

Nuclear norm $R(\mathbf{x}) = \sum_i \mathbf{S}_i$, let $\mathbf{y} = \mathbf{U}\mathbf{S}\mathbf{V}^T \in \mathbb{R}^{m \times n}$

$$\text{prox}_{\gamma R}(\mathbf{y}) = \mathbf{U}\mathcal{T}_{\gamma}(\mathbf{S})\mathbf{V}^T.$$

Quadratic perturbation $H(\mathbf{x}) = R(\mathbf{x}) + \frac{\alpha}{2} \|\mathbf{x}\|^2 + \langle \mathbf{x} | \mathbf{u} \rangle + \mathbf{b}$, $\alpha \geq 0$

$$\text{prox}_H(\mathbf{y}) = \text{prox}_{R/(\alpha+1)}\left(\frac{\mathbf{y}-\mathbf{u}}{\alpha+1}\right).$$

Quadratic perturbation $H(\mathbf{x}) = R(\mathbf{x}) + \frac{\alpha}{2} \|\mathbf{x}\|^2 + \langle \mathbf{x} \mid \mathbf{u} \rangle + \mathbf{b}$, $\alpha \geq 0$

$$\text{prox}_H(\mathbf{y}) = \text{prox}_{R/(\alpha+1)}\left(\frac{\mathbf{y}-\mathbf{u}}{\alpha+1}\right).$$

Translation $H(\mathbf{x}) = R(\mathbf{x} - \mathbf{z})$

$$\text{prox}_H(\mathbf{y}) = \mathbf{z} + \text{prox}_R(\mathbf{y} - \mathbf{z}).$$

Quadratic perturbation $H(\mathbf{x}) = R(\mathbf{x}) + \frac{\alpha}{2} \|\mathbf{x}\|^2 + \langle \mathbf{x} | \mathbf{u} \rangle + \mathbf{b}$, $\alpha \geq 0$

$$\text{prox}_H(\mathbf{y}) = \text{prox}_{R/(\alpha+1)}\left(\frac{\mathbf{y}-\mathbf{u}}{\alpha+1}\right).$$

Translation $H(\mathbf{x}) = R(\mathbf{x} - \mathbf{z})$

$$\text{prox}_H(\mathbf{y}) = \mathbf{z} + \text{prox}_R(\mathbf{y} - \mathbf{z}).$$

Scaling $H(\mathbf{x}) = R(\mathbf{x}/\rho)$

$$\text{prox}_H(\mathbf{y}) = \rho \text{prox}_{R/\rho^2}(\mathbf{y}/\rho).$$

Quadratic perturbation $H(\mathbf{x}) = R(\mathbf{x}) + \frac{\alpha}{2} \|\mathbf{x}\|^2 + \langle \mathbf{x} | \mathbf{u} \rangle + \mathbf{b}$, $\alpha \geq 0$

$$\text{prox}_H(\mathbf{y}) = \text{prox}_{R/(\alpha+1)}\left(\frac{\mathbf{y}-\mathbf{u}}{\alpha+1}\right).$$

Translation $H(\mathbf{x}) = R(\mathbf{x} - \mathbf{z})$

$$\text{prox}_H(\mathbf{y}) = \mathbf{z} + \text{prox}_R(\mathbf{y} - \mathbf{z}).$$

Scaling $H(\mathbf{x}) = R(\mathbf{x}/\rho)$

$$\text{prox}_H(\mathbf{y}) = \rho \text{prox}_{R/\rho^2}(\mathbf{y}/\rho).$$

Reflection $H(\mathbf{x}) = R(-\mathbf{x})$

$$\text{prox}_H(\mathbf{y}) = -\text{prox}_R(-\mathbf{y}).$$

Quadratic perturbation $H(\mathbf{x}) = R(\mathbf{x}) + \frac{\alpha}{2} \|\mathbf{x}\|^2 + \langle \mathbf{x} | \mathbf{u} \rangle + \mathbf{b}$, $\alpha \geq 0$

$$\text{prox}_H(\mathbf{y}) = \text{prox}_{R/(\alpha+1)}\left(\frac{\mathbf{y}-\mathbf{u}}{\alpha+1}\right).$$

Translation $H(\mathbf{x}) = R(\mathbf{x} - \mathbf{z})$

$$\text{prox}_H(\mathbf{y}) = \mathbf{z} + \text{prox}_R(\mathbf{y} - \mathbf{z}).$$

Scaling $H(\mathbf{x}) = R(\mathbf{x}/\rho)$

$$\text{prox}_H(\mathbf{y}) = \rho \text{prox}_{R/\rho^2}(\mathbf{y}/\rho).$$

Reflection $H(\mathbf{x}) = R(-\mathbf{x})$

$$\text{prox}_H(\mathbf{y}) = -\text{prox}_R(-\mathbf{y}).$$

Composition $H = R \circ \mathbf{K}$ with \mathbf{K} being bijective bounded linear mapping such that $\mathbf{K}^{-1} = \mathbf{K}^*$,

$$\text{prox}_H(\mathbf{y}) = \mathbf{K}^* \text{prox}_R(\mathbf{K}\mathbf{y}).$$

Problem - Non-smooth optimization

Let $R \in \Gamma_0(\mathbb{R}^n)$ and $F \in C_L^1(\mathbb{R}^n)$,

$$\min_{\mathbf{x} \in \mathbb{R}^n} R(\mathbf{x}) + F(\mathbf{x}).$$

Problem - Non-smooth optimization

Let $R \in \Gamma_0(\mathbb{R}^n)$ and $F \in C_L^1(\mathbb{R}^n)$,

$$\min_{\mathbf{x} \in \mathbb{R}^n} R(\mathbf{x}) + F(\mathbf{x}).$$

Algorithm - Proximal Gradient descent

initial: $\mathbf{x}^{(0)} \in \text{dom}(F)$;

repeat:

1. Choose step-size $\gamma_k > 0$
2. GD: $\mathbf{x}^{(k+1/2)} = \mathbf{x}^{(k)} - \gamma_k \nabla F(\mathbf{x}^{(k)})$
3. Projection: $\mathbf{x}^{(k+1)} = \text{prox}_{\gamma_k R}(\mathbf{x}^{(k+1/2)})$

until: stopping criterion is satisfied.

In a compact form

$$\mathbf{x}^{(k+1)} = \text{prox}_{\gamma_k R}(\mathbf{x}^{(k)} - \gamma_k \nabla F(\mathbf{x}^{(k)})).$$

Convergence analysis

Fixed-point iteration perspective



Problem - Non-smooth optimization problem

Let $F, R \in \Gamma_0(\mathbb{R}^n)$, consider

$$\min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \Phi(\mathbf{x}) \stackrel{\text{def}}{=} F(\mathbf{x}) + R(\mathbf{x}) \right\},$$

with

F: smooth differentiable with ∇F being L -Lipschitz continuous.

R: non-smooth with proximal mapping easy to compute.

Problem - Non-smooth optimization problem

Let $F, R \in \Gamma_0(\mathbb{R}^n)$, consider

$$\min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \Phi(\mathbf{x}) \stackrel{\text{def}}{=} F(\mathbf{x}) + R(\mathbf{x}) \right\},$$

with

F: smooth differentiable with ∇F being L -Lipschitz continuous.

R: non-smooth with proximal mapping easy to compute.

Proposition - Optimality condition

Suppose $\text{zer}(\nabla F + \partial R)$ is non-empty, and let $\mathbf{x}^* \in \text{zer}(\nabla F + \partial R)$. Then

$$\mathbf{0} \in \nabla F(\mathbf{x}^*) + \partial R(\mathbf{x}^*).$$

Definition - Cocoercive operator

Let S be a non-empty subset of \mathbb{R}^n , $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\beta > 0$. Then \mathcal{B} is β -cocoercive if

$$(\forall \mathbf{x} \in S)(\mathbf{y} \in S) \quad \langle \mathbf{x} - \mathbf{y} \mid \mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y}) \rangle \geq \beta \|\mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y})\|^2.$$

Definition - Cocoercive operator

Let S be a non-empty subset of \mathbb{R}^n , $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\beta > 0$. Then \mathcal{B} is β -cocoercive if

$$(\forall \mathbf{x} \in S)(\mathbf{y} \in S) \quad \langle \mathbf{x} - \mathbf{y} \mid \mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y}) \rangle \geq \beta \|\mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y})\|^2.$$

Theorem - Cocoercivity and Lipschitz continuity

Let S be a non-empty subset of \mathbb{R}^n , $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\beta > 0$. If \mathcal{B} is β -cocoercive, then

$$(\forall \mathbf{x} \in S)(\mathbf{y} \in S) \quad \|\mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y})\| \leq \frac{1}{\beta} \|\mathbf{x} - \mathbf{y}\|.$$

- Cocoercivity implies Lipschitz continuity, the reverse in general is not true.

Definition - Cocoercive operator

Let S be a non-empty subset of \mathbb{R}^n , $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\beta > 0$. Then \mathcal{B} is β -cocoercive if

$$(\forall \mathbf{x} \in S)(\mathbf{y} \in S) \quad \langle \mathbf{x} - \mathbf{y} \mid \mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y}) \rangle \geq \beta \|\mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y})\|^2.$$

Theorem - Cocoercivity and Lipschitz continuity

Let S be a non-empty subset of \mathbb{R}^n , $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\beta > 0$. If \mathcal{B} is β -cocoercive, then

$$(\forall \mathbf{x} \in S)(\mathbf{y} \in S) \quad \|\mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y})\| \leq \frac{1}{\beta} \|\mathbf{x} - \mathbf{y}\|.$$

- Cocoercivity implies Lipschitz continuity, the reverse in general is not true.

Theorem - [Baillon-Haddad '77]

For $F \in C_L^1(\mathbb{R}^n)$, its gradient ∇F is $\frac{1}{L}$ -cocoercive.

Definition - Cocoercive operator

Let S be a non-empty subset of \mathbb{R}^n , $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\beta > 0$. Then \mathcal{B} is β -cocoercive if

$$(\forall \mathbf{x} \in S)(\mathbf{y} \in S) \quad \langle \mathbf{x} - \mathbf{y} \mid \mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y}) \rangle \geq \beta \|\mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{y})\|^2.$$

Proposition - Cocoercivity and non-expansiveness

Let $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be β -cocoercive for some $\beta > 0$, then

- $\beta\mathcal{B}$ is firmly non-expansive.
- $\text{Id} - \gamma\mathcal{B}$ is $\frac{\gamma}{2\beta}$ -averaged non-expansive for $\gamma \in]0, 2\beta[$.

Definition - Resolvent

Let $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ be monotone. The resolvent of \mathcal{A} is

$$\mathcal{J}_{\mathcal{A}} = (\mathbf{Id} + \mathcal{A})^{-1}.$$

Definition - Resolvent

Let $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ be monotone. The resolvent of \mathcal{A} is

$$\mathcal{J}_{\mathcal{A}} = (\mathbf{Id} + \mathcal{A})^{-1}.$$

Example - Proximal mapping

Let $R \in \Gamma_0(\mathbb{R}^n)$ and $\gamma > 0$. Then

$$\mathcal{J}_{\gamma \partial R} = \text{prox}_{\gamma R}.$$

Definition - Resolvent

Let $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ be monotone. The resolvent of \mathcal{A} is

$$\mathcal{J}_{\mathcal{A}} = (\mathbf{Id} + \mathcal{A})^{-1}.$$

Theorem - Monotonicity and firmly non-expansiveness

Let S be a nonempty subset of \mathbb{R}^n , let $\mathcal{F} : S \rightarrow \mathbb{R}^n$ and set $\mathcal{A} = \mathcal{F}^{-1} - \mathbf{Id}$. Then the following holds

- $\mathcal{F} = \mathcal{J}_{\mathcal{A}}$.
- \mathcal{F} is *firmly non-expansive* if and only if \mathcal{A} is monotone.
- \mathcal{F} is *firmly non-expansive* and $S = \mathbb{R}^n$ if and only if \mathcal{A} is *maximally monotone*.

Definition - Resolvent

Let $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ be monotone. The resolvent of \mathcal{A} is

$$\mathcal{J}_{\mathcal{A}} = (\mathbf{Id} + \mathcal{A})^{-1}.$$

Corollary

Let $\mathcal{F} : \mathbb{R}^n \rightarrow \mathbb{R}^n$. Then \mathcal{F} is firmly non-expansive if and only if it is the resolvent of a maximally monotone operator $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$.

Definition - Monotone inclusion problem

Let $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ be maximal monotone and $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be β -cocoercive for some $\beta > 0$. Then monotone inclusion problem associated to $\mathcal{A} + \mathcal{B}$ reads

$$\text{find } \mathbf{x} \in \mathbb{R}^n \text{ such that } \mathbf{0} \in \mathcal{A}(\mathbf{x}) + \mathcal{B}(\mathbf{x}).$$

Definition - Monotone inclusion problem

Let $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ be maximal monotone and $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be β -cocoercive for some $\beta > 0$. Then monotone inclusion problem associated to $\mathcal{A} + \mathcal{B}$ reads

$$\text{find } \mathbf{x} \in \mathbb{R}^n \quad \text{such that} \quad \mathbf{0} \in \mathcal{A}(\mathbf{x}) + \mathcal{B}(\mathbf{x}).$$

Definition - Forward-Backward splitting

Let $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ be maximal monotone and $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$. Let \mathbf{x} be such that $\mathbf{0} \in \mathcal{A}(\mathbf{x}) + \mathcal{B}(\mathbf{x})$ and $\gamma > 0$. Then

$$\begin{aligned} \mathbf{0} \in \mathcal{A}(\mathbf{x}) + \mathcal{B}(\mathbf{x}) &\iff \mathbf{0} \in \gamma\mathcal{A}(\mathbf{x}) + \gamma\mathcal{B}(\mathbf{x}) \\ &\iff -\gamma\mathcal{B}(\mathbf{x}) \in \gamma\mathcal{A}(\mathbf{x}) \\ &\iff \mathbf{x} - \gamma\mathcal{B}(\mathbf{x}) \in \mathbf{x} + \gamma\mathcal{A}(\mathbf{x}) \\ &\iff \mathbf{x} = (\mathbf{Id} + \gamma\mathcal{A})^{-1}(\mathbf{Id} - \gamma\mathcal{B})(\mathbf{x}) \end{aligned}$$

Definition - Monotone inclusion problem

Let $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ be maximal monotone and $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be β -cocoercive for some $\beta > 0$. Then monotone inclusion problem associated to $\mathcal{A} + \mathcal{B}$ reads

$$\text{find } \mathbf{x} \in \mathbb{R}^n \text{ such that } \mathbf{0} \in \mathcal{A}(\mathbf{x}) + \mathcal{B}(\mathbf{x}).$$

Proposition - Fixed-point operator

Let $\mathcal{A} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ be maximal monotone and $\mathcal{B} : \mathbb{R}^n \rightarrow \mathbb{R}^n$. Set

$$\mathcal{F}_{\mathcal{A}, \mathcal{B}} = \mathcal{J}_{\mathcal{A}} \circ (\mathbf{Id} - \mathcal{B}).$$

Suppose \mathcal{B} is β -cocoercive for some $\beta > 0$ and that $\gamma \in]0, 2\beta[$. Let

$$\alpha = \frac{2\beta}{4\beta - \gamma},$$

then $\mathcal{F}_{\gamma\mathcal{A}, \gamma\mathcal{B}}$ is α -averaged non-expansive.

Theorem - Convergence with constant step-size

For proximal gradient descent, let $R \in \Gamma_0(\mathbb{R}^n)$ and $F \in C_L^1(\mathbb{R}^n)$. Let

$$\gamma_k \equiv \gamma \in]0, 2/L[.$$

Then

- $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ converges to a point \mathbf{x}^* in $\text{zer}(\partial R + \nabla F)$.

Theorem - Convergence speed

With the above convergence result,

- Sequence

$$\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| = o\left(\frac{1}{\sqrt{k}}\right).$$

- Objective function

$$(R + F)(\mathbf{x}^{(k)}) - (R + F)(\mathbf{x}^*) = o\left(\frac{1}{k}\right).$$

- Patrick L. Combettes, and Valérie R. Wajs. “Signal recovery by proximal forward-backward splitting.” *Multiscale modeling & simulation* 4.4 (2005): 1168-1200.
- Cesare Molinari, Jingwei Liang, and Jalal Fadili. “Convergence rates of Forward–Douglas–Rachford splitting method.” *Journal of Optimization Theory and Applications* 182.2 (2019): 606-639.
- Jingwei Liang, Jalal Fadili, and Gabriel Peyré. “Convergence rates with inexact non-expansive operators.” *Mathematical Programming* 159.1 (2016): 403-434.
- A. Beck: *First-order methods in optimization*, Vol. 25. SIAM, 2017.
- H. H. Bauschke and P. L. Combettes: *Convex analysis and monotone operator theory in Hilbert spaces*, Vol. 408. New York: Springer, 2011.