

An Introduction to Non-smooth Optimization

Lecture 06 - Several Acceleration Schemes

Jingwei LIANG

Institute of Natural Sciences, Shanghai Jiao Tong University

Email: optimization.sjtu@gmail.com

Office: Room 355, No. 6 Science Building

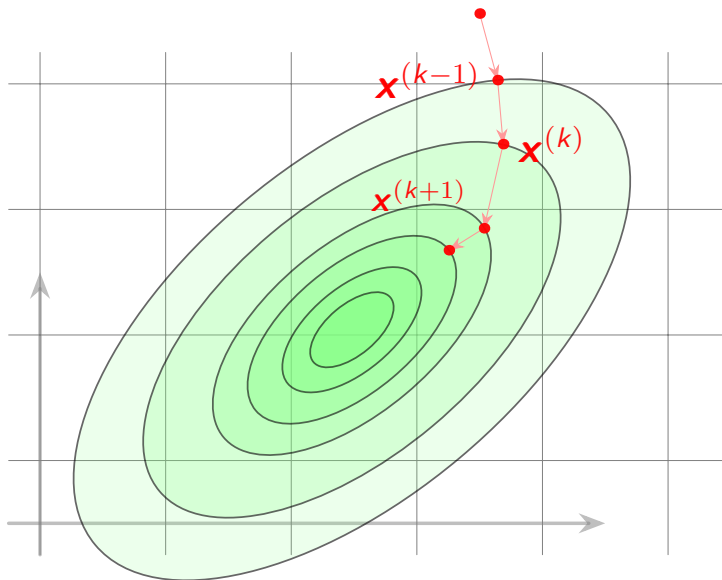


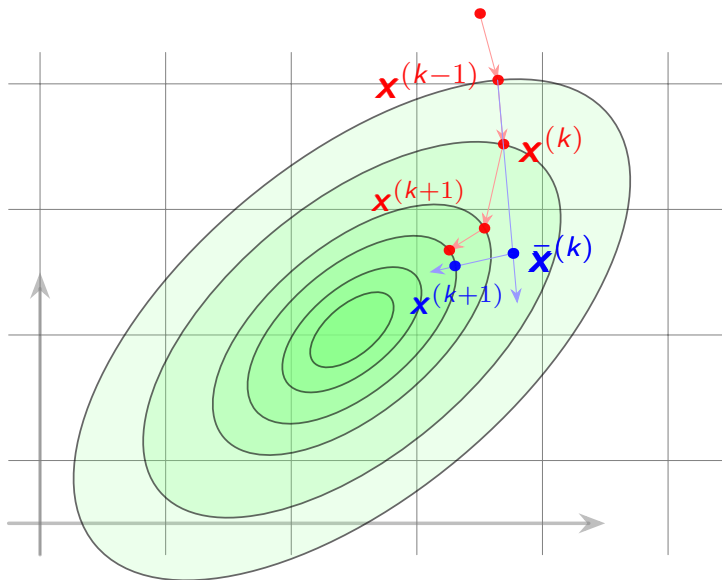
饮水思源 · 爱国荣校

Outline

- 1 Accelerate gradient descent
- 2 Accelerate proximal gradient descent
- 3 Accelerate ADMM
- 4 Accelerate fixed-point iteration







Algorithm - Heavy-ball method [Polyak '64]

Choose $\mathbf{x}^{(0)} \in \text{dom}(F)$, let $a > 0$ and $\gamma \in]0, 2/L[$

$$\begin{aligned}\mathbf{y}^{(k)} &= \mathbf{x}^{(k)} + a(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}), \\ \mathbf{x}^{(k+1)} &= \mathbf{y}^{(k)} - \gamma \nabla F(\mathbf{x}^{(k)}).\end{aligned}$$

Algorithm - Heavy-ball method [Polyak '64]

Choose $\mathbf{x}^{(0)} \in \text{dom}(F)$, let $a > 0$ and $\gamma \in]0, 2/L[$

$$\begin{aligned}\mathbf{y}^{(k)} &= \mathbf{x}^{(k)} + a(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}), \\ \mathbf{x}^{(k+1)} &= \mathbf{y}^{(k)} - \gamma \nabla F(\mathbf{x}^{(k)}).\end{aligned}$$

- $\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}$ is called the inertial term or momentum term.
- a is called the inertial parameter.
- In general, no convergence rate for $F \in C_L^1$. Local rate if moreover F is twice-differentiable and strongly convex.

Theorem - Optimal rate

Let \mathbf{x}^* be a (local) minimizer of F such that $\mu \mathbf{ld} \preceq \nabla^2 F(\mathbf{x}^*) \preceq L \mathbf{ld}$ and choose a, γ with $a \in [0, 1[, \gamma \in]0, 2(1+a)/L[$. There exists $\underline{\rho} < 1$ such that if $\underline{\rho} < \rho < 1$ and if $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}$ are close enough to \mathbf{x}^* , one has

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq C\rho^k.$$

Moreover, if

$$a = \left(\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^2, \quad \gamma = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2} \quad \text{then} \quad \underline{\rho} = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}.$$

Theorem - Optimal rate

Let \mathbf{x}^* be a (local) minimizer of F such that $\mu \mathbf{Id} \preceq \nabla^2 F(\mathbf{x}^*) \preceq L \mathbf{Id}$ and choose a, γ with $a \in [0, 1[, \gamma \in]0, 2(1+a)/L[$. There exists $\underline{\rho} < 1$ such that if $\underline{\rho} < \rho < 1$ and if $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}$ are close enough to \mathbf{x}^* , one has

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq C\rho^k.$$

Moreover, if

$$a = \left(\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^2, \quad \gamma = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2} \quad \text{then} \quad \underline{\rho} = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}.$$

- Starting points need to be close enough to \mathbf{x}^*
- Almost the optimal rate can be achieved by gradient method (or first-order method)
- Gradient descent

$$\underline{\rho} = \frac{L - \mu}{L + \mu}.$$

Algorithm - Nesterov's acceleration scheme [Nesterov '83]

Choose $\mathbf{x}^{(0)} \in \text{dom}(F)$ and $\mathbf{y}^{(0)} = \mathbf{x}^{(0)}$; Let $\phi_0 \in]0, 1[$ and $q = \mu/L$

$$\phi_{k+1}^2 = (1 - \phi_{k+1})\phi_k^2 + q\phi_{k+1}$$

$$a_k = \frac{\phi_k(1 - \phi_k)}{\phi_k^2 + \phi_{k+1}}$$

$$\mathbf{y}^{(k)} = \mathbf{x}^{(k)} + a_k(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$$

$$\mathbf{x}^{(k+1)} = \mathbf{y}^{(k)} - \frac{1}{L}\nabla F(\mathbf{y}^{(k)})$$

Algorithm - Nesterov's acceleration scheme [Nesterov '83]

Choose $\mathbf{x}^{(0)} \in \text{dom}(F)$ and $\mathbf{y}^{(0)} = \mathbf{x}^{(0)}$; Let $\phi_0 \in]0, 1[$ and $q = \mu/L$

$$\phi_{k+1}^2 = (1 - \phi_{k+1})\phi_k^2 + q\phi_{k+1}$$

$$a_k = \frac{\phi_k(1 - \phi_k)}{\phi_k^2 + \phi_{k+1}}$$

$$\mathbf{y}^{(k)} = \mathbf{x}^{(k)} + a_k(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$$

$$\mathbf{x}^{(k+1)} = \mathbf{y}^{(k)} - \frac{1}{L}\nabla F(\mathbf{y}^{(k)})$$

Parameter choices

■ $F \in C_L^1$: $\phi_0 = 1$,

$$q = 0, \quad \phi_k \approx \frac{2}{k+1} \rightarrow 0 \quad \text{and} \quad a_k \approx \frac{1 - \phi_k}{1 + \phi_k} \rightarrow 1.$$

■ $F \in S_{\mu,L}^1$: $\phi_0 = \sqrt{\mu/L}$

$$q = \sqrt{\frac{\mu}{L}}, \quad \phi_k \equiv \sqrt{\frac{\mu}{L}} \quad \text{and} \quad a_k \equiv \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}.$$

Accelerate proximal gradient descent

FISTA

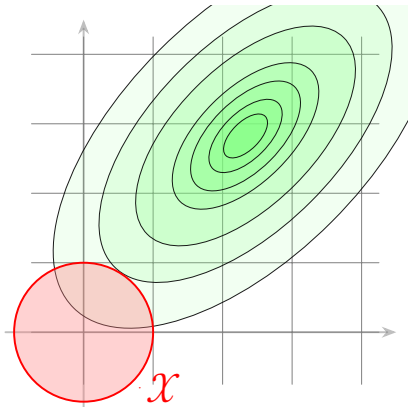


饮水思源 · 爱国荣校

Problem - Unconstrained smooth optimization

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

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



Algorithm - Gradient descent

Choose $\mathbf{x}^{(0)} \in \text{dom}(F)$ and $\gamma \in]0, 2/L[$

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

Algorithm - FISTA [Beck & Teboulle '09]

Choose $\mathbf{x}^{(0)} \in \text{dom}(F)$ and $\mathbf{y}^{(0)} = \mathbf{x}^{(0)}$; Let $t_0 = 1$ and $\gamma = 1/L$

$$t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}$$

$$a_k = \frac{t_{k-1} - 1}{t_k}$$

$$\mathbf{y}^{(k)} = \mathbf{x}^{(k)} + a_k(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$$

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

- A special case of inertial proximal gradient descent.
- Inertial parameters

$$t_k \approx \frac{k+1}{2} \quad \text{and} \quad a_k \rightarrow 1.$$

Algorithm - FISTA [Beck & Teboulle '09]

Choose $\mathbf{x}^{(0)} \in \text{dom}(F)$ and $\mathbf{y}^{(0)} = \mathbf{x}^{(0)}$; Let $t_0 = 1$ and $\gamma = 1/L$

$$t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}$$

$$a_k = \frac{t_{k-1} - 1}{t_k}$$

$$\mathbf{y}^{(k)} = \mathbf{x}^{(k)} + a_k(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$$

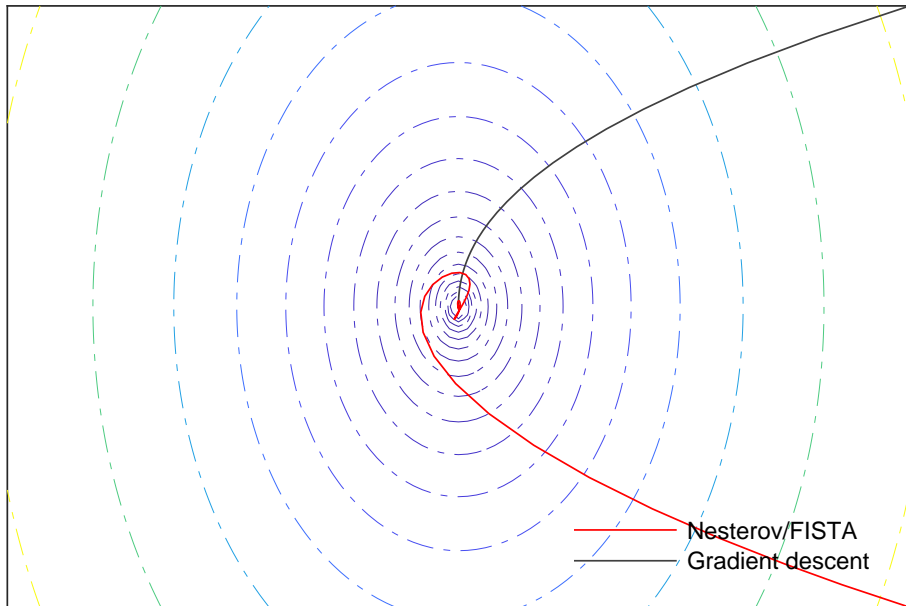
$$\mathbf{x}^{(k+1)} = \text{prox}_{\gamma R}(\mathbf{y}^{(k)} - \gamma \nabla F(\mathbf{y}^{(k)}))$$

Theorem - Convergence rate

Let $\mathbf{x}^* \in \text{Argmin}(F + R)$,

$$(F + R)(\mathbf{x}^{(k)}) - (F + R)(\mathbf{x}^*) \leq \frac{L \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2}{2(k+1)^2}.$$

Restarting FISTA



Why FISTA oscillates

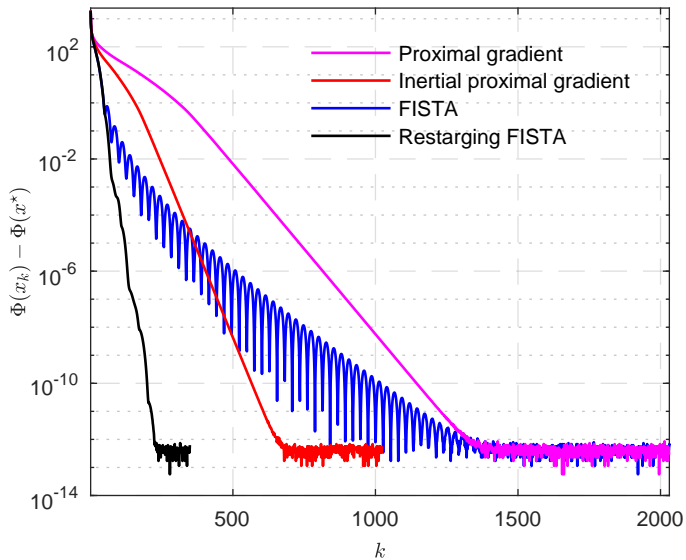
- for LSE, leading eigenvalue of the system is complex.
- over extrapolation, momentum beats gradient.

Algorithm - Restarting FISTA [O'Donoghue & Candès '12]

repeat:

1. Run FISTA iteration
2. If $\langle \mathbf{y}^{(k)} - \mathbf{x}^{(k+1)} | \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \rangle > 0$: $t_k = 1, \mathbf{y}^{(k)} = \mathbf{x}^{(k)}$.

Restarting FISTA



Accelerate ADMM

Strong convexity



Problem

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

such that $\mathbf{Ax} + \mathbf{By} = \mathbf{f}$.

Algorithm - ADMM [Gabay, Mercier, Glowinski, Marrocco '76]

$$\begin{aligned}\mathbf{x}^{(k+1)} &\in \operatorname{Argmin}_{\mathbf{x}} F(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{By}^{(k)} - \mathbf{f} + \mathbf{u}^{(k)} / \rho\|^2, \\ \mathbf{y}^{(k+1)} &\in \operatorname{Argmin}_{\mathbf{y}} R(\mathbf{y}) + \frac{\rho}{2} \|\mathbf{Ax}^{(k+1)} + \mathbf{By} - \mathbf{f} + \mathbf{u}^{(k)} / \rho\|^2, \\ \mathbf{u}^{(k+1)} &= \mathbf{u}^{(k)} + \rho(\mathbf{Ax}^{(k+1)} + \mathbf{By}^{(k+1)} - \mathbf{f}).\end{aligned}$$

Assumption

- Both F and R are strongly convex.

Algorithm - Fast ADMM [Goldstein et al '14]

Let $\mathbf{y}^{(0)} \in \mathbb{R}^n$, $\bar{\mathbf{y}}^{(0)} = \mathbf{y}^{(0)}$ and $\mathbf{u}^{(0)} \in \mathbb{R}^p$, $\bar{\mathbf{u}}^{(0)} = \mathbf{u}^{(0)}$; Let $\rho > 0$ and $t_0 = 1$:

$$\mathbf{x}^{(k+1)} \in \operatorname{Argmin}_{\mathbf{x}} F(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{B}\bar{\mathbf{y}}^{(k)} - \mathbf{f} + \bar{\mathbf{u}}^{(k)} / \rho\|^2,$$

$$\mathbf{y}^{(k+1)} \in \operatorname{Argmin}_{\mathbf{y}} R(\mathbf{y}) + \frac{\rho}{2} \|\mathbf{Ax}^{(k+1)} + \mathbf{By} - \mathbf{f} + \bar{\mathbf{u}}^{(k)} / \rho\|^2,$$

$$\mathbf{u}^{(k+1)} = \bar{\mathbf{u}}^{(k)} + \rho(\mathbf{Ax}^{(k+1)} + \mathbf{By}^{(k+1)} - \mathbf{f})$$

$$t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}$$

$$\bar{\mathbf{y}}^{(k+1)} = \mathbf{y}^{(k)} + \frac{t_{k-1} - 1}{t_k} (\mathbf{y}^{(k)} - \mathbf{y}^{(k-1)})$$

$$\bar{\mathbf{u}}^{(k+1)} = \mathbf{u}^{(k)} + \frac{t_{k-1} - 1}{t_k} (\mathbf{u}^{(k)} - \mathbf{u}^{(k-1)})$$

Accelerate fixed-point iteration

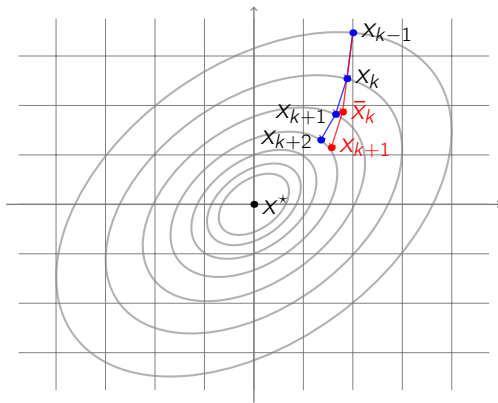
Inertial and over-relaxation



饮水思源 · 爱国荣校

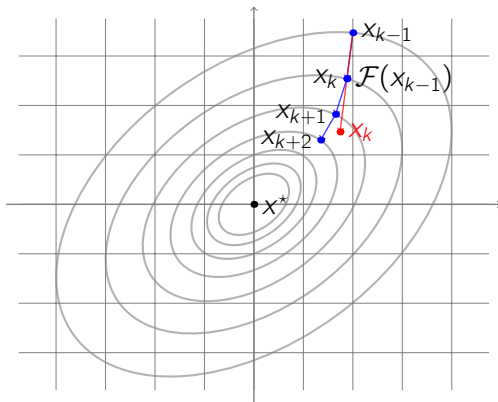
Algorithm - Inertial technique [Polyak '64, Nesterov '83, Beck & Teboulle '09]

$$\begin{cases} \bar{\mathbf{x}}^{(k)} = \mathbf{x}^{(k)} + a_k(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}), \\ \mathbf{x}^{(k+1)} = \mathcal{F}(\bar{\mathbf{x}}^{(k)}). \end{cases}$$



Algorithm - Successive over-relaxation [Richardson '1911, Young '50]

$$\mathbf{x}^{(k+1)} = (1 - \lambda_k)\mathbf{x}^{(k)} + \lambda_k\mathcal{F}(\mathbf{x}^{(k)}) \xrightarrow{a_k = \lambda_k - 1} \begin{cases} \bar{\mathbf{x}}^{(k)} = \mathbf{x}^{(k)} + a_k(\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k-1)}), \\ \mathbf{x}^{(k+1)} = \mathcal{F}(\bar{\mathbf{x}}^{(k)}). \end{cases}$$



Problem - Sum of two Γ_0 functions

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

Let $\gamma > 0$

$$\mathcal{F}_{\text{DR}} \stackrel{\text{def}}{=} \frac{1}{2}(\mathbf{Id} + (2\text{prox}_{\gamma R} - \mathbf{Id})(2\text{prox}_{\gamma F} - \mathbf{Id})).$$

Problem - Sum of two Γ_0 functions

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

Let $\gamma > 0$

$$\mathcal{F}_{\text{DR}} \stackrel{\text{def}}{=} \frac{1}{2}(\text{Id} + (2\text{prox}_{\gamma R} - \text{Id})(2\text{prox}_{\gamma F} - \text{Id})).$$

Douglas-Rachford splitting [[Douglas & Rachford '56](#)]

$$\mathbf{z}^{(k+1)} = \mathcal{F}_{\text{DR}}(\mathbf{z}^{(k)}),$$

- Sequence $o(1/\sqrt{k})$, objective **NA**.

Problem - Sum of two Γ_0 functions

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

Let $\gamma > 0$

$$\mathcal{F}_{\text{DR}} \stackrel{\text{def}}{=} \frac{1}{2}(\mathbf{Id} + (2\text{prox}_{\gamma R} - \mathbf{Id})(2\text{prox}_{\gamma F} - \mathbf{Id})).$$

Inertial Douglas–Rachford [Boţ, Csetnek & Hendrich '15]

$$\begin{aligned}\bar{\mathbf{z}}_k &= \mathbf{z}^{(k)} + a_k(\mathbf{z}^{(k)} - \mathbf{z}^{(k-1)}), \\ \mathbf{z}^{(k+1)} &= \mathcal{F}_{\text{DR}}(\bar{\mathbf{z}}_k).\end{aligned}$$

- No rates available, **may fail to provide acceleration.**

Alert: acceleration NOT guaranteed!



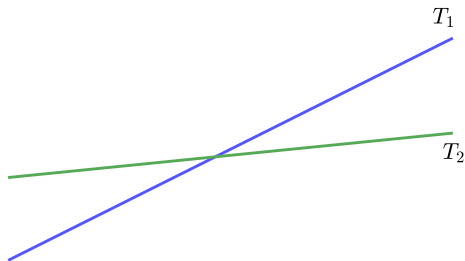
Problem - Feasibility problem in \mathbb{R}^2

Let $T_1, T_2 \subset \mathbb{R}^2$ be two subspaces such that $T_1 \cap T_2 \neq \emptyset$,

Find $x \in \mathbb{R}^2$ such that $x \in T_1 \cap T_2$.

Define

$$\mathcal{F}_{\text{DR}} \stackrel{\text{def}}{=} \frac{1}{2}(\mathbf{Id} + (2\mathcal{P}_{T_1} - \mathbf{Id})(2\mathcal{P}_{T_2} - \mathbf{Id})).$$



Alert: acceleration NOT guaranteed!



Problem - Feasibility problem in \mathbb{R}^2

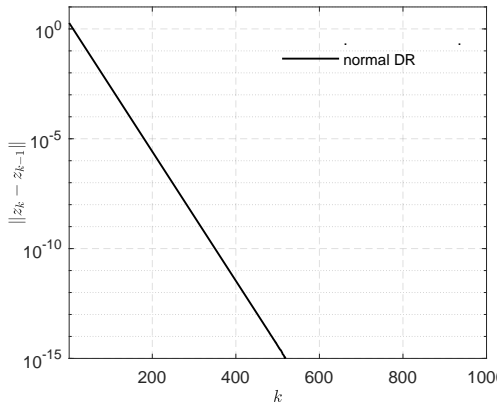
Let $T_1, T_2 \subset \mathbb{R}^2$ be two subspaces such that $T_1 \cap T_2 \neq \emptyset$,

Find $x \in \mathbb{R}^2$ such that $x \in T_1 \cap T_2$.

Douglas-Rachford:

$$\bar{z}_k = z^{(k)},$$

$$z^{(k+1)} = \mathcal{F}_{\text{DR}}(\bar{z}_k).$$



Alert: acceleration NOT guaranteed!



Problem - Feasibility problem in \mathbb{R}^2

Let $T_1, T_2 \subset \mathbb{R}^2$ be two subspaces such that $T_1 \cap T_2 \neq \emptyset$,

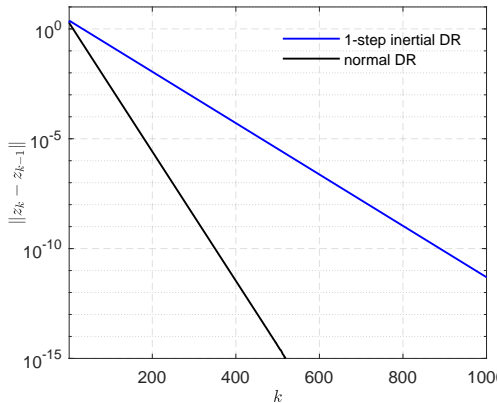
Find $x \in \mathbb{R}^2$ such that $x \in T_1 \cap T_2$.

Inertial Douglas-Rachford:

$$\bar{\mathbf{z}}_k = \mathbf{z}^{(k)} + a(\mathbf{z}^{(k)} - \mathbf{z}^{(k-1)}),$$

$$\mathbf{z}^{(k+1)} = \mathcal{F}_{\text{DR}}(\bar{\mathbf{z}}_k).$$

- 1-step inertial: $a = 0.3$.



- B. Polyak. “Introduction to optimization”. Optimization Software, 1987.
- Y. Nesterov. “Introductory lectures on convex optimization: A basic course”. Vol. 87. Springer Science & Business Media, 2013.
- A. Beck and M. Teboulle. “A fast iterative shrinkage-thresholding algorithm for linear inverse problems”. SIAM Journal on Imaging Sciences, 2(1):183–202, 2009.
- B. O’Donoghue and E. J. Candés. “Adaptive restart for accelerated gradient schemes”. Foundations of Computational Mathematics, pages 1–18, 2012.
- T. Goldstein, B. O’Donoghue, S. Setzer, R. Baraniuk. “Fast alternating direction optimization methods”. SIAM Journal on Imaging Sciences, pages 1588–1623, 2014.
- 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.