

On the Acceleration of First-order Methods

and a New One

Jingwei LIANG

Institute of Natural Sciences, Shanghai Jiao Tong University

Joint work with: **Clarice Poon (Bath UK)**

Outline

- ① Introduction
- ② Trajectory of first-order methods
- ③ Adaptive acceleration for FoM
- ④ Relation with previous work
- ⑤ Numerical experiments
- ⑥ Conclusions

Example - Constrained non-smooth optimization

$$\min_{x \in \mathbb{R}^n} F(x) + \sum_{i=1}^r R_i(K_i x),$$

where

F : smooth differentiable with L -Lipschitz continuous gradient.

R_i : proper and lower semi-continuous.

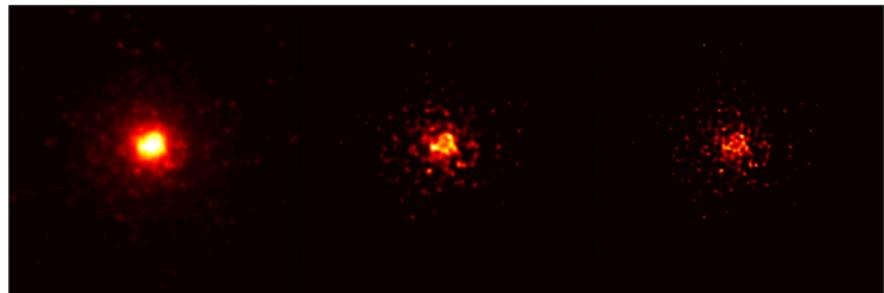
K_i : bounded linear mapping.

Applications: signal/image processing, inverse problems, machine learning, data science, statistics...

Challenges: non-smooth, (non-convex), composite, high dimension...

Hope: well structured.

Example: blind image deconvolution

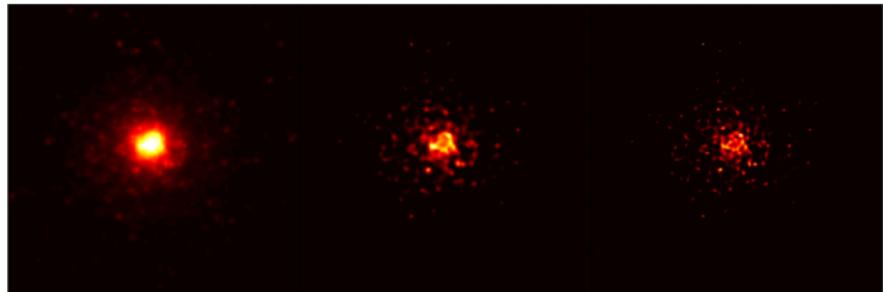


Forward model

$$A = y \odot x + \varepsilon.$$

Figure: NGC224 by Hubble Space Telescope from Wikipedia.

Example: blind image deconvolution



Forward model

$$A = y \odot x + \varepsilon.$$

Figure: NGC224 by Hubble Space Telescope from Wikipedia.

Example - blind image deconvolution

$$\min_{x \in \mathbb{R}^{p \times r}, y \in \mathbb{R}^{s \times s}} \frac{1}{2} \|A - y \odot x\|^2 + \lambda R(x) \quad \text{s.t.} \quad 0 \preceq x \preceq 1, 0 \preceq y \preceq 1, \|y\|_1 = 1.$$

Example: sparse non-negative matrix factorization



$$\begin{matrix} \text{Portrait of a person} \\ = \end{matrix} \begin{matrix} \text{A large sparse matrix} \\ \times \\ \text{A small sparse matrix} \end{matrix}$$

Decomposition model

$$A = yx + \varepsilon.$$

Example: sparse non-negative matrix factorization



Decomposition model

$$A = yx + \varepsilon.$$

Example - sparse NMF

$$\min_{x \in \mathbb{R}^{p \times r}, y \in \mathbb{R}^{r \times q}} \frac{1}{2} \|A - xy\|^2 \quad \text{s.t.} \quad x, y \geq 0, \quad \|x_i\|_0 \leq \kappa_x, \quad i = 1, \dots, r.$$



Gradient descent [Cauchy '1847]

$$\min_x F(x)$$

where F is convex smooth differentiable with ∇F being L -Lipschitz.



Gradient descent [Cauchy '1847]

$$\min_x F(x)$$

where F is convex smooth differentiable with ∇F being L -Lipschitz.

EXplicit Euler scheme:

$$x_{k+1} = x_k - \gamma_k \nabla F(x_k), \quad \gamma_k \in]0, 2/L[.$$



Gradient descent [Cauchy '1847]

$$\min_x F(x)$$

where F is convex smooth differentiable with ∇F being L -Lipschitz.

EXplicit Euler scheme:

$$x_{k+1} = x_k - \gamma_k \nabla F(x_k), \quad \gamma_k \in]0, 2/L[.$$

Proximal point algorithm [Rockafellar '76]

$$\min_x R(x)$$

with R being proper closed convex.



Gradient descent [Cauchy '1847]

$$\min_x F(x)$$

where F is convex smooth differentiable with ∇F being L -Lipschitz.

EXplicit Euler scheme:

$$x_{k+1} = x_k - \gamma_k \nabla F(x_k), \quad \gamma_k \in]0, 2/L[.$$

Proximal point algorithm [Rockafellar '76]

$$\min_x R(x)$$

with R being proper closed convex. Define “proximity operator” by

$$\text{prox}_{\gamma R}(v) \stackrel{\text{def}}{=} \operatorname{argmin}_x \gamma R(x) + \frac{1}{2} \|x - v\|^2.$$



Gradient descent [Cauchy '1847]

$$\min_x F(x)$$

where F is convex smooth differentiable with ∇F being L -Lipschitz.

EXplicit Euler scheme:

$$x_{k+1} = x_k - \gamma_k \nabla F(x_k), \quad \gamma_k \in]0, 2/L[.$$

Proximal point algorithm [Rockafellar '76]

$$\min_x R(x)$$

with R being proper closed convex. Define “proximity operator” by

$$\text{prox}_{\gamma R}(v) \stackrel{\text{def}}{=} \operatorname{argmin}_x \gamma R(x) + \frac{1}{2} \|x - v\|^2.$$

IMplicit Euler scheme:

$$x_{k+1} = \text{prox}_{\gamma_k R}(x_k), \quad \gamma_k > 0.$$

First-order methods



A rich class of first-order methods

$F + R$ Forward–Backward splitting [Lions & Mercier '79]

$R_1 + R_2$ Douglas–Rachford splitting [Douglas & Rachford '56; Lions & Mercier '79]

ADMM [Glowinski & Marrocco '75; Gabay & Mercier '76]...

$F + R(K \cdot)$ Primal–Dual splitting methods [Arrow, Hurwicz & Uzawa '58; Esser, Zhang & Chan '10; Chambolle & Pock '11]

$F + \sum_i R_i$ Generalized Forward–Backward splitting [Raguet, Fadili & Peyré '13]

— ...

Origins from numerical PDE back to 1950s, now ubiquitous in signal/image processing, inverse problems, data science, statistics, machine learning...



A rich class of first-order methods

$F + R$ Forward–Backward splitting [Lions & Mercier '79]

$R_1 + R_2$ Douglas–Rachford splitting [Douglas & Rachford '56; Lions & Mercier '79]

ADMM [Glowinski & Marrocco '75; Gabay & Mercier '76]...

$F + R(K \cdot)$ Primal–Dual splitting methods [Arrow, Hurwicz & Uzawa '58; Esser, Zhang & Chan '10; Chambolle & Pock '11]

$F + \sum_i R_i$ Generalized Forward–Backward splitting [Raguet, Fadili & Peyré '13]

— ...

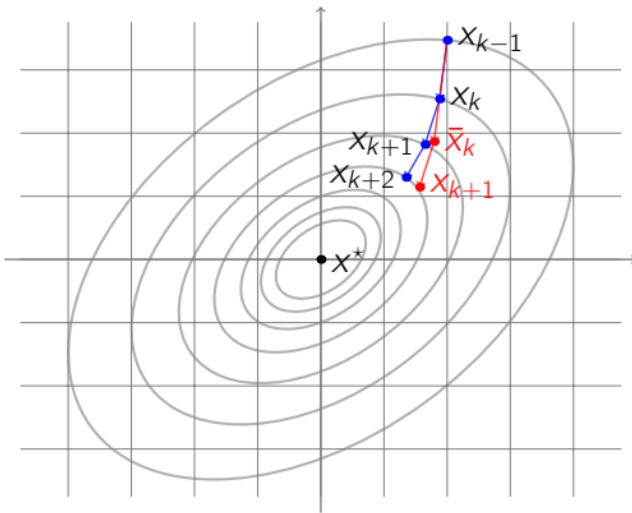
Fixed-point formulation

Let $\mathcal{F} : \mathcal{H} \rightarrow \mathcal{H}$ be non-expansive with $\text{fix}(\mathcal{F}) \stackrel{\text{def}}{=} \{z \in \mathcal{H} \mid z = \mathcal{F}(z)\} \neq \emptyset$,

$$z_{k+1} = \mathcal{F}(z_k).$$

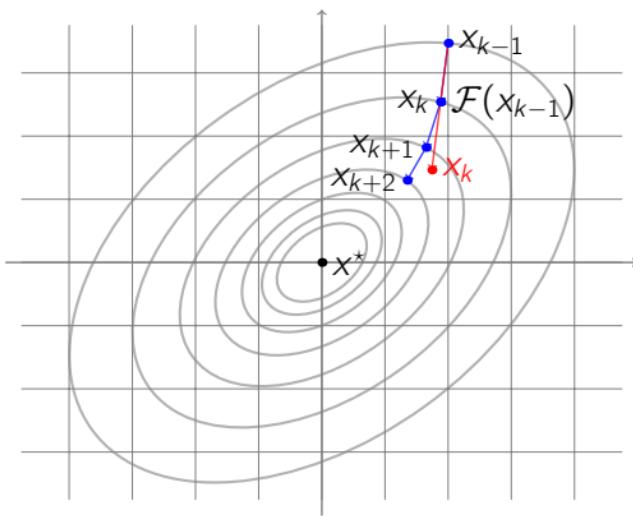
Definition - Inertial technique [Polyak '64, Nesterov '83, Beck & Teboulle '09]

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

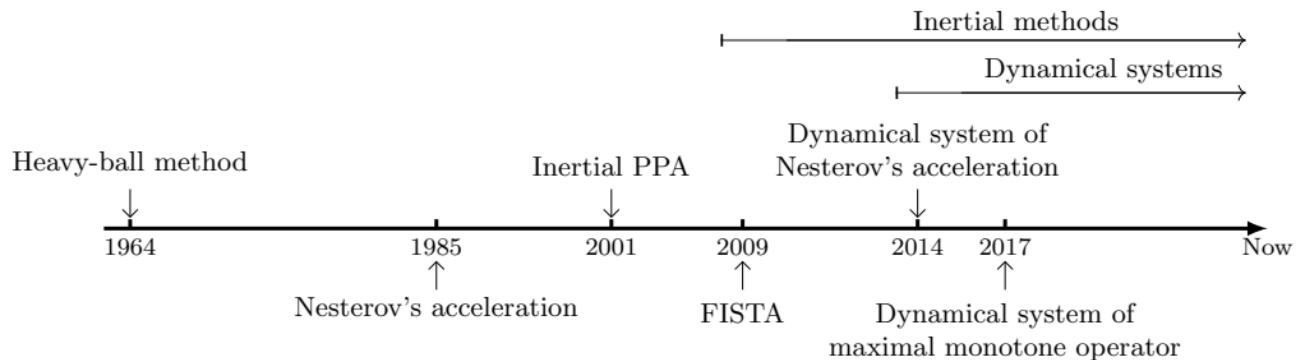


Definition - Successive over-relaxation [Richardson '1911, Young '50]

$$x_{k+1} = (1 - \lambda_k)x_k + \lambda_k \mathcal{F}(x_k) \xrightarrow{a_k = \lambda_k - 1} \begin{cases} \bar{x}_k = x_k + a_k(x_k - \bar{x}_{k-1}), \\ x_{k+1} = \mathcal{F}(\bar{x}_k). \end{cases}$$



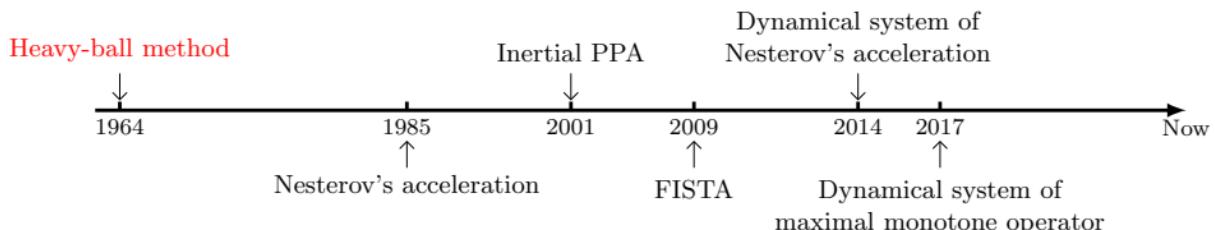
A brief history of inertial acceleration



A brief timeline of inertial acceleration.

A glimpse of references: [Polyak '64], [Polyak '87], [Nesterov '07], [Nesterov '83], [Nesterov '04], [Alvarez '00], [Attouch, Goudou & Redont '00], [Alvarez & Attouch '01], [H Attouch, Peypouquet & Redont '14], [Attouch & Peypouquet '16], [Attouch & Peypouquet '16], [Attouch & Peypouquet '19], [Adly & Attouch '20], [Moudafi & Oliny '03], [Maingé '08], [Beck & Teboulle '09], [O'Donoghue '12], [Su, Boyd & Candès '14], [Su, Boyd & Candès '16], [Chambolle & Dossal '15], [Aujol & Dossal '15], [Lorenz & Pock '14], [Bot, Csetnek & Hendrich '15], [Bot & Csetnek '16] [Liang, Fadili and Gabriél '17], [Apidopoulos, Aujol & Dossal '20], [França, Robinson and Vidal '18a], [França, Robinson and Vidal '18b], [Calatroni & Chambolle '19], [Chambolle & Tovey '21] and many others...

A brief history of inertial acceleration



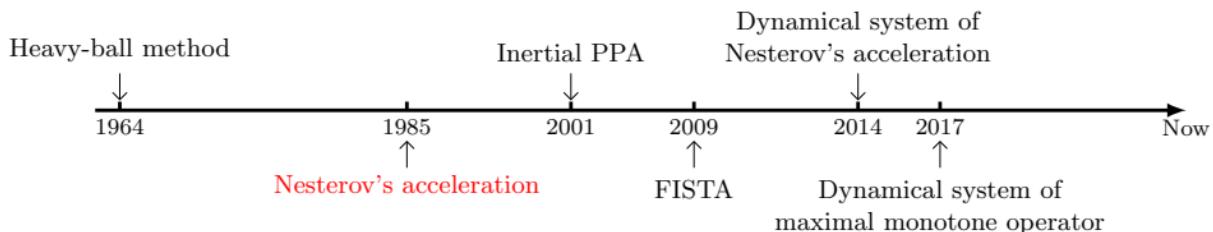
[Polyak '64] Consider minimizing $F \in C_L^1(\mathbb{R}^n)$

$$\begin{cases} \bar{x}_k = x_k + a_k(x_k - x_{k-1}), \\ x_{k+1} = \bar{x}_k - \gamma \nabla F(x_k). \end{cases}$$

- If $F \in C_{L,\mu}^2(\mathbb{R}^n)$, the optimal rate can be achieved $\rho^\star = \frac{1 - \sqrt{\mu/L}}{1 + \sqrt{\mu/L}}$.
- Heavy-ball dynamics with friction:

$$\ddot{x}(t) + \gamma \dot{x}(t) = -\nabla F(x(t)).$$

A brief history of inertial acceleration

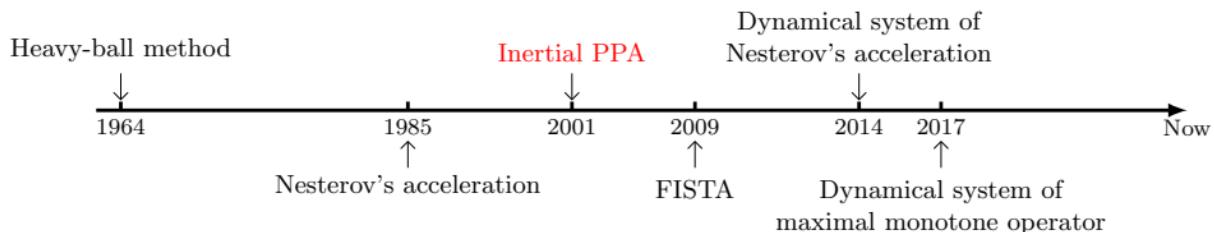


[Nesterov '83] Consider minimizing $F \in C_L^1(\mathbb{R}^n)$: $\gamma \leq 1/L$

$$\begin{cases} \bar{x}_k = x_k + \frac{k-1}{k+2}(x_k - x_{k-1}), \\ x_{k+1} = \bar{x}_k - \gamma \nabla F(\bar{x}_k). \end{cases}$$

- Convex case: $o(1/k) \rightarrow o(1/k^2)$.
- If $F \in C_{L,\mu}^1(\mathbb{R}^n)$, $1 - \frac{\mu}{L} \rightarrow 1 - \sqrt{\frac{\mu}{L}}$.

A brief history of inertial acceleration



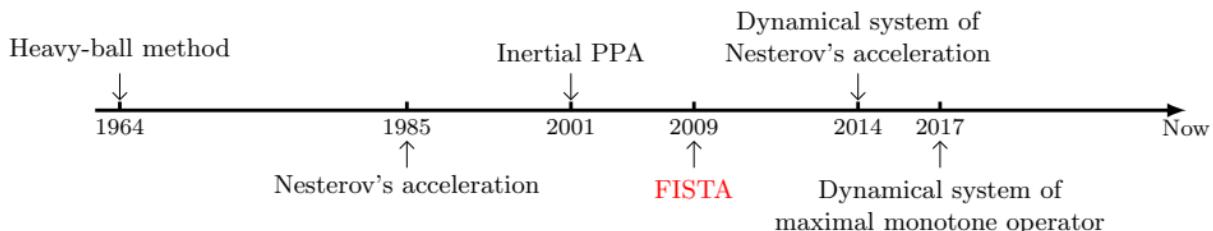
[Alvarez & Attouch '01] Consider maximal monotone inclusion problem $0 \in A(x)$

$$\ddot{x}(t) + \gamma \dot{x}(t) \in -A(x(t)).$$

- No co-coercivity.
- Guaranteed convergence for $a_k < \frac{1}{3}$ with

$$x_{k+1} = (\text{Id} + \lambda_k A)^{-1}(x_k + a_k(x_k - x_{k-1})).$$

A brief history of inertial acceleration



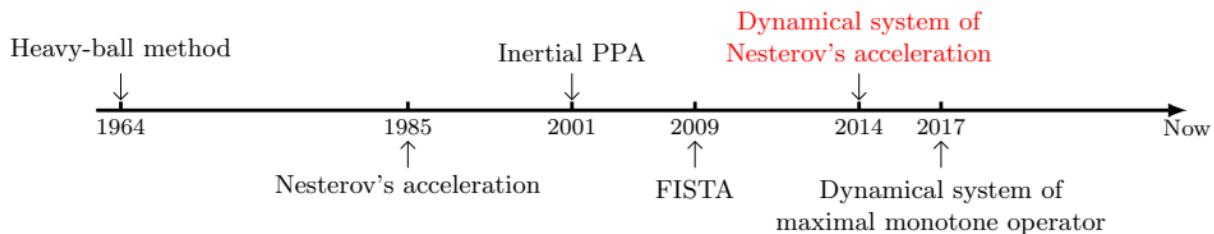
[Beck & Teboulle '09][Chambolle & Dossal '15] $F \in C_L^1(\mathbb{R}^n)$, $R \in \Gamma_0(\mathbb{R}^n)$:

$$\gamma \in [0, 1/L], d > 2$$

$$\begin{aligned}\bar{x}_k &= x_k + \frac{k-1}{k+d}(x_k - x_{k-1}), \\ x_{k+1} &= \text{prox}_{\gamma R}(\bar{x}_k - \gamma \nabla F(\bar{x}_k)).\end{aligned}$$

- Rate of convergence: sequence $o(1/\sqrt{k}) \rightarrow o(1/k)$, objective $o(1/k) \rightarrow o(1/k^2)$.
- Convergence of sequence.

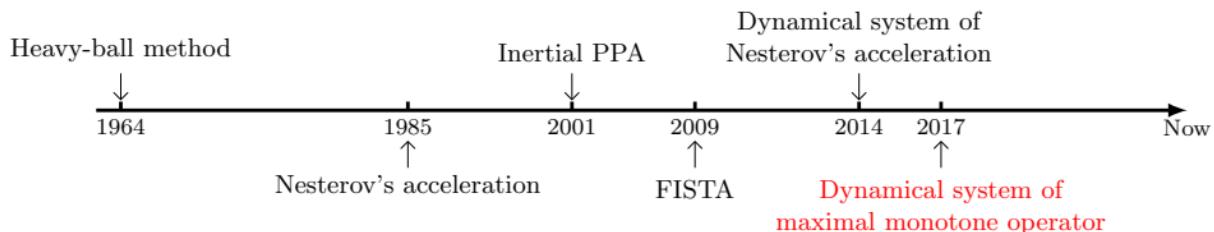
A brief history of inertial acceleration



[Su, Boyd & Candès '14] A differential equation for Nesterov's acceleration

$$\ddot{x}(t) + \frac{\alpha}{t} \dot{x}(t) = -\nabla F(x(t)).$$

A brief history of inertial acceleration



[Attouch & Peypouquet '19] Yosida's regularization

$$A_\lambda = \frac{1}{\lambda} (\text{Id} - (\text{Id} + \lambda A)^{-1}).$$

Damped inertial dynamics

$$\ddot{x}(t) + \gamma(t)\dot{x}(t) = -A_{\lambda(t)}(x(t)).$$

- $\gamma(t) = \frac{\alpha}{t}$ and $\lambda(t)\gamma^2(t) > 1$.
- Sequence rate of convergence: $o(1/\sqrt{k}) \rightarrow o(1/k)$.

Are we blessed with guaranteed acceleration?



Problem -

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

Let $\gamma > 0$

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



Problem -

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

Let $\gamma > 0$

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

Douglas–Rachford splitting [Douglas & Rachford '56]

$$z_{k+1} = \mathcal{F}_{\text{DR}}(z_k),$$

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

Are we blessed with guaranteed acceleration?



Problem -

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

Let $\gamma > 0$

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

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

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

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

Are we blessed with guaranteed acceleration?



Problem -

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

Let $\gamma > 0$

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

Regularized inertial Douglas–Rachford: let $\mathcal{J}_A = (\text{Id} + A)^{-1} = \mathcal{F}_{\text{DR}}$, $\alpha > 2$ and $\beta > 0$

$$\begin{aligned} y_k &= x_k + \frac{k-1}{k+\alpha}(x_k - x_{k-1}), \\ x_{k+1} &= \mathcal{J}_{\beta A_{\gamma_k}}(y_k). \end{aligned}$$

Let $\gamma_k = (1 + \epsilon) \frac{\beta}{\alpha^2} k^2$ with $\epsilon > \frac{2}{\alpha-2}$,

- $\|x_k - x_{k-1}\| = o(1/k)$, $x_k \rightarrow x^* \in \text{zer}(A)$.

Are we blessed with guaranteed acceleration?



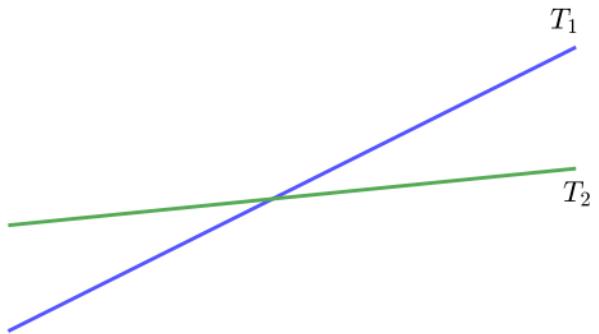
Problem - Feasibility problem in \mathbb{R}^2

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

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



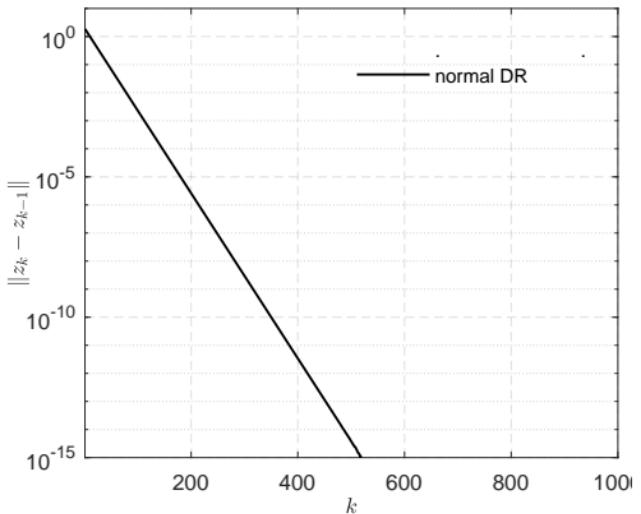
Problem - Feasibility problem in \mathbb{R}^2

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:

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



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

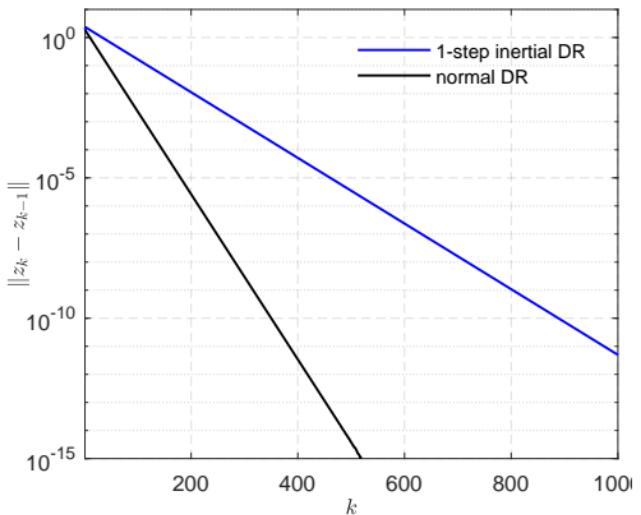
Let $T_1, T_2 \subset \mathbb{R}^2$ be two subspaces such that $T_1 \cap T_2 \neq \emptyset$,

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

Inertial Douglas–Rachford:

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

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



Are we blessed with guaranteed acceleration?



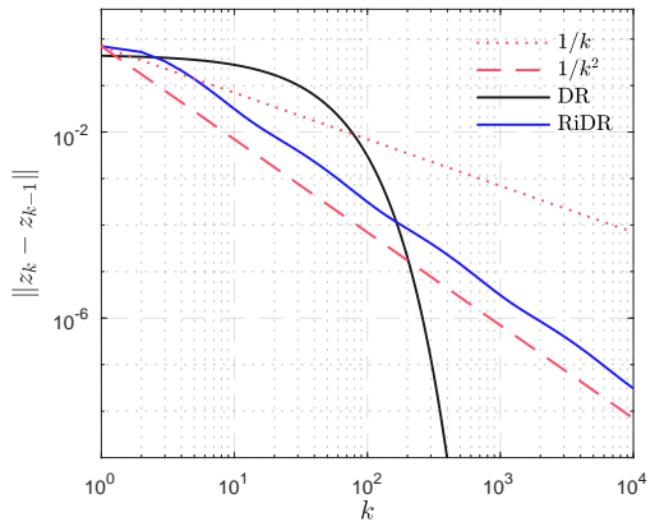
Problem - Feasibility problem in \mathbb{R}^2

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 $x \in T_1 \cap T_2$.

Regularized inertial Douglas–Rachford:

$$y_k = x_k + \frac{k-1}{k+\alpha}(x_k - x_{k-1}),$$
$$x_{k+1} = \mathcal{J}_{\beta A_{\gamma_k}}(y_k).$$



Are we blessed with guaranteed acceleration?



Consider minimizing

$$\min_{x \in \mathbb{R}^n} R(x) + J(z) \quad \text{s.t.} \quad z = Ax.$$

Proposition - Dynamics of ADMM [França et al '18a]

Let $V(x) = R(x) + J(Ax)$,

$$(A^T A) \left(\ddot{x}(t) + \frac{\alpha}{t} \dot{x}(t) \right) = -\nabla V(x(t)).$$

- Smoothness are required by both R and J .
- Improve the objective rate of convergence from $O(1/t)$ to $O(1/t^2)$.
- With non-smooth objective, Yosida regularization can be applied [França et al '18b].
- Not result available in discrete setting...

Are we blessed with guaranteed acceleration?



Some remarks

- Nesterov/FISTA provide optimal convergence rate.
- Generalization of inertial technique to first-order methods, or in general fixed-point iteration, is achievable:
 - Guaranteed sequence convergence.
 - **NO** acceleration guarantees. Unless stronger assumptions are imposed, *e.g.* strong convexity or Lipschitz smoothness.
- For a given method, *e.g.* Douglas–Rachford, the outcome its inertial/SOR versions is problem and **parameters** dependent.

A general acceleration framework with acceleration guarantees is missing!

Acceleration of First-order Methods

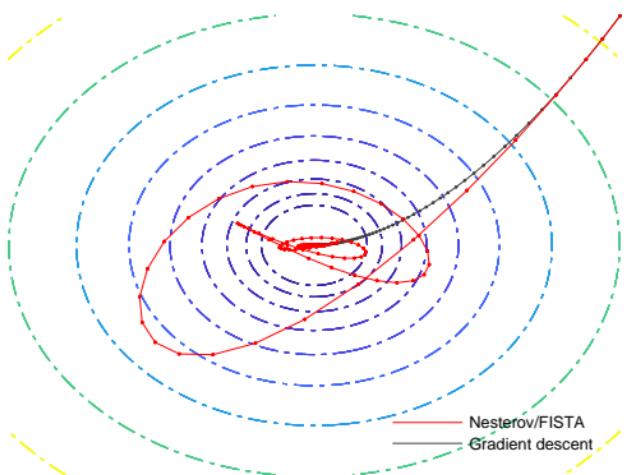
Trajectory of first-order methods

Straight line, logarithmic/elliptical spirals



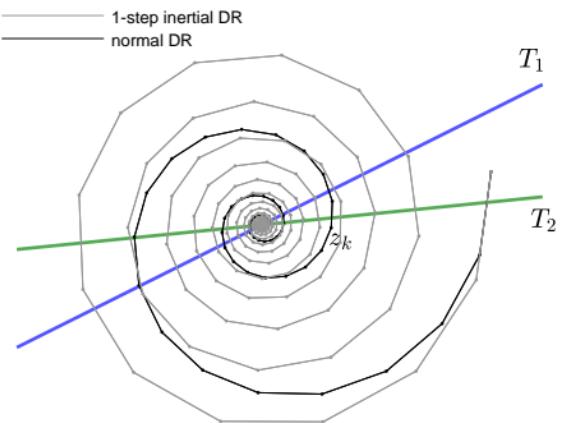
饮水思源 · 爱国荣校

What makes the difference?



Trajectory of $\{x_k\}_{k \in \mathbb{N}}$.

Projection of elliptical spiral in \mathbb{R}^4 onto \mathbb{R}^2 .

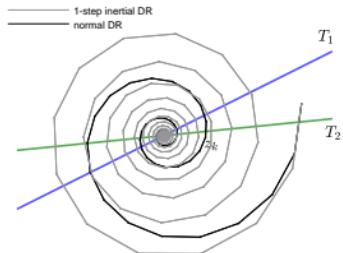
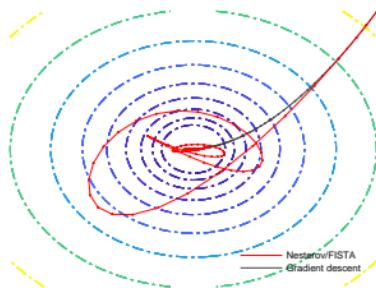
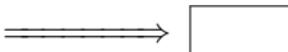


Trajectory of $\{z_k\}_{k \in \mathbb{N}}$.

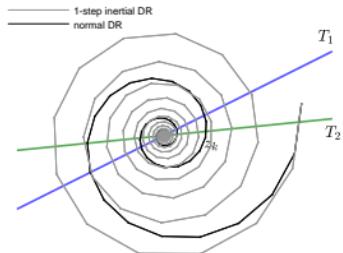
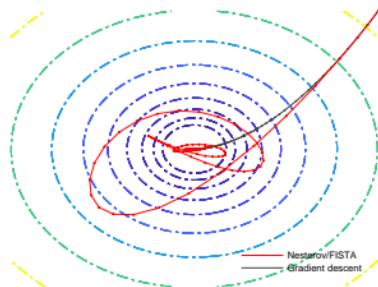
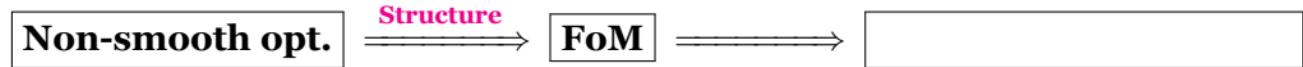
What makes the difference?



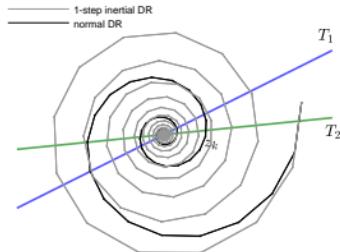
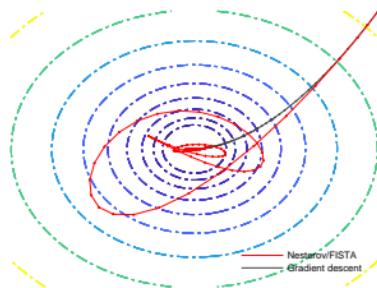
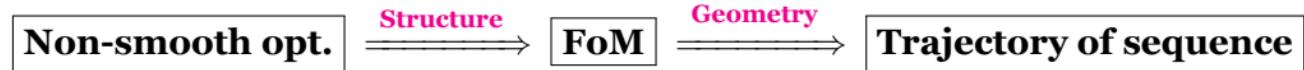
Non-smooth opt.



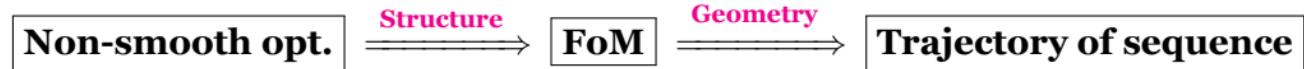
What makes the difference?



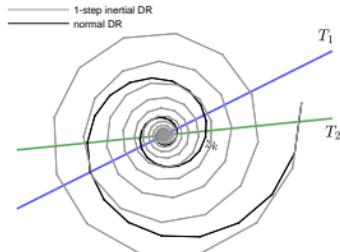
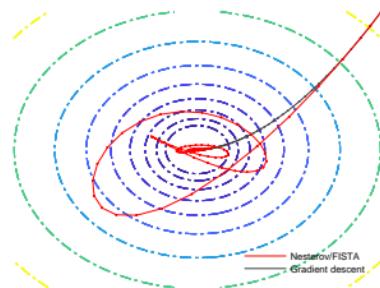
What makes the difference?



What makes the difference?



However, FoM are **non-linear** in general...



Definition - Partly smooth function [Lewis '03]

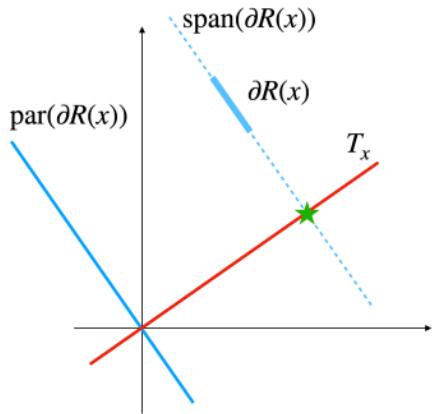
R is *partly smooth at x relative to a set \mathcal{M}_x containing x* if $\partial R(x) \neq \emptyset$ and

Smoothness \mathcal{M}_x is a C^2 -manifold, $R|_{\mathcal{M}_x}$ is C^2 near x .

Sharpness Tangent space $\mathcal{T}_{\mathcal{M}_x}(x)$ is $T_x \stackrel{\text{def}}{=} \text{par}(\partial R(x))^\perp$.

Continuity $\partial R : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ is continuous along \mathcal{M}_x near x .

$\text{par}(C)$: sub-space parallel to C , where $C \subset \mathbb{R}^n$ is a non-empty convex set.



Examples:

- $\ell_1, \ell_{1,2}, \ell_\infty$ -norm
- Nuclear norm
- Total variation
- Partly smooth sets...

Definition - Partly smooth function [Lewis '03]

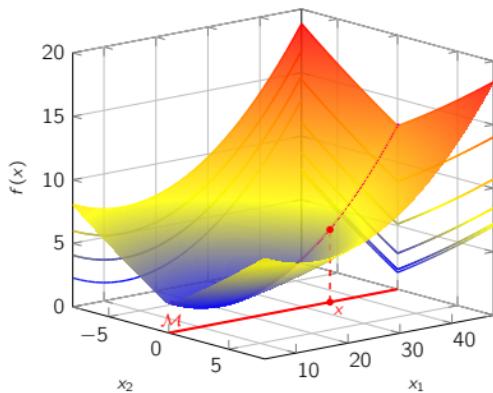
R is *partly smooth at x relative to a set \mathcal{M}_x containing x* if $\partial R(x) \neq \emptyset$ and

Smoothness \mathcal{M}_x is a C^2 -manifold, $R|_{\mathcal{M}_x}$ is C^2 near x .

Sharpness Tangent space $\mathcal{T}_{\mathcal{M}_x}(x)$ is $T_x \stackrel{\text{def}}{=} \text{par}(\partial R(x))^\perp$.

Continuity $\partial R : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ is continuous along \mathcal{M}_x near x .

$\text{par}(C)$: sub-space parallel to C , where $C \subset \mathbb{R}^n$ is a non-empty convex set.



Examples:

- $\ell_1, \ell_{1,2}, \ell_\infty$ -norm
- Nuclear norm
- Total variation
- Partly smooth sets...

Trajectory of first-order methods



Framework for analyzing the local trajectory of FoM

First-order method (non-linear)



Convergence & Non-degeneracy: finite identification of \mathcal{M}



Local linearization along \mathcal{M} : matrix M (linear)



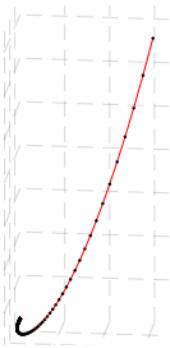
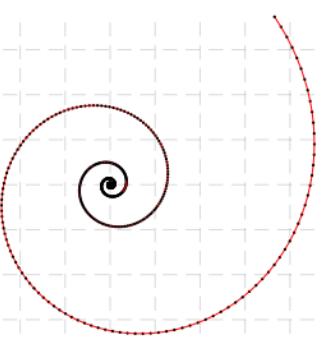
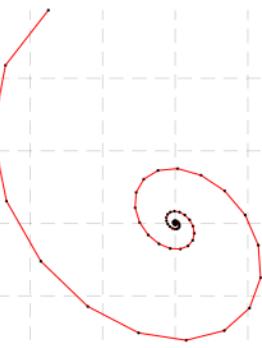
Spectral properties of M



Local trajectory

Trajectory of first-order methods



Forward–Backward	Douglas–Rachford/ADMM	Primal–Dual
<p>M is similar to a symmetric matrix with real eigenvalues in $] -1, 1]$.</p> 	<p>Both functions are polyhedral, M is normal with complex eigenvalues of the form $\cos(\theta)e^{\pm i\theta}$.</p> 	<p>Both functions are polyhedral, M is block diagonalizable.</p> 

NB: For DR/ADMM, if smoothness or (local) strong convexity is posed, **straight-line** trajectory can be obtained under proper parameters.

Acceleration of First-order Methods

Adaptive acceleration for FoM

Trajectory following linear prediction



饮水思源 · 爱国荣校

Linear prediction: illustration



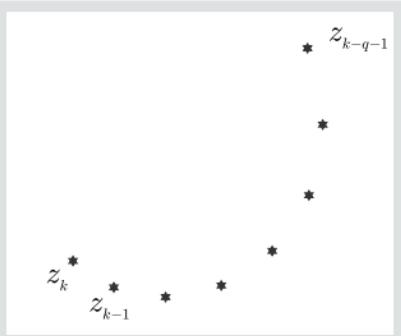
Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, how to predict z_{k+1} ?

Define $\{v_{k-j} \stackrel{\text{def}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Fit v_k using past directions v_{k-1}, \dots, v_{k-q} :

$$c_k \stackrel{\text{def}}{=} \operatorname{argmin}_{c \in \mathbb{R}^q} \left\| \sum_{j=1}^q c_j v_{k-j} - v_k \right\|^2.$$

-
-
-



Linear prediction: illustration



Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, how to predict z_{k+1} ?

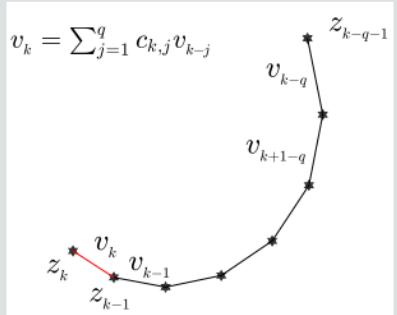
Define $\{v_{k-j} \stackrel{\text{def}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Fit v_k using past directions v_{k-1}, \dots, v_{k-q} :

$$c_k \stackrel{\text{def}}{=} \operatorname{argmin}_{c \in \mathbb{R}^q} \left\| \sum_{j=1}^q c_j v_{k-j} - v_k \right\|^2.$$

- Suppose z_{k+1} is given, we can compute c_{k+1} as above...
- Approximation

$$v_{k+1} = \sum_j \color{red}{c_{k+1,j}} v_{k+1-j} \approx \sum_j \color{red}{c_{k,j}} v_{k+1-j}.$$



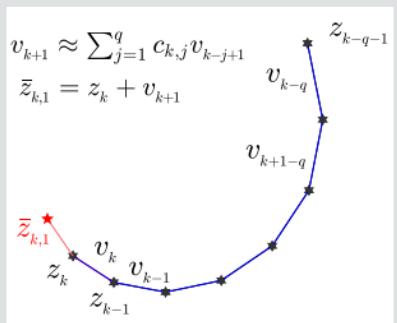
Linear prediction: illustration



Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, how to predict z_{k+1} ?

Define $\{v_{k-j} \stackrel{\text{def}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

■ Approximation $z_{k+1} \approx \bar{z}_{k,1} \stackrel{\text{def}}{=} z_k + \sum_{j=1}^q c_{k,j} v_{k-j+1}$.



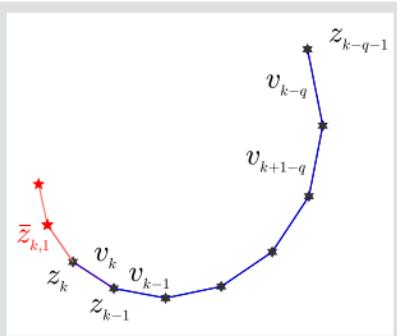
Linear prediction: illustration



Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, how to predict z_{k+1} ?

Define $\{v_{k-j} \stackrel{\text{def}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Approximation $z_{k+1} \approx \bar{z}_{k,1} \stackrel{\text{def}}{=} z_k + \sum_{j=1}^q c_{k,j} v_{k-j+1}$.
- Repeat on $\{z_{k-j}\}_{j=0}^q \cup \{\bar{z}_{k,1}\}$ and so on...



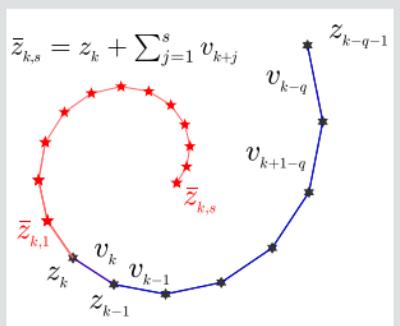
Linear prediction: illustration



Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, how to predict z_{k+1} ?

Define $\{v_{k-j} \stackrel{\text{def}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Approximation $z_{k+1} \approx \bar{z}_{k,1} \stackrel{\text{def}}{=} z_k + \sum_{j=1}^q c_{k,j} v_{k-j+1}$.
- Repeat on $\{z_{k-j}\}_{j=0}^q \cup \{\bar{z}_{k,1}\}$ and so on...
- Repeat s times: $z_{k+s} \approx \bar{z}_{k,s} = z_k + \sum_{j=1}^s v_{k+j}$.



Linear prediction: illustration

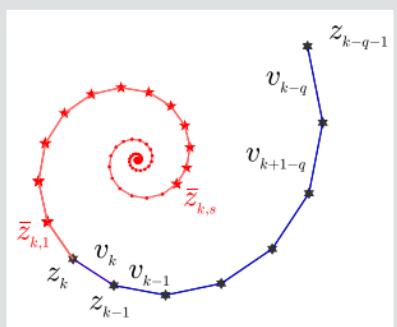


Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, how to predict z_{k+1} ?

Define $\{v_{k-j} \stackrel{\text{def}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Approximation $z_{k+1} \approx \bar{z}_{k,1} \stackrel{\text{def}}{=} z_k + \sum_{j=1}^q c_{k,j} v_{k-j+1}$.
- Repeat on $\{z_{k-j}\}_{j=0}^q \cup \{\bar{z}_{k,1}\}$ and so on...
- Repeat s times: $z_{k+s} \approx \bar{z}_{k,s} = z_k + \sum_{j=1}^s v_{k+j}$.
- Let $s \rightarrow +\infty$, if converges

$$z^* \approx \bar{z}_{k,+\infty} = z_k + \sum_{j=1}^{+\infty} v_{k+j}.$$



Linear prediction: illustration

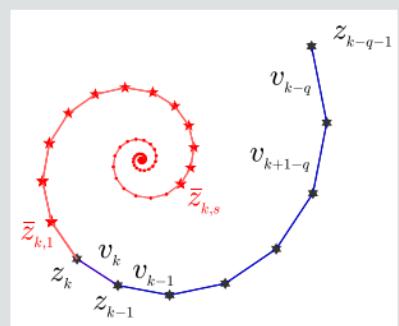


Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, how to predict z_{k+1} ?

Define $\{v_{k-j} \stackrel{\text{def}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Approximation $z_{k+1} \approx \bar{z}_{k,1} \stackrel{\text{def}}{=} z_k + \sum_{j=1}^q c_{k,j} v_{k-j+1}$.
- Repeat on $\{z_{k-j}\}_{j=0}^q \cup \{\bar{z}_{k,1}\}$ and so on...
- Repeat s times: $z_{k+s} \approx \bar{z}_{k,s} = z_k + \sum_{j=1}^s v_{k+j}$.
- Let $s \rightarrow +\infty$, if converges

$$z^* \approx \bar{z}_{k,+\infty} = z_k + \sum_{j=1}^{+\infty} v_{k+j}.$$



The ***s*-step extrapolation** is $\bar{z}_{k,s} = z_k + \mathcal{E}_{s,q,k}$, where $\mathcal{E}_{s,q,k} = \sum_{j=1}^q \hat{c}_j v_{k-j+1}$ and

$$\hat{c} \stackrel{\text{def}}{=} \left(\sum_{j=1}^s H(c^k)^j \right)_{(:,1)} \quad \text{with} \quad H(c^k) \stackrel{\text{def}}{=} \begin{bmatrix} c_1 & 1 & & \\ c_2 & & 1 & \\ \vdots & & & \ddots \\ c_{q-1} & 0 & & 1 \\ c_q & & \dots & 0 \end{bmatrix}.$$

Given first-order method

$$z_{k+1} = \mathcal{F}(z_k).$$

Algorithm - A²FoM via linear prediction

Let $s \geq 1, q \geq 1$ be integers. Let $z_0 \in \mathbb{R}^n$ and $\bar{z}_0 = z_0$, set $D = 0 \in \mathbb{R}^{n \times (q+1)}$:

- For $k \geq 1$:

$$z_k = \mathcal{F}(\bar{z}_{k-1}),$$

$$v_k = z_k - z_{k-1},$$

$$D = [v_k, D(:, 1 : q)].$$

- If $\text{mod}(k, q+2) = 0$: compute c and H_c ,

If $\rho(H_c) < 1$: $\bar{z}_k = z_k + V_k \left(\sum_{i=1}^s H_c^i \right)_{(:,1)};$

else: $\bar{z}_k = z_k.$

If $\text{mod}(k, q+2) \neq 0$: $\bar{z}_k = z_k.$



Remarks

- Every $(q + 2)$ -iteration we apply 1-step LP.
- Only apply the linear prediction when $\rho(H_c) < 1$.
- Extra memory cost $n \times (q + 1)$ (the difference vector matrix). Usually $q \leq 10$.
- Extra computation cost, $q^2 n$ from V_{k-1}^+ .
-
-
-



Remarks

- Every $(q + 2)$ -iteration we apply 1-step LP.
- Only apply the linear prediction when $\rho(H_c) < 1$.
- Extra memory cost $n \times (q + 1)$ (the difference vector matrix). Usually $q \leq 10$.
- Extra computation cost, $q^2 n$ from V_{k-1}^+ .
- Conditional convergence can be obtained by treating LP as **perturbation error**,

$$z_{k+1} = \mathcal{F}(z_k + \epsilon_k).$$

- Weighted LP

$$\bar{z}_k = z_k + \textcolor{red}{a_k} V_k \left(\sum_{i=1}^s H_c^i \right)_{(:,1)},$$

with a_k updated online.





Remarks

- Every $(q + 2)$ -iteration we apply 1-step LP.
- Only apply the linear prediction when $\rho(H_c) < 1$.
- Extra memory cost $n \times (q + 1)$ (the difference vector matrix). Usually $q \leq 10$.
- Extra computation cost, $q^2 n$ from V_{k-1}^+ .
- Conditional convergence can be obtained by treating LP as **perturbation error**,

$$z_{k+1} = \mathcal{F}(z_k + \epsilon_k).$$

- Weighted LP

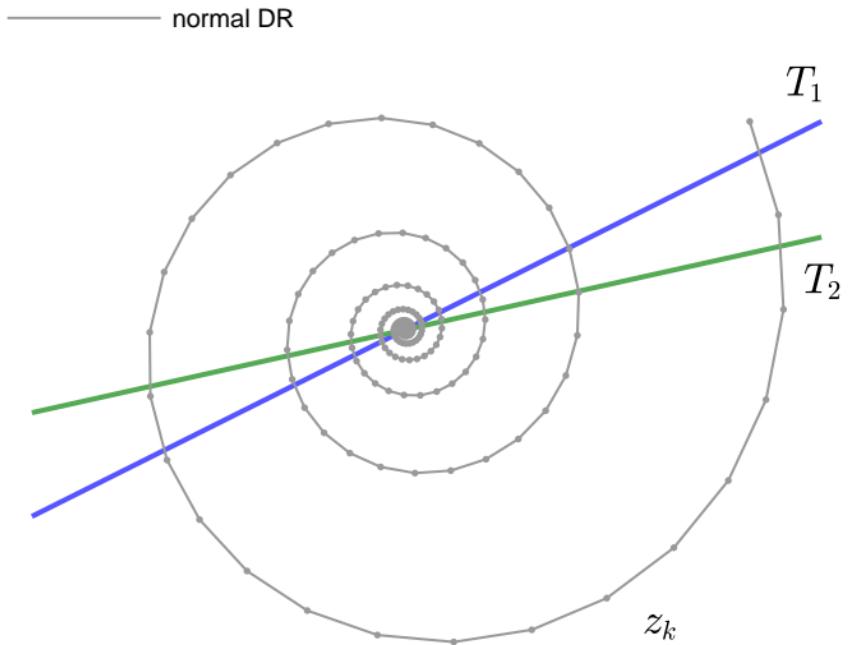
$$\bar{z}_k = z_k + \textcolor{red}{a_k} V_k \left(\sum_{i=1}^s H_c^i \right)_{(:,1)},$$

with a_k updated online.

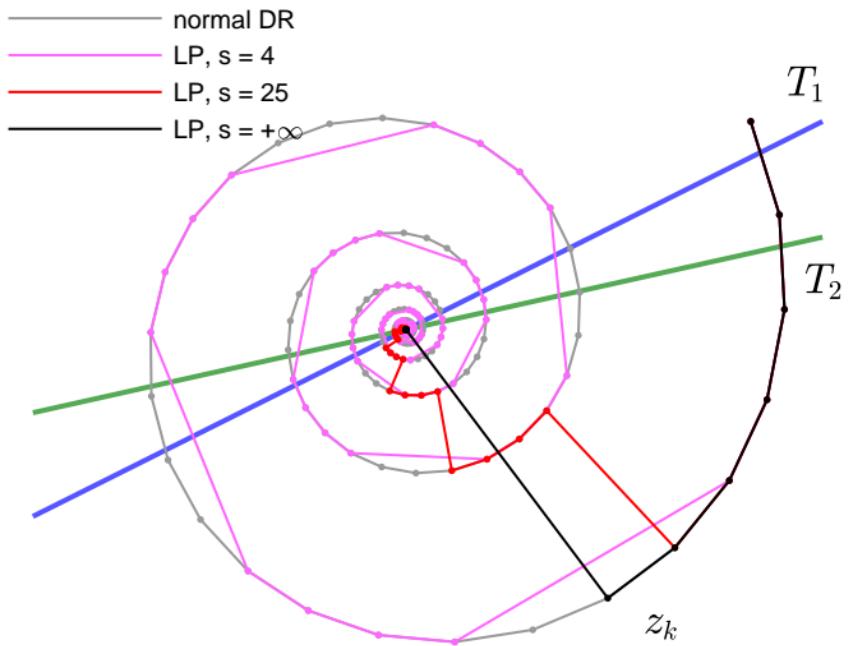
- If $\rho(H_c) < 1$, the Neumann series is convergent

$$\sum_{i=0}^{+\infty} H_c^i = (\text{Id} - H_c)^{-1}.$$

Example: Douglas–Rachford continue



Example: Douglas–Rachford continue



Acceleration of First-order Methods

Relation with previous work

Polynomial extrapolation, regularized non-linear acceleration

Convergence acceleration



Given a sequence $\{z_k\}_{k \in \mathbb{N}}$ which converges to z^* . Can we generate another sequence $\{\bar{z}_k\}_{k \in \mathbb{N}}$ such that $\|\bar{z}_k - z^*\| = o(\|z_k - z^*\|)$?

Convergence acceleration



Given a sequence $\{z_k\}_{k \in \mathbb{N}}$ which converges to z^* . Can we generate another sequence $\{\bar{z}_k\}_{k \in \mathbb{N}}$ such that $\|\bar{z}_k - z^*\| = o(\|z_k - z^*\|)$?

This is called **convergence acceleration** and is well-established in numerical analysis:

1927 Aitkin's Δ -process.

1965 Andersen's acceleration.

1970's Vector extrapolation techniques such as minimal polynomial extrapolation (MPE) and reduced rank extrapolation (RRE) [Sidi '17].

Recent Regularized non-linear acceleration (RNA) is a regularised version of RRE introduced by [Scieur, D'Aspremont, Bach '16].



Polynomial extrapolation [Cabay & Jackson '76]

Consider $z_{k+1} = Mz_k + d$ with $\rho(M) < 1$ such that $z_k \rightarrow z^*$:

- $z_k - z^* = M(z_{k-1} - z^*) = M^k(z_0 - z^*)$,
- If $P(\lambda) = \sum_{j=0}^q c_j \lambda^j$ is the minimal polynomial of M w.r.t. $z_0 - z^*$, that is

$$P(M)(z_0 - z^*) = \sum_{j=0}^q c_j M^j(z_0 - z^*) = 0.$$

then $z^* = \frac{\sum_{j=0}^q c_j z_j}{\sum_j c_j}$.

- The coefficients c can be computed without knowledge of z^* :

$$V_q c_{(0:q-1)} = -v_{q+1}$$

where $V_q = [v_1 | v_2 | \cdots | v_q]$ and $v_j = z_j - z_{j-1}$.

Vector extrapolation methods with applications (SIAM, 2017) by Avram Sidi.

Minimal polynomial extrapolation (MPE)

Let $z_0 = \bar{z}$:

S.1 Generate points $\{z_j\}_{j=0}^{q+1}$ and let $v_j = z_j - z_{j-1}$.

S.2 Let $c \in \mathbb{R}^{q+1}$ be s.t. $c_q = 1$ and $V_q c_{(0:q-1)} = -v_{q+1}$ where $V_q = [v_1 | \cdots | v_q]$.
For $j \in [0, q-1]$, $\tilde{c}_j \stackrel{\text{def}}{=} c_j / (\sum_{j=0}^q c_j)$.

S.3 $\bar{z} \stackrel{\text{def}}{=} \sum_{j=0}^q \tilde{c}_j z_j$.

Vector extrapolation techniques



Vector extrapolation methods with applications (SIAM, 2017) by Avram Sidi.

Minimal polynomial extrapolation (MPE)

Let $z_0 = \bar{z}$:

S.1 Generate points $\{z_j\}_{j=0}^{q+1}$ and let $v_j = z_j - z_{j-1}$.

S.2 Let $c \in \mathbb{R}^{q+1}$ be s.t. $c_q = 1$ and $V_q c_{(0:q-1)} = -v_{q+1}$ where $V_q = [v_1 | \cdots | v_q]$.
For $j \in [0, q-1]$, $\tilde{c}_j \stackrel{\text{def}}{=} c_j / (\sum_{j=0}^q c_j)$.

S.3 $\bar{z} \stackrel{\text{def}}{=} \sum_{j=0}^q \tilde{c}_j z_j$.

Reduced rank extrapolation (RRE)

[Andersen '65; Kaniel & Stein '74; Eddy '79; Mešina '77] Replace step **S.2** by

$$\tilde{c} \in \operatorname{argmin}_c \|V_{q+1} c\| \quad \text{subject to } \mathbf{1}^T c = 1.$$

Vector extrapolation techniques



Vector extrapolation methods with applications (SIAM, 2017) by Avram Sidi.

Minimal polynomial extrapolation (MPE)

Let $z_0 = \bar{z}$:

S.1 Generate points $\{z_j\}_{j=0}^{q+1}$ and let $v_j = z_j - z_{j-1}$.

S.2 Let $c \in \mathbb{R}^{q+1}$ be s.t. $c_q = 1$ and $V_q c_{(0:q-1)} = -v_{q+1}$ where $V_q = [v_1 | \cdots | v_q]$.
For $j \in [0, q-1]$, $\tilde{c}_j \stackrel{\text{def}}{=} c_j / (\sum_{j=0}^q c_j)$.

S.3 $\bar{z} \stackrel{\text{def}}{=} \sum_{j=0}^q \tilde{c}_j z_j$.

Reduced rank extrapolation (RRE)

[Andersen '65; Kaniel & Stein '74; Eddy '79; Mešina '77] Replace step **S.2** by

$$\tilde{c} \in \operatorname{argmin}_c \|V_{q+1} c\| \quad \text{subject to } \mathbf{1}^T c = 1.$$

LP is equivalent to MPE with **S.3** replaced by $\bar{z} \stackrel{\text{def}}{=} \sum_{j=0}^q \tilde{c}_j z_{j+1}$.

Relationship between MPE and LP



Our derivation

- is based on sequence **trajectory**.
- motivates checking $\rho(H_c) < 1$.

Relationship between MPE and LP

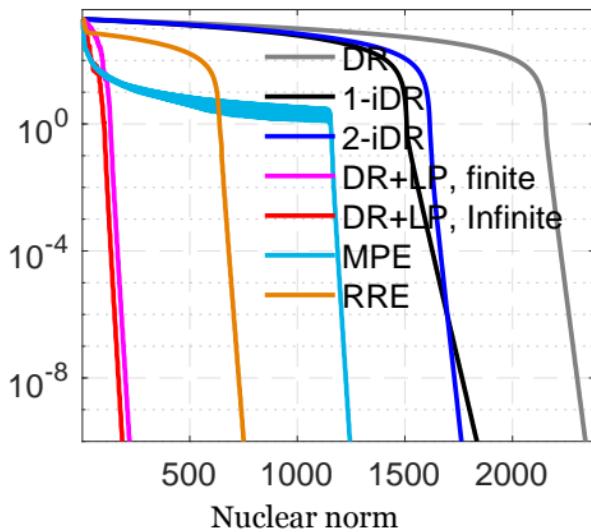
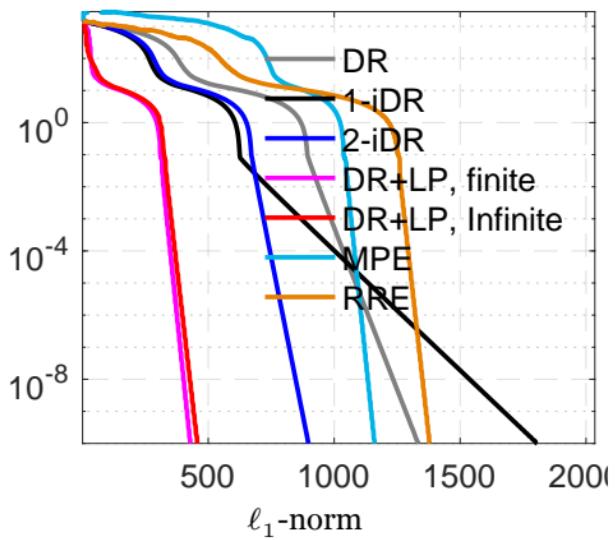


Our derivation

- is based on sequence **trajectory**.
- motivates checking $\rho(H_c) < 1$.

Example: Solve with DR

$$\min_x R(x) \text{ subject to } Ax = b.$$



Acceleration guarantees



When $z_{k+1} - z_k = M(z_k - z_{k-1})$,

$$\|\bar{z}_{k,s} - z^*\| \leq \|z_{k+s} - z^*\| + B\epsilon_k$$

where $\epsilon_k = \|V_{k-1}c - v_k\|$ and $B \stackrel{\text{def}}{=} \sum_{\ell=1}^s \|M^\ell\| |\sum_{i=0}^{s-\ell} (H_c^i)_{(1,1)}|$.

When $z_{k+1} - z_k = M(z_k - z_{k-1})$,

$$\|\bar{z}_{k,s} - z^*\| \leq \|z_{k+s} - z^*\| + B\epsilon_k$$

where $\epsilon_k = \|V_{k-1}c - v_k\|$ and $B \stackrel{\text{def}}{=} \sum_{\ell=1}^s \|M^\ell\| |\sum_{i=0}^{s-\ell} (H_c^i)_{(1,1)}|$.

Asymptotic bound ($k \rightarrow \infty$):

$$\epsilon_k = O(|\lambda_{q+1}|^k)$$

where λ_{q+1} is the $(q+1)^{th}$ largest eigenvalue. Without extrapolation, we just have $O(|\lambda_1|^k)$.

Acceleration guarantees



When $z_{k+1} - z_k = M(z_k - z_{k-1})$,

$$\|\bar{z}_{k,s} - z^*\| \leq \|z_{k+s} - z^*\| + B\epsilon_k$$

where $\epsilon_k = \|V_{k-1}c - v_k\|$ and $B \stackrel{\text{def}}{=} \sum_{\ell=1}^s \|M^\ell\| |\sum_{i=0}^{s-\ell} (H_c^i)_{(1,1)}|$.

Asymptotic bound ($k \rightarrow \infty$):

$$\epsilon_k = O(|\lambda_{q+1}|^k)$$

where λ_{q+1} is the $(q+1)^{th}$ largest eigenvalue. Without extrapolation, we just have $O(|\lambda_1|^k)$.

Non-asymptotic bound: If $\Sigma(M) \subset [\alpha, \beta]$ with $-1 < \alpha < \beta < 1$, then

$$B\epsilon_k \leq K\beta^{k-q} \left(\frac{\sqrt{\eta}-1}{\sqrt{\eta}+1} \right)^q, \quad \text{where } \eta = \frac{1-\alpha}{1-\beta}$$

When $z_{k+1} - z_k = M(z_k - z_{k-1})$,

$$\|\bar{z}_{k,s} - z^*\| \leq \|z_{k+s} - z^*\| + B\epsilon_k$$

where $\epsilon_k = \|V_{k-1}c - v_k\|$ and $B \stackrel{\text{def}}{=} \sum_{\ell=1}^s \|M^\ell\| |\sum_{i=0}^{s-\ell} (H_c^i)_{(1,1)}|$.

Asymptotic bound ($k \rightarrow \infty$):

$$\epsilon_k = O(|\lambda_{q+1}|^k)$$

where λ_{q+1} is the $(q+1)^{th}$ largest eigenvalue. Without extrapolation, we just have $O(|\lambda_1|^k)$.

Non-asymptotic bound: If $\Sigma(M) \subset [\alpha, \beta]$ with $-1 < \alpha < \beta < 1$, then

$$B\epsilon_k \leq K\beta^{k-q} \left(\frac{\sqrt{\eta} - 1}{\sqrt{\eta} + 1} \right)^q, \quad \text{where } \eta = \frac{1 - \alpha}{1 - \beta}$$

- Similar error bounds also hold for MPE and RRE [Sidi '98].

When $z_{k+1} - z_k = M(z_k - z_{k-1})$,

$$\|\bar{z}_{k,s} - z^*\| \leq \|z_{k+s} - z^*\| + B\epsilon_k$$

where $\epsilon_k = \|V_{k-1}c - v_k\|$ and $B \stackrel{\text{def}}{=} \sum_{\ell=1}^s \|M^\ell\| |\sum_{i=0}^{s-\ell} (H_c^i)_{(1,1)}|$.

Asymptotic bound ($k \rightarrow \infty$):

$$\epsilon_k = O(|\lambda_{q+1}|^k)$$

where λ_{q+1} is the $(q+1)^{th}$ largest eigenvalue. Without extrapolation, we just have $O(|\lambda_1|^k)$.

Non-asymptotic bound: If $\Sigma(M) \subset [\alpha, \beta]$ with $-1 < \alpha < \beta < 1$, then

$$B\epsilon_k \leq K\beta^{k-q} \left(\frac{\sqrt{\eta} - 1}{\sqrt{\eta} + 1} \right)^q, \quad \text{where } \eta = \frac{1-\alpha}{1-\beta}$$

- Similar error bounds also hold for MPE and RRE [Sidi '98].
- For PD, DR with polyhedral functions, guaranteed acceleration with $q = 2$.

Regularised non-linear acceleration (RNA): regularised RRE

[Scieur, D'Aspremont, Bach '16] studied RRE for the case of

$$z_{k+1} - z^* = A(z_k - z^*) + O(\|z_k - z^*\|^2)$$

where A is **symmetric** with $0 \preceq A \preceq \sigma \text{Id}$, $\sigma < 1$.

Regularised non-linear acceleration (RNA): regularised RRE

[Scieur, D'Aspremont, Bach '16] studied RRE for the case of

$$z_{k+1} - z^* = A(z_k - z^*) + O(\|z_k - z^*\|^2)$$

where A is **symmetric** with $0 \preceq A \preceq \sigma \text{Id}$, $\sigma < 1$.

To deal with the possible ill-conditioning of V_q , regularise with $\lambda > 0$:

$$\tilde{c} \in \operatorname{Argmin}_c \|c^T(V_q^T V_q + \lambda \text{Id})c\| \quad \text{subject to } \mathbf{1}^T c = 1.$$

Regularised non-linear acceleration (RNA): regularised RRE

[Scieur, D'Aspremont, Bach '16] studied RRE for the case of

$$z_{k+1} - z^* = A(z_k - z^*) + O(\|z_k - z^*\|^2)$$

where A is **symmetric** with $0 \preceq A \preceq \sigma \text{Id}$, $\sigma < 1$.

To deal with the possible ill-conditioning of V_q , regularise with $\lambda > 0$:

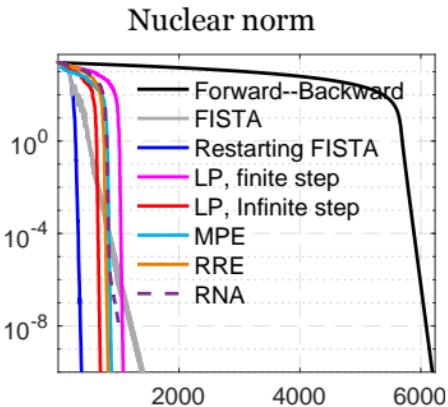
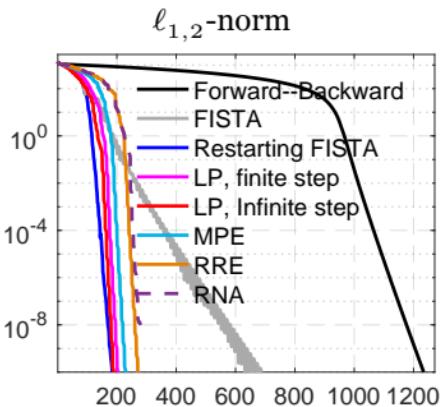
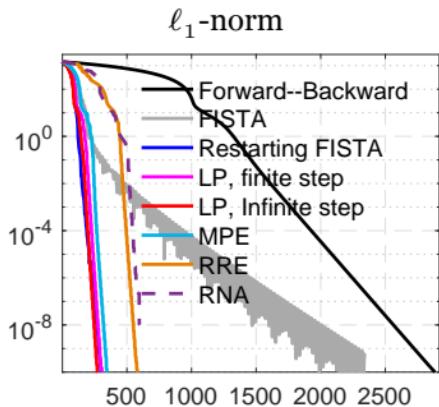
$$\tilde{c} \in \operatorname{Argmin}_c \|c^T(V_q^T V_q + \lambda \text{Id})c\| \quad \text{subject to } \mathbf{1}^T c = 1.$$

- In practice, grid search on optimal $\lambda \in [\lambda_{\min}, \lambda_{\max}]$. Potentially many evaluations of the objective.
- The angle between $z_k - z_{k-1}$ and $z_{k+1} - z_k$ converges to zero, intuitively, this is the regime where standard inertial works well...

LASSO-type problems with FB



NB: Recall the straight-line trajectory of FB.



Note the performance of restarted-FISTA [O'Donoghue & Candès '12]!

Acceleration of First-order Methods

Numerical experiments

Some more results



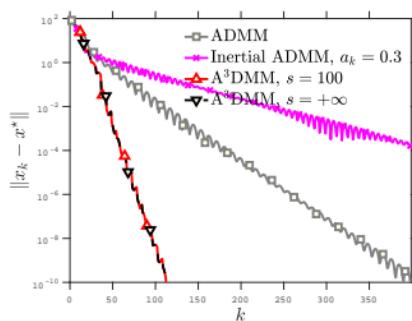
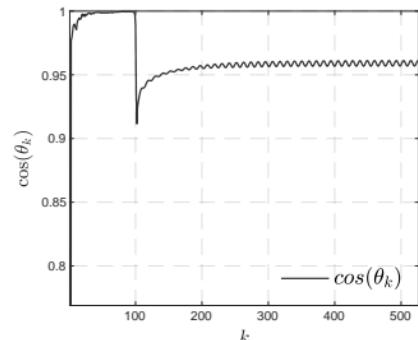
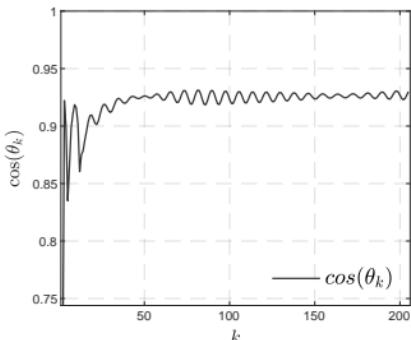
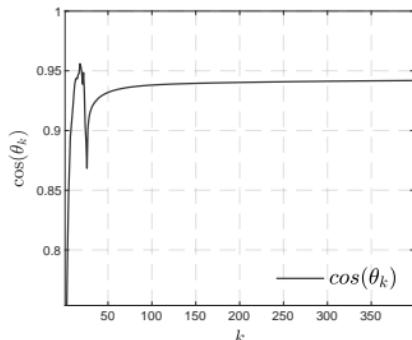
饮水思源 · 爱国荣校

Experiment: 2 non-smooth terms

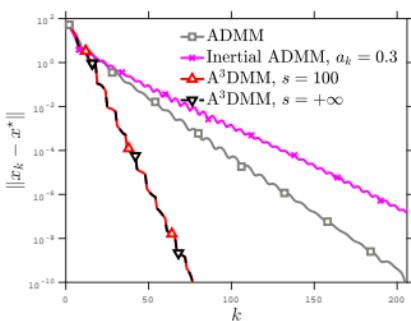


Basis pursuit type problem with $\Omega \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n : Kx = f\}$:

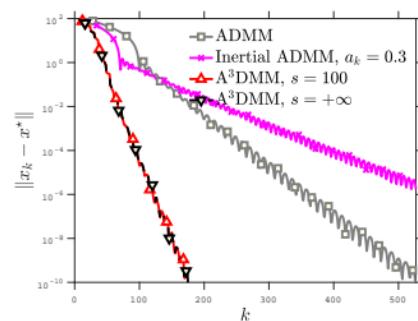
$$\min_{x, y \in \mathbb{R}^n} R(x) + \iota_\Omega(y) \quad \text{such that} \quad x - y = 0.$$



ℓ_1 -norm

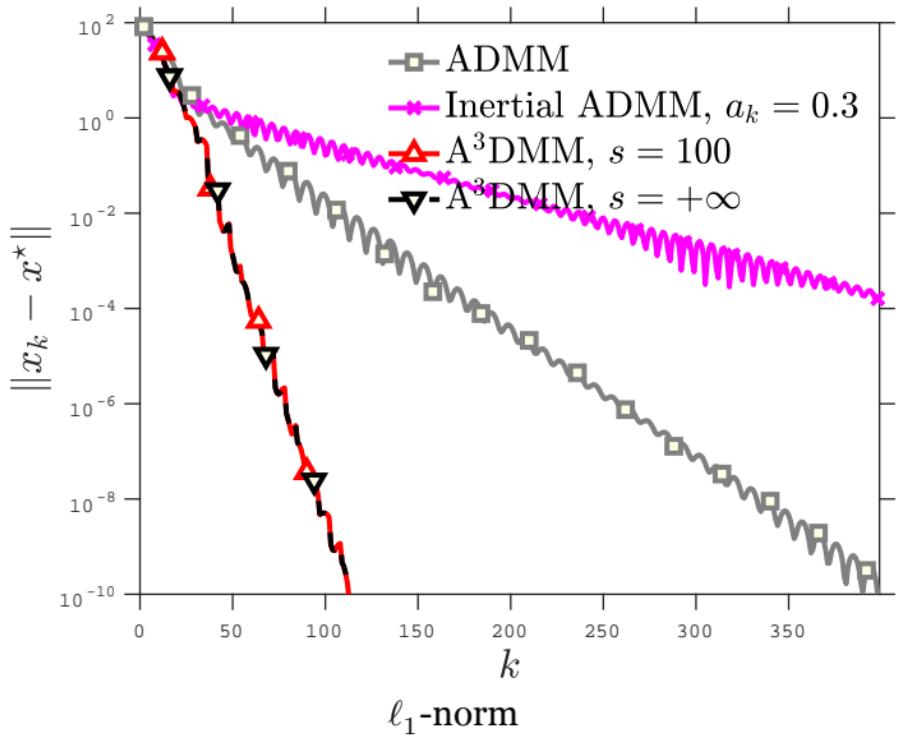


$\ell_{1,2}$ -norm



Nuclear norm

Experiment: 2 non-smooth terms



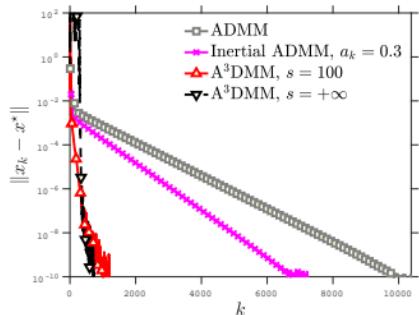
Inertial ADMM is **slower** than ADMM as eventual trajectory is a spiral.

Experiment: LASSO

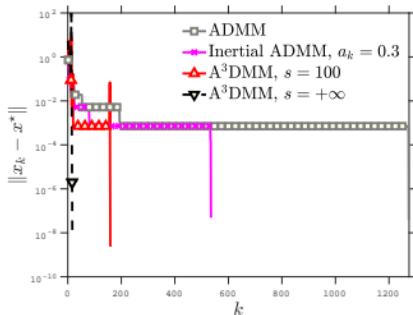


The LASSO problem

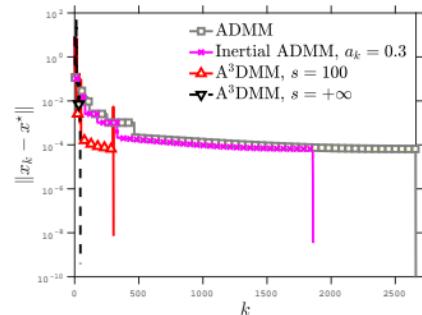
$$\min_{x,y \in \mathbb{R}^n} R(x) + \frac{1}{2} \|Ky - f\|^2 \quad \text{such that} \quad x - y = 0.$$



covtype

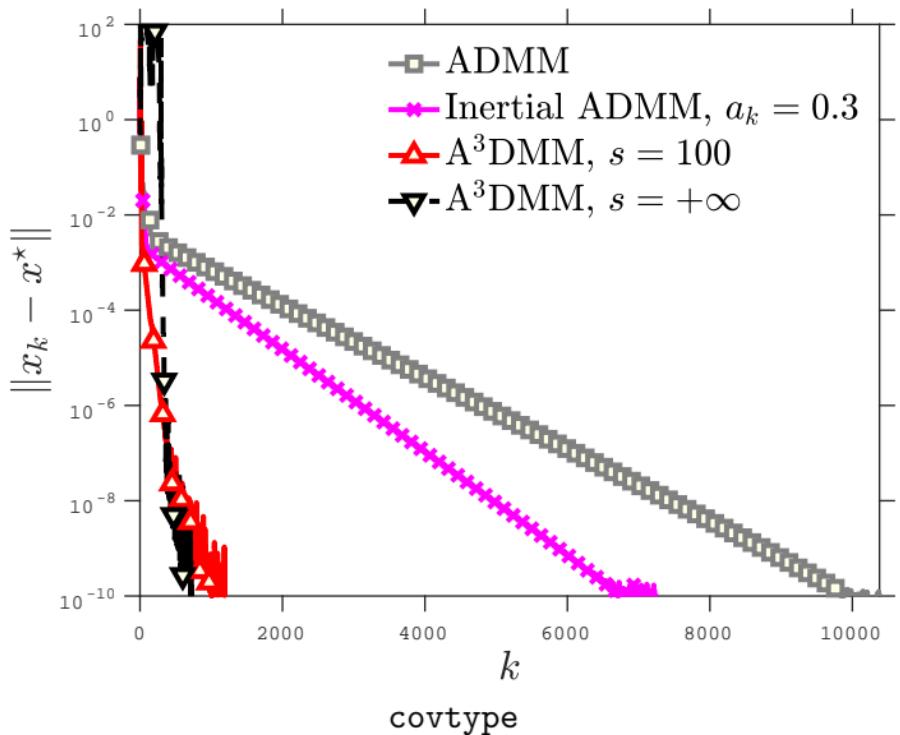


ijcnn1



phishing

Experiment: LASSO



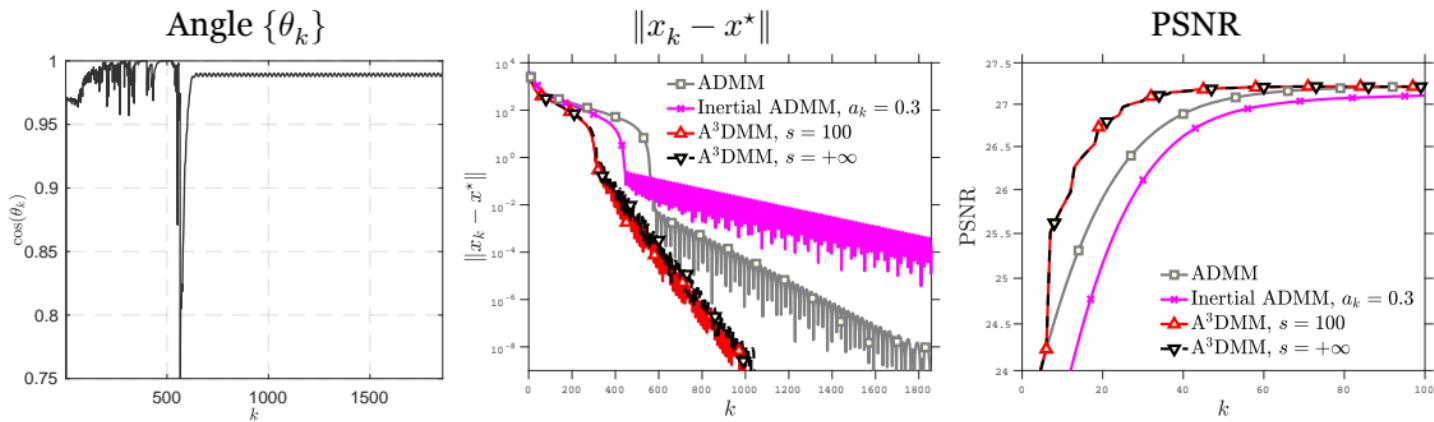
Inertial ADMM does accelerate, but A^3DMM is significantly faster.

Experiment: Total variation based image inpainting



Let $\Omega \stackrel{\text{def}}{=} \{x \in \mathbb{R}^{n \times n} : P_{\mathcal{D}}(x) = f\}$, $P_{\mathcal{D}}$ randomly sets 50% pixels to zero and consider

$$\min_{x \in \mathbb{R}^{n \times n}} \|y\|_1 + \iota_{\Omega}(x) \quad \text{such that} \quad \nabla x - y = 0.$$



- Both functions are polyhedral, trajectory is a spiral.
- Inertial ADMM is **slower** than ADMM.

Experiment: Total variation based image inpainting



Original image



ADMM, PSNR = 26.5448



Inertial ADMM, PSNR = 26.1096



Corrupted image



A^3 DMM $s = 100$, PSNR = 27.0402



A^3 DMM $s = +\infty$, PSNR = 27.0402



Trajectory of FoM For fixed-point sequence $\{z_k\}_{k \in \mathbb{N}}$

- Linearization of FoM.
- Locally different FoM demonstrate **distinct trajectories**: straight line, (logarithmic and elliptic) spirals.

Trajectory of FoM For fixed-point sequence $\{z_k\}_{k \in \mathbb{N}}$

- Linearization of FoM.
- Locally different FoM demonstrate **distinct trajectories**: straight line, (logarithmic and elliptic) spirals.

An adaptive acceleration for FoM

- Though motivated by local trajectory, linear prediction works **globally**.
- For polyhedral functions, guaranteed **local acceleration** for DR/PD using 4 past points.
- For FB, the trajectory is eventually a straight line and one can guarantee local acceleration by extrapolating from 3 past points under our framework, but one could just apply **restarted FISTA**...

Trajectory of FoM For fixed-point sequence $\{z_k\}_{k \in \mathbb{N}}$

- Linearization of FoM.
- Locally different FoM demonstrate **distinct trajectories**: straight line, (logarithmic and elliptic) spirals.

An adaptive acceleration for FoM

- Though motivated by local trajectory, linear prediction works **globally**.
- For polyhedral functions, guaranteed **local acceleration** for DR/PD using 4 past points.
- For FB, the trajectory is eventually a straight line and one can guarantee local acceleration by extrapolating from 3 past points under our framework, but one could just apply **restarted FISTA**...

Reference

- Trajectory of Alternating Direction Method of Multipliers and Adaptive Acceleration, NeurIPS19.
- Geometry of First-Order Methods and Adaptive Acceleration, arXiv:2003.03910.

Many thanks for your time!

<https://jliang993.github.io/>