

Trajectory of Alternating Direction Method of Multipliers and Adaptive Acceleration

Clarice Poon (University of Bath)

Jingwei Liang (University of Cambridge)

Question: *How should one accelerate the convergence of ADMM?*

Constrained and composite optimisation problem:

$$\min_{x \in \mathbb{R}^n, y \in \mathbb{R}^m} R(x) + J(y) \quad \text{such that} \quad Ax + By = b \quad (\mathcal{P})$$

under basic assumptions

- R, J are proper, convex, lower semi-continuous functions.
- $A : \mathbb{R}^n \rightarrow \mathbb{R}^p$ and $B : \mathbb{R}^m \rightarrow \mathbb{R}^p$ are injective linear operators.
- $\text{ri}(\text{dom}(R) \cap \text{dom}(J)) \neq \emptyset$ and the set of minimizers is non-empty.

Given a fixed point sequence $z_{k+1} = \mathcal{F}(z_k)$, accelerate by

$$\bar{z}_k = z_k + a_k(z_k - z_{k-1}), \quad a_k > 0,$$

$$z_{k+1} = \mathcal{F}(\bar{z}_k).$$

Given a fixed point sequence $z_{k+1} = \mathcal{F}(z_k)$, accelerate by

$$\begin{aligned}\bar{z}_k &= z_k + a_k(z_k - z_{k-1}), \quad a_k > 0, \\ z_{k+1} &= \mathcal{F}(\bar{z}_k).\end{aligned}$$

Inertial is well-studied for algorithms such as gradient descent and Forward-Backward.

Improves the objective convergence rate from $\mathcal{O}(k^{-1})$ to $\mathcal{O}(k^{-2})$.

[Heavy-Ball/Nesterov accelerated gradient/FISTA]

Given a fixed point sequence $z_{k+1} = \mathcal{F}(z_k)$, accelerate by

$$\begin{aligned}\bar{z}_k &= z_k + a_k(z_k - z_{k-1}), \quad a_k > 0, \\ z_{k+1} &= \mathcal{F}(\bar{z}_k).\end{aligned}$$

Inertial is well-studied for algorithms such as gradient descent and Forward-Backward.

Improves the objective convergence rate from $\mathcal{O}(k^{-1})$ to $\mathcal{O}(k^{-2})$.

[Heavy-Ball/Nesterov accelerated gradient/FISTA]

Most works on inertial-ADMM impose extra assumptions (e.g. smoothness, uniform convexity).

Given a fixed point sequence $z_{k+1} = \mathcal{F}(z_k)$, accelerate by

$$\begin{aligned}\bar{z}_k &= z_k + a_k(z_k - z_{k-1}), \quad a_k > 0, \\ z_{k+1} &= \mathcal{F}(\bar{z}_k).\end{aligned}$$

Inertial is well-studied for algorithms such as gradient descent and Forward-Backward.

Improves the objective convergence rate from $\mathcal{O}(k^{-1})$ to $\mathcal{O}(k^{-2})$.

[Heavy-Ball/Nesterov accelerated gradient/FISTA]

Most works on inertial-ADMM impose extra assumptions (e.g. smoothness, uniform convexity).

The performance of inertial-ADMM in general is less clear.

1. We study the local trajectory of a sequence generated by ADMM under the framework of partial smoothness.

1. We study the local trajectory of a sequence generated by ADMM under the framework of partial smoothness.

Based on this trajectory analysis:

2. We obtain insight into when inertial will work and fail.

1. We study the local trajectory of a sequence generated by ADMM under the framework of partial smoothness.

Based on this trajectory analysis:

2. We obtain insight into when inertial will work and fail.

3. We develop an acceleration scheme with local acceleration rates.

Augmented Lagrangian: For $\gamma > 0$ and Lagrangian multiplier $\psi \in \mathbb{R}^p$

$$\mathcal{L}(x, y, \psi) \stackrel{\text{def.}}{=} R(x) + J(y) + \langle \psi, Ax + By - b \rangle + \frac{\gamma}{2} \|Ax + By - b\|_2^2.$$

The ADMM iterations:

$$x_k = \operatorname{argmin}_{x \in \mathbb{R}^n} R(x) + \frac{\gamma}{2} \|Ax + By_{k-1} - b + \frac{1}{\gamma} \psi_{k-1}\|_2^2,$$

$$y_k = \operatorname{argmin}_{y \in \mathbb{R}^m} J(y) + \frac{\gamma}{2} \|Ax_k + By - b + \frac{1}{\gamma} \psi_{k-1}\|_2^2,$$

$$\psi_k = \psi_{k-1} + \gamma(Ax_k + By_k - b).$$

Augmented Lagrangian: For $\gamma > 0$ and Lagrangian multiplier $\psi \in \mathbb{R}^p$

$$\mathcal{L}(x, y, \psi) \stackrel{\text{def.}}{=} R(x) + J(y) + \langle \psi, Ax + By - b \rangle + \frac{\gamma}{2} \|Ax + By - b\|_2^2.$$

The ADMM iterations:

$$x_k = \operatorname{argmin}_{x \in \mathbb{R}^n} R(x) + \frac{\gamma}{2} \|Ax + By_{k-1} - b + \frac{1}{\gamma} \psi_{k-1}\|_2^2,$$

$$y_k = \operatorname{argmin}_{y \in \mathbb{R}^m} J(y) + \frac{\gamma}{2} \|Ax_k + By - b + \frac{1}{\gamma} \psi_{k-1}\|_2^2,$$

$$\psi_k = \psi_{k-1} + \gamma(Ax_k + By_k - b).$$

Define $z_k \stackrel{\text{def.}}{=} \psi_{k-1} + \gamma Ax_k$.

Augmented Lagrangian: For $\gamma > 0$ and Lagrangian multiplier $\psi \in \mathbb{R}^p$

$$\mathcal{L}(x, y, \psi) \stackrel{\text{def.}}{=} R(x) + J(y) + \langle \psi, Ax + By - b \rangle + \frac{\gamma}{2} \|Ax + By - b\|_2^2.$$

The ADMM iterations:

$$x_k = \operatorname{argmin}_{x \in \mathbb{R}^n} R(x) + \frac{\gamma}{2} \|Ax - \frac{1}{\gamma}(z_{k-1} - 2\psi_{k-1})\|^2,$$

$$z_k = \psi_{k-1} + \gamma Ax_k,$$

$$y_k = \operatorname{argmin}_{y \in \mathbb{R}^m} J(y) + \frac{\gamma}{2} \|By + \frac{1}{\gamma}(z_k - \gamma b)\|^2,$$

$$\psi_k = z_k + \gamma(By_k - b).$$

Then, $z_k = \mathcal{F}(z_{k-1})$ for some fixed point operator \mathcal{F}^\dagger .

[†] Due to the equivalence between ADMM and Douglas-Rachford splitting [Gabay '83].

Augmented Lagrangian: For $\gamma > 0$ and Lagrangian multiplier $\psi \in \mathbb{R}^p$

$$\mathcal{L}(x, y, \psi) \stackrel{\text{def.}}{=} R(x) + J(y) + \langle \psi, Ax + By - b \rangle + \frac{\gamma}{2} \|Ax + By - b\|_2^2.$$

The ADMM iterations:

$$x_k = \operatorname{argmin}_{x \in \mathbb{R}^n} R(x) + \frac{\gamma}{2} \|Ax - \frac{1}{\gamma}(z_{k-1} - 2\psi_{k-1})\|^2,$$

$$z_k = \psi_{k-1} + \gamma Ax_k,$$

$$y_k = \operatorname{argmin}_{y \in \mathbb{R}^m} J(y) + \frac{\gamma}{2} \|By + \frac{1}{\gamma}(z_k - \gamma b)\|^2,$$

$$\psi_k = z_k + \gamma(By_k - b).$$

We will analyse the behaviour of $\{z_k\}_k$.

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

Smoothness:

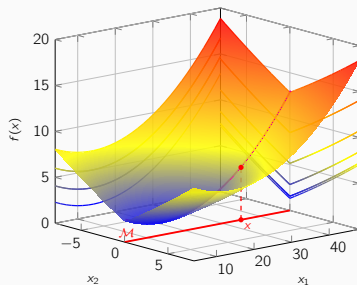
\mathcal{M} is a C^2 -manifold, $R|_{\mathcal{M}}$ is C^2 near x .

Sharpness:

Tangent space $\mathcal{T}_{\mathcal{M}}(x)$ is $\text{par}(\partial R(x))^{\perp}$.

Continuity:

∂R is continuous along \mathcal{M} near x .



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

$\text{PSF}_x(\mathcal{M}_x)$: function that is partly smooth at x relative to \mathcal{M}_x .

Examples: ℓ_1 , $\ell_{1,2}$, ℓ_{∞} -norm, nuclear norm, total variation.

If $R \in \text{PSF}_{x^*}(\mathcal{M}_{x^*}^R)$ and $J \in \text{PSF}_{y^*}(\mathcal{M}_{y^*}^J)$, then under **non-degeneracy conditions** around x^* and y^* :

Manifold identification and local linearisation [Liang, Fadili & Peyré '16]:

There exists $K \in \mathbb{N}$ and a matrix M_{ADMM} such that for all $k \geq K$,

- $x_k \in \mathcal{M}_{x^*}^R$ and $y_k \in \mathcal{M}_{y^*}^J$.
- $z_k - z^* = M_{\text{ADMM}}(z_{k-1} - z^*) + o(\|z_{k-1} - z^*\|)$.

If $R \in \text{PSF}_{x^*}(\mathcal{M}_{x^*}^R)$ and $J \in \text{PSF}_{y^*}(\mathcal{M}_{y^*}^J)$, then under **non-degeneracy conditions** around x^* and y^* :

Manifold identification and local linearisation [Liang, Fadili & Peyré '16]:

There exists $K \in \mathbb{N}$ and a matrix M_{ADMM} such that for all $k \geq K$,

- $x_k \in \mathcal{M}_{x^*}^R$ and $y_k \in \mathcal{M}_{y^*}^J$.
- $z_k - z^* = M_{\text{ADMM}}(z_{k-1} - z^*) + o(\|z_{k-1} - z^*\|)$.

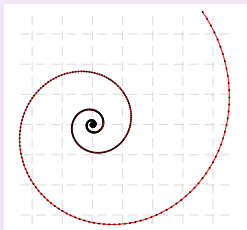
The behaviour of z_k is eventually **regular**.

Partial smoothness and sequence trajectory

Let $v_k \stackrel{\text{def.}}{=} z_k - z_{k-1}$ and $\theta_k = \angle(v_k, v_{k-1})$.

Two non-smooth terms

R and J are locally polyhedral around x^* and y^* .



Spiral trajectory:

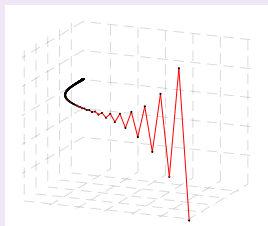
$$\cos(\theta_k) = \cos(\alpha) + \mathcal{O}(\eta^{2k})$$

with $\eta < 1, \alpha > 0$.

M_{ADMM} has **complex** eigenvalues

At least one smooth term

A is an invertible square matrix and R is locally \mathcal{C}^2 around x^* .



Straight line trajectory:

$\cos(\theta_k) \rightarrow 1$ when

$$\gamma > \|(A^\top A)^{-\frac{1}{2}} \nabla^2 R(x^*) (A^\top A)^{-\frac{1}{2}}\|.$$

M_{ADMM} has all **real** eigenvalues

One **inertial**-ADMM iteration:

$$x_k = \operatorname{argmin}_{x \in \mathbb{R}^n} R(x) + \frac{\gamma}{2} \|Ax - \frac{1}{\gamma}(\bar{z}_{k-1} - 2\psi_{k-1})\|^2,$$

$$z_k = \psi_{k-1} + \gamma Ax_k,$$

$$\bar{z}_k = z_k + a_k(z_k - z_{k-1}),$$

$$y_k = \operatorname{argmin}_{y \in \mathbb{R}^m} J(y) + \frac{\gamma}{2} \|By + \frac{1}{\gamma}(\bar{z}_k - \gamma b)\|^2,$$

$$\psi_k = \bar{z}_k + \gamma(By_k - b).$$

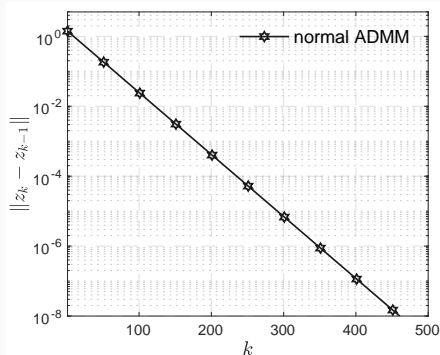
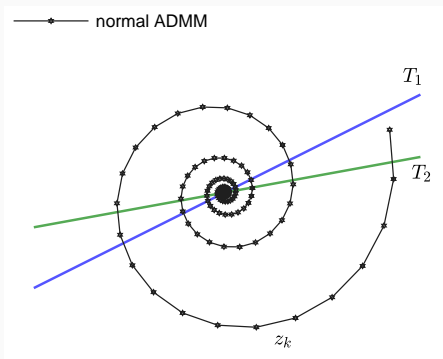
Intuition: inertial-ADMM accelerates if z_k is moving along a straight path...

Failure of inertial-ADMM

Find $x \in T_1 \cap T_2$. Solve using ADMM

$$\min_{x,y} \iota_{T_1}(x) + \iota_{T_2}(y) \quad \text{such that} \quad x - y = 0.$$

Consider $z_k \stackrel{\text{def.}}{=} \psi_{k-1} + \gamma x_k$. **Standard ADMM:**

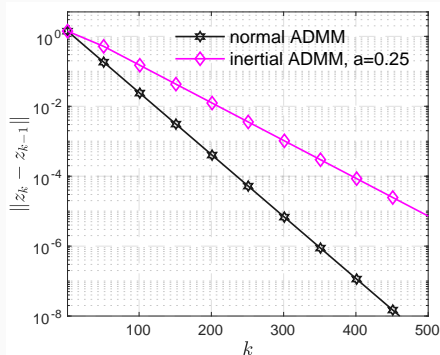
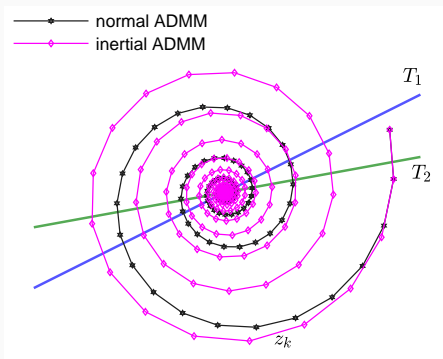


Failure of inertial-ADMM

Find $x \in T_1 \cap T_2$. Solve using ADMM

$$\min_{x,y} \iota_{T_1}(x) + \iota_{T_2}(y) \quad \text{such that} \quad x - y = 0.$$

Consider $z_k \stackrel{\text{def.}}{=} \psi_{k-1} + \gamma x_k$. **Inertial-ADMM with $a = 0.25$:**

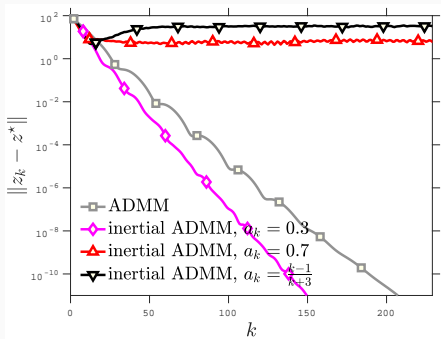


Failure of inertial-ADMM

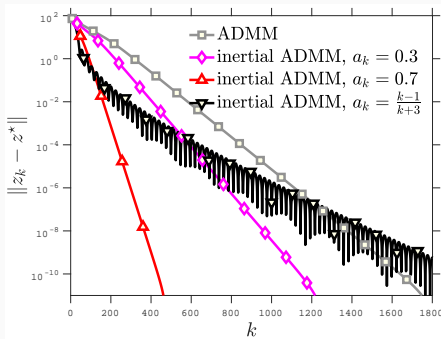
LASSO example:

$$\min_{x,y \in \mathbb{R}^n} \mu \|x\|_1 + \frac{1}{2} \|Ky - f\|_2^2 \quad \text{such that} \quad x - y = 0.$$

$$\gamma = \|K\|^2/10$$



$$\gamma = \|K\|^2 + 0.1$$



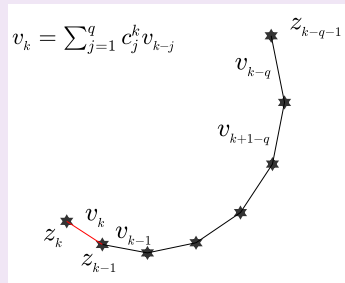
Eventual trajectory:

- Straight line when $\gamma > \|K\|^2$
- M_{ADMM} may have complex leading eigenvalue if $\gamma \leq \|K\|^2$.

Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, define $\{v_{k-j} \stackrel{\text{def.}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Fit the past directions v_{k-1}, \dots, v_{k-q} to the latest direction v_k :

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



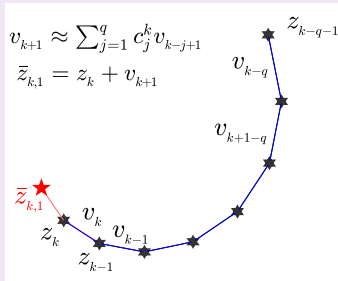
Adaptive acceleration for ADMM (A³DMM)

Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, define $\{v_{k-j} \stackrel{\text{def.}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Fit the past directions v_{k-1}, \dots, v_{k-q} to the latest direction v_k :

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

- Let $\bar{z}_{k,1} \stackrel{\text{def.}}{=} z_k + \sum_{j=1}^q c_j^k v_{k-j+1}$.



Adaptive acceleration for ADMM (A³DMM)

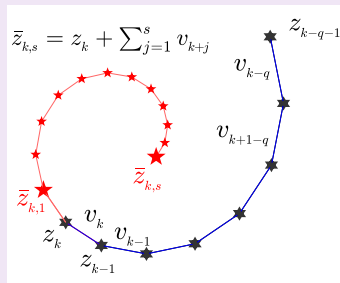
Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, define $\{v_{k-j} \stackrel{\text{def.}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Fit the past directions v_{k-1}, \dots, v_{k-q} to the latest direction v_k :

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

- Let $\bar{z}_{k,1} \stackrel{\text{def.}}{=} z_k + \sum_{j=1}^q c_j^k v_{k-j+1}$.

Repeat on $\{z_{k-j}\}_{j=0}^q \cup \{\bar{z}_{k,1}\}$ and so on.



Adaptive acceleration for ADMM (A³DMM)

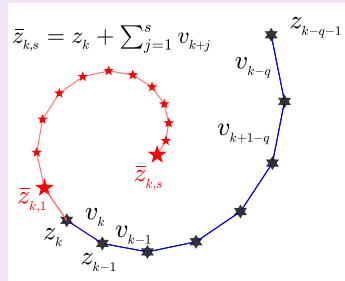
Idea: Given past points $\{z_{k-j}\}_{j=0}^{q+1}$, define $\{v_{k-j} \stackrel{\text{def.}}{=} z_{k-j} - z_{k-j-1}\}_{j=0}^q$.

- Fit the past directions v_{k-1}, \dots, v_{k-q} to the latest direction v_k :

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

- Let $\bar{z}_{k,1} \stackrel{\text{def.}}{=} z_k + \sum_{j=1}^q c_j^k v_{k-j+1}$.

Repeat on $\{z_{k-j}\}_{j=0}^q \cup \{\bar{z}_{k,1}\}$ and so on.



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

$$\hat{c} \stackrel{\text{def.}}{=} \left(\sum_{j=1}^s H(c^k)^j \right)_{(:,1)} \quad \text{with} \quad H(c^k) \stackrel{\text{def.}}{=} \begin{bmatrix} c^k & \left| \frac{\text{Id}_{q-1}}{0_{1,q-1}} \right. \end{bmatrix}.$$

Initial: Let $s \geq 1, q \geq 1$. Let $\bar{z}_0 = z_0 \in \mathbb{R}^p$ and $V_0 = 0_{p \times (q+1)}$.

Repeat: For $k \geq 1$

$$y_k = \operatorname{argmin}_{y \in \mathbb{R}^m} J(y) + \frac{\gamma}{2} \|By + \frac{1}{\gamma} (\bar{z}_{k-1} - \gamma b)\|^2,$$

$$\psi_k = \bar{z}_{k-1} + \gamma(By_k - b),$$

$$x_k = \operatorname{argmin}_{x \in \mathbb{R}^n} R(x) + \frac{\gamma}{2} \|Ax - \frac{1}{\gamma} (\bar{z}_{k-1} - 2\psi_k)\|^2,$$

$$z_k = \psi_k + \gamma Ax_k,$$

$$v_k = z_k - z_{k-1} \quad \text{and} \quad V_k = [v_k, V_k(:, 1:q)].$$

If $\operatorname{mod}(k, q+2) = 0$: Compute coefficients c^k and let $C_k \stackrel{\text{def.}}{=} H(c^k)$

If $\rho(C_k) < 1$: $\bar{z}_k = z_k + a_k \mathcal{E}_{s,q,k}$; else: $\bar{z}_k = z_k$.

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

Global convergence is guaranteed for appropriate choice of a_k .

Local acceleration depends on $\varepsilon_k \stackrel{\text{def.}}{=} \min_c \|V_{k-1}c - v_k\|$.

- If M_{ADMM} is diagonalisable, then $\varepsilon_k = \mathcal{O}(|\lambda_{q+1}|^k)$ where λ_{q+1} is the $(q+1)^{\text{th}}$ largest eigenvalue.
- Guaranteed local acceleration for $q = 2$ if R and J are polyhedral.

Related to vector extrapolation techniques from the 1960's.

[Aitken '27, Wynn '62, Andersen '65...]

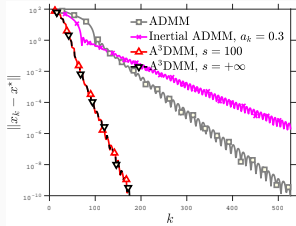
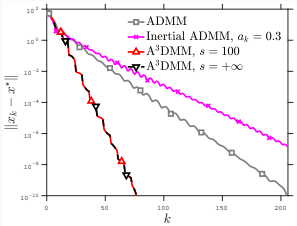
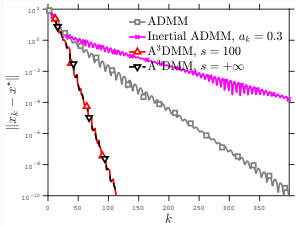
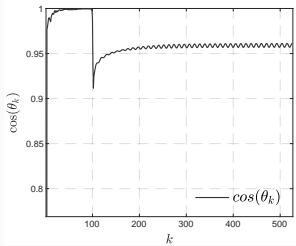
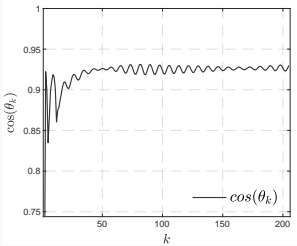
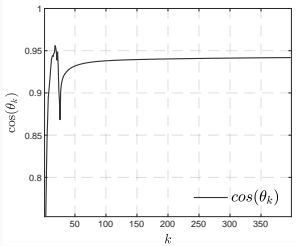
Implementation:

- Typically set $q \leq 10$.
- Extra memory cost of $p \times (q + 1)$ (storing V_k).
- Extra computation cost of $q^2 p$ every $(q + 2)$ iterations.
- One could also extrapolate $\{x_k, y_k\}$ simultaneously. But this would require extra storage of past directions.

Experiment: 2 non-smooth terms

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

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

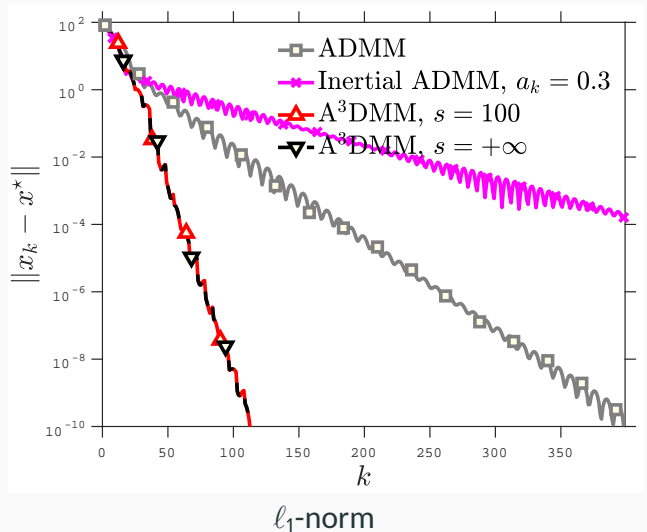


ℓ_1 -norm

$\ell_{1,2}$ -norm

Nuclear norm

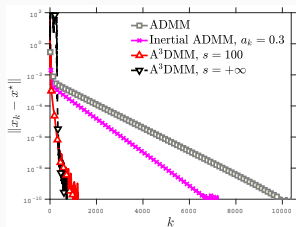
Experiment: 2 non-smooth terms



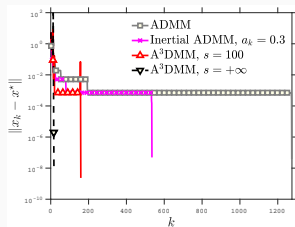
Inertial ADMM is **slower** than ADMM as eventual trajectory is a spiral.

The LASSO problem

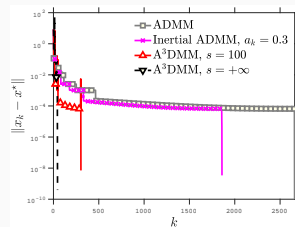
$$\min_{x,y \in \mathbb{R}^n} R(x) + \frac{1}{2} \|Ky - f\|^2 \quad \text{such that} \quad x - y = 0.$$



covtype

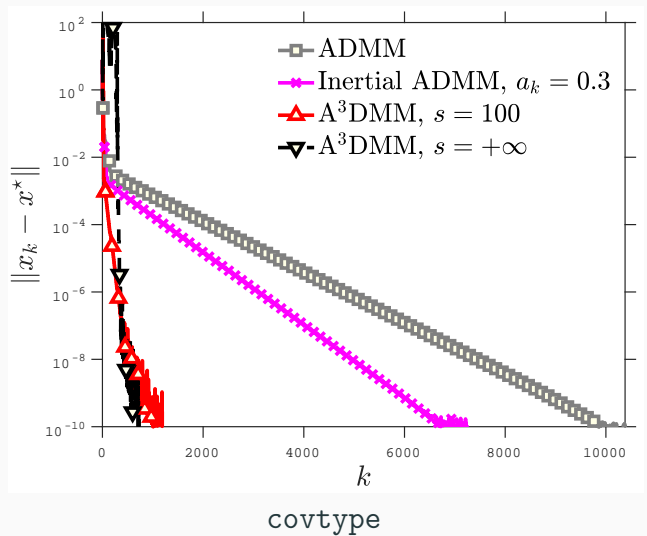


ijcnn1



phishing

Experiment: LASSO

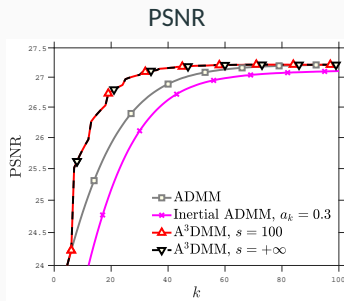
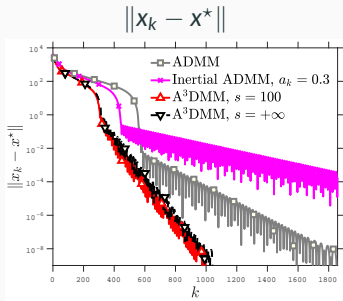
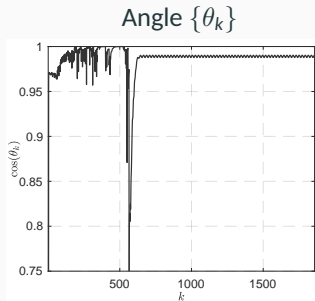


Inertial ADMM does accelerate, but A³DMM is significantly faster.

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



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

Trajectory of ADMM For sequence $\{z_k\}_{k \in \mathbb{N}}$

- When both R and J are locally polyhedral around the fixed point, $\{z_k\}_{k \in \mathbb{N}}$ eventually moves along a **spiral**.
- When at least one of R or J is smooth, the trajectory of $\{z_k\}_{k \in \mathbb{N}}$ depends on γ and can be either a spiral or a **straight line**.

Trajectory of ADMM For sequence $\{z_k\}_{k \in \mathbb{N}}$

- When both R and J are locally polyhedral around the fixed point, $\{z_k\}_{k \in \mathbb{N}}$ eventually moves along a **spiral**.
- When at least one of R or J is smooth, the trajectory of $\{z_k\}_{k \in \mathbb{N}}$ depends on γ and can be either a spiral or a **straight line**.

An adaptive acceleration for ADMM

- The different trajectory behaviour of ADMM can lead to the **failure** of the inertial technique.
- We propose an acceleration strategy based on the idea of following the sequence trajectory.

Trajectory of ADMM For sequence $\{z_k\}_{k \in \mathbb{N}}$

- When both R and J are locally polyhedral around the fixed point, $\{z_k\}_{k \in \mathbb{N}}$ eventually moves along a **spiral**.
- When at least one of R or J is smooth, the trajectory of $\{z_k\}_{k \in \mathbb{N}}$ depends on γ and can be either a spiral or a **straight line**.

An adaptive acceleration for ADMM

- The different trajectory behaviour of ADMM can lead to the **failure** of the inertial technique.
- We propose an acceleration strategy based on the idea of following the sequence trajectory.

Poster: East Exhibition Hall B+C #115!

