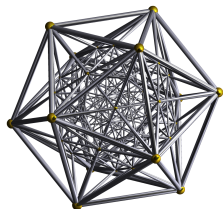


# ICASSP 2022 Short Course One on Low-Dimensional Models for High-Dimensional Data

## Lecture 2: Scalable Convex Optimization Methods for Low-Dimensional Structures

**Sam Buchanan, Yi Ma, Qing Qu  
John Wright, Yuqian Zhang, Zhihui Zhu**

May 24, 2022



# Recap for Recovery of Low-dimensional Structures

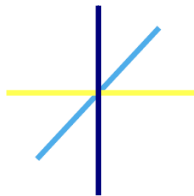
**Parallel developments for sparse vectors and low-rank matrices.**

Sparse v.s. Low-rank	Sparse Vector	Low-rank Matrix
Low-dimensionality of	individual signal $\mathbf{x}$	a set of signals $\mathbf{X}$
Compressive sensing	$\mathbf{y} = \mathbf{A}\mathbf{x}$	$\mathbf{Y} = \mathcal{A}(\mathbf{X})$
Low-dim measure	$\ell^0$ norm $\ \mathbf{x}\ _0$	$\text{rank}(\mathbf{X})$
Convex surrogate	$\ell^1$ norm $\ \mathbf{x}\ _1$	nuclear norm $\ \mathbf{X}\ _*$
Success conditions (RIP)	$\delta_{2k}(\mathbf{A}) \geq \sqrt{2} - 1$	$\delta_{4r}(\mathbf{A}) \geq \sqrt{2} - 1$
Random measurements	$m = O(k \log(n/k))$	$m = O(nr)$
Stable/Inexact recovery	$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{z}$	$\mathbf{Y} = \mathcal{A}(\mathbf{X}) + \mathbf{Z}$
Phase transition at	Stat. dim. of descent cone: $m^* = \delta(\mathbf{D})$	

- ① Motivating Examples for Recovery of Low-Dim Models
- ② (Accelerated) Proximal Methods
- ③ Alternating Direction Methods of Multipliers (ADMM)
- ④ Summary & Extensions

## Sparse Recovery

$$y \in \mathbb{R}^m \begin{matrix} \color{blue}{\boxed{\phantom{y}}} \\ \color{yellow}{\boxed{\phantom{y}}} \\ \color{red}{\boxed{\phantom{y}}} \\ \color{cyan}{\boxed{\phantom{y}}} \\ \color{blue}{\boxed{\phantom{y}}} \\ \color{red}{\boxed{\phantom{y}}} \\ \color{cyan}{\boxed{\phantom{y}}} \end{matrix} = \begin{matrix} \color{yellow}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} \\ \color{yellow}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} \\ \color{yellow}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} \\ \color{yellow}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} \\ \color{yellow}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} \\ \color{yellow}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} \\ \color{yellow}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} & \color{red}{\boxed{\phantom{A}}} & \color{cyan}{\boxed{\phantom{A}}} & \color{blue}{\boxed{\phantom{A}}} \end{matrix} \begin{matrix} \color{red}{\boxed{\phantom{x}}} \\ \color{cyan}{\boxed{\phantom{x}}} \\ \color{blue}{\boxed{\phantom{x}}} \\ \color{red}{\boxed{\phantom{x}}} \\ \color{cyan}{\boxed{\phantom{x}}} \\ \color{blue}{\boxed{\phantom{x}}} \\ \color{red}{\boxed{\phantom{x}}} \\ \color{cyan}{\boxed{\phantom{x}}} \\ \color{blue}{\boxed{\phantom{x}}} \\ \color{red}{\boxed{\phantom{x}}} \\ \color{cyan}{\boxed{\phantom{x}}} \\ \color{blue}{\boxed{\phantom{x}}} \end{matrix} \quad x \in \mathbb{R}^n$$



Recovering a sparse signal  $x_o$  from:

$$\begin{matrix} \mathbf{y} \\ \text{observation} \end{matrix} = \mathbf{A} \begin{matrix} \mathbf{x}_o \\ \text{unknown} \end{matrix}, \quad (1)$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is a linear map.

# Sparse Recovery

- **Basis pursuit (BP):**

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \mathbf{y} = \mathbf{A}\mathbf{x}.$$

# Sparse Recovery

- **Basis pursuit (BP):**

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \mathbf{y} = \mathbf{A}\mathbf{x}.$$

- **Basis pursuit denoising (BPDN):**

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \delta,$$

which is equivalent to the lasso problem with properly chosen  $\lambda$ :

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

# Low-Rank Matrix Recovery

## Recommendation Ratings:

Users

$$\begin{bmatrix}
 5 & 3 & \dots & ? \\
 ? & 2 & \dots & 4 \\
 \vdots & \vdots & \ddots & \vdots \\
 5 & ? & \dots & ?
 \end{bmatrix}
 = \mathcal{P}_\Omega \left( \begin{bmatrix}
 5 & 3 & \dots & 5 \\
 4 & 2 & \dots & 4 \\
 \vdots & \vdots & \ddots & \vdots \\
 5 & 5 & \dots & 3
 \end{bmatrix} \right)$$

Items

Observed (Incomplete) Ratings  $\mathbf{Y}$

Complete Ratings  $\mathbf{X}$

We observe:

$$\mathbf{Y} = \mathcal{P}_\Omega \left[ \mathbf{X} \right],$$

Observed ratings

Complete ratings

where  $\Omega \doteq \{(i, j) \mid \text{user } i \text{ has rated product } j\}$ .

# Low-Rank Matrix Recovery

## Recommendation Ratings:

Users

$$\begin{bmatrix}
 5 & 3 & \dots & ? \\
 ? & 2 & \dots & 4 \\
 \vdots & \vdots & \ddots & \vdots \\
 5 & ? & \dots & ?
 \end{bmatrix}
 = \mathcal{P}_\Omega
 \left(
 \begin{bmatrix}
 5 & 3 & \dots & 5 \\
 4 & 2 & \dots & 4 \\
 \vdots & \vdots & \ddots & \vdots \\
 5 & 5 & \dots & 3
 \end{bmatrix}
 \right)$$

Items

Observed (Incomplete) Ratings  $\mathbf{Y}$

Complete Ratings  $\mathbf{X}$



# Low-Rank Matrix Recovery

## Recommendation Ratings:

$$\begin{array}{c} \text{Users} \\ \begin{matrix} \text{User 1} \\ \text{User 2} \\ \vdots \\ \text{User } m \end{matrix} \end{array} \begin{bmatrix} 5 & 3 & \dots & ? \\ ? & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & ? & \dots & ? \end{bmatrix} = \mathcal{P}_\Omega \left( \begin{array}{c} \begin{bmatrix} 5 & 3 & \dots & 5 \\ 4 & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & 5 & \dots & 3 \end{bmatrix} \\ \text{Complete Ratings } \mathbf{X} \end{array} \right)$$

Items  
Observed (Incomplete) Ratings  $\mathbf{Y}$

## Recovering a low-rank matrix $\mathbf{X}_o$ :

$$\begin{array}{c} \mathbf{y} \\ \text{observation} \end{array} = \mathcal{A} \begin{bmatrix} \mathbf{X}_o \\ \text{unknown} \end{bmatrix}, \quad (2)$$

where,  $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$  is a linear map.

# Low-Rank Matrix Recovery

- **Nuclear norm minimization:**

$$\min_{\mathbf{X}} \|\mathbf{X}\|_*, \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}[\mathbf{X}].$$

# Low-Rank Matrix Recovery

- **Nuclear norm minimization:**

$$\min_{\mathbf{X}} \|\mathbf{X}\|_*, \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}[\mathbf{X}].$$

- **Nuclear norm minimization under noise:**

$$\min_{\mathbf{X}} \|\mathbf{X}\|_*, \quad \text{s.t.} \quad \|\mathbf{y} - \mathcal{A}[\mathbf{X}]\|_F \leq \delta.$$

# Low-Rank Matrix Recovery

- **Nuclear norm minimization:**

$$\min_{\mathbf{X}} \|\mathbf{X}\|_*, \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}[\mathbf{X}].$$

- **Nuclear norm minimization under noise:**

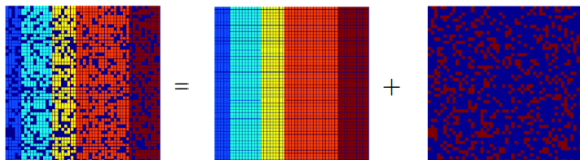
$$\min_{\mathbf{X}} \|\mathbf{X}\|_*, \quad \text{s.t.} \quad \|\mathbf{y} - \mathcal{A}[\mathbf{X}]\|_F \leq \delta.$$

Similarly, the problem is equivalent to

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}[\mathbf{X}]\|_F^2 + \lambda \|\mathbf{X}\|_*,$$

for properly chosen  $\lambda > 0$ .

# Low-rank & Sparse Decomposition



Robustly recover a low-rank matrix  $L_o$  from:

$$\underset{\text{observation}}{\mathbf{Y}} = \underset{\text{unknown low-rank}}{\mathbf{L}_o} + \underset{\text{unknown sparse}}{\mathbf{S}_o}, \quad (3)$$

where  $\mathbf{L}_o \in \mathbb{R}^{n_1 \times n_2}$  is low-rank, and  $\mathbf{S}_o \in \mathbb{R}^{n_1 \times n_2}$  is sparse.

# Low-rank & Sparse Decomposition

- **Principal component pursuit (PCP):**

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1, \quad \text{s.t.} \quad \mathbf{Y} = \mathbf{L} + \mathbf{S}.$$

# Low-rank & Sparse Decomposition

- **Principal component pursuit (PCP):**

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1, \quad \text{s.t. } \mathbf{Y} = \mathbf{L} + \mathbf{S}.$$

- **Stable principal component pursuit (Stable PCP):**

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2$$

# Convex Nonsmooth Problems

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = \underbrace{f(\mathbf{x})}_{\text{smooth}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth}}$$

- **Basis pursuit denoising:**

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2, \quad g(\mathbf{x}) = \lambda \|\mathbf{x}\|_1.$$

- **Stable low-rank matrix recovery:**

$$f(\mathbf{X}) = \frac{1}{2} \|\mathbf{y} - \mathcal{A}[\mathbf{X}]\|_F^2, \quad g(\mathbf{X}) = \lambda \|\mathbf{X}\|_*.$$

- **Stable PCP:**

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



# Optimization Challenges for Structured Data Recovery

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \doteq \underbrace{f(\mathbf{x})}_{\text{smooth convex}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth convex}}. \quad (4)$$

- **Challenge of Scale:** scale algorithms to when  $n$  is very large.

$$\text{Second order methods} \implies \text{First order methods...} \quad (5)$$

- **Nonsmoothness:** first order methods are slow for nonsmooth.

$$\mathcal{O}(1/\sqrt{k}) \implies \mathcal{O}(1/k) \implies \mathcal{O}(1/k^2) \implies \mathcal{O}(e^{-\alpha k}) \quad (6)$$

- **Equality Constraints:** augmented Lagrange multiplier (ALM).
- **Separable Structures:** alternating direction of multipliers method (ADMM).

# Gradient Descent for Smooth Functions, [Cauchy, 1847]

For minimizing a smooth convex function (App. B):

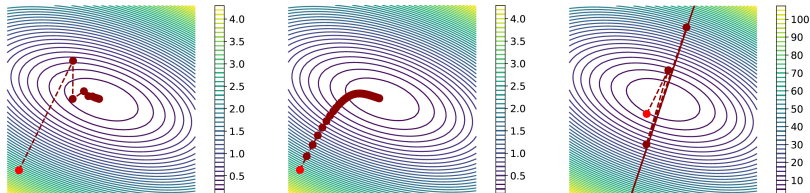
$$\min f(\mathbf{x}), \quad \mathbf{x} \in \mathcal{C} \text{ (a convex set),} \quad (7)$$

conduct **local gradient descent search** (App. D):

$$\mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{x}_k - \tau_k \nabla f(\mathbf{x}_k)), \quad (8)$$



where a rule of thumb:  $\tau_k \approx 1/L$ , where  $L$  the Lipschitz constant.



- figure courtesy of Prof. Carlos Fernandez of NYU.

# (Projected) Subgradient Methods

$$\min_{\mathbf{x}} \underset{\text{nonsmooth}}{F(\mathbf{x})}, \quad \text{s.t.} \quad \underset{\text{convex constraint}}{\mathbf{x} \in \mathcal{C}}.$$

- The loss function  $F(\cdot)$  is nonsmooth: **cannot** apply gradient descent, as  $\nabla F(\mathbf{x}_0)$  might not exist.

# (Projected) Subgradient Methods

$$\min_{\mathbf{x}} \underset{\text{nonsmooth}}{F(\mathbf{x})}, \quad \text{s.t.} \quad \underset{\text{convex constraint}}{\mathbf{x} \in \mathcal{C}}.$$

- The loss function  $F(\cdot)$  is nonsmooth: **cannot** apply gradient descent, as  $\nabla F(\mathbf{x}_0)$  might not exist.
- Similar to GD, a natural choice is the “subgradient-based method”:

$$\mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{x}_k - \tau_k \cdot \mathbf{g}_k),$$

# (Projected) Subgradient Methods

$$\min_{\mathbf{x}} F(\mathbf{x}) \text{ , s.t. } \mathbf{x} \in \mathcal{C} \text{ .}$$

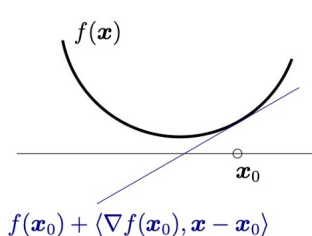
**nonsmooth**
**convex constraint**

- The loss function  $F(\cdot)$  is nonsmooth: **cannot** apply gradient descent, as  $\nabla F(\mathbf{x}_0)$  might not exist.
- Similar to GD, a natural choice is the “subgradient-based method”:

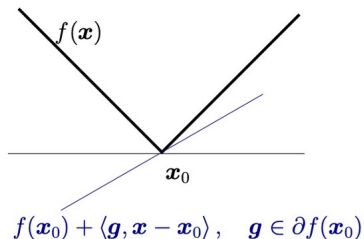
$$\mathbf{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{x}_k - \tau_k \cdot \mathbf{g}_k) \text{ ,}$$

- Here,  $\mathbf{g}_k \in \partial F(\mathbf{x}_k)$  is any subgradient of  $F(\cdot)$  at  $\mathbf{x}_k$ ;  
 $\mathcal{P}_{\mathcal{C}}(\cdot)$  is the projection onto the set  $\mathcal{C}$ .

# Subgradient & Subdifferential



**differentiable**



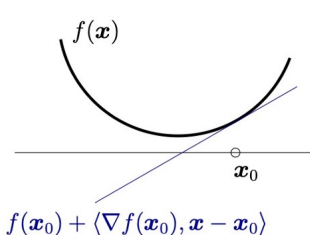
**nondifferentiable**

## Definition (Subgradient)

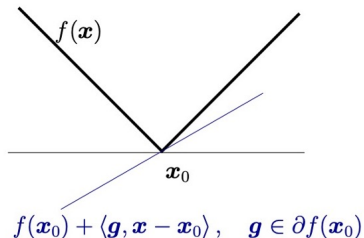
Let  $F : \mathbb{R}^n \mapsto \mathbb{R}$  be *convex*. A subgradient of  $F$  at  $x_0$  is any  $g$  satisfying

$$F(x) \geq F(x_0) + \langle g, x - x_0 \rangle, \quad \forall x \in \mathbb{R}^n.$$

# Subgradient & Subdifferential



**differentiable**



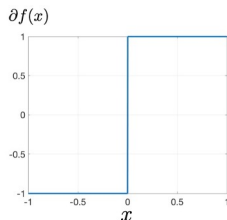
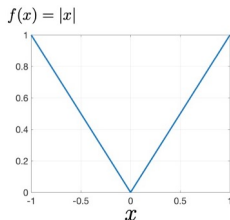
**nondifferentiable**

## Definition (Subdifferential)

Subdifferential is the set of *all* subgradients of  $F(\cdot)$  at  $x_0$ :

$$\partial F(x_0) := \{g \mid F(x) \geq F(x_0) + \langle g, x - x_0 \rangle, \forall x \in \mathbb{R}^n\}.$$

# Examples of Subgradient

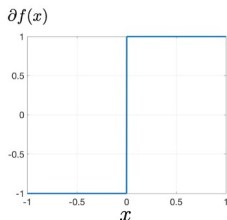
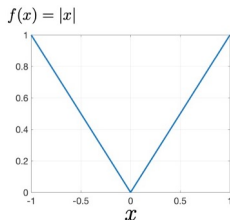


- Absolute value function  $F(x) = |x|$ :

$$\partial F(x) = \begin{cases} \{1\} & x > 0, \\ [-1, 1] & x = 0, \\ \{-1\} & x < 0. \end{cases}$$



# Examples of Subgradient



- Absolute value function  $F(x) = |x|$ :

$$\partial F(x) = \begin{cases} \{1\} & x > 0, \\ [-1, 1] & x = 0, \\ \{-1\} & x < 0. \end{cases}$$

- $\ell_1$ -norm  $F(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$ :

$$\partial F(\mathbf{x}) = \partial F(x_1) \times \cdots \times \partial F(x_n), \quad \mathbf{x} \in \mathbb{R}^n.$$

# Slow Convergence of the Subgradient Method

Suppose  $\mathcal{C} = \mathbb{R}^n$ , and the subgradient method

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \cdot \mathbf{g}_k,$$

with  $\tau_k$  being the stepsize and  $\mathbf{g}_k \in \partial F(\mathbf{x}_k)$  is a subgradient of  $F$  at  $\mathbf{x}_k$ .

# Slow Convergence of the Subgradient Method

Suppose  $\mathcal{C} = \mathbb{R}^n$ , and the subgradient method

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \cdot \mathbf{g}_k,$$

with  $\tau_k$  being the stepsize and  $\mathbf{g}_k \in \partial F(\mathbf{x}_k)$  is a subgradient of  $F$  at  $\mathbf{x}_k$ .

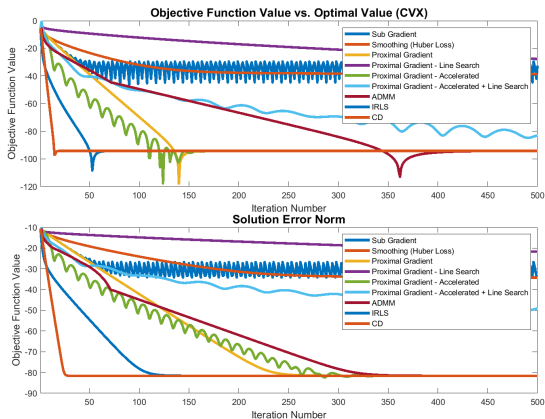
**Theorem ( $\mathcal{O}(1/\sqrt{k})$  convergence of subgradient methods)**

*Suppose  $F : \mathbb{R}^n \mapsto \mathbb{R}$  is  $L$ -Lipschitz, and the optimal function value is  $F_\star$ . Choose the stepsize satisfying  $\tau_k = \frac{F_k - F_\star}{\|\mathbf{g}_k\|_2^2}$ , then we have the convergence rate*

$$F_{best,k} - F_\star \leq \frac{RL}{\sqrt{k}},$$

*where  $F_{best,k} = \min_{1 \leq i \leq k} F(\mathbf{x}_i)$  and  $R \geq \|\mathbf{x}_0 - \mathbf{x}_\star\|_2$ .*

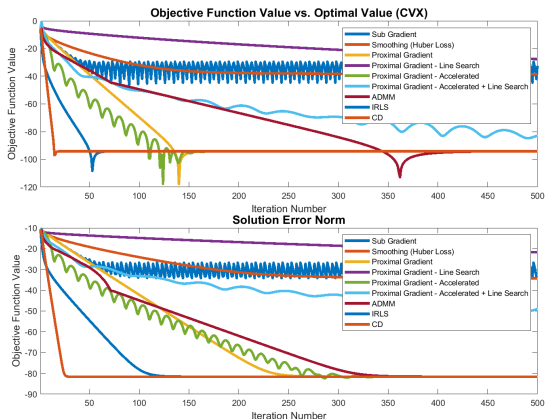
# Comparison of Convergence<sup>1</sup>



$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

<sup>1</sup> <https://github.com/RoyiAvital/Projects/tree/master/Optimization/LsL1SolversAnalysis>

# Comparison of Convergence<sup>2</sup>



- For recovery of low-dimensional models, **generic solvers are slow** (e.g., subgradient method).

<sup>2</sup><https://github.com/RoyiAvital/Projects/tree/master/Optimization/LsL1SolversAnalysis>

- 1 Motivating Examples for Recovery of Low-Dim Models
- 2 (Accelerated) Proximal Methods
- 3 Alternating Direction Methods of Multipliers (ADMM)
- 4 Summary & Extensions

# Composite Nonsmooth Problems

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = \underbrace{f(\mathbf{x})}_{\text{smooth}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth}}$$

- The function  $f(\cdot) : \mathbb{R}^n \mapsto \mathbb{R}$  is convex, continuously differentiable, and  $L$ -smooth with

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

- $g(\cdot) : \mathbb{R}^n \mapsto \mathbb{R}$  is convex but possibly *nonsmooth*.

# Composite Nonsmooth Problems

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = \underbrace{f(\mathbf{x})}_{\text{smooth}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth}}$$

- **Basis pursuit denoising:**

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2, \quad g(\mathbf{x}) = \lambda \|\mathbf{x}\|_1.$$

- **Stable low-rank matrix recovery:**

$$f(\mathbf{X}) = \frac{1}{2} \|\mathbf{y} - \mathcal{A}[\mathbf{X}]\|_F^2, \quad g(\mathbf{X}) = \lambda \|\mathbf{X}\|_*.$$

- **Stable PCP:**

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



# Gradient Descent, [Cauchy, 1847]

For minimizing a smooth convex function (App. B):

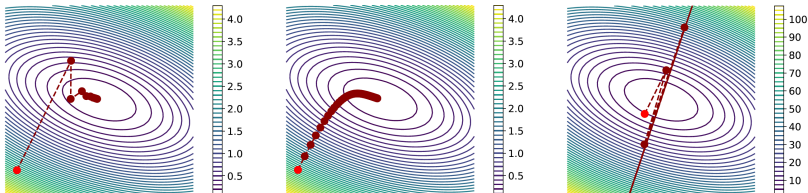
$$\min f(\mathbf{x}), \quad \mathbf{x} \in \mathcal{C} \text{ (a convex set),}$$

conduct **local gradient descent search** (App. D):

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma_k \nabla f(\mathbf{x}_k),$$



where a rule of thumb:  $\gamma \approx 1/L$ , where  $L$  the Lipschitz constant (why?).



- figure courtesy of Prof. Carlos Fernandez of NYU.

# Gradient Descent

For  $f(\mathbf{x})$  has  $L$ -Lipschitz continuous gradients if

$$\|\nabla f(\mathbf{x}') - \nabla f(\mathbf{x})\|_2 \leq L\|\mathbf{x}' - \mathbf{x}\|_2, \quad \forall \mathbf{x}', \mathbf{x} \in \mathbb{R}^n. \quad (9)$$

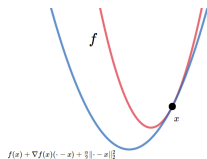
# Gradient Descent

For  $f(\mathbf{x})$  has  $L$ -Lipschitz continuous gradients if

$$\|\nabla f(\mathbf{x}') - \nabla f(\mathbf{x})\|_2 \leq L\|\mathbf{x}' - \mathbf{x}\|_2, \quad \forall \mathbf{x}', \mathbf{x} \in \mathbb{R}^n. \quad (9)$$

This gives a matching **quadratic upper bound**:

$$\begin{aligned} f(\mathbf{x}') &\leq \hat{f}(\mathbf{x}', \mathbf{x}) \\ &\doteq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2 \\ &= \frac{L}{2} \|\mathbf{x}' - (\mathbf{x} - \tau \nabla f(\mathbf{x}))\|_2^2 + h(\mathbf{x}). \end{aligned}$$



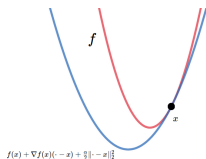
## Gradient Descent

For  $f(\mathbf{x})$  has  $L$ -Lipschitz continuous gradients if

$$\|\nabla f(\mathbf{x}') - \nabla f(\mathbf{x})\|_2 \leq L\|\mathbf{x}' - \mathbf{x}\|_2, \quad \forall \mathbf{x}', \mathbf{x} \in \mathbb{R}^n. \quad (9)$$

This gives a matching **quadratic upper bound**:

$$\begin{aligned} f(\mathbf{x}') &\leq \hat{f}(\mathbf{x}', \mathbf{x}) \\ &\doteq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2 \\ &= \frac{L}{2} \|\mathbf{x}' - (\mathbf{x} - \tau \nabla f(\mathbf{x}))\|_2^2 + h(\mathbf{x}). \end{aligned}$$



Take a step to the **minimizer of this bound**:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}'} \hat{f}(\mathbf{x}', \mathbf{x}_k) = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k).$$

**Fact:** this gives a convergence rate of  $O(1/k)$ .

## Proximal Gradient Descent

The same (local) strategy for a convex function with a nonsmooth term:

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = \underbrace{f(\mathbf{x})}_{\text{smooth}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth}}$$

**Upper bound:**

$$\begin{aligned} \hat{F}(\mathbf{x}, \mathbf{x}_k) &= f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 + g(\mathbf{x}) \\ &= \frac{L}{2} \|\mathbf{x} - (\mathbf{x}_k - \tau \nabla f(\mathbf{x}_k))\|_2^2 + g(\mathbf{x}) + h(\mathbf{x}_k). \end{aligned}$$

## Proximal Gradient Descent

The same (local) strategy for a convex function with a nonsmooth term:

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = \underbrace{f(\mathbf{x})}_{\text{smooth}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth}}$$

**Upper bound:**

$$\begin{aligned} \hat{F}(\mathbf{x}, \mathbf{x}_k) &= f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2 + g(\mathbf{x}) \\ &= \frac{L}{2} \|\mathbf{x} - (\mathbf{x}_k - \tau \nabla f(\mathbf{x}_k))\|_2^2 + g(\mathbf{x}) + h(\mathbf{x}_k). \end{aligned}$$

**A step to the minimizer of the bound  $\hat{F}(\mathbf{x}, \mathbf{x}_k)$ :**

$$\begin{aligned} \mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \frac{L}{2} \left\| \mathbf{x} - \underbrace{\left( \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k) \right)}_{\mathbf{w}_k} \right\|_2^2 + g(\mathbf{x}) \\ &= \boxed{\arg \min_{\mathbf{x}} \frac{L}{2} \|\mathbf{x} - \mathbf{w}_k\|_2^2 + g(\mathbf{x})}.} \end{aligned}$$

# Proximal Operator

## Definition (Proximal Operator)

The proximal operator of a convex function  $g(\cdot) : \mathbb{R}^n \mapsto \mathbb{R}$  is

$$\text{prox}_g(\mathbf{w}) := \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{w}\|_2^2 + g(\mathbf{x}) \right\}$$

# Proximal Operator

## Definition (Proximal Operator)

The proximal operator of a convex function  $g(\cdot) : \mathbb{R}^n \mapsto \mathbb{R}$  is

$$\text{prox}_g(\mathbf{w}) := \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{w}\|_2^2 + g(\mathbf{x}) \right\}$$

Thus, the proximal iteration can be written as

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

**For many structured low-dim problems, the proximal mapping has closed-form and can be computed efficiently!**



# Proximal Operator

## Definition (Proximal Operator)

The proximal operator of a convex function  $g(\cdot) : \mathbb{R}^n \mapsto \mathbb{R}$  is

$$\text{prox}_g(\mathbf{w}) := \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{w}\|_2^2 + g(\mathbf{x}) \right\}$$

### Example: $\ell_1$ -norm

- $g(\mathbf{x}) = t \|\mathbf{x}\|_1$ ,  $\text{prox}_g(\cdot)$  is the **soft-thresholding operator**  $\text{soft}_t(\cdot)$ :

$$[\text{prox}_g(\mathbf{w})]_i = [\text{soft}_t(\mathbf{w})]_i = \begin{cases} w_i - t & w_i \geq t \\ 0 & |w_i| \leq t \\ w_i + t & w_i \leq -t \end{cases}$$

# Proximal Operator

How to prove it?

- **Subgradient characterization:** for any convex function  $F : \mathbb{R}^n \mapsto \mathbb{R}$ ,

$$F(\mathbf{x}_\star) = \min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \iff \mathbf{0} \in \partial F(\mathbf{x}_\star)$$

# Proximal Operator

How to prove it?

- **Subgradient characterization:** for any convex function  $F : \mathbb{R}^n \mapsto \mathbb{R}$ ,

$$F(\mathbf{x}_\star) = \min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \iff \mathbf{0} \in \partial F(\mathbf{x}_\star)$$

- **Proof ideas of the proximal operator for  $g(\mathbf{x}) = t \|\mathbf{x}\|_1$ :**

The objective function reaches minimum when the subdifferential of

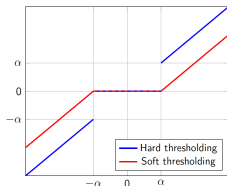
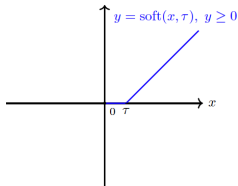
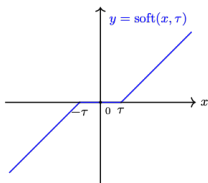
$$F(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \mathbf{w}\|_2^2 + t \|\mathbf{x}\|_1$$

contains zero, that is

$$\mathbf{0} \in (\mathbf{x} - \mathbf{w}) + t\partial\|\mathbf{x}\|_1 = \begin{cases} x_i - w_i + t, & x_i > 0 \\ -w_i + t[-1, 1], & x_i = 0 \\ x_i - w_i - t, & x_i < 0 \end{cases}, \quad i = 1, \dots, n.$$

# Proximal Operator

## Thresholding:



## Example: $\ell_1$ -norm

- $g(\mathbf{x}) = t \|\mathbf{x}\|_1$ ,  $\text{prox}_g(\cdot)$  is the **soft-thresholding operator**:

$$[\text{prox}_g(\mathbf{w})]_i = [\text{soft}_t(\mathbf{w})]_i = \begin{cases} w_i - t & w_i \geq t \\ 0 & |w_i| \leq t \\ w_i + t & w_i \leq -t \end{cases}$$

# Proximal Operator

## Definition (Proximal Operator)

The proximal operator of a convex function  $g(\cdot) : \mathbb{R}^{n \times m} \mapsto \mathbb{R}$  is

$$\text{prox}_g(\mathbf{W}) := \min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \left\{ \frac{1}{2} \|\mathbf{X} - \mathbf{W}\|_F^2 + g(\mathbf{X}) \right\}$$

### Example: nuclear norm

- $g(\mathbf{W}) = t \|\mathbf{W}\|_*$  with  $\mathbf{W} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ , then  $\text{prox}_g(\cdot)$  is the singular value thresholding operator:

$$\text{prox}_g(\mathbf{W}) = \mathbf{U} \text{soft}_t(\mathbf{\Sigma}) \mathbf{V}^\top.$$

# Proximal Gradient Algorithm

## Proximal Gradient (PG)

**Problem Class:**  $\min_{\mathbf{x}} F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$

$f, g : \mathbb{R}^n \rightarrow \mathbb{R}$  convex,  $\nabla f$   $L$ -Lipschitz and  $g$  nonsmooth.

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

Repeat:

$$\begin{aligned}\mathbf{w}_k &\leftarrow \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k), \\ \mathbf{x}_{k+1} &\leftarrow \text{prox}_{g/L}[\mathbf{w}_k].\end{aligned}$$

**Convergence Guarantee:**

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

# Example: Proximal Gradient for Basis Pursuit Denoising

## Iterative soft-thresholding algorithm (ISTA):

- 1: **Problem:**  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$ , given  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .
- 2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $L \geq \lambda_{\max}(\mathbf{A}^* \mathbf{A})$ .
- 3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**
- 4:      $\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^* (\mathbf{A}\mathbf{x}_k - \mathbf{y})$ .
- 5:      $\mathbf{x}_{k+1} \leftarrow \text{soft}_{\lambda/L}(\mathbf{w}_k)$ .
- 6: **end for**
- 7: **Output:**  $\mathbf{x}_\star \leftarrow \mathbf{x}_K$ .

---

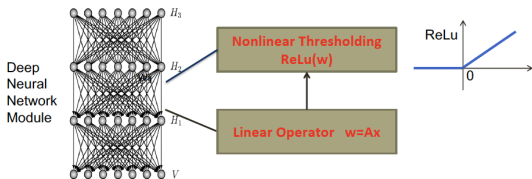
<sup>3</sup>Learning Fast Approximations of Sparse Coding, Karol Gregor and Yann LeCun, ICML 2010. Also known as the Learned ISTA (LISTA).

# Example: Proximal Gradient for Basis Pursuit Denoising

## Iterative soft-thresholding algorithm (ISTA):

- 1: **Problem:**  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$ , given  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .
- 2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $L \geq \lambda_{\max}(\mathbf{A}^* \mathbf{A})$ .
- 3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**
- 4:      $\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^* (\mathbf{A}\mathbf{x}_k - \mathbf{y})$ .
- 5:      $\mathbf{x}_{k+1} \leftarrow \text{soft}_{\lambda/L}(\mathbf{w}_k)$ .
- 6: **end for**
- 7: **Output:**  $\mathbf{x}_\star \leftarrow \mathbf{x}_K$ .

The unrolled iterations resemble a deep neural network!<sup>3</sup>



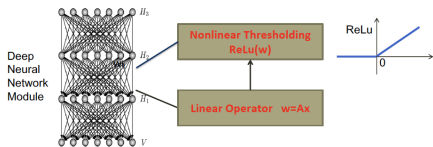
<sup>3</sup>Learning Fast Approximations of Sparse Coding, Karol Gregor and Yann LeCun, ICML 2010. Also known as the Learned ISTA (LISTA).



## From ISTA to Learned ISTA (LISTA)

- 1: **Problem:**  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$ , given  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .
- 2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $L \geq \lambda_{\max}(\mathbf{A}^* \mathbf{A})$ .
- 3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**
- 4:      $\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^* (\mathbf{A}\mathbf{x}_k - \mathbf{y})$ .
- 5:      $\mathbf{x}_{k+1} \leftarrow \text{soft}_{\lambda/L}(\mathbf{w}_k)$ .
- 6: **end for**
- 7: **Output:**  $\mathbf{x}_* \leftarrow \mathbf{x}_K$ .

The unrolled iterations resemble a deep neural network!



We can optimize the optimization path of ISTA using supervised learning<sup>4</sup>:

$$\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^* (\mathbf{A}\mathbf{x}_k - \mathbf{y}) = \underbrace{\left( \mathbf{I} - \frac{1}{L} \mathbf{A}^* \mathbf{A} \right)}_{\text{learnable parameter } \mathbf{S}} \mathbf{x}_k + \underbrace{\frac{1}{L} \mathbf{A}^* \mathbf{A} \mathbf{y}}_{\text{learnable parameter } \mathbf{b}}$$

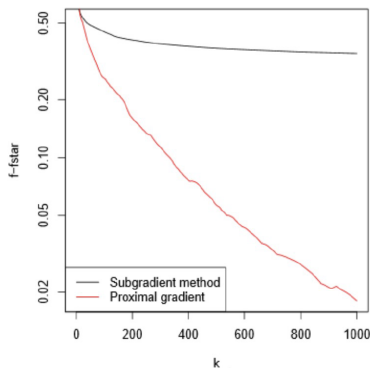
<sup>4</sup>Learning Fast Approximations of Sparse Coding, Karol Gregor and Yann LeCun, ICML 2010.

# Example: Proximal Gradient for Basis Pursuit Denoising

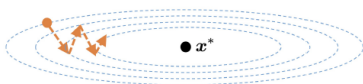
## Iterative soft-thresholding algorithm (ISTA):

- 1: **Problem:**  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$ , given  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .
- 2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $L \geq \lambda_{\max}(\mathbf{A}^* \mathbf{A})$ .
- 3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**
- 4:      $\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^* (\mathbf{A}\mathbf{x}_k - \mathbf{y})$ .
- 5:      $\mathbf{x}_{k+1} \leftarrow \text{soft}_{\lambda/L}(\mathbf{w}_k)$ .
- 6: **end for**
- 7: **Output:**  $\mathbf{x}_* \leftarrow \mathbf{x}_K$ .

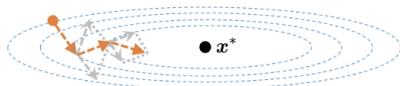
## Proximal Gradient vs. Subgradient Method.



# Can We Further Accelerate Convergence?



gradient descent



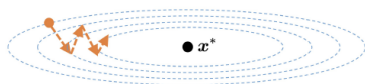
heavy-ball method

Recall gradient descent for smooth  $\min_{\mathbf{x}} f(\mathbf{x})$ :

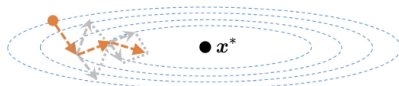
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k).$$



# Can We Further Accelerate Convergence?



gradient descent



heavy-ball method

Recall gradient descent for smooth  $\min_{\mathbf{x}} f(\mathbf{x})$ :

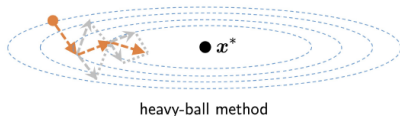
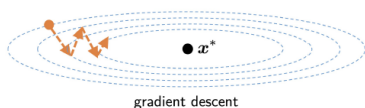
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k).$$

The **heavy ball method** [Polyak, 1964]

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k) + \underbrace{\beta(\mathbf{x}_k - \mathbf{x}_{k-1})}_{\text{momentum}}.$$



# Can We Further Accelerate Convergence?



Recall gradient descent for smooth  $\min_{\mathbf{x}} f(\mathbf{x})$ :

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k).$$

The **heavy ball method** [Polyak, 1964]

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k) + \underbrace{\beta(\mathbf{x}_k - \mathbf{x}_{k-1})}_{\text{momentum}}.$$



It is also called the *momentum method*:

- Basis for popular ADAM for train deep neural networks.
- Worst convergence rate is still  $O(1/k)$ , yet best possible is  $O(1/k^2)$ .

# Accelerated Gradient Descent [Nesterov, 1983]

Generate an auxiliary point  $\mathbf{p}_{k+1}$  of the form:

$$\mathbf{p}_{k+1} \doteq \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1}).$$

Move from  $\mathbf{x}_k$  to  $\mathbf{p}_{k+1}$ , and gradient descend from it:

$$\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \alpha \underbrace{\nabla f(\mathbf{p}_{k+1})}_{\text{a stroke of genius}}.$$



# Accelerated Gradient Descent [Nesterov, 1983]

Generate an auxiliary point  $\mathbf{p}_{k+1}$  of the form:

$$\mathbf{p}_{k+1} \doteq \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1}).$$

Move from  $\mathbf{x}_k$  to  $\mathbf{p}_{k+1}$ , and gradient descend from it:

$$\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \alpha \underbrace{\nabla f(\mathbf{p}_{k+1})}_{\text{a stroke of genius}}.$$



The weights  $\alpha$  and  $\{\beta_{k+1}\}$  are carefully chosen:

$$t_1 = 1, \quad t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}, \quad \beta_{k+1} = \frac{t_k - 1}{t_{k+1}}, \quad \alpha = 1/L.$$

# Accelerated Gradient Descent [Nesterov, 1983]

Generate an auxiliary point  $\mathbf{p}_{k+1}$  of the form:

$$\mathbf{p}_{k+1} \doteq \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1}).$$

Move from  $\mathbf{x}_k$  to  $\mathbf{p}_{k+1}$ , and gradient descend from it:

$$\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \alpha \underbrace{\nabla f(\mathbf{p}_{k+1})}_{\text{a stroke of genius}}.$$



The weights  $\alpha$  and  $\{\beta_{k+1}\}$  are carefully chosen:

$$t_1 = 1, \quad t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}, \quad \beta_{k+1} = \frac{t_k - 1}{t_{k+1}}, \quad \alpha = 1/L.$$

- We may not always have  $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k)$ .
- Achieve optimal convergence rate  $O(1/k^2)$  among 1st order methods.



## Accelerated Proximal Gradient for Nonsmooth Problems

**Accelerated Proximal Gradient (APG)****Problem Class:**  $\min_{\mathbf{x}} F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$ , $f, g$  convex, with  $\nabla f$   $L$ -Lipschitz and  $g$  **nonsmooth**.**Basic Iteration:** set  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\mathbf{p}_1 = \mathbf{x}_1 \leftarrow \mathbf{x}_0$ , and  $t_1 \leftarrow 1$ .Repeat for  $k = 1, 2, \dots, K$ :

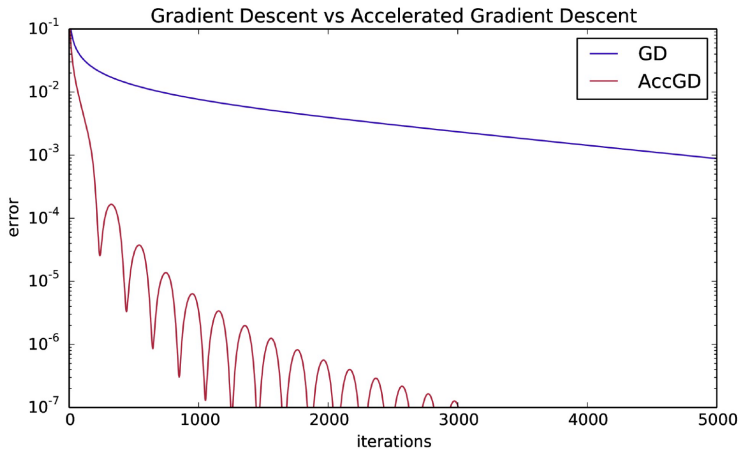
$$t_{k+1} \leftarrow \frac{1 + \sqrt{1 + 4t_k^2}}{2}, \quad \beta_{k+1} \leftarrow \frac{t_k - 1}{t_{k+1}}.$$

$$\mathbf{p}_{k+1} \leftarrow \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1}).$$

$$\mathbf{x}_{k+1} \leftarrow \underbrace{\text{prox}_{g/L} \left[ \mathbf{p}_{k+1} - \frac{1}{L} \nabla f(\mathbf{p}_{k+1}) \right]}_{\text{proximal gradient}}.$$

**Convergence Guarantee:** $F(\mathbf{x}_k) - F(\mathbf{x}_\star)$  converges at a rate of  $O(1/k^2)$ .

## Proximal Gradient versus Accelerated Proximal Gradient

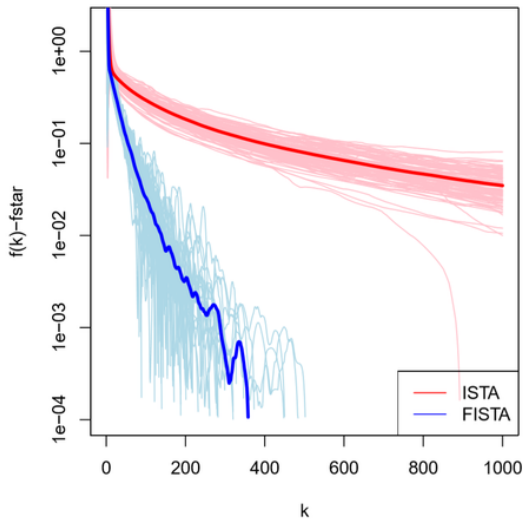


# Example I: APG for Basis Pursuit Denoising

## FISTA: Accelerated Proximal Gradient (APG) for LASSO

- 1: **Problem:**  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$ , given  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .
- 2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\mathbf{p}_1 = \mathbf{x}_1 \leftarrow \mathbf{x}_0$ , and  $t_1 \leftarrow 1$ , and  $L \geq \lambda_{\max}(\mathbf{A}^* \mathbf{A})$ .
- 3: **for** ( $k = 1, 2, \dots, K - 1$ ) **do**
- 4:    $t_{k+1} \leftarrow \frac{1 + \sqrt{1 + 4t_k^2}}{2}$ ;  $\beta_{k+1} \leftarrow \frac{t_k - 1}{t_{k+1}}$ .
- 5:    $\mathbf{p}_{k+1} \leftarrow \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1})$ .
- 6:    $\mathbf{w}_{k+1} \leftarrow \mathbf{p}_{k+1} - \frac{1}{L} \mathbf{A}^*(\mathbf{A}\mathbf{p}_{k+1} - \mathbf{y})$ .
- 7:    $\mathbf{x}_{k+1} \leftarrow \text{soft}[\mathbf{w}_{k+1}, \lambda/L]$ .
- 8: **end for**
- 9: **Output:**  $\mathbf{x}_* \leftarrow \mathbf{x}_K$ .

## ISTA vs. FISTA



## Example II: APG for Stable PCP

### Accelerated Proximal Gradient (APG) for Stable PCP

- 1: **Problem:**  $\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \frac{\mu}{2} \|\mathbf{Y} - \mathbf{L} - \mathbf{S}\|_F^2$ , given  $\mathbf{Y}$ .
- 2: **Input:**  $\mathbf{L}_0, \mathbf{S}_0 \in \mathbb{R}^{m \times n}$ ,  $\mathbf{P}_1^S = \mathbf{S}_1 \leftarrow \mathbf{S}_0$ ,  $\mathbf{P}_1^L = \mathbf{L}_1 \leftarrow \mathbf{L}_0$ ,  $t_1 \leftarrow 1$ .
- 3: **for** ( $k = 1, 2, \dots, K - 1$ ) **do**
- 4:    $t_{k+1} \leftarrow \frac{1 + \sqrt{1 + 4t_k^2}}{2}$ ,  $\beta_{k+1} \leftarrow \frac{t_{k-1}}{t_{k+1}}$ .
- 5:    $\mathbf{P}_{k+1}^L \leftarrow \mathbf{L}_k + \beta_{k+1}(\mathbf{L}_k - \mathbf{L}_{k-1})$ ;  $\mathbf{P}_{k+1}^S \leftarrow \mathbf{S}_k + \beta_{k+1}(\mathbf{S}_k - \mathbf{S}_{k-1})$ .
- 6:    $\mathbf{W}_{k+1} \leftarrow \mathbf{Y} - \mathbf{P}_{k+1}^S$  and compute SVD:  $\mathbf{W}_{k+1} = \mathbf{U}_{k+1} \mathbf{\Sigma}_{k+1} \mathbf{V}_{k+1}^*$ .
- 7:    $\mathbf{L}_{k+1} \leftarrow \mathbf{U}_{k+1} \text{soft}[\mathbf{\Sigma}_{k+1}, 1/\mu] \mathbf{V}_{k+1}^*$ ;  $\mathbf{S}_{k+1} \leftarrow \text{soft}[(\mathbf{Y} - \mathbf{P}_{k+1}^L), \lambda/\mu]$ .
- 8: **end for**
- 9: **Output:**  $\mathbf{L}_* \leftarrow \mathbf{L}_K$ ;  $\mathbf{S}_* \leftarrow \mathbf{S}_K$ .

# GD for Strongly Convex Problems

**A troubling fact though:** Not supposed to be this fast!

**Reason?** Consider minimizing a  $L$ -**Lipschitz continuous** function

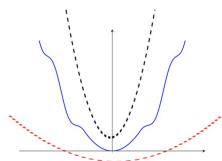
$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^n. \quad (10)$$

Assume  $f(\mathbf{x})$  is  $\mu$ -**strongly convex**:

$$f(\mathbf{x}') \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2. \quad (11)$$

This implies (assuming  $f$  is twice differentiable):

$$\mathbf{0} \prec \mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}.$$



# Convergence of GD for Strongly Convex Problems

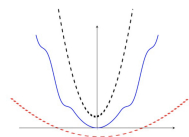
**Theorem (see Appendix D).**

$f(\mathbf{x})$ :  $\mu$ -strongly convex and  $L$ -Lipschitz continuous.

For gradient descent with a step size  $t = \frac{2}{L+\mu}$ , we have:

$$\|\mathbf{x}_k - \mathbf{x}_\star\|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|\mathbf{x}_0 - \mathbf{x}_\star\|_2, \quad (12)$$

where  $\kappa = L/\mu$  and  $\mathbf{x}_\star$  is the minimizer.

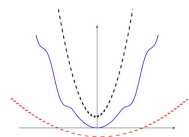


# Convergence of GD for Strongly Convex Problems

## Theorem (see Appendix D).

$f(\mathbf{x})$ :  $\mu$ -strongly convex and  $L$ -Lipschitz continuous.

For gradient descent with a step size  $t = \frac{2}{L+\mu}$ , we have:



$$\|\mathbf{x}_k - \mathbf{x}_\star\|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|\mathbf{x}_0 - \mathbf{x}_\star\|_2, \quad (12)$$

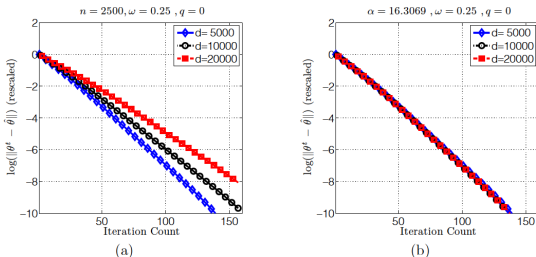
where  $\kappa = L/\mu$  and  $\mathbf{x}_\star$  is the minimizer.

## Convergence Rates for Gradient Descent:

- ①  $f$  non-smooth:  $O(1/\sqrt{k})$ .
- ②  $f$  differentiable:  $O(1/k)$ .
- ③  $f$  smooth,  $\nabla f$  Lipschitz:  $O(1/k^2)$ .
- ④  $f$  strongly convex:  $O(e^{-\alpha k})$ .



# Convergence of Restricted Strong Convex Problems



**Figure 1.** Convergence rates of projected gradient descent in application to Lasso programs ( $\ell_1$ -constrained least-squares). Each panel shows the log optimization error  $\log \|\theta^t - \hat{\theta}\|$  versus the iteration number  $t$ . Panel (a) shows three curves, corresponding to dimensions  $d \in \{5000, 10000, 20000\}$ , sparsity  $s = \lceil \sqrt{d} \rceil$ , and all with the same sample size  $n = 2500$ . All cases show geometric convergence, but the rate for larger problems becomes progressively slower. (b) For an appropriately rescaled sample size ( $\alpha = \frac{n}{s \log d}$ ), all three convergence rates should be roughly the same, as predicted by the theory.

- **Fact:** Structured signal recovery problems such as LASSO and PCP satisfy **restricted strong convexity**.
- Hence, gradient descent enjoys **globally linear convergence** up to the statistical precision of the model.<sup>5</sup>

<sup>5</sup>Fast global convergence of gradient methods for high-dimensional statistical recovery, Agarwal, Negahban, Wainwright, NIPS 2010.

- 1 Motivating Examples for Recovery of Low-Dim Models
- 2 (Accelerated) Proximal Methods
- 3 Alternating Direction Methods of Multipliers (ADMM)
- 4 Summary & Extensions

# Optimization Challenges for Structured Data Recovery

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \doteq \underbrace{f(\mathbf{x})}_{\text{smooth convex}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth convex}}. \quad (13)$$

- **Challenge of Scale:** scale algorithms to when  $n$  is very large.

$$\text{Second order methods} \implies \text{First order methods...} \quad (14)$$

- **Nonsmoothness:** first order methods are slow for nonsmooth.

$$\mathcal{O}(1/\sqrt{k}) \implies \mathcal{O}(1/k) \implies \mathcal{O}(1/k^2) \implies \mathcal{O}(e^{-\alpha k}) \quad (15)$$

- **Equality Constraints:** augmented Lagrange multiplier (ALM).
- **Separable Structures:** alternating direction of multipliers method (ADMM).

# Equality Constrained Problems with Separable Structures

Let us consider the **two-block equality constrained problem**:

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}), \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}.$$

- $g : \mathbb{R}^n \mapsto \mathbb{R}$  and  $h : \mathbb{R}^n \mapsto \mathbb{R}$  are (probably nonsmooth) convex functions.
- $\mathbf{A}$  and  $\mathbf{B}$  are matrices and  $\mathbf{y} \in \text{range}([\mathbf{A} \mid \mathbf{B}])$ , so that the problem is feasible.

# Constrained Nonsmooth Problem (Examples)

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}), \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}.$$

- **Basis pursuit denoising** (let  $\mathbf{z} = \mathbf{x}$ ):

$$\begin{aligned} & \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1 \\ \iff & \min_{\mathbf{x}, \mathbf{z}} \underbrace{\frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2}_{g(\mathbf{x})} + \underbrace{\lambda \|\mathbf{z}\|_1}_{h(\mathbf{z})}, \quad \text{s.t.} \quad \mathbf{x} - \mathbf{z} = \mathbf{0}. \end{aligned}$$

# Constrained Nonsmooth Problem (Examples)

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}), \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}.$$

- **Basis pursuit denoising** (let  $\mathbf{z} = \mathbf{x}$ ):

$$\begin{aligned} & \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1 \\ \iff & \min_{\mathbf{x}, \mathbf{z}} \underbrace{\frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2}_{g(\mathbf{x})} + \underbrace{\lambda \|\mathbf{z}\|_1}_{h(\mathbf{z})}, \quad \text{s.t.} \quad \mathbf{x} - \mathbf{z} = \mathbf{0}. \end{aligned}$$

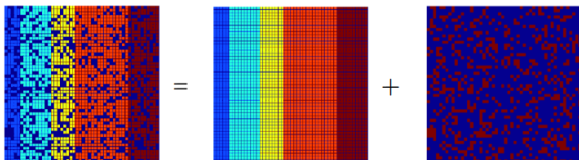
- **Stable low-rank matrix recovery** (let  $\mathbf{Z} = \mathbf{X}$ ):

$$\begin{aligned} & \min_{\mathbf{X}} \frac{1}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2 + \lambda \|\mathbf{X}\|_* \\ \iff & \min_{\mathbf{X}, \mathbf{Z}} \underbrace{\frac{1}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2}_{g(\mathbf{X})} + \underbrace{\lambda \|\mathbf{Z}\|_*}_{h(\mathbf{Z})}, \quad \text{s.t.} \quad \mathbf{X} - \mathbf{Z} = \mathbf{0}. \end{aligned}$$

# Examples: Constrained Nonsmooth Problem

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}), \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}.$$

- **Robust PCA**



$$\min_{\mathbf{L}, \mathbf{S}} \underbrace{\|\mathbf{L}\|_*}_{g(\mathbf{L})} + \lambda \underbrace{\|\mathbf{S}\|_1}_{h(\mathbf{S})}, \quad \text{s.t.} \quad \mathbf{L} + \mathbf{S} = \mathbf{Y}.$$

# Linear Equality Constrained Optimization

Let us first consider a simpler **one-block constrained problem**:

$$\min_{\mathbf{x}} g(\mathbf{x}) \quad \text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{y}, \quad (16)$$

where

- $g : \mathbb{R}^n \rightarrow \mathbb{R}$  is a (probably nonsmooth) convex function,
- $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{y} \in \text{range}(\mathbf{A})$  (so that the problem is feasible).



# Linear Equality Constrained Optimization

Let us first consider a simpler **one-block constrained problem**:

$$\min_{\mathbf{x}} g(\mathbf{x}) \quad \text{subject to} \quad \mathbf{Ax} = \mathbf{y}, \quad (16)$$

where

- $g : \mathbb{R}^n \rightarrow \mathbb{R}$  is a (probably nonsmooth) convex function,
- $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{y} \in \text{range}(\mathbf{A})$  (so that the problem is feasible).

**A Natural Attempt:** solve the unconstrained by penalizing the constraint:

$$\hat{\mathbf{x}}(\mu) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 \quad \text{for a large } \mu. \quad (17)$$

- **Pros:** As  $\mu \rightarrow +\infty$ ,  $\hat{\mathbf{x}}(\mu) \rightarrow \mathbf{x}_\star$  (the “continuation method”).
- **Cons:** The rate of convergence depends on  $L = \mu \|\mathbf{A}\|_2^2$ .

# A More Principled Approach via Lagrangian

## Definition (The Lagrange Duality)

The *Lagrangian* function of the constrained problem (16):

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) \doteq g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle,$$

where  $\boldsymbol{\lambda} \in \mathbb{R}^m$  is a vector of *Lagrange multipliers*. This gives a *dual function*:

$$d(\boldsymbol{\lambda}) \doteq \inf_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle.$$

# A More Principled Approach via Lagrangian

## Definition (The Lagrange Duality)

The *Lagrangian* function of the constrained problem (16):

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) \doteq g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle,$$

where  $\boldsymbol{\lambda} \in \mathbb{R}^m$  is a vector of *Lagrange multipliers*. This gives a *dual function*:

$$d(\boldsymbol{\lambda}) \doteq \inf_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle.$$

**Fact** (credited to Lagrange):  $\exists \boldsymbol{\lambda}_*$  such that the optimal solution  $(\mathbf{x}_*, \boldsymbol{\lambda}_*)$  is a saddle point of the Lagrangian:

$$\sup_{\boldsymbol{\lambda}} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \sup_{\boldsymbol{\lambda}} \inf_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle = \sup_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda}).$$

# Dual Ascent Algorithm for the Lagrangian

**Fact:** If

$$\mathbf{x}'(\boldsymbol{\lambda}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle,$$

then  $\mathbf{Ax}'(\boldsymbol{\lambda}) - \mathbf{y}$  is a gradient  $\nabla d(\boldsymbol{\lambda})$  of the concave dual function  $d(\boldsymbol{\lambda})$  at  $\boldsymbol{\lambda}$ .

# Dual Ascent Algorithm for the Lagrangian

**Fact:** If

$$\mathbf{x}'(\boldsymbol{\lambda}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle,$$

then  $\mathbf{A}\mathbf{x}'(\boldsymbol{\lambda}) - \mathbf{y}$  is a gradient  $\nabla d(\boldsymbol{\lambda})$  of the concave dual function  $d(\boldsymbol{\lambda})$  at  $\boldsymbol{\lambda}$ .

**A Natural Attempt** to find the saddle point  $(\mathbf{x}_*, \boldsymbol{\lambda}_*)$  is via *dual ascent*:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k), \quad (18)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + t_{k+1}(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}). \quad (19)$$

- For certain problem classes, this converges to the optimal  $(\mathbf{x}_*, \boldsymbol{\lambda}_*)$ .
- However, unfortunately it fails for problems in our settings.

## An Example of Failure

Consider the basis pursuit problem:

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{y}.$$

We have:

$$d(\boldsymbol{\lambda}) = \inf_{\mathbf{x}} \|\mathbf{x}\|_1 + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}).$$

The *dual ascent* algorithm:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + t_{k+1} (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}).$$

- For certain problem classes, dual ascent yields efficient convergent algorithms to an optimal primal-dual solution  $(\mathbf{x}_*, \boldsymbol{\lambda}_*)$ .
- However, it may *fail* for problems in structured signal recovery.

## An Example of Failure

Consider the basis pursuit problem:

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{y}.$$

One can show that the dual function

$$d(\boldsymbol{\lambda}) = \inf_{\mathbf{x}} \|\mathbf{x}\|_1 + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle = \begin{cases} -\infty & \|\mathbf{A}^\top \boldsymbol{\lambda}\|_\infty > 1 \\ -\langle \boldsymbol{\lambda}, \mathbf{y} \rangle & \|\mathbf{A}^\top \boldsymbol{\lambda}\|_\infty \leq 1 \end{cases}$$

Whenever the dual ascent step (19) happens to produce a  $\boldsymbol{\lambda}$  such that  $\|\mathbf{A}^* \boldsymbol{\lambda}\|_\infty > 1$ , the algorithm will break down.

## An Example of Failure

Consider the basis pursuit problem:

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{y}.$$

One can show that the dual function

$$d(\boldsymbol{\lambda}) = \inf_{\mathbf{x}} \|\mathbf{x}\|_1 + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle = \begin{cases} -\infty & \|\mathbf{A}^\top \boldsymbol{\lambda}\|_\infty > 1 \\ -\langle \boldsymbol{\lambda}, \mathbf{y} \rangle & \|\mathbf{A}^\top \boldsymbol{\lambda}\|_\infty \leq 1 \end{cases}$$

Whenever the dual ascent step (19) happens to produce a  $\boldsymbol{\lambda}$  such that  $\|\mathbf{A}^\top \boldsymbol{\lambda}\|_\infty > 1$ , the algorithm will break down.

**The reason is  $g(\mathbf{x}) = \|\mathbf{x}\|_1$  here is not “strongly” convex enough to dominate the linear term  $\langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} \rangle$ .**



## Remedy: Augmented Lagrangian

### Definition (Augmented Lagrangian Function [Hestenes'69, Powell'69])

The augmented Lagrangian function is defined as

$$\mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda}) := g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2,$$

where  $\mu > 0$  is a penalty parameter.

The augmented Lagrangian can be regarded as the Lagrangian for

$$\min_{\mathbf{x}} \underbrace{g(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2}_{\text{strongly convex}}, \quad \text{s.t.} \quad \mathbf{Ax} = \mathbf{y},$$

which has the same optimal solution as the original un-penalized problem.

## Augmented Lagrange Multiplier

Apply dual ascent to  $\mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda})$  with a particular step size  $t_{k+1} = \mu$ ,

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda}_k), \quad (20)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}). \quad (21)$$

**Fact:**  $\mathbf{x}_{k+1}$  always minimizes the unaugmented Lagrangian  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_{k+1})$  at  $\boldsymbol{\lambda} \equiv \boldsymbol{\lambda}_{k+1}$ , because:

$$\begin{aligned} \mathbf{0} &\in \partial \mathcal{L}_\mu(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_k), \\ &= \partial g(\mathbf{x}_{k+1}) + \mathbf{A}^* \boldsymbol{\lambda}_k + \mu \mathbf{A}^* (\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y}), \\ &= \partial g(\mathbf{x}_{k+1}) + \mathbf{A}^* \boldsymbol{\lambda}_{k+1}, \\ &= \partial \mathcal{L}(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_{k+1}). \end{aligned}$$

**$\boldsymbol{\lambda}_{k+1}$  is always feasible, no bad behaviors!**

# Augmented Lagrange Multiplier

## Augmented Lagrange Multiplier (ALM)

**Problem Class:**  $\min_{\mathbf{x}} g(\mathbf{x})$  subject to  $\mathbf{Ax} = \mathbf{y}$ .  
with  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  convex and coercive,  $\mathbf{y} \in \text{range}(\mathbf{A})$ .

**Basic Iteration:** set

$$\mathcal{L}_{\mu}(\mathbf{x}, \boldsymbol{\lambda}) = g(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2.$$

Repeat:

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_{\mu}(\mathbf{x}, \boldsymbol{\lambda}_k),$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu (\mathbf{Ax}_{k+1} - \mathbf{y}).$$

**Convergence Guarantee:**

$\{\mathbf{x}_k\}$  converges to an optimal solution at a rate  $O(1/k)$ .

## Example: ALM for Basis Pursuit

### Augmented Lagrange Multiplier (ALM) for BP

- 1: **Problem:**  $\min_{\mathbf{x}} \|\mathbf{x}\|_1$  subject to  $\mathbf{y} = \mathbf{A}\mathbf{x}$ ,  
 given  $\mathbf{y} \in \mathbb{R}^m$  and  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . The augmented Lagrangian is:

$$\mathcal{L}_\mu(\mathbf{x}, \boldsymbol{\lambda}) = \|\mathbf{x}\|_1 + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2.$$

- 2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\boldsymbol{\lambda}_0 \in \mathbb{R}^m$ , and  $\beta > 1$ .  
 3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**  
 4:    $\mathbf{x}_{k+1} \leftarrow \arg \min \mathcal{L}_{\mu_k}(\mathbf{x}, \boldsymbol{\lambda}_k)$  using APG.  
 5:    $\boldsymbol{\lambda}_{k+1} \leftarrow \boldsymbol{\lambda}_k + \mu_k(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{y})$ .  
 6:    $\mu_{k+1} \leftarrow \min\{\beta\mu_k, \mu_{\max}\}$ .  
 7: **end for**  
 8: **Output:**  $\mathbf{x}_\star \leftarrow \mathbf{x}_K$ .

# Application of ALM to the Two-Block Problem

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}), \quad \text{s.t.} \quad \mathbf{Ax} + \mathbf{Bz} = \mathbf{y}.$$

- Form the augmented Lagrangian

$$\begin{aligned} \mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) &= g(\mathbf{x}) + h(\mathbf{z}) + \langle \mathbf{Ax} + \mathbf{Bz} - \mathbf{y}, \boldsymbol{\lambda} \rangle \\ &\quad + \frac{\mu}{2} \|\mathbf{Ax} + \mathbf{Bz} - \mathbf{y}\|_2^2 \end{aligned}$$

- Solve the problem via

$$(\mathbf{x}_{k+1}, \mathbf{z}_{k+1}) \in \arg \min_{\mathbf{x}, \mathbf{z}} \mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}_k), \quad (\text{could be expensive})$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu_k (\mathbf{Ax}_{k+1} + \mathbf{Bz}_{k+1} - \mathbf{y}).$$

# Application of ALM to the Two-Block Problem

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}), \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}.$$

- Form the augmented Lagrangian

$$\begin{aligned} \mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) &= g(\mathbf{x}) + h(\mathbf{z}) + \langle \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{y}, \boldsymbol{\lambda} \rangle \\ &\quad + \frac{\mu}{2} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{y}\|_2^2 \end{aligned}$$

- Solve the problem via

$$(\mathbf{x}_{k+1}, \mathbf{z}_{k+1}) \in \arg \min_{\mathbf{x}, \mathbf{z}} \mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}_k), \quad (\text{could be expensive})$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu_k (\mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}\mathbf{z}_{k+1} - \mathbf{y}).$$

**The primal subproblem for  $\mathbf{x}$  does not have closed-form and could be expensive to solve.**

# Alternating Directions Method of Multipliers (ADMM)

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}), \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}.$$

**Remedy for solving**  $\min_{\mathbf{x}, \mathbf{z}} \mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}_k)$ :

- fix  $\mathbf{z}$  and  $\boldsymbol{\lambda}$ , minimize  $\mathbf{x}$ :

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_\mu(\mathbf{x}, \mathbf{z}_k, \boldsymbol{\lambda}_k),$$

- fix  $\mathbf{x}$  and  $\boldsymbol{\lambda}$ , minimize  $\mathbf{z}$ :

$$\mathbf{z}_{k+1} \in \arg \min_{\mathbf{z}} \mathcal{L}_\mu(\mathbf{x}_{k+1}, \mathbf{z}, \boldsymbol{\lambda}_k),$$

- fix  $\mathbf{x}$  and  $\mathbf{z}$ , take a dual ascent step on  $\boldsymbol{\lambda}$ :

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu_k \nabla_{\boldsymbol{\lambda}} \mathcal{L}_\mu(\mathbf{x}_{k+1}, \mathbf{z}_{k+1}, \boldsymbol{\lambda}).$$

## Solutions for Subproblems via Proximal Operators

- fix  $z$  and  $\lambda$ , minimize  $x$ :

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \left\{ g(\mathbf{x}) + \frac{\mu}{2} \left\| \mathbf{A}\mathbf{x} + \mathbf{B}z_k - \mathbf{y} + \frac{1}{\mu} \boldsymbol{\lambda}_k \right\|_2^2 \right\},$$

- fix  $x$  and  $\lambda$ , minimize  $z$ :

$$z_{k+1} = \arg \min_z \left\{ h(z) + \frac{\mu}{2} \left\| \mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}z - \mathbf{y} + \frac{1}{\mu} \boldsymbol{\lambda}_k \right\|_2^2 \right\},$$

- fix  $x$  and  $z$ , take a dual ascent step on  $\lambda$ :

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu_k (\mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}z_{k+1} - \mathbf{y}).$$



## Solutions for Subproblems via Proximal Operators

- fix  $z$  and  $\lambda$ , minimize  $x$ :

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \left\{ g(\mathbf{x}) + \frac{\mu}{2} \left\| \mathbf{A}\mathbf{x} + \mathbf{B}z_k - \mathbf{y} + \frac{1}{\mu} \boldsymbol{\lambda}_k \right\|_2^2 \right\},$$

- fix  $x$  and  $\lambda$ , minimize  $z$ :

$$\mathbf{z}_{k+1} = \arg \min_z \left\{ h(z) + \frac{\mu}{2} \left\| \mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}z - \mathbf{y} + \frac{1}{\mu} \boldsymbol{\lambda}_k \right\|_2^2 \right\},$$

- fix  $x$  and  $z$ , take a dual ascent step on  $\lambda$ :

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu_k (\mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}z_{k+1} - \mathbf{y}).$$

**The solution of each subproblem is the proximal operator, which has closed-form solution for structured  $g$  and  $h$ !**

## Optimization with Separable Structures

The augmented Lagrangian  $\mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda})$  is:

$$\mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = g(\mathbf{x}) + h(\mathbf{z}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} + \mathbf{Bz} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{Ax} + \mathbf{Bz} - \mathbf{y}\|_2^2.$$

The **alternating directions method of multipliers** (ADMM) conducts a simple, alternating iteration:

$$\mathbf{z}_{k+1} \in \arg \min_{\mathbf{z}} \mathcal{L}_\mu(\mathbf{x}_k, \mathbf{z}, \boldsymbol{\lambda}_k), \quad (22)$$

$$\mathbf{x}_{k+1} \in \arg \min_{\mathbf{x}} \mathcal{L}_\mu(\mathbf{x}, \mathbf{z}_{k+1}, \boldsymbol{\lambda}_k), \quad (23)$$

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \mu (\mathbf{Ax}_{k+1} + \mathbf{Bz}_{k+1} - \mathbf{y}). \quad (24)$$

This is also known as the *Gauss-Seidel iteration*.

**ADMM converges at a rate of  $O(1/k)$ . (proof no picnic<sup>6</sup>)**

<sup>6</sup>On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators. J. Eckstein and D. Bertsekas, 1992.

## Example I: Basis Pursuit

**Basis pursuit denoising** (let  $z = x$ ):

$$\begin{aligned} & \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1 \\ \iff & \min_{\mathbf{x}, \mathbf{z}} \underbrace{\frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2}_{g(\mathbf{x})} + \underbrace{\lambda \|\mathbf{z}\|_1}_{h(\mathbf{z})}, \quad \text{s.t.} \quad \mathbf{x} - \mathbf{z} = \mathbf{0}. \end{aligned}$$

- Form the augmented Lagrangian:

$$\mathcal{L}_\mu(\mathbf{x}, \mathbf{z}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{z}\|_1 + \langle \boldsymbol{\beta}, \mathbf{x} - \mathbf{z} \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{z}\|_2^2.$$

## Example I: Basis Pursuit Denoising

- Fix  $\mathbf{z}_k$  and  $\beta_k$ , find  $\mathbf{x}_{k+1}$  via

$$\begin{aligned}\mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \mathcal{L}_{\mu}(\mathbf{x}, \mathbf{z}_k, \beta_k) \\ &= \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \frac{\mu}{2} \left\| \mathbf{x} - \mathbf{z}_k + \frac{1}{\mu} \beta_k \right\|_2^2 \right\}.\end{aligned}$$

- Fix  $\mathbf{x}_{k+1}$  and  $\beta_k$ , find  $\mathbf{z}_{k+1}$  via

$$\begin{aligned}\mathbf{z}_{k+1} &= \arg \min_{\mathbf{z}} \mathcal{L}_{\mu}(\mathbf{x}_{k+1}, \mathbf{z}, \beta_k) \\ &= \arg \min_{\mathbf{z}} \left\{ \lambda \|\mathbf{z}\|_1 + \frac{\mu}{2} \left\| \mathbf{x}_{k+1} - \mathbf{z} + \frac{1}{\mu} \beta_k \right\|_2^2 \right\}.\end{aligned}$$

- Fix  $\mathbf{x}_{k+1}$  and  $\mathbf{z}_{k+1}$ , take a dual ascent step on  $\beta$ :

$$\beta_{k+1} = \beta_k + \mu_k (\mathbf{x}_{k+1} - \mathbf{z}_{k+1}).$$

## Example I: Basis Pursuit

- Fix  $\mathbf{z}_k$  and  $\boldsymbol{\beta}_k$ , find  $\mathbf{x}_{k+1}$  via

$$\mathbf{x}_{k+1} = (\mathbf{A}^T \mathbf{A} + \mu \mathbf{I})^{-1} (\mathbf{A}^T \mathbf{y} + \mu \mathbf{z}_k - \boldsymbol{\beta}_k)$$

- Fix  $\mathbf{x}_{k+1}$  and  $\boldsymbol{\beta}_k$ , find  $\mathbf{z}_{k+1}$  via

$$\mathbf{z}_{k+1} = \text{prox}_{\lambda \mu^{-1} \|\cdot\|_1} \left( \mathbf{x}_{k+1} + \frac{1}{\mu} \boldsymbol{\beta}_k \right)$$

- Fix  $\mathbf{x}_{k+1}$  and  $\mathbf{z}_{k+1}$ , take a dual ascent step on  $\boldsymbol{\beta}$ :

$$\boldsymbol{\beta}_{k+1} = \boldsymbol{\beta}_k + \mu_k (\mathbf{x}_{k+1} - \mathbf{z}_{k+1}).$$

## Example II: Robust PCA

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1, \quad \text{s.t.} \quad \mathbf{L} + \mathbf{S} = \mathbf{Y}.$$

- Form the augmented Lagrangian:

$$\begin{aligned} \mathcal{L}_\mu(\mathbf{L}, \mathbf{S}, \mathbf{\Lambda}) &= \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 + \langle \mathbf{\Lambda}, \mathbf{L} + \mathbf{S} - \mathbf{Y} \rangle \\ &\quad + \frac{\mu}{2} \|\mathbf{L} + \mathbf{S} - \mathbf{Y}\|_F^2. \end{aligned}$$

## Example II: Robust PCA

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1, \quad \text{s.t.} \quad \mathbf{L} + \mathbf{S} = \mathbf{Y}.$$

- Fix  $\mathbf{S}_k$  and  $\mathbf{\Lambda}_k$ , find  $\mathbf{L}_{k+1}$  via

$$\begin{aligned} \mathbf{L}_{k+1} &= \arg \min_{\mathbf{L}} \mathcal{L}_\mu(\mathbf{L}, \mathbf{S}_k, \mathbf{\Lambda}_k) \\ &= \arg \min_{\mathbf{L}} \left\{ \|\mathbf{L}\|_* + \frac{\mu}{2} \left\| \mathbf{L} + \mathbf{S}_k - \mathbf{Y} + \frac{1}{\mu} \mathbf{\Lambda}_k \right\|_F^2 \right\} \\ &= \text{prox}_{\mu^{-1} \|\cdot\|_*} (\mathbf{Y} - \mathbf{S}_k - \mu^{-1} \mathbf{\Lambda}_k). \end{aligned}$$

## Example II: Robust PCA

$$\min_{\mathbf{L}, \mathbf{S}} \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1, \quad \text{s.t.} \quad \mathbf{L} + \mathbf{S} = \mathbf{Y}.$$

- Fix  $\mathbf{L}_{k+1}$  and  $\mathbf{\Lambda}_k$ , find  $\mathbf{S}_{k+1}$  via

$$\begin{aligned} \mathbf{S}_{k+1} &= \arg \min_{\mathbf{S}} \mathcal{L}_\mu(\mathbf{L}_{k+1}, \mathbf{S}, \mathbf{\Lambda}_k) \\ &= \arg \min_{\mathbf{S}} \left\{ \lambda \|\mathbf{S}\|_1 + \frac{\mu}{2} \left\| \mathbf{L}_{k+1} + \mathbf{S} - \mathbf{Y} + \frac{1}{\mu} \mathbf{\Lambda}_k \right\|_F^2 \right\} \\ &= \text{prox}_{\lambda\mu^{-1}\|\cdot\|_1} \left( \mathbf{Y} - \mathbf{L}_{k+1} - \mu^{-1} \mathbf{\Lambda}_k \right). \end{aligned}$$



## Example II: Robust PCA

- Fix  $\mathbf{S}_k$  and  $\mathbf{\Lambda}_k$ , find  $\mathbf{L}_{k+1}$  via

$$\mathbf{L}_{k+1} = \text{prox}_{\mu^{-1}\|\cdot\|_*}(\mathbf{Y} - \mathbf{S}_k - \mu^{-1}\mathbf{\Lambda}_k).$$

- Fix  $\mathbf{L}_{k+1}$  and  $\mathbf{\Lambda}_k$ , find  $\mathbf{S}_{k+1}$  via

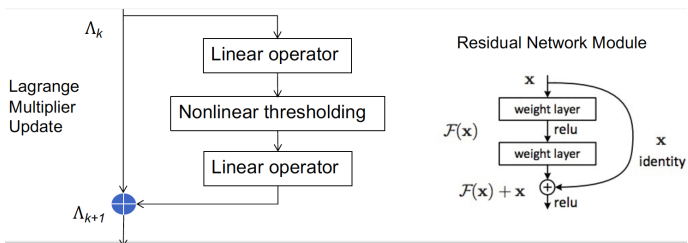
$$\mathbf{S}_{k+1} = \text{prox}_{\lambda\mu^{-1}\|\cdot\|_1}(\mathbf{Y} - \mathbf{L}_{k+1} - \mu^{-1}\mathbf{\Lambda}_k).$$

- Fix  $\mathbf{L}_{k+1}$  and  $\mathbf{S}_{k+1}$ , take a dual ascent step on  $\mathbf{\Lambda}$ :

$$\mathbf{\Lambda}_{k+1} = \mathbf{\Lambda}_k + \mu_k(\mathbf{L}_{k+1} + \mathbf{S}_{k+1} - \mathbf{Y}).$$

## Example II: Robust PCA

- 1: **Problem:**  $\min_{L, S} \mathcal{L}_\mu(L, S, \Lambda)$ , given  $Y$ ,  $\lambda, \mu > 0$ .
- 2: **Input:**  $L_0, S_0, \Lambda_0 \in \mathbb{R}^{m \times n}$ .
- 3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**
- 4:    $L_{k+1} \leftarrow \text{prox}_{\mu^{-1}\|\cdot\|_*} [Y - S_k - \mu^{-1}\Lambda_k]$ .
- 5:    $S_{k+1} \leftarrow \text{prox}_{\lambda\mu^{-1}\|\cdot\|_1} [Y - L_{k+1} - \mu^{-1}\Lambda_k]$ .
- 6:    $\Lambda_{k+1} \leftarrow \Lambda_k + \mu(L_{k+1} + S_{k+1} - Y)$ .
- 7: **end for**
- 8: **Output:**  $L_* \leftarrow L_K; S_* \leftarrow S_K$ .



# Summary for ADMM

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}), \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}.$$

- ADMM is easy to implement and use, and scalable for large-scale problem.
- ADMM is slow to converge to high accuracy (with  $O(1/k)$  convergence rate).

- 1 Motivating Examples for Recovery of Low-Dim Models
- 2 (Accelerated) Proximal Methods
- 3 Alternating Direction Methods of Multipliers (ADMM)
- 4 Summary & Extensions

# Optimization Challenges for Structured Data Recovery

- **Challenge of Scale:** scale algorithms to when  $n$  is very large.

Second order methods  $\implies$  First order methods...

- **Nonsmoothness:** first order methods are slow for nonsmooth.

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \doteq \underbrace{f(\mathbf{x})}_{\text{smooth convex}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth convex}},$$

that we deal with proximal gradient method:

$$\underbrace{\mathcal{O}(1/\sqrt{k})}_{\text{subgradient}} \implies \underbrace{\mathcal{O}(1/k)}_{\text{proximal gradient}} \implies \underbrace{\mathcal{O}(1/k^2)}_{\text{APG}} \implies \underbrace{\mathcal{O}(e^{-\alpha k})}_{\text{RSC}}$$

# Optimization Challenges for Structured Data Recovery

- **Challenge of Scale:** scale algorithms to when  $n$  is very large.

Second order methods  $\implies$  First order methods...

- **Nonsmoothness:** first order methods are slow for nonsmooth.

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \doteq \underbrace{f(\mathbf{x})}_{\text{smooth convex}} + \underbrace{g(\mathbf{x})}_{\text{nonsmooth convex}},$$

that we deal with proximal gradient method:

$$\underbrace{\mathcal{O}(1/\sqrt{k})}_{\text{subgradient}} \implies \underbrace{\mathcal{O}(1/k)}_{\text{proximal gradient}} \implies \underbrace{\mathcal{O}(1/k^2)}_{\text{APG}} \implies \underbrace{\mathcal{O}(e^{-\alpha k})}_{\text{RSC}}$$

- **Equality Constraints with Separable Structures:** ADMM

$$\min_{\mathbf{x}, \mathbf{z}} g(\mathbf{x}) + h(\mathbf{z}), \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{y}.$$

## Other Ideas for Better Scalability

Typical optimization problem:  $\min_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m h_i(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^n.$

**Complexity = per iteration cost  $\times$  # of iterations.**

- **Block Coordinate Descent** reduces dependency on the dimension  $n$ :

$$O(n) \rightarrow O(n^{1/2}).$$

- **Stochastic Gradient Descent** (with variance reduction) reduces dependency on sample size  $m$ :

$$O(m) \rightarrow O(m^{1/2}).$$

- **Acceleration Schemes** reduce the number of iterations  $k$ :

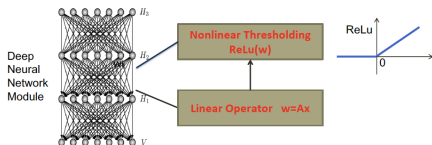
$$O(\epsilon^{-2}) \rightarrow O(\epsilon^{-1/2}).$$

Nonconvex programs are a different story... (later, Lecture 3).

## Algorithmic Unrolling Beyond LISTA

- 1: **Problem:**  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$ , given  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .
- 2: **Input:**  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $L \geq \lambda_{\max}(\mathbf{A}^* \mathbf{A})$ .
- 3: **for** ( $k = 0, 1, 2, \dots, K - 1$ ) **do**
- 4:      $\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^* (\mathbf{A}\mathbf{x}_k - \mathbf{y})$ .
- 5:      $\mathbf{x}_{k+1} \leftarrow \text{soft}_{\lambda/L}(\mathbf{w}_k)$ .
- 6: **end for**
- 7: **Output:**  $\mathbf{x}_* \leftarrow \mathbf{x}_K$ .

The unrolled iterations resemble a deep neural network!



We can optimize the optimization path of ISTA using supervised learning<sup>7</sup>:

$$\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^* (\mathbf{A}\mathbf{x}_k - \mathbf{y}) = \underbrace{\left( \mathbf{I} - \frac{1}{L} \mathbf{A}^* \mathbf{A} \right)}_{\text{learnable parameter } \mathbf{S}} \mathbf{x}_k + \underbrace{\frac{1}{L} \mathbf{A}^* \mathbf{A} \mathbf{y}}_{\text{learnable parameter } \mathbf{b}}$$

<sup>7</sup>Learning Fast Approximations of Sparse Coding, Karol Gregor and Yann LeCun, ICML 2010.



## Summary

**Next lecture:** Learning Low-dimensional Models via Nonconvex Optimization.

# Thank You! Questions?