

The background image shows the exterior of a classical-style building with large columns and a pedimented roof. A vertical banner hangs from one of the columns, featuring a portrait of a person and the text "CATHERINE COPIE 100 NIMES ROAD JUNE 11 - SEPTEMBER 11, 2016". Above the entrance, the words "ALUMNI MEMORIAL" are visible. To the right of the building, the University of Michigan logo (a yellow "M" on a blue square) is displayed, along with the text "UNIVERSITY OF MICHIGAN".

# EECS 559 Optimization Methods for SIPML

Lecture 8 – Accelerated Proximal Gradient Methods

Instructor: Prof. Qing Qu ([qingqu@umich.edu](mailto:qingqu@umich.edu))

# Lecture Agenda

- Accelerated Proximal Gradient Method
- Proximal Gradient Homotopy Method

# Lecture Agenda

- Accelerated Proximal Gradient Method
- Proximal Gradient Homotopy Method

# Unconstrained Minimization

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x})$$

- Suppose  $f(\boldsymbol{x})$  belongs the class of convex and smooth with its gradient  $L$ -Lipschitz:  $f \in \mathcal{F}_{L,R}$  with

$$\mathcal{F}_{L,R} := \{f : \mathcal{B}(\mathbf{0}, R) \mapsto \mathbb{R} \mid \|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{x}')\|_2 \leq L \|\boldsymbol{x} - \boldsymbol{x}'\|_2\}$$

- Solve the problem via *iterative optimization methods*, producing a sequence of points

$$\boldsymbol{x}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_k, \dots$$

starting from an initialization  $\boldsymbol{x}_0$ .

# Suboptimal Convergence Rate

- Now suppose the iterates  $\{\boldsymbol{x}_i\}_{i \geq 0}$  are generated from a *black-box model*:

$$\boldsymbol{x}_{k+1} = \varphi_k \left( \{\boldsymbol{x}_i\}_{i=0}^k, \{f(\boldsymbol{x}_i)\}_{i=0}^k, \{\nabla f(\boldsymbol{x}_i)\}_{i=0}^k \right)$$

- We ask how well an algorithm does on the *worst-function* in  $\mathcal{F}_{L,R}$  ?

$$\sup_{f \in \mathcal{F}_{L,R}, \boldsymbol{x}_0} \left\{ f(\boldsymbol{x}_k) - \inf_{\boldsymbol{x}} f(\boldsymbol{x}) \right\}$$

# Suboptimal Convergence Rate

- The gradient descent method achieves

$$\sup_{f \in \mathcal{F}_{L,R}, \mathbf{x}_0} \left\{ f(\mathbf{x}_k) - \inf_{\mathbf{x}} f(\mathbf{x}) \right\} \leq \frac{CLR^2}{k}$$

- The best lower bound proved by Nesterov'83:

$$\sup_{f \in \mathcal{F}_{L,R}, \mathbf{x}_0} \left\{ f(\mathbf{x}_k) - \inf_{\mathbf{x}} f(\mathbf{x}) \right\} \geq \frac{cLR^2}{k^2}$$

# Suboptimal Convergence Rate

	Method	Iterate	Convergence
Smooth	Gradient Descent (GD)	$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau \nabla f(\mathbf{x}_k)$	$O(1/k)$
	Accelerated GD	$\mathbf{p}_{k+1} = \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1})$ $\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \tau \nabla f(\mathbf{p}_{k+1})$	$O(1/k^2)$
Nonsmooth	Proximal Gradient (PG)	$\mathbf{x}_{k+1} = \mathbf{x}_k - \mu G_\mu(\mathbf{x}_k)$	$O(1/k)$
	Accelerated PG?	??	$O(1/k^2)??$

# Suboptimal Convergence Rate

	Method	Iterate	Convergence
Smooth	Gradient Descent (GD)	$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau \nabla f(\mathbf{x}_k)$	$O(1/k)$
	Accelerated GD	$\mathbf{p}_{k+1} = \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1})$ $\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \tau \nabla f(\mathbf{p}_{k+1})$	$O(1/k^2)$
Nonsmooth	Proximal Gradient (PG)	$\mathbf{x}_{k+1} = \mathbf{x}_k - \mu G_\mu(\mathbf{x}_k)$	$O(1/k)$
	Accelerated PG?	$\mathbf{p}_{k+1} = \mathbf{x}_k + \beta_{k+1} \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}),$ $\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \mu G_\mu(\mathbf{p}_{k+1})$	$O(1/k^2)??$

# Nonsmooth Composite Problems

$$\min_{\boldsymbol{x}} F(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})$$

- $f$ : convex, continuously differentiable, and  $L$ -smooth,

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

- $g$ : convex but may *not* be differentiable.

# Accelerated Proximal Gradient (APG)

- Accelerated Gradient Descent

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

$$\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \tau \nabla f(\mathbf{p}_{k+1})$$

- Accelerated Proximal Gradient

$$\mathbf{p}_{k+1} = \mathbf{x}_k + \beta_{k+1} \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}),$$

$$\begin{aligned}\mathbf{x}_{k+1} &= \text{prox}_{\mu g}(\mathbf{p}_{k+1} - \mu \nabla f(\mathbf{p}_{k+1})) \\ &= \mathbf{p}_{k+1} - \mu G_\mu(\mathbf{p}_{k+1})\end{aligned}$$

# Accelerated Proximal Gradient (APG)

$$\min_{\boldsymbol{x}} F(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x}).$$

$$\boldsymbol{p}_{k+1} = \boldsymbol{x}_k + \beta_{k+1} \cdot (\boldsymbol{x}_k - \boldsymbol{x}_{k-1}),$$

$$\begin{aligned}\boldsymbol{x}_{k+1} &= \text{prox}_{\mu g}(\boldsymbol{p}_{k+1} - \mu \nabla f(\boldsymbol{p}_{k+1})) \\ &= \boldsymbol{p}_{k+1} - \mu G_\mu(\boldsymbol{p}_{k+1})\end{aligned}$$

$$\tau_1 = 1, \quad \tau_{k+1} \frac{1 + \sqrt{1 + 4\tau_k^2}}{2}, \quad \beta_{k+1} = \frac{\tau_k - 1}{\tau_{k+1}}.$$

# Convergence of APG

## Theorem (Convergence of APG)

Let the sequence  $\{\mathbf{x}_k\}_{k \geq 1}$  be generated by the APG for the convex composite function  $F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$  with  $\mu = 1/L$ , where  $f$  is  $L$ -smooth. Let  $\mathbf{x}_\star$  be a minimizer of  $F(\mathbf{x})$ , then for any  $k \geq 1$ ,

$$F(\mathbf{x}_k) - F(\mathbf{x}_\star) \leq \frac{2L \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2}{(k+1)^2}.$$

# Suboptimal Convergence Rate

	Method	Iterate	Convergence
Smooth	Gradient Descent (GD)	$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau \nabla f(\mathbf{x}_k)$	$O(1/k)$
	Accelerated GD	$\mathbf{p}_{k+1} = \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1})$ $\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \tau \nabla f(\mathbf{p}_{k+1})$	$O(1/k^2)$
Nonsmooth	Proximal Gradient (PG)	$\mathbf{x}_{k+1} = \mathbf{x}_k - \mu G_\mu(\mathbf{x}_k)$	$O(1/k)$
	Accelerated PG?	$\mathbf{p}_{k+1} = \mathbf{x}_k + \beta_{k+1} \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}),$ $\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \mu G_\mu(\mathbf{p}_{k+1})$	$O(1/k^2)$

# Example: Solving Lasso

---

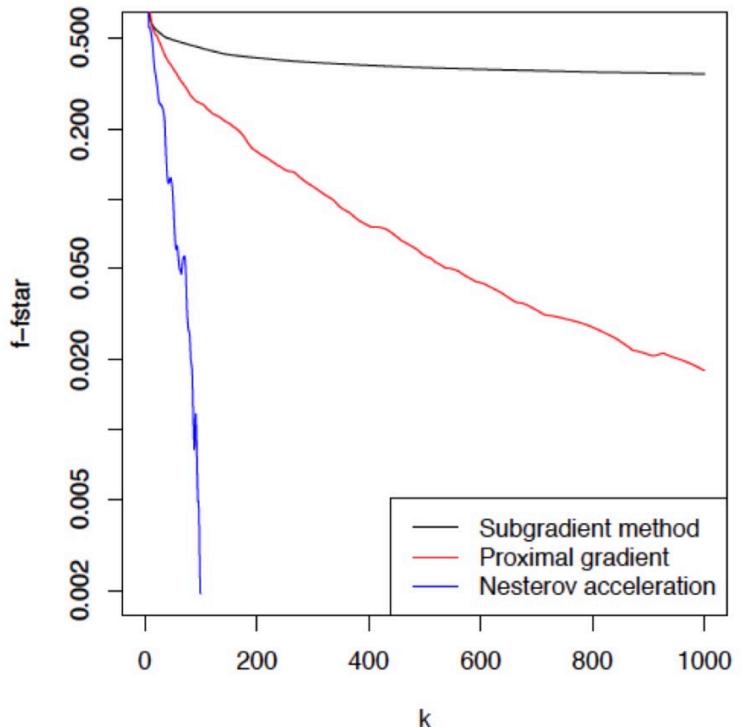
## Algorithm 8.3 Accelerated Proximal Gradient (APG) for BPDN

---

- 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:**  $\mathbf{x}_* \leftarrow \mathbf{x}_K$ .
-

# Comparison on Lasso

$$\min_{\mathbf{x}} F(\mathbf{x}) = \underbrace{\frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2}_{f(\mathbf{x})} + \underbrace{\lambda \|\mathbf{x}\|_1}_{g(\mathbf{x})}$$



# Example: Solving Stable PCP

---

**Algorithm 8.4** 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  $\mathbf{Y}$ ,  $\lambda, \mu > 0$ .
  - 2: **Input:**  $\mathbf{L}_0 \in \mathbb{R}^{m \times n}$ ,  $\mathbf{S}_0 \in \mathbb{R}^{m \times n}$ ,  $\mathbf{P}_1^S = \mathbf{S}_1 \leftarrow \mathbf{S}_0$ ,  $\mathbf{P}_1^L = \mathbf{L}_1 \leftarrow \mathbf{L}_0$ ,  $t_1 \leftarrow 1$ .
  - 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:    $\mathbf{P}_{k+1}^L \leftarrow \mathbf{L}_k + \beta_{k+1}(\mathbf{L}_k - \mathbf{L}_{k-1})$ .
  - 6:    $\mathbf{P}_{k+1}^S \leftarrow \mathbf{S}_k + \beta_{k+1}(\mathbf{S}_k - \mathbf{S}_{k-1})$ .
  - 7:    $\mathbf{W}_{k+1} \leftarrow \mathbf{Y} - \mathbf{P}_{k+1}^S$  and compute the SVD:  $\mathbf{W}_{k+1} = \mathbf{U}_{k+1} \boldsymbol{\Sigma}_{k+1} \mathbf{V}_{k+1}^*$ .
  - 8:    $\mathbf{L}_{k+1} \leftarrow \mathbf{U}_{k+1} \text{soft}(\boldsymbol{\Sigma}_{k+1}, 1/\mu) \mathbf{V}_{k+1}^*$ .
  - 9:    $\mathbf{S}_{k+1} \leftarrow \text{soft}((\mathbf{Y} - \mathbf{P}_{k+1}^L), \lambda/\mu)$ .
  - 10: **end for**
  - 11: **Output:**  $\mathbf{L}_* \leftarrow \mathbf{L}_K$ ;  $\mathbf{S}_* \leftarrow \mathbf{S}_K$ .
-

# Convergence for Strongly Convex Function

**Theorem. (Convergence for  $\gamma$ -Strongly Convex Function)**

Suppose  $f$  is  $L$ -smooth and  $\gamma$ -strongly convex. Let  $\mu = 1/L$  and  $\kappa = L/\gamma$ , then for any  $k \geq 1$ ,

$$F(\mathbf{x}_k) - F(\mathbf{x}_*) \leq \left(1 - \frac{1}{\sqrt{\kappa}}\right)^k \left( F(\mathbf{x}_0) - F(\mathbf{x}_*) + \frac{\mu \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{2} \right)$$

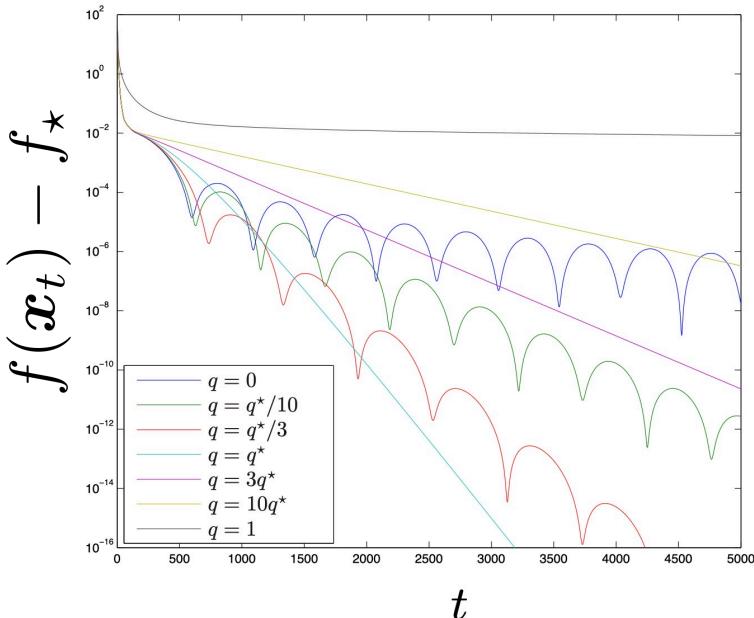
$$\begin{aligned} \mathbf{p}_{k+1} &= \mathbf{x}_k + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}), \quad \beta_k = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}, \\ \mathbf{x}_{k+1} &= \text{prox}_{\mu g} (\mathbf{p}_{k+1} - \mu \nabla f(\mathbf{p}_{k+1})). \end{aligned}$$

# A Practical Issue

- Fast convergence requires knowledge of estimating  $\kappa = L/\gamma$ , strong convexity parameter  $\gamma$  could be challenging in practice.
- For sparse recovery, the APG method is also called the Fast Iterative Shrinkage-Thresholding Algorithm (FISTA).
- A common observation: *ripples/bumps* in the traces of cost values.

# Rippling Behavior

Numerical example:  $p_{k+1} = x_k + \frac{1-\sqrt{q}}{1+\sqrt{q}} (x_k - x_{k-1})$ ,  $q^* = 1/\kappa$



- Period of ripples is often proportional to  $\kappa^{1/2} = \sqrt{L/\gamma}$
- when  $q < q^*$ : we overestimate momentum — rippling behavior
- when  $q > q^*$ : we underestimate momentum — slow convergence

# Adaptive Restart (O'Donoghue, Candes'12)

When a certain criterion is met, *restart* running APG with

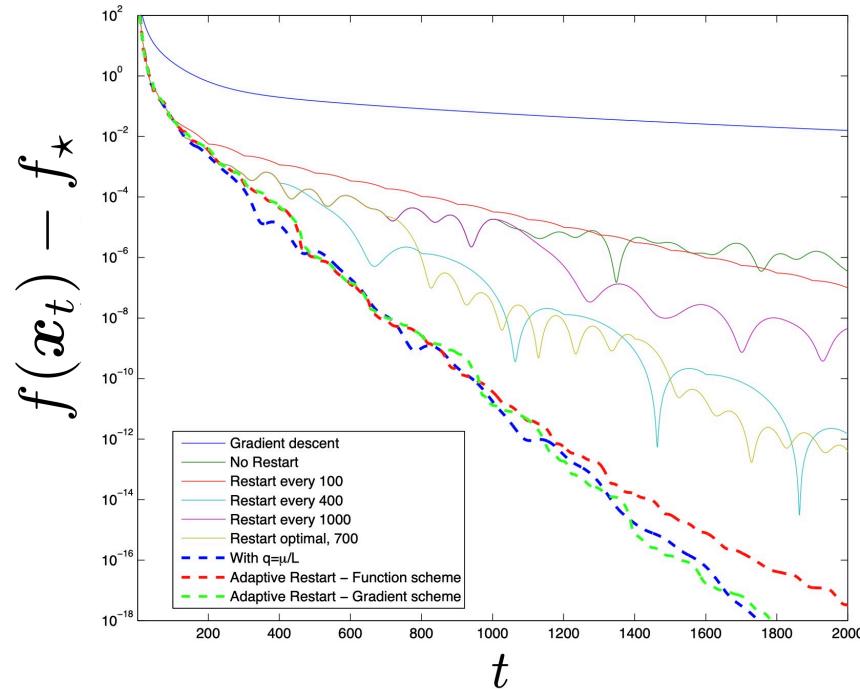
$$\mathbf{p}_0 \leftarrow \mathbf{x}_k$$

$$\mathbf{x}_0 \leftarrow \mathbf{x}_k$$

$$t_0 \leftarrow 1$$

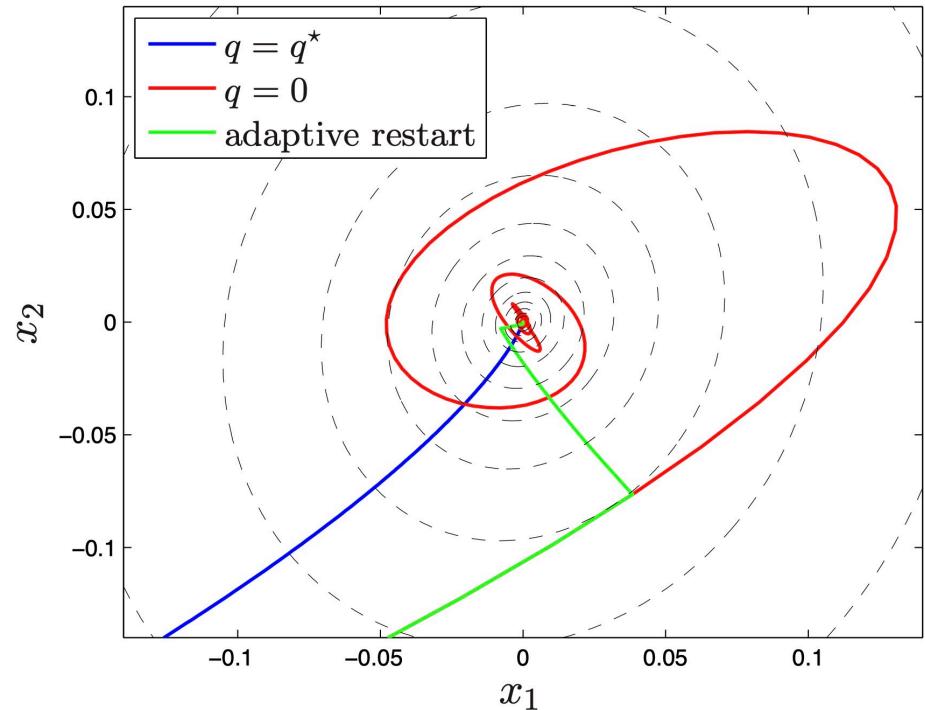
- Take the current iterate as a new start pointing;
- Erase all memory of previous iterates and reset the momentum back to zero.

# Adaptive Restart (O'Donoghue, Candes'12)



- **function scheme:**  
restart if  $f(\mathbf{x}_k) > f(\mathbf{x}_{k-1})$
- **gradient scheme:**  
restart if  $\langle \nabla f(\mathbf{p}_k), \mathbf{x}_k - \mathbf{x}_{k-1} \rangle > 0$

# Adaptive Restart (O'Donoghue, Candes'12)



- With overestimated momentum (e.g.,  $q = 0$ ), one sees spiraling trajectory.
- Adaptive restart helps mitigate this issue.

# Summary of Proximal Methods

	Method	Iterate	Convergence
Smooth	Gradient Descent (GD)	$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau \nabla f(\mathbf{x}_k)$	$O(1/k)$
	Accelerated GD	$\mathbf{p}_{k+1} = \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1})$ $\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \tau \nabla f(\mathbf{p}_{k+1})$	$O(1/k^2)$
Nonsmooth	Proximal Gradient (PG)	$\mathbf{x}_{k+1} = \mathbf{x}_k - \mu G_\mu(\mathbf{x}_k)$	$O(1/k)$
	Accelerated PG?	$\mathbf{p}_{k+1} = \mathbf{x}_k + \beta_{k+1}(\mathbf{x}_k - \mathbf{x}_{k-1})$ $\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \tau G_\mu(\mathbf{p}_{k+1})$	$O(1/k^2)$

# Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma.  
**(Chapter 8.3)**
- *Proximal Algorithms, Foundations & Trends in Optimization.*  
Neal Parikh, Stephen Boyd, 2014.  
[https://web.stanford.edu/~boyd/papers/prox\\_algs.html](https://web.stanford.edu/~boyd/papers/prox_algs.html)
- *A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems.* Amir Beck and Marc Teboulle, 2009.  
<https://pubs.siam.org/doi/abs/10.1137/080716542?mobileUi=0>

# Lecture Agenda

- Accelerated Proximal Gradient Method
- Proximal Gradient Homotopy Method

# Nonsmooth Composite Problems

$$\min_{\mathbf{x}} F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$$

- $f$ : convex, continuously differentiable, and  $L$ -smooth,
- $g$ : convex but may *not* be differentiable.

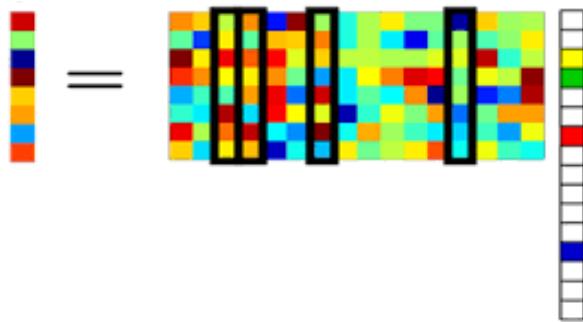
We know that

- APG achieves  $O(1/k^2)$  convergence

*Can we have faster convergence of proximal method by exploiting extra data structures?*

# Sparse Recovery

$$\mathbf{y} \in \mathbb{R}^m \quad \mathbf{A} \in \mathbb{R}^{m \times n} \quad \mathbf{x}_\star \in \mathbb{R}^n$$



**Problem:** given the measurement  $\mathbf{y} \in \mathbb{R}^m$  and sensing matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  ( $m \ll n$ ) recover the sparsest  $\mathbf{x} \in \mathbb{R}^n$  from

$$\mathbf{y} = \mathbf{A} \cdot \mathbf{x}_\star + \mathbf{n}.$$

# Stable Sparse Recovery via Lasso

Basis pursuit denoising (BPDN)

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

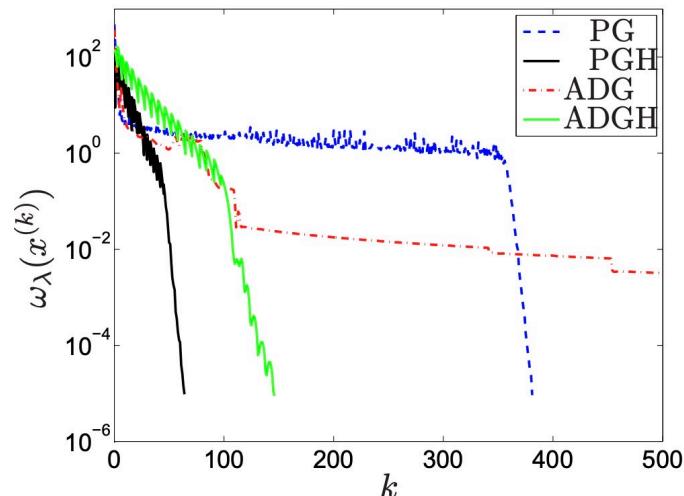
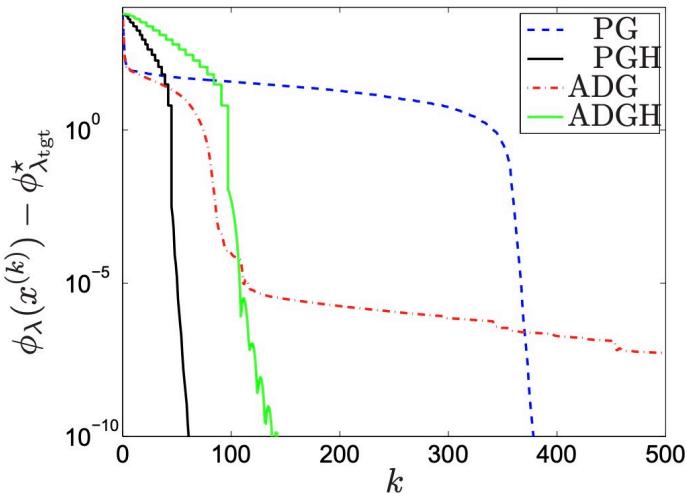
equivalent to

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2 + \lambda \cdot \|\boldsymbol{x}\|_1,$$

for some properly chosen  $\lambda > 0$ . The problem is termed as Lasso (least absolute shrinkage and selection operator).

# Solving Lasso

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \cdot \|\mathbf{x}\|_1,$$



- PG: proximal gradient
- PGH: proximal gradient-homotopy
- ADG: accelerated proximal gradient
- ADGH: accelerated proximal gradient - homotopy

# Solving Lasso: Linear Convergence?

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

- In theory, for convex (but *not* strongly convex) problems, the convergence rate for APG is  $O(1/k^2)$ .
- For solving Lasso, why is the final stage convergence *linear*?

**Intuition:** in the final stage, the solution is *sparse*, the problem satisfies certain *restricted strong convexity property*.

# Restricted Strong Convexity and Smoothness

**Theorem.** Let  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$  and suppose  $\mathbf{A}$  satisfies REC. Suppose  $\mathbf{x}$  and  $\mathbf{x}'$  are two *sparse vectors* such that  $|\text{supp}(\mathbf{x}) \cup \text{supp}(\mathbf{x}')| \leq s$  for some integer  $s < m$ . Then the following hold:

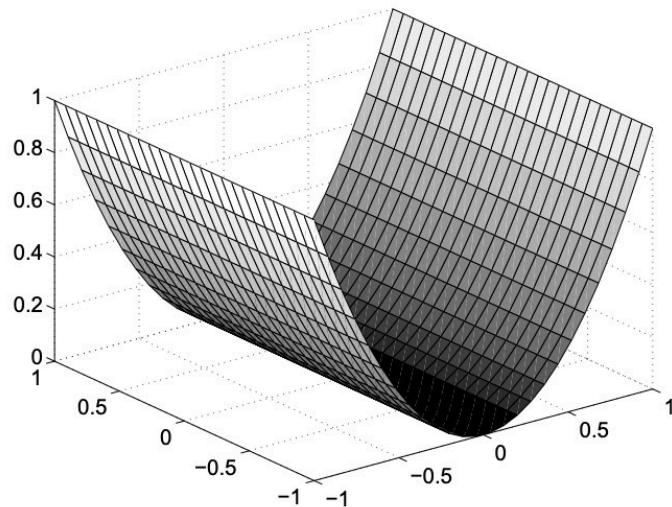
- **Restricted smoothness:**

$$f(\mathbf{x}') \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{\rho_+(\mathbf{A}, s)}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2,$$

- **Restricted strong convexity:**

$$f(\mathbf{x}') \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{\rho_-(\mathbf{A}, s)}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2.$$

# Restricted Strong Convexity and Smoothness



- When  $\mathbf{A}$  satisfies REC, the function  $F(\mathbf{x})$  is restricted strongly convex along *sparse* directions.
- Restricted strong convexity implies *linear* convergence.

# Restricted Eigenvalue Condition

**Definition. (Restricted eigenvalue condition, REC)**

Given an integer  $s > 0$ , we say that  $\mathbf{A}$  satisfies the REC at the sparsity level  $s$  if there exist positive constants  $\rho_+(\mathbf{A}, s)$  and  $\rho_-(\mathbf{A}, s)$  such that

$$\rho_+(\mathbf{A}, s) = \sup_{\|\mathbf{x}\|_2=1, \|\mathbf{x}\|_0 \leq s} \|\mathbf{Ax}\|_2^2,$$

$$\rho_-(\mathbf{A}, s) = \inf_{\|\mathbf{x}\|_2=1, \|\mathbf{x}\|_0 \leq s} \|\mathbf{Ax}\|_2^2.$$

- If  $\mathbf{A}$  satisfies RIP, then  $\rho_+ = 1 + \delta$  and  $\rho_- = 1 - \delta$ .

# Restricted Isometry Property (RIP)

**Definition. (Restricted isometry property, RIP)**

Given an integer  $s > 0$ , we say that  $\mathbf{A}$  satisfies the RIP at sparsity level  $s$  if there exist a positive constant  $\delta > 0$  such that

$$(1 - \delta) \cdot \|\mathbf{x}\|_2^2 \leq \|\mathbf{Ax}\|_2^2 \leq (1 + \delta) \cdot \|\mathbf{x}\|_2^2$$

for all  $\mathbf{x}$  with  $\|\mathbf{x}\|_0 \leq s$ .

**Example:** when  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is i.i.d. Gaussian  $A_{ij} \sim_{i.i.d.} \mathcal{N}(0, 1)$  and  $m \geq \Omega(\frac{n}{\delta} \log n)$ , it satisfies  $\delta$ -RIP with high probability.

# Idea: Homotopy Continuation of $\lambda$

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \cdot \|\mathbf{x}\|_1,$$

Adaptively adjust the  $\lambda$  to make  $\{\mathbf{x}_k\}_{k \geq 1}$  *all sparse* along the solution path, so that the restricted strong convexity is satisfied.

- Pick a large initial  $\lambda_0 = \|\mathbf{A}^\top \mathbf{y}\|_\infty$  so that  $\mathbf{x}_0$  is sparse.
- *Shrink* the penalty  $\lambda$  to the desired  $\lambda_{tgt}$  *gradually*, making sure all  $\{\mathbf{x}_k\}_{k \geq 1}$  are always *sparse*.

# Proximal Gradient Homotopy Method

---

**Algorithm 1** Proximal Gradient Homotopy Method

---

**Input:**  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{y} \in \mathbb{R}^m$ ,  $\lambda_{tgt} > 0$ ,  $\varepsilon > 0$ ,  $L_{\min} > 0$

**Parameters:**  $\eta \in (0, 1)$ ,  $\delta \in (0, 1)$

**Initialize:**  $\lambda_0 \leftarrow \|\mathbf{A}^\top \mathbf{y}\|_\infty$ ,  $\varepsilon_0 \leftarrow \delta \lambda_0$ ,  $\hat{\mathbf{x}}_0 \leftarrow \mathbf{0}$ ,  $M_0 \leftarrow L_{\min}$

$N \leftarrow \lfloor \log(\lambda_0 / \lambda_{tgt}) / \log(1/\eta) \rfloor$

**for**  $K = 0, 1, \dots, N - 1$  **do**

$\{\hat{\mathbf{x}}_{K+1}, M_{K+1}\} \leftarrow \text{ProxGrad}(\lambda_K, \varepsilon_K, \hat{\mathbf{x}}_K, M_K)$

$\lambda_{K+1} \leftarrow \eta \lambda_K$ ,  $\varepsilon_{K+1} \leftarrow \delta \lambda_{K+1}$

**end for**

$\{\hat{\mathbf{x}}_{tgt}, M_{tgt}\} \leftarrow \text{ProxGrad}(\lambda_{tgt}, \varepsilon, \hat{\mathbf{x}}_N, M_N)$

**Return:**  $\hat{\mathbf{x}}_{tgt}$

---

# Linear Convergence

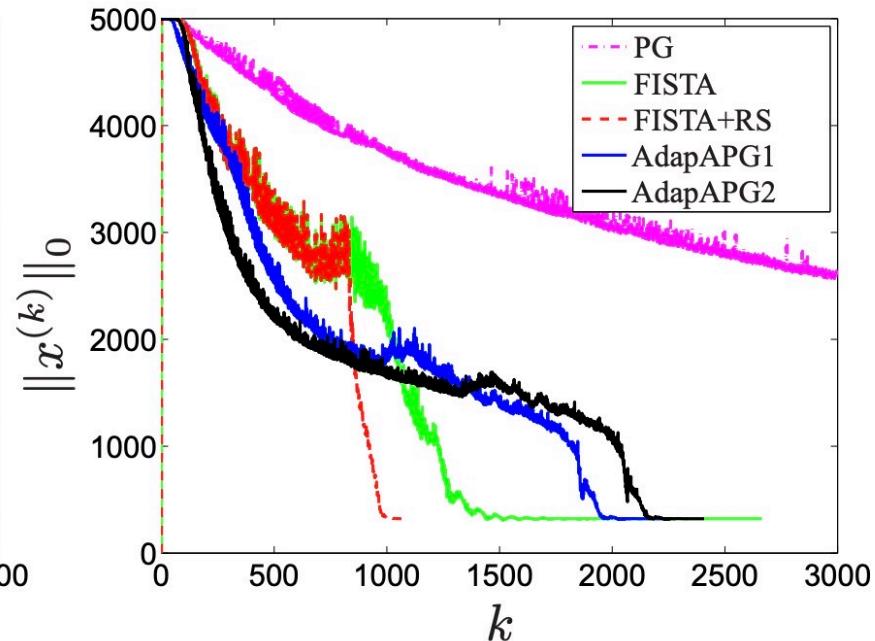
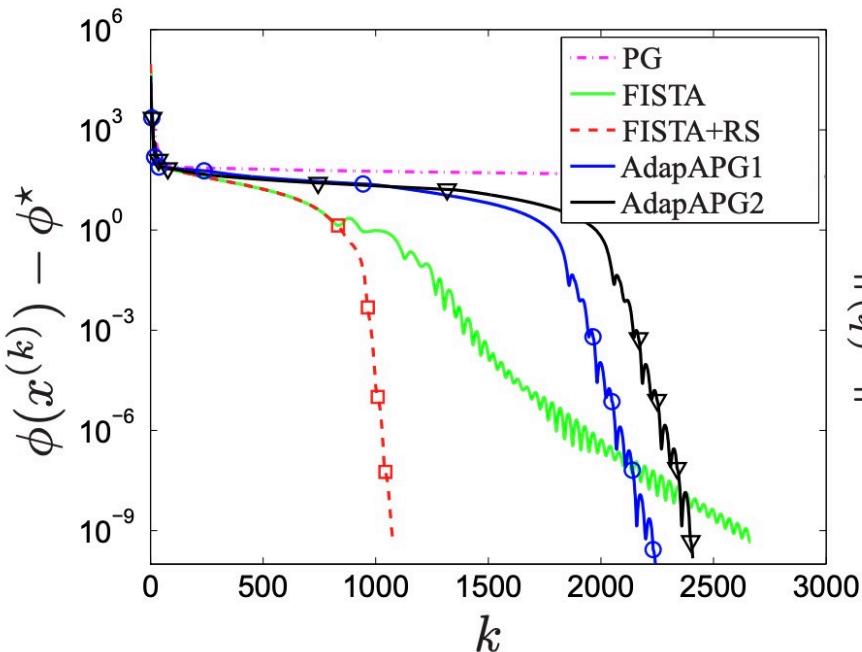
## Theorem. (Informal)

Suppose  $A \in \mathbb{R}^{m \times n}$  satisfies REC. Under mild assumptions and with properly chosen  $\delta$  and  $\eta$ , the function value  $F(x_k)$  of the proximal gradient homotopy method converges with

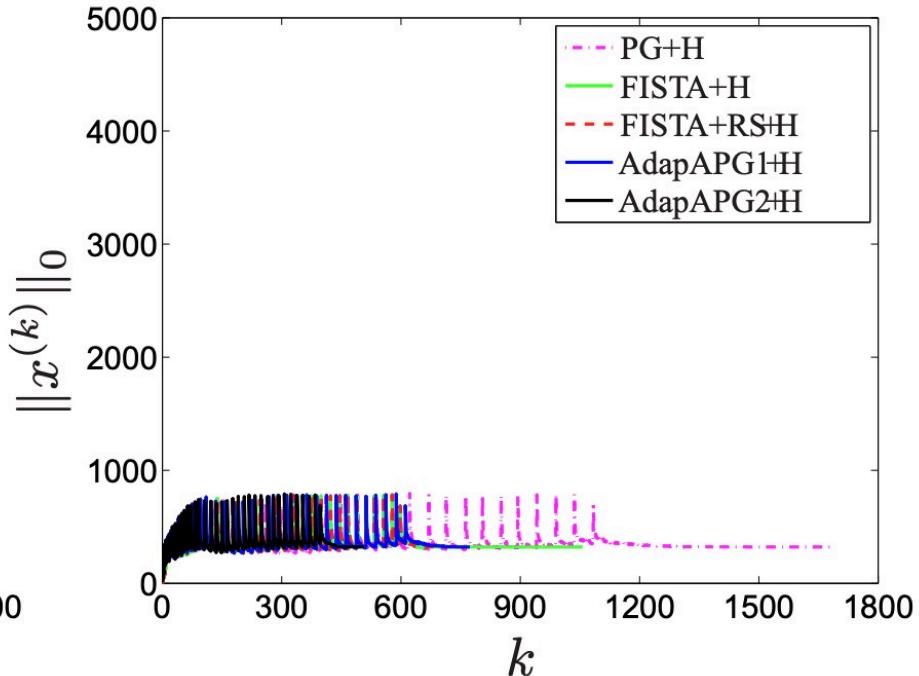
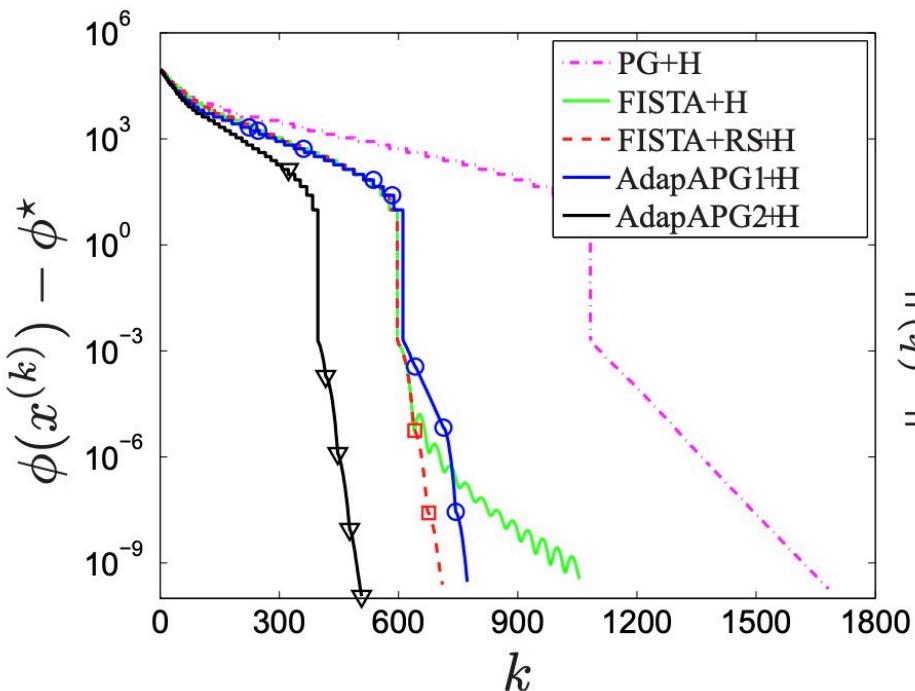
$$\#\text{Iter} \leq O(\kappa \log(\lambda_0/\varepsilon)),$$

where  $\varepsilon > 0$  is the precision and  $\kappa := \rho_+(A, s)/\rho_-(A, s)$ .

# Sublinear Convergence



# Linear Convergence with Homotopy



# Further Readings

- *Fixed-Point Continuation for  $\ell_1$ -Minimization: Methodology and Convergence.* SIAM Journal on Optimization, Elaine T. Hale, Wotao Yin, and Yin Zhang, 2008.  
<https://locus.siam.org/doi/abs/10.1137/070698920>
- *A Proximal-Gradient Homotopy Method for the Sparse Least-Squares Problem.* Lin Xiao, Tong Zhang, 2012.  
<https://arxiv.org/abs/1203.3002>
- *An Adaptive Accelerated Proximal Gradient Method and its Homotopy Continuation for Sparse Optimization.* Qihang Lin, Xiao Lin, 2014. [https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/apg\\_homotopy.pdf](https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/apg_homotopy.pdf)