# Computational Principles for High-dim Data Analysis (Lecture Thirteen)

#### Yi Ma

EECS Department, UC Berkeley

October 12, 2021





EECS208, Fall 2021

## Unconstrained Convex Optimization for Structured Data Recovery

1 Challenges and Opportunities

2 Proximal Gradient Methods

3 Accelerated Proximal Gradient Methods

"Since the fabric of the universe is most perfect and the work of a most wise Creator, nothing at all takes place in the universe in which some rule of maximum or minimum does not appear."

- Leonhard Euler

### Optimization Problems for Structured Data Recovery

**Sparse Vector Recovery:** recover a sparse  $x_o$  from  $y = Ax_o \in \mathbb{R}^m$  or  $y = Ax_o + z \in \mathbb{R}^m$  via convex programs:

• **Basis Pursuit** (BP):

$$\min_{\boldsymbol{x}} \|\boldsymbol{x}\|_1 \quad \text{subject to} \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}. \tag{1}$$

• LASSO:  

$$\min_{\boldsymbol{x}} \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{A} \boldsymbol{x}\|_{2}^{2} + \lambda \|\boldsymbol{x}\|_{1}.$$
(2)

## Optimization Problems for Structured Data Recovery

Matrix Completion or Recovery: recover a low-rank  $L_o$  from incomplete  $Y = \mathcal{P}_{\Omega}[X_o]$  or corrupted  $Y = L_o + S_o \in \mathbb{R}^{m \times n}$  via convex programs:

Matrix Completion:

$$\min \|\boldsymbol{X}\|_*$$
 subject to  $\mathcal{P}_{\Omega}[\boldsymbol{X}] = \boldsymbol{Y}.$  (3)

• Principal Component Pursuit (PCP):

$$\min_{\boldsymbol{L},\boldsymbol{S}} \|\boldsymbol{L}\|_* + \lambda \|\boldsymbol{S}\|_1 \quad \text{subject to} \quad \boldsymbol{L} + \boldsymbol{S} = \boldsymbol{Y}. \tag{4}$$

• Stable PCP:

$$\min_{\boldsymbol{L},\boldsymbol{S}} \|\boldsymbol{L}\|_{*} + \lambda \|\boldsymbol{S}\|_{1} + \frac{\mu}{2} \|\boldsymbol{Y} - \boldsymbol{L} - \boldsymbol{S}\|_{F}^{2}.$$
 (5)

(4 何) トイヨト イヨト

Optimization Challenges for Structured Data Recovery



• Challenge of Scale: scale algorithms to when n is very large.

Second order methods  $\implies$  First order methods... (7)

• Nonsmoothness: first order methods are slow for nonsmooth.

$$O(1/\sqrt{k}) \implies O(1/k) \implies O(1/k^2)...$$
 (8)

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

< □ > < 同 > < 三 > < 三 >

# Gradient Descent [Cauchy, 1847]

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

$$\min f(\boldsymbol{x}), \quad \boldsymbol{x} \in \mathsf{C} \text{ (a convex set)},$$
 (9)

conduct local gradient descent search (App. D):

 $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \gamma_k \nabla f(\boldsymbol{x}_k), \quad (10)$ 

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



figure courtesy of prof. Carlos Fernandez of NYU.



3

< □ > < 同 > < 回 > < 回 > < 回 >

### Gradient Descent

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

 $\|\nabla f(\boldsymbol{x}') - \nabla f(\boldsymbol{x})\|_2 \le L \|\boldsymbol{x}' - \boldsymbol{x}\|_2, \quad \forall \boldsymbol{x}', \boldsymbol{x} \in \mathbb{R}^n.$ (11)

This gives a matching quadratic upper bound:

$$\begin{split} f(\boldsymbol{x}') &\leq \hat{f}(\boldsymbol{x}', \boldsymbol{x}) \\ &\doteq f(\boldsymbol{x}) + \langle \nabla f(\boldsymbol{x}), \boldsymbol{x}' - \boldsymbol{x} \rangle + \frac{L}{2} \left\| \boldsymbol{x}' - \boldsymbol{x} \right\|_2^2 \\ &= \frac{L}{2} \left\| \boldsymbol{x}' - (\boldsymbol{x} - \frac{1}{L} \nabla f(\boldsymbol{x})) \right\|_2^2 + h(\boldsymbol{x}). \end{split}$$



Take a step to the minimizer of this bound:

$$x_{k+1} = \arg\min_{x'} \hat{f}(x', x_k) = x_k - \frac{1}{L} \nabla f(x_k).$$
 (12)

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

EECS208, Fall 2021

## Proximal Gradient Descent

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

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

Upper bound:

$$\hat{F}(\boldsymbol{x}, \boldsymbol{x}_{k}) = f(\boldsymbol{x}_{k}) + \langle \nabla f(\boldsymbol{x}_{k}), \boldsymbol{x} - \boldsymbol{x}_{k} \rangle + \frac{L}{2} \|\boldsymbol{x} - \boldsymbol{x}_{k}\|_{2}^{2} + g(\boldsymbol{x}) (14)$$

$$= \frac{L}{2} \|\boldsymbol{x} - (\boldsymbol{x}_{k} - \frac{1}{L} \nabla f(\boldsymbol{x}_{k}))\|_{2}^{2} + g(\boldsymbol{x}) + h(\boldsymbol{x}_{k}). \quad (15)$$

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

$$\boldsymbol{x}_{k+1} = \arg\min_{\boldsymbol{x}} \frac{L}{2} \| \boldsymbol{x} - \underbrace{(\boldsymbol{x}_k - \frac{1}{L} \nabla f(\boldsymbol{x}_k))}_{\boldsymbol{w}_k} \|_2^2 + g(\boldsymbol{x}) \quad (16)$$
$$= \arg\min_{\boldsymbol{x}} g(\boldsymbol{x}) + \frac{L}{2} \| \boldsymbol{x} - \boldsymbol{w}_k \|_2^2. \quad (17)$$

## **Proximal Operators**

#### Definition (Proximal Operator)

The proximal operator of a convex function g is

$$\operatorname{prox}_{g}[\boldsymbol{w}] \doteq \arg\min_{\boldsymbol{x}} \left\{ g(\boldsymbol{x}) + \frac{1}{2} \|\boldsymbol{x} - \boldsymbol{w}\|_{2}^{2} \right\}.$$
(18)

Iteration (17) can be written as:

$$\boldsymbol{x}_{k+1} = \operatorname{prox}_{g/L}[\boldsymbol{w}_k]. \tag{19}$$

For many convex functions g:

 $prox_a[w]$  has a closed form or can be computed efficiently.

< □ > < □ > < □ > < □ > < □ > < □ >

# **Proximal Operators**

#### Proposition

Proximal operators for the  $\ell^1$  norm and nuclear norm are given by:

• Let  $g(x) = \lambda ||x||_1$  be the  $\ell^1$  norm. Then  $prox_g[w]$  is the soft-thresholding function applied element-wise:

$$(\operatorname{prox}_{g}[w])_{i} = \operatorname{soft}\{w_{i}, \lambda\} \doteq \operatorname{sign}(w_{i}) \max(|w_{i}| - \lambda, 0).$$

2 Let  $g(X) = \lambda ||X||_*$  be the matrix nuclear norm. Then  $prox_g[W]$  is the singular-value soft thresholding function:

$$prox_{g}[\boldsymbol{W}] = \boldsymbol{U}soft\{\boldsymbol{\Sigma}, \lambda\}\boldsymbol{V}^{*},$$

where  $(U, \Sigma, V)$  are the SVD of W. In other words,  $prox_g[W]$  applies component-wise soft thresholding on the singular values of W.

イロト 不得 トイヨト イヨト 二日

### **Proximal Operators**

**Proof ideas:** The objective function reaches minimum when the subdifferential of  $\lambda \|x\|_1 + \frac{1}{2} \|x - w\|_2^2$  contains zero,

$$0 \in (\boldsymbol{x} - \boldsymbol{w}) + \lambda \partial \|\boldsymbol{x}\|_1 = \begin{cases} x_i - w_i + \lambda, & x_i > 0\\ -w_i + \lambda[-1, 1], & x_i = 0, \\ x_i - w_i - \lambda, & x_i < 0 \end{cases}$$

#### Thresholding:



# Proximal Gradient Algorithm

#### Proximal Gradient (PG)

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

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

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

$$oldsymbol{w}_k \leftarrow oldsymbol{x}_k - rac{1}{L} 
abla f(oldsymbol{x}_k), \ oldsymbol{x}_{k+1} \leftarrow \operatorname{prox}_{g/L}[oldsymbol{w}_k].$$

#### **Convergence Guarantee:**

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

- 4 回 ト 4 ヨ ト 4 ヨ ト

# Proximal Gradient for LASSO

Iterative soft-thresholding algorithm (ISTA):

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

The unrolled iterations resemble a deep neural network!<sup>1</sup>

EECS208, Fall 2021

## The Heavy Ball Method [Polyak, 1964]

Gradient descent:

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \alpha \nabla f(\boldsymbol{x}_k). \tag{20}$$

The **heavy ball method** (a.k.a the *momentum method*):

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \alpha \nabla f(\boldsymbol{x}_k) + \underbrace{\beta(\boldsymbol{x}_k - \boldsymbol{x}_{k-1})}_{\text{momentum}}.$$
 (21)



- Basis for popular ADAM for train deep neural networks.
- Worst convergence rate is still O(1/k).



EECS208, Fall 2021

## Accelerated Gradient Descent [Nesterov, 1983]

Generate an auxiliary point  $p_{k+1}$  of the form:

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

Move from  $x_k$  to  $p_{k+1}$ , and gradient descend from it:

 $\boldsymbol{x}_{k+1} = \boldsymbol{p}_{k+1} - \alpha \underbrace{\nabla f(\boldsymbol{p}_{k+1})}_{\text{stroke of genius}}$  (22)



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

$$t_1 = 1, \quad t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}, \quad \beta_{k+1} = \frac{t_k - 1}{t_{k+1}}, \quad \alpha = 1/L.$$
 (23)

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

Image: A math a math

## Accelerated Gradient Descent [Nesterov, 1983] Accelerated Proximal Gradient (APG)

### **Problem Class:** $\min_{\boldsymbol{x}} F(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x}),$ f, g convex, with $\nabla f$ L-Lipschitz and g nonsmooth.

**Basic Iteration:** set  $x_0 \in \mathbb{R}^n$ ,  $p_1 = x_1 \leftarrow x_0$ , and  $t_1 \leftarrow 1$ . Repeat for  $k = 1, 2, \dots, K$ :

$$t_{k+1} \leftarrow \frac{1 + \sqrt{1 + 4t_k^2}}{2}, \quad \beta_{k+1} \leftarrow \frac{t_k - 1}{t_{k+1}}.$$
$$p_{k+1} \leftarrow x_k + \beta_{k+1} (x_k - x_{k-1}).$$
$$x_{k+1} \leftarrow \operatorname{prox}_{g/L} \left[ \underbrace{p_{k+1} - \frac{1}{L} \nabla f(p_{k+1})}_{\text{proximal gradient}} \right].$$

#### **Convergence Guarantee:**

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

### GD versus Accelerated GD



Image courtesy of Prof. Qing Qu of Univ. Michigan.

## APG for LASSO

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

1: **Problem:** 
$$\min_{\boldsymbol{x}} \frac{1}{2} \| \boldsymbol{y} - \boldsymbol{A} \boldsymbol{x} \|_{2}^{2} + \lambda \| \boldsymbol{x} \|_{1}$$
, given  $\boldsymbol{y} \in \mathbb{R}^{m}$ ,  $\boldsymbol{A} \in \mathbb{R}^{m \times n}$ .  
2: **Input:**  $\boldsymbol{x}_{0} \in \mathbb{R}^{n}$ ,  $\boldsymbol{p}_{1} = \boldsymbol{x}_{1} \leftarrow \boldsymbol{x}_{0}$ , and  $t_{1} \leftarrow 1$ , and  $L \geq \lambda_{\max}(\boldsymbol{A}^{*}\boldsymbol{A})$ .  
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:  $\boldsymbol{p}_{k+1} \leftarrow \boldsymbol{x}_{k} + \beta_{k+1}(\boldsymbol{x}_{k} - \boldsymbol{x}_{k-1})$ .  
6:  $\boldsymbol{w}_{k+1} \leftarrow \boldsymbol{p}_{k+1} - \frac{1}{L}\boldsymbol{A}^{*}(\boldsymbol{A}\boldsymbol{p}_{k+1} - \boldsymbol{y})$ .  
7:  $\boldsymbol{x}_{k+1} \leftarrow \operatorname{soft}[\boldsymbol{w}_{k+1}, \lambda/L]$ .  
8: end for  
9: **Output:**  $\boldsymbol{x}_{\star} \leftarrow \boldsymbol{x}_{K}$ .

3

< □ > < 同 > < 回 > < 回 > < 回 >

### APG for Stable PCP

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

1: **Problem:**  $\min_{L,S} \|L\|_* + \lambda \|S\|_1 + \frac{\mu}{2} \|Y - L - S\|_F^2$ , given Y. 2: **Input:**  $L_0, S_0 \in \mathbb{R}^{m \times n}, P_1^S = S_1 \leftarrow S_0, P_1^L = L_1 \leftarrow L_0, t_1 \leftarrow 1$ . 3: for (k = 1, 2, ..., 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:  $P_{k+1}^L \leftarrow L_k + \beta_{k+1} (L_k - L_{k-1}); P_{k+1}^S \leftarrow S_k + \beta_{k+1} (S_k - S_{k-1})$ . 6:  $W_{k+1} \leftarrow Y - P_{k+1}^S$  and compute SVD:  $W_{k+1} = U_{k+1} \Sigma_{k+1} V_{k+1}^*$ . 7:  $L_{k+1} \leftarrow U_{k+1} \text{soft} [\Sigma_{k+1}, 1/\mu] V_{k+1}^*; S_{k+1} \leftarrow \text{soft} [(Y - P_{k+1}^L), \lambda/\mu]$ . 8: end for 9: **Output:**  $L_* \leftarrow L_K; S_* \leftarrow S_K$ .

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

# Algorithm: A Little Lesson from History

Comparison from chronological development of algorithms for solving the PCP problem: **the older the algorithm, the more efficient!** 

GOOD NEWS: Scalable first-order gradient-descent algorithms:

- Proximal Gradient [Osher, Mao, Dong, Yin '09,Wright et. al.'09, Cai et. al.'09].
- Accelerated Proximal Gradient [Nesterov '83, Beck and Teboulle '09]:
- Augmented Lagrange Multiplier [Hestenes '69, Powell '69]:
- Alternating Direction Method of Multipliers [Gabay and Mercier '76].

For a 1000x1000 matrix of rank 50, with 10% (100,000) entries randomly corrupted: min  $||A||_* + \lambda ||E||_1$  subj A + E = D.

Algorithms	Accuracy	Rank	E  _0	# iterations	time (sec)	10,000 times speedup!
IT	5.99e-006	50	101,268	8,550	119,370.3	
DUAL	8.65e-006	50	100,024	822	1,855.4	
APG	5.85e-006	50	100,347	134	1,468.9	
APG <sub>P</sub>	5.91e-006	50	100,347	134	82.7	
EALM <sub>P</sub>	2.07e-007	50	100,014	34	37.5	
IALM <sub>P</sub>	3.83e-007	50	99,996	23	11.8	Ļ

Ma (EECS Department, UC Berkeley)

# GD for Strongly Convex Problems

Consider minimizing a L-Lipschitz continuous function

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}), \quad \boldsymbol{x} \in \mathbb{R}^n.$$
(24)

Assume f(x) is  $\mu$ -strongly convex:

$$f((\boldsymbol{x}') \ge f(\boldsymbol{x}) + \langle \nabla f(\boldsymbol{x}), \boldsymbol{x}' - \boldsymbol{x} \rangle + \frac{\mu}{2} \|\boldsymbol{x}' - \boldsymbol{x}\|_2^2.$$
(25)

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(x):  $\mu$ -strongly convex and L-Lipschitz continuous. For gradient descent with a step size  $t = \frac{2}{L+\mu}$ , we have:

$$\|\boldsymbol{x}_k - \boldsymbol{x}_\star\|_2 \le \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|\boldsymbol{x}_0 - \boldsymbol{x}_\star\|_2,$$
 (26)

where  $\kappa = L/\mu$  and  $\boldsymbol{x_{\star}}$  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:  $O(e^{-\alpha k})$ .

(4個) (4回) (4回) (日)

### Convergence of Restricted Strong Convex Problems

**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>2</sup>



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  $||d^- \hat{\theta}||$  versus the iteration number t. Panel (a) shows three curves, corresponding to dimensions  $d \in \{5000, 1000, 20000\}$ , 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  $(\alpha = \frac{1}{s \log d})$ , all three convergence rates should be roughly the same, as predicted by the theory.

<sup>2</sup>Fast global convergence of gradient methods for high-dimensional statistical recovery, Agarwal, Negahban, Wainwright, NIPS 2010.

Ma (EECS Department, UC Berkeley)

EECS208, Fall 202

Assignments

- Reading: Section 8.1 8.3 of Chapter 8. Appendix B, C, and D.
- Programming Homework #3.

э

-