

# Lecture 10: Issues of Large Scale Optimization

- Motivating examples: compressive sensing, matrix completion, etc
- Sparsity-promoting  $L_1$  regularization, nonsmoothness
- Quick introduction to subgradient and subdifferential calculus
- Linear convergence without strong convexity?

# Large Scale Convex Optimization

$$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in X \end{array} \quad (1)$$

where  $f(\cdot)$  is convex differentiable,  $X \subseteq \mathbb{R}^n$  is a *nonempty* convex set.

- If  $X$  is well-represented, and  $f(\mathbf{x})$ ,  $\nabla f(\mathbf{x})$ ,  $\nabla^2 f(\mathbf{x})$  are easily computed, then (1) is efficiently solvable (e.g., via interior point methods in polynomial time).
- Second order methods (e.g., interior point methods) require solving a linear system of size  $n \times n$ , which can be slow when  $n$  is large (e.g.,  $> 10^5$ ).
- For large scale optimization problems, first order methods are preferred, since each iteration is cheaper, requiring only matrix-vector multiplications.
- First order methods typically preserve problem sparsity in large scale optimization problems.

## Example: Routing in a Data Network

- Consider a directed graph with nodes  $\{1, \dots, N\}$  and directed links  $\mathcal{A} \subseteq \{1, \dots, N\} \times \{1, \dots, N\}$  connecting the nodes.
- $M$  pairs of nodes (called origin-destination (OD) pairs) labeled from 1 to  $M$
- For each OD pair  $w \in \{1, \dots, M\}$ , a set of directed loop-free paths from the origin to the destination, with traffic arrival rate  $r_w > 0$ .
- Label the given paths from 1 to  $P$  and let  $\mathcal{P}_w \subseteq \{1, \dots, P\}$  be the set of paths from the origin to the destination of OD pair  $w$ .
- For each link  $(i, j) \in \mathcal{A}$ , we are given  $\bar{D}_{ij}(\mathbf{f}_{ij})$ , the delay on this link when the flow it carries is equal to  $\mathbf{f}_{ij}$ , so that the total delay is

$$\bar{D}(\cdots, \mathbf{f}_{ij}, \cdots)_{(i,j) \in \mathcal{A}} = \sum_{(i,j) \in \mathcal{A}} \bar{D}_{ij}(\mathbf{f}_{ij}).$$

## Routing in a Data Network, continued

- Let

$$X = \{ (\mathbf{x}_1, \dots, \mathbf{x}_P) \mid \mathbf{x}_1 \geq 0, \dots, \mathbf{x}_P \geq 0, \sum_{p \in \mathcal{P}_w} \mathbf{x}_p = r_w, w = 1, \dots, M \}.$$

- Let  $E$  denote the  $|\mathcal{A}| \times P$  link-path incidence matrix associated with the given paths (i.e., an entry of  $E$  is 1 if the link corresponding to its row is on the path corresponding to its column, and is 0 otherwise).
- Then, the optimal routing problem may be formulated as the following nonlinear program:

$$\text{minimize } D(\mathbf{x}) = \bar{D}(E\mathbf{x}) \quad \text{subject to } \mathbf{x} \in X.$$

- Assume  $\bar{D}_{ij}$  is continuous on an interval  $[0, C_{ij})$ , tends to  $\infty$  at  $C_{ij}$  (the “transmission capacity” of link  $(i, j)$ ), and is strictly convex and smooth.

## Routing in a Data Network, continued

- A well-known example is the M/M/1 queue delay function

$$\bar{D}_{ij}(f_{ij}) = f_{ij}/(C_{ij} - f_{ij}).$$

- There may be more than one optimal routing since, although  $\bar{D}$  is strictly convex,  $D$  need not be.
- Given the large network size, many OD pairs and the large number of routes, the routing problem is a large scale convex optimization problem, which may admit multiple optimal path flows.

## Example: MRI, Compressive Sensing

- In a MRI system, we have measurement in the form of

$$\mathbf{b} = \mathbf{A}\mathbf{x} + \mathbf{n}$$

where  $\mathbf{x} \in \mathbb{C}^n$  is the image of interest,  $\mathbf{n} \in \mathbb{C}^m$  is measurement noise, and  $\mathbf{A} \in \mathbb{C}^{m \times n}$  is the measurement matrix (partial Fourier Transform matrix)

- Large under-determined linear system  $m \ll n$ , to reduce data volume.



Figure 1: MRI, Sampling points, Image

## MRI and Compressive Sensing

- $\mathbf{x}$  is typically sparse  $\Rightarrow$  finding the sparsest solution
- $L_0$ -norm minimization

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \lambda \|\mathbf{x}\|_0 + \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

where  $\lambda > 0$  is a penalty parameter

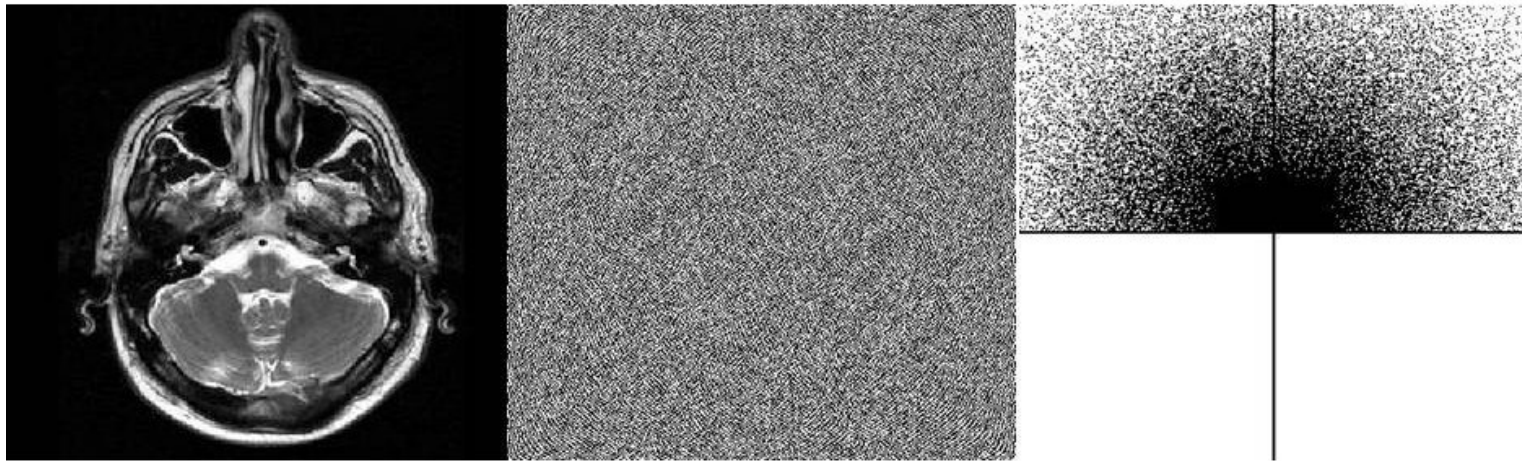
- Nonconvex, difficult to solve
- Compressive sensing:  $L_1$ -norm minimization

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \lambda \|\mathbf{x}\|_1 + \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

which is a large scale (non-smooth) convex optimization problem.

- Under suitable conditions (e.g.,  $\mathbf{A}$  is random),  $L_0$  minimization  $\Leftrightarrow L_1$  minimization

# MRI and Compressive Sensing



- Pick 25% coefficients at random (with bias)
- Reconstruct image from the 25% coefficients



## MRI and Compressive Sensing

Use 1/4 Fourier coefficients

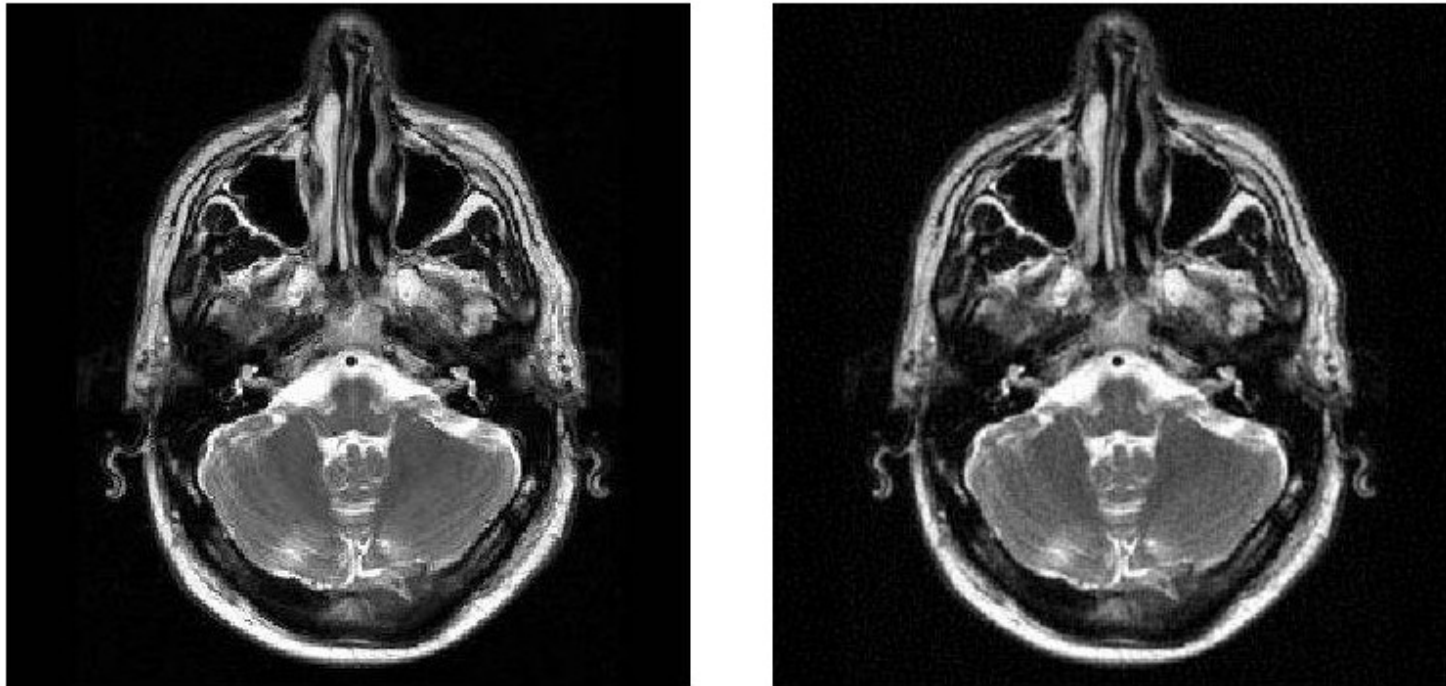


Figure 2: Original vs. Reconstructed Image (courtesy of Y. Zhang, Rice University)

## $L_1$ -minimization: Sparsity-Inducing

The optimality condition for

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \lambda \|\mathbf{x}\|_1 + \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

is  $-\mathbf{A}'(\mathbf{Ax}^* - \mathbf{b}) \in \lambda \partial \|\mathbf{x}^*\|_1$ , or

- $-\lambda \leq (\mathbf{A}'(\mathbf{Ax}^* - \mathbf{b}))_i \leq \lambda$
- if  $x_i^* > 0$ , then  $(\mathbf{A}'(\mathbf{Ax}^* - \mathbf{b}))_i = -\lambda$
- if  $x_i^* < 0$ , then  $(\mathbf{A}'(\mathbf{Ax}^* - \mathbf{b}))_i = \lambda$

Thus,  $L_1$  minimization is sparsity-inducing

$$-\lambda < (\mathbf{A}'(\mathbf{Ax}^* - \mathbf{b}))_i < \lambda \quad \text{implies} \quad x_i^* = 0.$$

$\Rightarrow$  The larger the value of  $\lambda$ , the sparser is the solution  $\mathbf{x}^*$ .

## Subgradient

A vector  $g$  is a **subgradient** of  $f$  (not necessarily convex) at  $x$  if

$$f(y) \geq f(x) + g'(y - x) \quad \text{for all } y$$

- subgradient gives affine global underestimator of  $f$
- if  $f$  is convex, it has at least one subgradient at every point in  $\text{relint dom } f$
- if  $f$  is convex and differentiable,  $\nabla f(x)$  is a subgradient of  $f$  at  $x$

**Example:**  $f = \max\{f_1, f_2\}$ , with  $f_1, f_2$  convex and differentiable

## Subdifferential

The set of all subgradients of  $f$  at  $x$  is called the **subdifferential** of  $f$  at  $x$ , written  $\partial f(x)$

- $\partial f(x)$  is a closed convex set
- $\partial f(x)$  nonempty (if  $f$  convex, and finite near  $x$ )
- $\partial f(x) = \{\nabla f(x)\}$  if  $f$  is differentiable at  $x$
- if  $\partial f(x) = \{g\}$ , then  $f$  is differentiable at  $x$  and  $g = \nabla f(x)$

**Example:**  $f(x) = |x_1| + |x_2|$

**Optimality condition:**  $x^*$  is a minimizer of  $\min f(x)$  iff  $\mathbf{0} \in \partial f(x^*)$ .

## Subdifferential Calculus

Assumption: all functions are finite near  $x$

- $\partial f(x) = \{\nabla f(x)\}$  if  $f$  is differentiable at  $x$
- **scaling:**  $\partial(\alpha f) = \alpha \partial f$  (if  $\alpha > 0$ )
- **addition:**  $\partial(f_1 + f_2) = \partial f_1 + \partial f_2$  (RHS is addition of sets)
- **affine transformation of variables:** if  $g(x) = f(Ax + b)$ , then  $\partial g(x) = A^T \partial f(Ax + b)$
- **pointwise maximum:** if  $f = \max_{i=1, \dots, m} f_i$ , then

$$\partial f(x) = \mathbf{Co} \bigcup \{\partial f_i(x) \mid f_i(x) = f(x)\},$$

i.e., convex hull of union of subdifferentials of 'active' functions at  $x$

## Compressive Sensing: Other Convex Formulations

Variations of formulations

- $\min_{\mathbf{x} \in \mathbb{C}^n} \|\mathbf{x}\|_1$ , subject to:  $\mathbf{Ax} = \mathbf{b}$
- $\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{x}\|_1$ , subject to:  $\mathbf{Ax} = \mathbf{b}; \mathbf{x} \geq 0$
- $\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{Ax} - \mathbf{b}\|_2^2$ , subject to:  $\|\mathbf{x}\|_1 \leq \nu$
- The interior point methods are ineffective. The classical first order methods are preferred, such as the coordinate descent, the alternating direction method of multipliers, LASSO - Least Absolute Shrinkage and Selection Operator
- Note  $\mathbf{A}$  is a fat matrix, thus not full column rank  $\Rightarrow$  existing (rate of) convergence theory does not apply.
- Non-smoothness is involved in the objective function, “causing problem” to the coordinate descent method??

# Coordinate Descent for Compressive Sensing

Consider the nonsmooth formulation:

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

- Coordinate descent algorithm (CD): iteratively and cyclically minimize wrt each variable
- Soft thresholding: let  $x^+ = \arg \min_{y \in \mathbb{R}} \lambda |y| + \frac{1}{2}(y - x)^2$ , then

$$x^+ = \begin{cases} x + \lambda, & x \leq -\lambda \\ 0, & -\lambda \leq x \leq \lambda \\ x - \lambda, & x \geq \lambda \end{cases}$$

- Simple, fast, sparsity-inducing ... how about convergence??

## Coordinate Descent Method

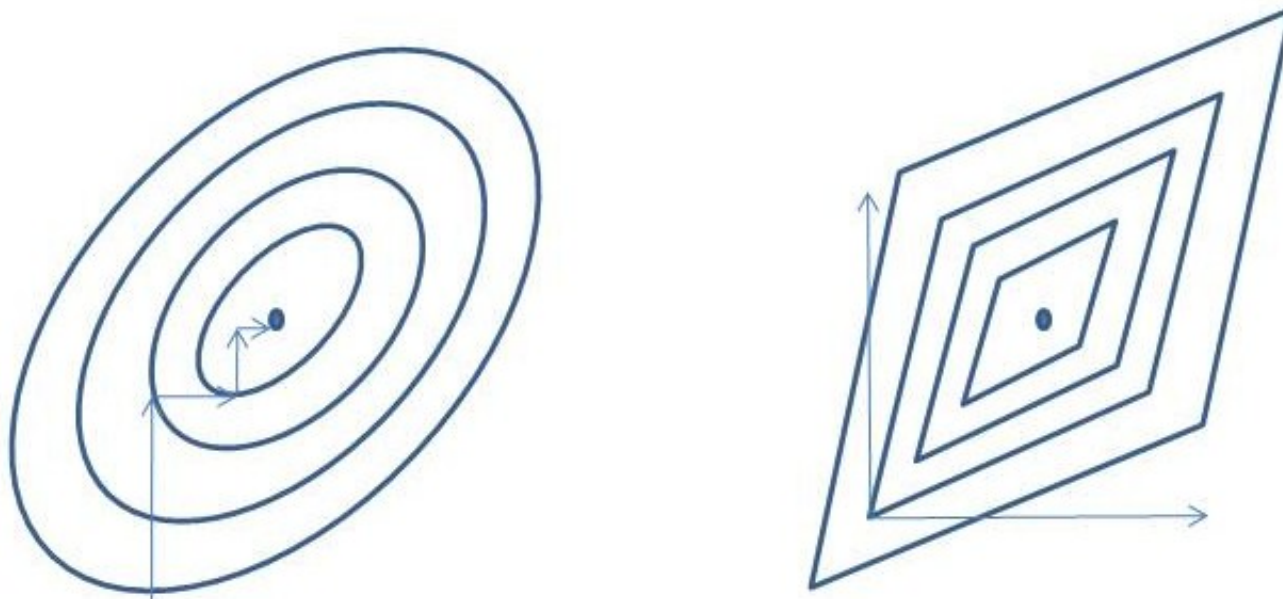


Figure 3: CD method for smooth/non-smooth minimization

Non-smoothness can cause the CD method to get stuck!



## Example: Matrix Completion

- Restore an image  $\mathbf{M} \in \mathbb{R}^{m \times n}$  using a subset  $E$  of pixel values
- $\mathbf{M}$  is a large matrix, but low rank (say rank  $r \ll \min\{m, n\}$ )
- Matrix completion: assume  $\mathbf{M}$  is rank  $r$ , let  $\mathbf{X} \in \mathbb{R}^{m \times r}$ ,  $\mathbf{Y} \in \mathbb{R}^{n \times r}$ , solve

$$\begin{aligned} & \text{minimize}_{\mathbf{X}, \mathbf{Y}, \mathbf{Z}} \quad \|\mathbf{Z} - \mathbf{X}\mathbf{Y}'\|_2^2 \\ & \text{subject to} \quad \mathbf{Z}_{ij} = \mathbf{M}_{ij}, \quad (i, j) \in E. \end{aligned} \tag{2}$$

- Can do coordinate descent by cyclically minimizing  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ .
- Nuclear norm minimization: let  $P_E$  denote projection to the entries in  $E$

$$\min_{\mathbf{Z}} \text{rank}(\mathbf{Z}) + \lambda \|P_E(\mathbf{Z} - \mathbf{M})\|^2 \Rightarrow \min_{\mathbf{Z}} \|\mathbf{Z}\|_* + \lambda \|P_E(\mathbf{Z} - \mathbf{M})\|^2.$$

$\Rightarrow$  relaxation to a large scale convex optimization problem

- Related to the Netflix problem:  $\mathbf{M}$  contains the movie rankings from customers

## Matrix Completion

- Under suitable conditions (e.g., randomly generated low rank matrix  $M$ ,  $|E| \geq O^*(nr)$ , ignoring log factors), the nuclear norm minimization leads to exact recovery of  $M$  w.h.p.. That is, matrix completion via convex optimization
- Nuclear norm minimization  $\Leftrightarrow$  Semidefinite programming, slow for large  $m, n$ .
- Solving the nonconvex formulation (2) using CD is much more efficient; also has exact recovery property?
- Special case:  $E = \{1, 2, \dots, m\} \times \{1, 2, \dots, n\}$ . Then the nonconvex formulation (2) becomes

$$\min_{\mathbf{X} \in \mathbb{R}^{m \times r}, \mathbf{Y} \in \mathbb{R}^{n \times r}} \|\mathbf{M} - \mathbf{X}\mathbf{Y}'\|^2 \quad (3)$$

whose optimal solution is given by the SVD of  $M$ .

- **Surprise(?):** For almost all initial values of  $\mathbf{X}, \mathbf{Y}$ , the CD method for (3) converges to the global optimal solution (i.e., the SVD solution) linearly.

## Box Constrained Convex QP

Consider a convex quadratic minimization problem over  $\mathbb{R}_+^n$

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \mathbf{x}' \mathbf{A} \mathbf{x} + \mathbf{b}' \mathbf{x} \\ \text{subject to} & \mathbf{x} \in X \end{array} \quad (4)$$

If  $X = \mathbb{R}^n$ , we can solve (4) by a simple coordinate descent (Gauss-Seidel type) method: write  $\mathbf{A} = \mathbf{A}_{\text{low}} + \mathbf{D} + \mathbf{A}_{\text{upp}} \succeq 0$  (assume  $\mathbf{D} \succ 0$ ) and iterate

$$\begin{aligned} (\mathbf{A}_{\text{low}} + \mathbf{D}) \mathbf{x}^{r+1} + \mathbf{A}_{\text{upp}} \mathbf{x}^r + \mathbf{b} &= 0 \\ \Rightarrow \mathbf{x}^{r+1} &= -(\mathbf{A}_{\text{low}} + \mathbf{D})^{-1} (\mathbf{A}_{\text{upp}} \mathbf{x}^r + \mathbf{b}). \end{aligned}$$

If  $X = \mathbb{R}_+^n$ , the coordinate descent (CD) algorithm becomes

$$\mathbf{x}^{r+1} = [\mathbf{x}^{r+1} - (\mathbf{A}_{\text{low}} + \mathbf{D}) \mathbf{x}^{r+1} - \mathbf{A}_{\text{upp}} \mathbf{x}^r - \mathbf{b}]_+$$

where  $[u]_+ = \max\{0, u\}$ ; related to SOR, block SOR.

## Convex QP

Consider a convex quadratic minimization problem over  $\mathbb{R}_+^n$

$$\begin{array}{ll} \text{minimize} & \frac{1}{2}\mathbf{x}'\mathbf{A}\mathbf{x} + \mathbf{b}'\mathbf{x} \\ \text{subject to} & \mathbf{x} \in X \end{array}$$

- CD has many applications in image processing, MRI, ... Convergence was known for symmetric  $\mathbf{A} \succ 0$ , with unique optimal solution
- Convergence for the case symmetric  $\mathbf{A} \succeq 0$  was unresolved for many years ..., unbounded optimal solution set

$$X^* = \{\mathbf{x}^* \mid \mathbf{x}^* = \text{proj}_X[\mathbf{x}^* - \alpha \nabla f(\mathbf{x}^*)]\}, \quad \alpha > 0$$

with  $\text{proj}_X[\cdot] = \text{projection to } X$ .

- Studied by Hildreth, Cryer, Mangasarian, Pang, ... in various contexts since 1950's.

## Gradient Projection Method

Given  $\mathbf{x}^0 \in X$ , generate

$$\mathbf{x}^{r+1} = \text{proj}_X[\mathbf{x}^r - \alpha^r \nabla f(\mathbf{x}^r)], \quad r = 0, 1, 2, \dots$$

where  $\alpha^r > 0$  is the stepsize (chosen as constant or by Armijo-like rule).

**Convergence for the non-degenerate case**  $\mathbf{A} \succ 0$ :  $\exists$  a unique solution  $\mathbf{x}^*$  satisfying  $\mathbf{x}^* = \text{proj}_X[\mathbf{x}^* - \alpha^r \nabla f(\mathbf{x}^*)]$ , so

$$\begin{aligned} \|\mathbf{x}^{r+1} - \mathbf{x}^*\| &= \|\text{proj}_X[\mathbf{x}^r - \alpha^r \nabla f(\mathbf{x}^r)] - \text{proj}_X[\mathbf{x}^* - \alpha^r \nabla f(\mathbf{x}^*)]\| \\ &\leq \|(\mathbf{x}^r - \mathbf{x}^*) - \alpha^r (\nabla f(\mathbf{x}^r) - \nabla f(\mathbf{x}^*))\| \\ &= \|(\mathbf{I} - \alpha^r \mathbf{A})(\mathbf{x}^r - \mathbf{x}^*)\| \\ &\leq \left( \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} \right) \|\mathbf{x}^r - \mathbf{x}^*\| = \left( 1 - \frac{2}{\kappa + 1} \right) \|\mathbf{x}^r - \mathbf{x}^*\|. \end{aligned}$$

Convergence rate depends on  $\kappa = \lambda_{\max}/\lambda_{\min}$ , but *independent* of dimension.  
Requires  $O(1)\kappa \ln(1/\epsilon)$  to find an  $\epsilon$ -relative optimal solution.