

On Convergence of Projected Gradient Descent for Minimizing a Large-Scale Quadratic over the Unit Sphere

Trung Vu, Raviv Raich, and Xiao Fu

School of EECS, Oregon State University, Corvallis, OR 97331-5501 USA

{vutru, raich, xiao.fu}@oregonstate.edu

Motivation

- Eigen-decomposition problem
 - repeatedly solving a sequence of problems of form

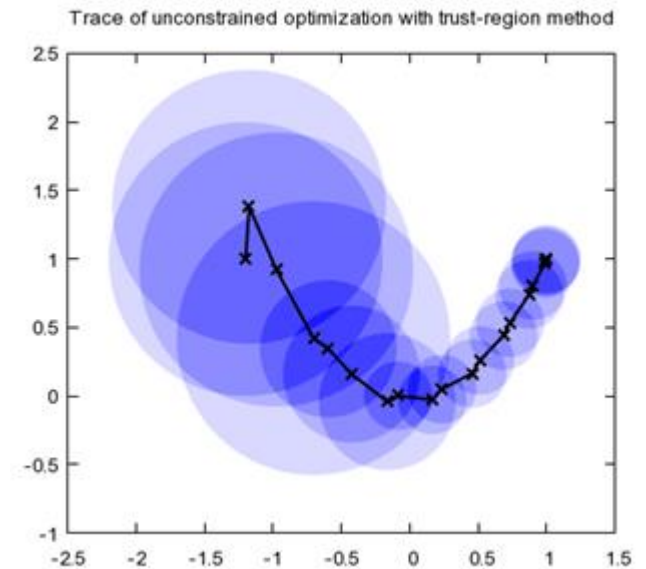
$$\max_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^T \mathbf{A} \mathbf{x} \quad \text{s.t.} \quad \mathbf{x}^T \mathbf{x} = 1$$

- Trust-region subproblem
 - using the quadratic model to approximate the original objective function

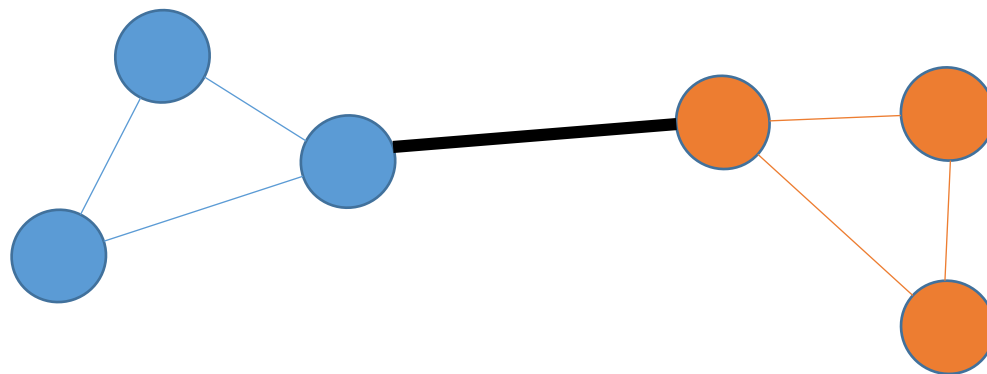
$$\begin{aligned} \min_{\mathbf{x}} & f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{x} + \mathbf{x}^T \nabla^2 f(\mathbf{x}_k) \mathbf{x} \\ \text{s.t.} & \quad \|\mathbf{x}\| \leq \Delta_k \end{aligned}$$

Algorithm Power Method

- 1: Initialize $\mathbf{x}^{(0)}$ such that $[\mathbf{V} \mathbf{x}]_1 \neq 0$
 - 2: for $t = 0, 1, \dots$ do
 - 3: $\mathbf{z}^{(t+1)} = \mathbf{A} \mathbf{x}^{(t)}$
 - 4: $\mathbf{x}^{(t+1)} = \frac{\mathbf{z}^{(t+1)}}{\|\mathbf{z}^{(t+1)}\|}$
-



Graph Partitioning as Constrained Quadratic Optimization



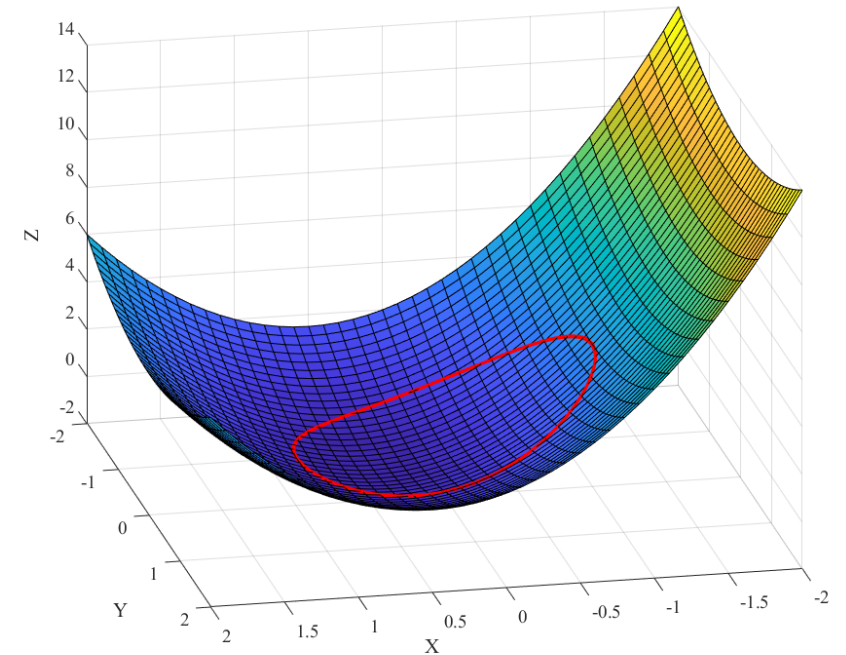
- *Bipartition*: cut a weighted, undirected graph into two subgraphs
 - roughly equals in size
 - the total weight of the cut edges is smallest
- Express the weight of a cut as a quadratic function of binary variables

$$\begin{array}{ccc} \boxed{\begin{array}{l} \min_{\mathbf{x} \in \{-1, +1\}^n} \mathbf{x}^T \mathbf{L} \mathbf{x} \\ \text{s.t. } \mathbf{1}^T \mathbf{x} = 0 \end{array}} & \xrightarrow{\text{Relaxed}} & \boxed{\begin{array}{l} \min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^T \mathbf{L} \mathbf{x} \\ \text{s.t. } \mathbf{1}^T \mathbf{x} = 0 \text{ and } \mathbf{x}^T \mathbf{x} = 1 \end{array}} & \xrightarrow{\text{Regularized}} & \boxed{\begin{array}{l} \min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^T \mathbf{L} \mathbf{x} + \lambda \mathbf{1}^T \mathbf{x} \\ \text{s.t. } \|\mathbf{x}\|_2 = 1 \end{array}} \end{array}$$

Problem Formulation

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x} \quad \text{subject to } \|\mathbf{x}\|^2 = 1$$

$$\min_{x_1, x_2} \frac{1}{2} (4x_1^2 + x_2^2) - 2x_1 \quad \text{s.t. } x_1^2 + x_2^2 = 1$$



- $\|\cdot\|$ is the Euclidean norm
- $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric but not necessarily positive semidefinite
- Non-convex objective function with a non-convex constraint

Solution Properties

$$\mathcal{L}(x, \gamma) = \frac{1}{2}x^T Ax - b^T x - \frac{1}{2}\gamma(\|x\|^2 - 1)$$

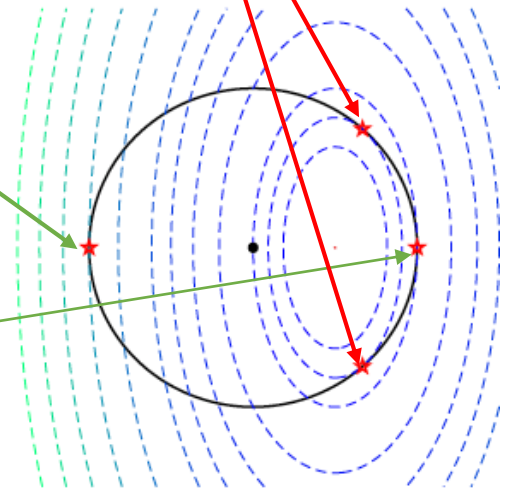
- Stationary points $\begin{cases} x_* \in \mathcal{S}^{n-1} \\ Ax_* - b = \gamma(x_*) \cdot x_* \end{cases}$
- Local minimum $\begin{cases} P_{x_*}^\perp = I - x_*x_*^T \\ \gamma(x_*) \leq \lambda_{n-1}(P_{x_*}^\perp AP_{x_*}^\perp) \end{cases}$
- Global minimum $\gamma(x_*) \leq \lambda_n(A) \leq \lambda_{n-1}(P_{x_*}^\perp AP_{x_*}^\perp)$

Global maximum
 $\gamma = 6, \lambda_{n-1} = 1$

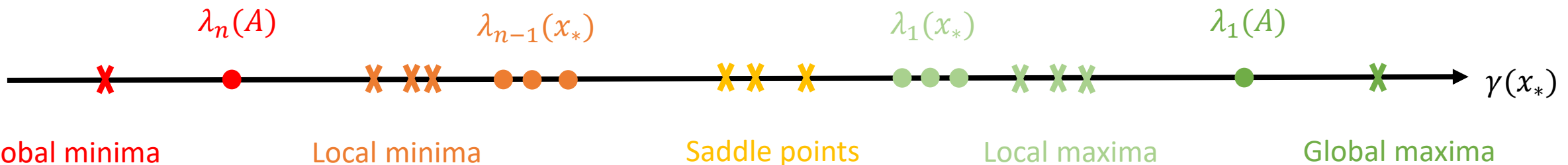
Local/Global minima
 $\gamma = 1, \lambda_{n-1} = 8/3$

$$\min_{x_1, x_2} \frac{1}{2}(4x_1^2 + x_2^2) - 2x_1 \quad \text{s.t. } x_1^2 + x_2^2 = 1$$

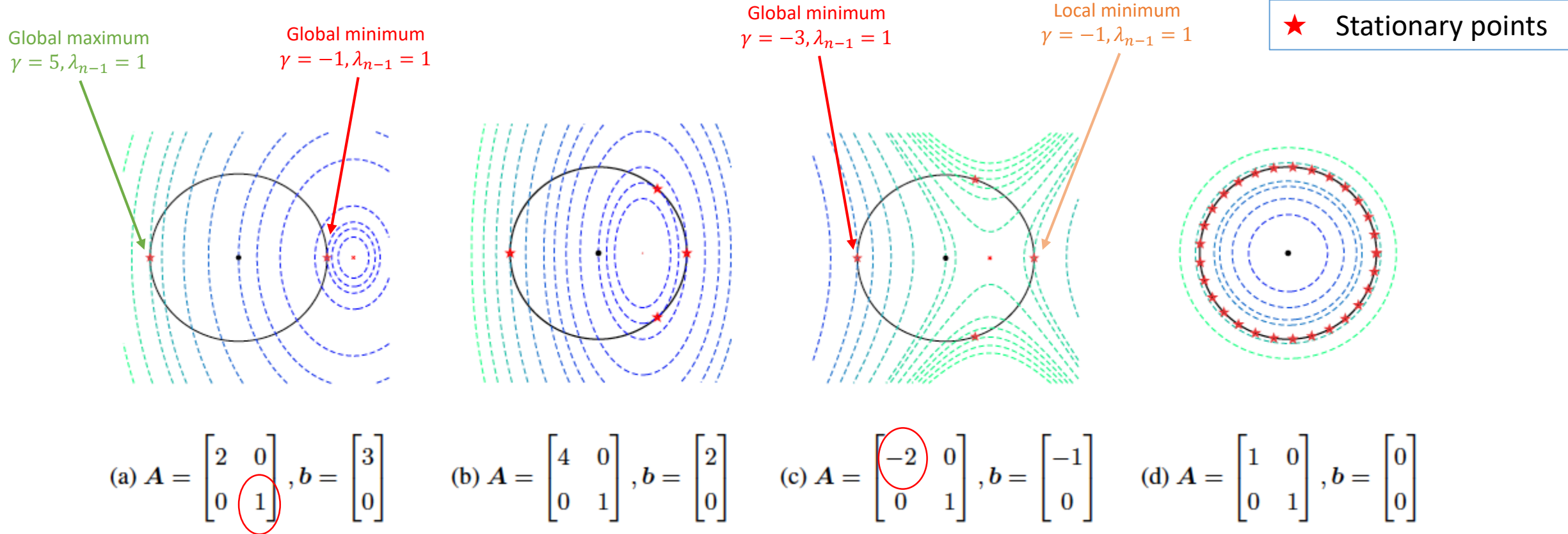
Local maximum
 $\gamma = 2, \lambda_{n-1} = 1$



$$A = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}, b = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$



2D-Examples



Related Work

- Quadratic Constrained Quadratic Program with only one constraint (QCQP-1)
 - semidefinite relaxation (SDR)
 - Lagrangian relaxation
 - $O(n^2)$!
- Trust-region subproblem
 - solvable when n is small (by matrix factorizations)
 - for large n , iterative methods are considered
 - Parameterized eigenvalue problem [Sorensen'97]
 - Sequential Subspace Method (SSM) [Hager'01]

Projected Gradient Descent

Algorithm Projected Gradient Descent (PGD)

- 1: Initialize $\mathbf{x}^{(0)} \in \mathcal{S}^{n-1}$
 - 2: for $t = 0, 1, \dots$ do
 - 3: $\mathbf{z}^{(t+1)} = \mathbf{x}^{(t)} - \alpha(\mathbf{A}\mathbf{x}^{(t)} - \mathbf{b})$
 - 4: $\mathbf{x}^{(t+1)} = \mathcal{P}_{\mathcal{S}^{n-1}}(\mathbf{z}^{(t+1)})$
-

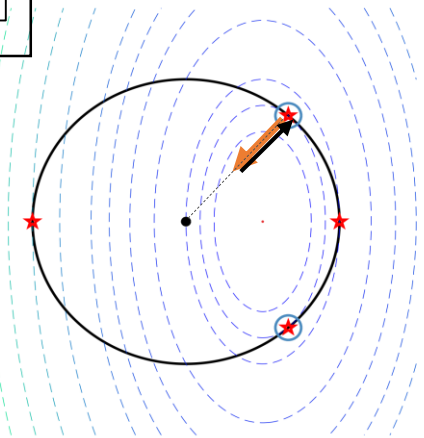
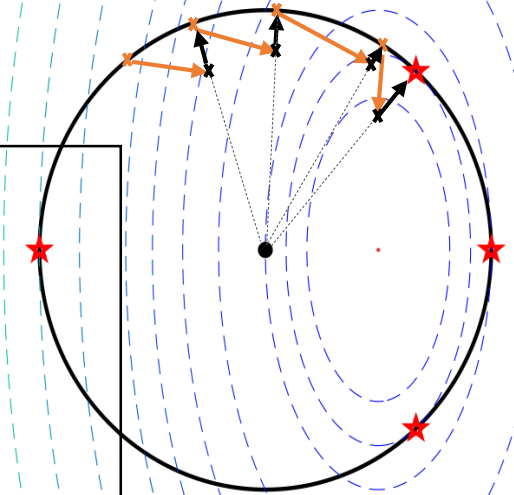
$$\begin{aligned} \max_{\mathbf{x} \in \mathbb{R}^n} \quad & \mathbf{x}^T \mathbf{A} \mathbf{x} \\ \text{s.t.} \quad & \mathbf{x}^T \mathbf{x} = 1 \end{aligned}$$

Power method $\alpha \rightarrow \infty$

$$\mathbf{z}^{(t+1)} = \mathbf{x}^{(t)} - \infty(\mathbf{A}\mathbf{x}^{(t)} - \mathbf{b}) \sim \mathbf{A}\mathbf{x}^{(t)}$$

- Fixed point

$$\mathcal{P}_{\mathcal{S}^{n-1}}(\bar{\mathbf{x}} - \alpha(\mathbf{A}\bar{\mathbf{x}} - \mathbf{b})) = \bar{\mathbf{x}} \iff \begin{cases} \bar{\mathbf{x}} \in \mathcal{S}^{n-1} \\ \gamma(\bar{\mathbf{x}}) < \frac{1}{\alpha} \\ \mathbf{A}\bar{\mathbf{x}} - \mathbf{b} = \gamma(\bar{\mathbf{x}}) \cdot \bar{\mathbf{x}} \end{cases}$$

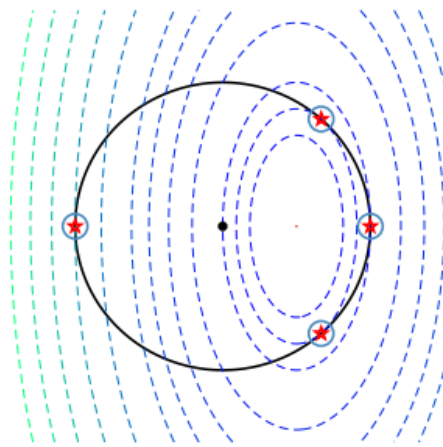


Stationary Point versus Fixed Point

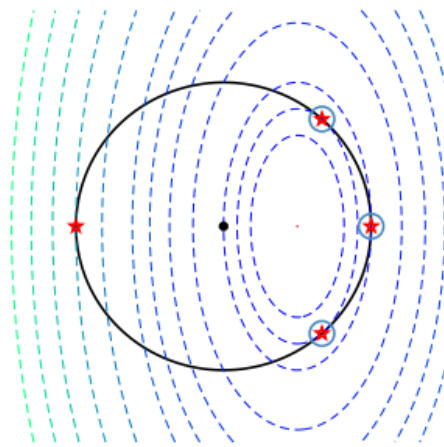
$$\min_{x_1, x_2} \frac{1}{2}(4x_1^2 + x_2^2) - 2x_1 \quad \text{s.t. } x_1^2 + x_2^2 = 1$$

$$\gamma(\bar{\mathbf{x}}) < \frac{1}{\alpha}$$

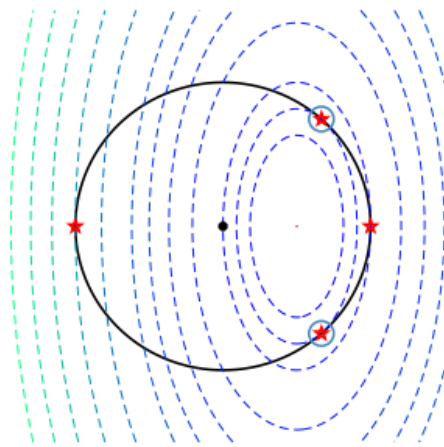
★	Stationary point
○	Fixed point



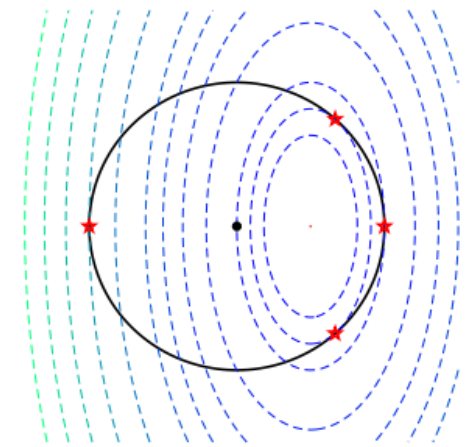
(a) $0 < \alpha < 1/6$



(b) $1/6 < \alpha < 1/2$



(c) $1/2 < \alpha < 1$



(d) $\alpha > 1$



Convergence Analysis

- Taylor Series Expansion of the Projection

$$\mathcal{P}_{\mathcal{S}^{n-1}}(\mathbf{x} + \boldsymbol{\delta}) = \frac{\mathbf{x} + \boldsymbol{\delta}}{\|\mathbf{x} + \boldsymbol{\delta}\|} = \mathcal{P}_{\mathcal{S}^{n-1}}(\mathbf{x}) + \frac{1}{\|\mathbf{x}\|} \left(\mathbf{I} - \frac{\mathbf{x}\mathbf{x}^T}{\|\mathbf{x}\|^2} \right) \boldsymbol{\delta} + O(\|\boldsymbol{\delta}\|^2)$$

- Recursion on the error vector

$$\begin{aligned} \boldsymbol{\delta}^{(t+1)} &= \mathbf{x}^{(t+1)} - \mathbf{x}_* = \mathcal{P}_{\mathcal{S}^{n-1}}(\mathbf{x}^{(t)} - \alpha(\mathbf{A}\mathbf{x}^{(t)} - \mathbf{b})) - \mathbf{x}_* \\ &= \mathcal{P}_{\mathcal{S}^{n-1}}(\mathbf{x}_* - \alpha(\mathbf{A}\mathbf{x}_* - \mathbf{b}) + (\mathbf{I} - \alpha\mathbf{A})\boldsymbol{\delta}^{(t)}) - \mathbf{x}_* \\ &= \frac{1}{1 - \alpha\gamma(\mathbf{x}_*)} (\mathbf{I} - \mathbf{x}_*\mathbf{x}_*^T)(\mathbf{I} - \alpha\mathbf{A})\boldsymbol{\delta}^{(t)} + O(\|\boldsymbol{\delta}^{(t)}\|^2) \end{aligned}$$

ρ_α

If $\rho_\alpha < 1$ and $\boldsymbol{\delta}^{(0)}$ is sufficiently small, the error series behaves similar to a geometric series

Rate of Convergence

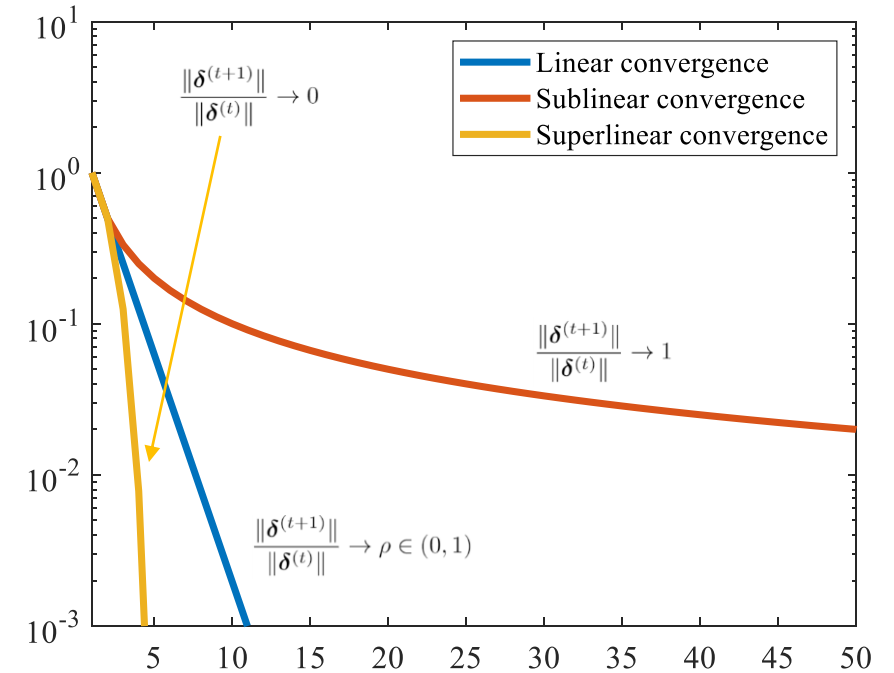
- **Our results:**

- PGD converges linearly locally to any strict local minimum with appropriate choice of α
- The asymptotic rate of convergence is given by

$$\rho_\alpha(\mathbf{x}_*) = \max_{1 \leq i \leq n-1} \frac{|1 - \alpha \lambda_i(\mathbf{P}_{\mathbf{x}_*}^\perp \mathbf{A} \mathbf{P}_{\mathbf{x}_*}^\perp)|}{1 - \alpha \gamma(\mathbf{x}_*)} < 1$$

- Optimizing over the step size yields faster convergence

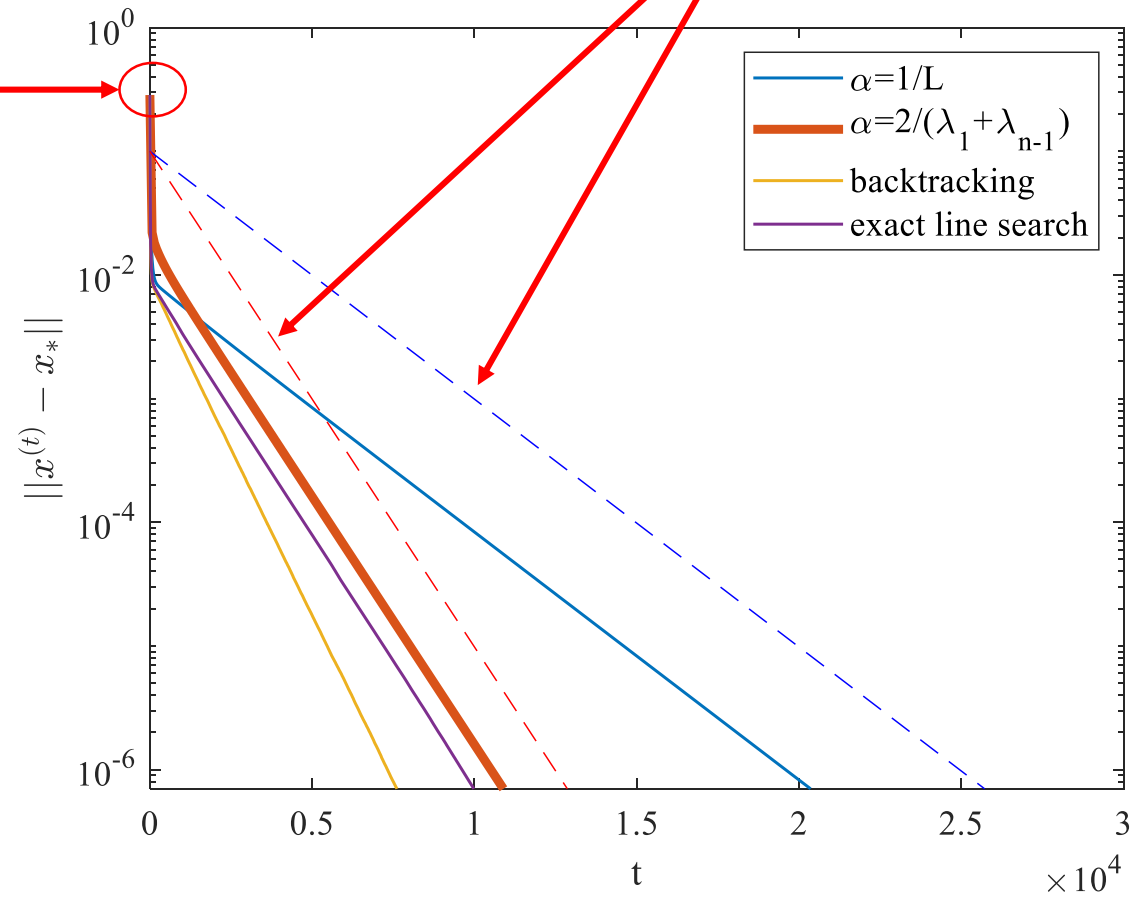
$$\alpha_* = \operatorname{argmin}_{\alpha > 0} \max_{1 \leq i \leq n-1} \frac{|1 - \alpha \lambda_i(\mathbf{P}_{\mathbf{x}_*}^\perp \mathbf{A} \mathbf{P}_{\mathbf{x}_*}^\perp)|}{1 - \alpha \gamma(\mathbf{x}_*)} = \frac{2}{\lambda_1(\mathbf{P}_{\mathbf{x}_*}^\perp \mathbf{A} \mathbf{P}_{\mathbf{x}_*}^\perp) + \lambda_{n-1}(\mathbf{P}_{\mathbf{x}_*}^\perp \mathbf{A} \mathbf{P}_{\mathbf{x}_*}^\perp)}$$



Numerical Evaluation

Theoretical analysis of $\rho_\alpha(x_*)^t$

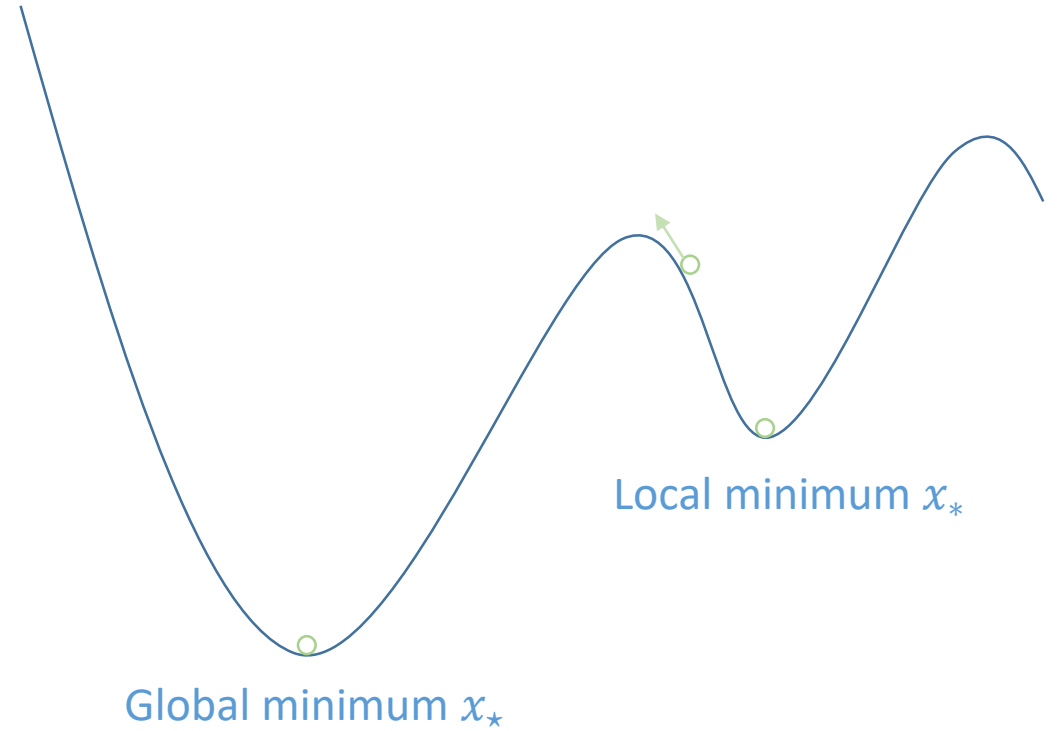
Initialize $x^{(0)}$ near the local minimum x_*



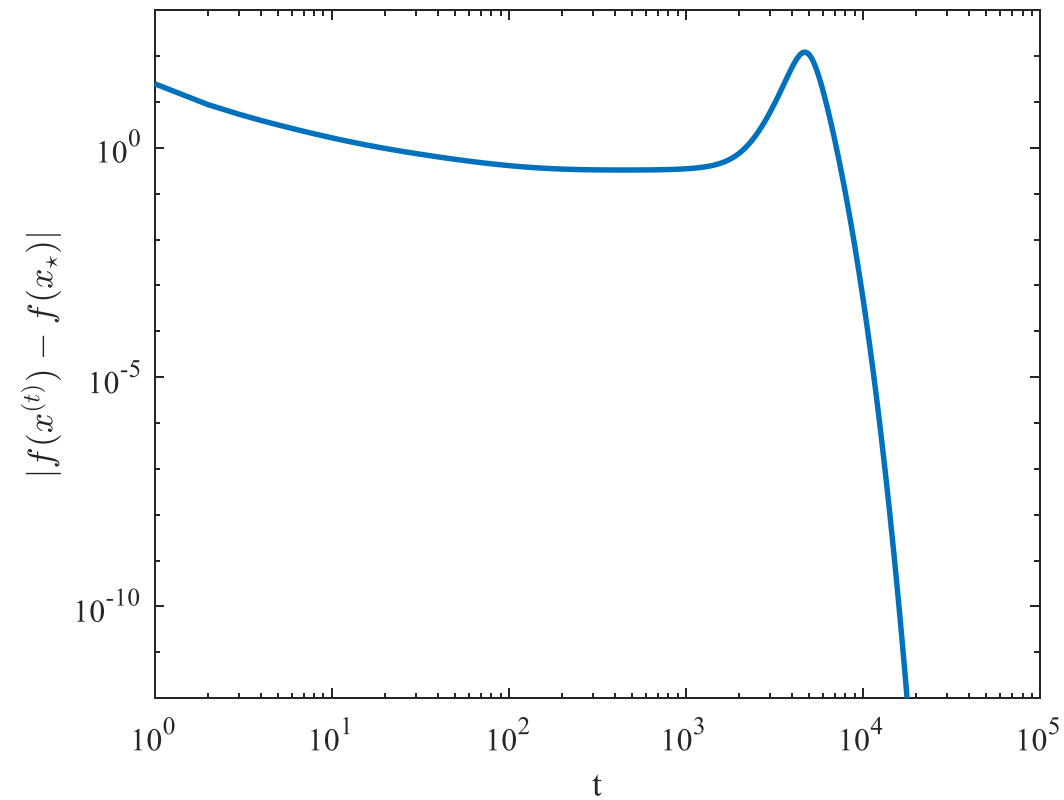
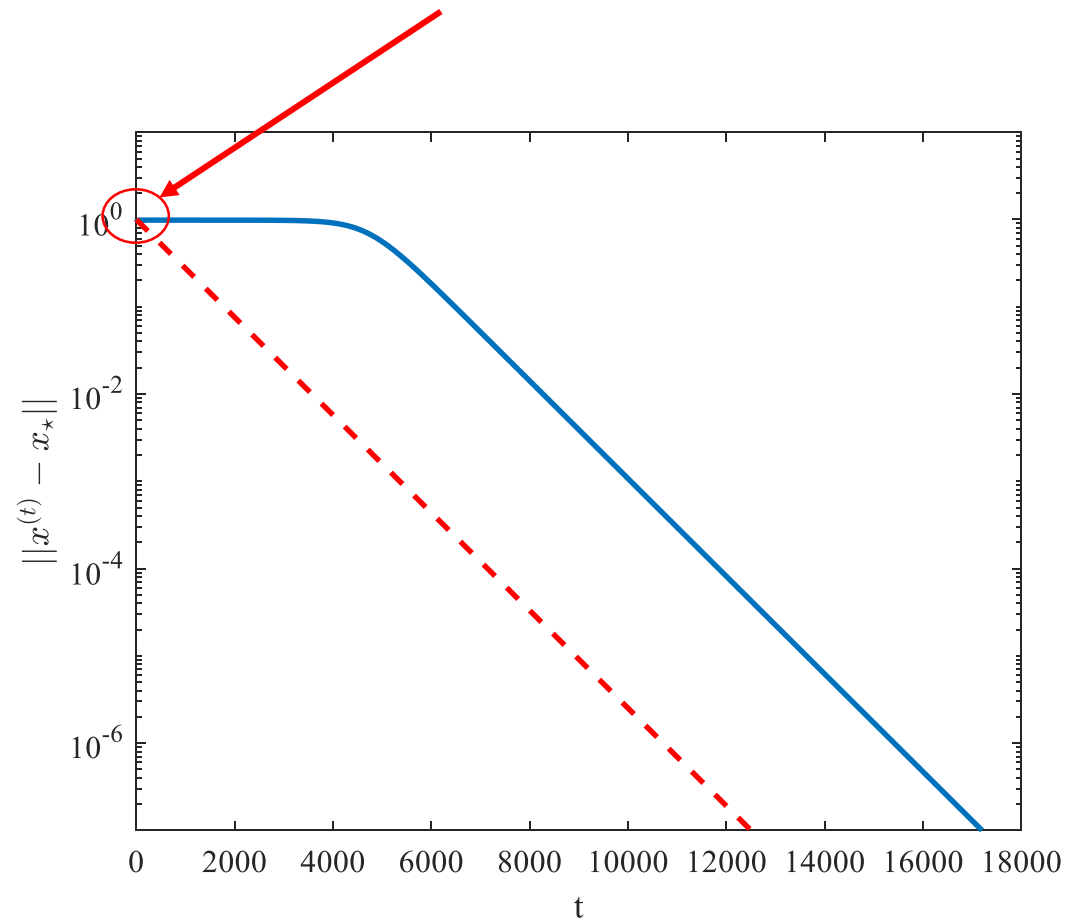
Escape from Local Minima

- For $\alpha > \frac{2}{\lambda_1(x_*) + \gamma(x_*)}$, the error series *w.r.t* x_* tends to diverge since $\rho_\alpha > 1$
- Define $g(\alpha, x_*) = \alpha(\lambda_1(x_*) + \gamma(x_*))$
- *Conjecture:*

Assume there exists sufficiently large α satisfying $g(\alpha, x_*) < 2$ for any global minimum x_* and $g(\alpha, x_*) \geq 2$ for any strict local minimum x_* . Then PGD with step size α converges to one of the optimal solutions x_* at an asymptotic geometric rate of $\rho_\alpha(x_*)$.



Initialize $x^{(0)}$ near the local minimum x_*



Conclusion and Future Works

- Conclusion
 - showed PGD converges linearly to a strict local minimum in its neighborhood
 - provided the closed-form expression for asymptotic convergence rate
 - identified ways of achieving optimal rate of convergence near the optimum
- Future works
 - minimizing a quadratic over an ellipsoid
 - acceleration of gradient projection using momentum
 - analysis of convergence to a continuum of optima

THANK YOU!

References

1. Danny C Sorensen, “Minimization of a large-scale quadratic function subject to a spherical constraint,” *SIAM Journal on Optimization*, vol. 7, no. 1, pp. 141–161, 1997.
2. William W Hager, “Minimizing a quadratic over a sphere,” *SIAM Journal on Optimization*, vol. 12, no.1, pp. 188–208, 2001.