

# Sharp convergence guarantees for iterative algorithms in random (nonconvex) optimization

or

Can you predict *if* and *how fast*  
your model-fitting algorithm will converge?

Ashwin Pananjady

Schools of Industrial and Systems Engineering and  
Electrical and Computer Engineering, Georgia Tech



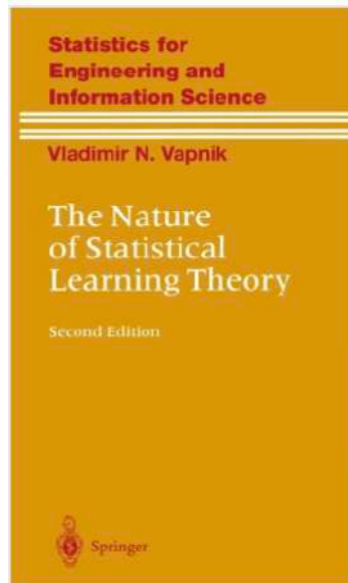
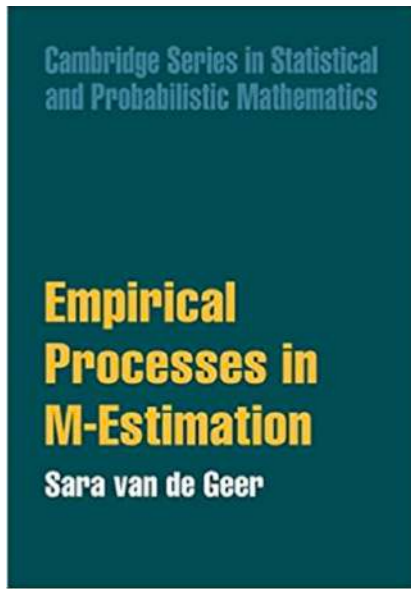
Kabir  
Chandrasekher  
(Stanford)



Christos  
Thrampoulidis  
(UBC)



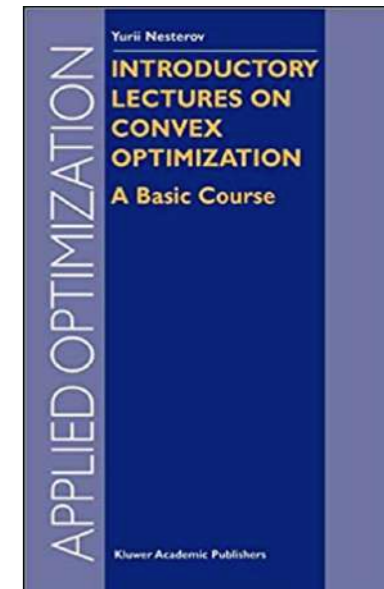
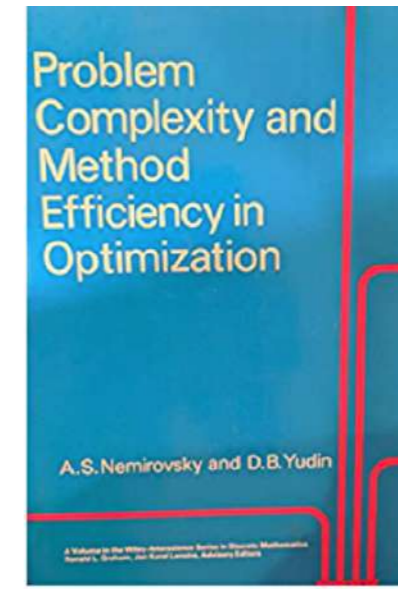
Mengqi  
Lou  
(Georgia Tech)



### M-estimation in supervised learning

$$\mathbb{R}^d \quad \mathbb{R}$$
$$(\mathbf{x}_i, y_i)_{i=1}^n$$

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell_i(f_{\theta}(\mathbf{x}_i), y_i)$$



- ▶ Optimal solution  $\hat{\theta}$  has nice properties under noise model, e.g.

$$\mathbb{E}[y_i | \mathbf{x}_i] = f_{\theta^*}(\mathbf{x}_i)$$

- ▶ General purpose: Rate of estimation measured in terms of sample size and model geometry around  $\theta^*$ .
- ▶ Optimality guarantees: Minimax lower bounds and information theory

- ▶ Iterative algorithms to converge to  $\hat{\theta}$

- ▶ *Efficiency* of method: #iterations to get  $\varepsilon$ -close (typically upper bounds under convexity/smoothness notions)
- ▶ Optimality guarantees: Oracle complexity lower bounds over worst-case family of loss functions

.....

Largely successful in convex optimization approaches to M-estimation:  
Decoupling statistics from optimization

$$(\mathbf{x}_i, y_i)_{i=1}^n \quad \hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell_i(f_{\boldsymbol{\theta}}(\mathbf{x}_i), y_i)$$

Q1

Do worst case efficiency estimates reflect performance in model-fitting problems with random data?

Q2

Can we provide exact efficiency estimates in random ensembles of optimization problems?

Q3

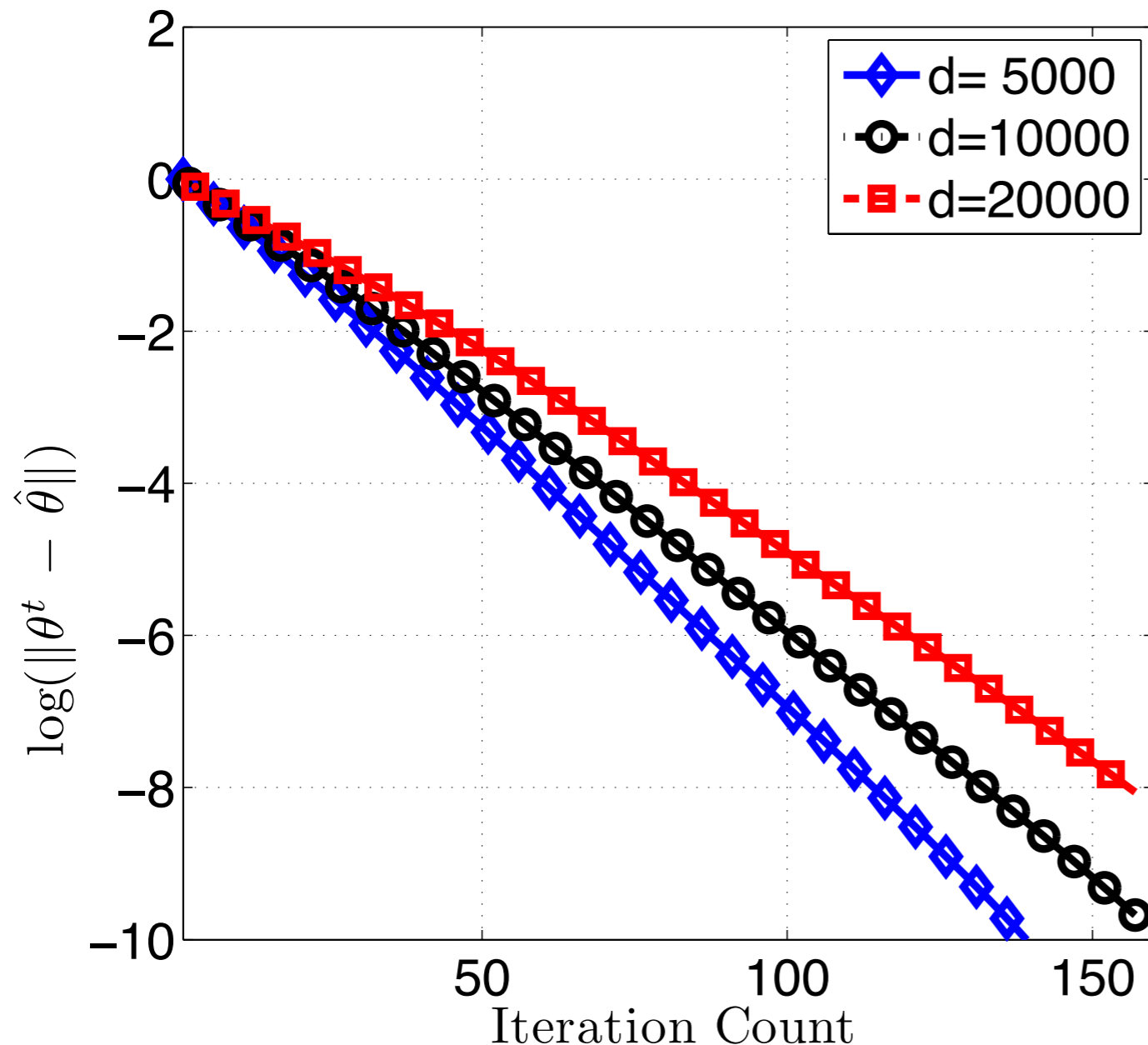
Is there hope of developing a parallel understanding for nonconvex problems?

# Illustration 1: Projected gradient descent on L1 constrained quadratic programming with standard Gaussian covariates

$n = 2500$  Sparsity  $k = \sqrt{d}$

$$y_i = \langle \mathbf{x}_i, \boldsymbol{\theta}^* \rangle + \epsilon_i$$

$$\|\boldsymbol{\theta}^*\|_0 = k, \quad \|\boldsymbol{\theta}^*\|_2 = 1$$



► Worst-case efficiency guarantees pessimistic

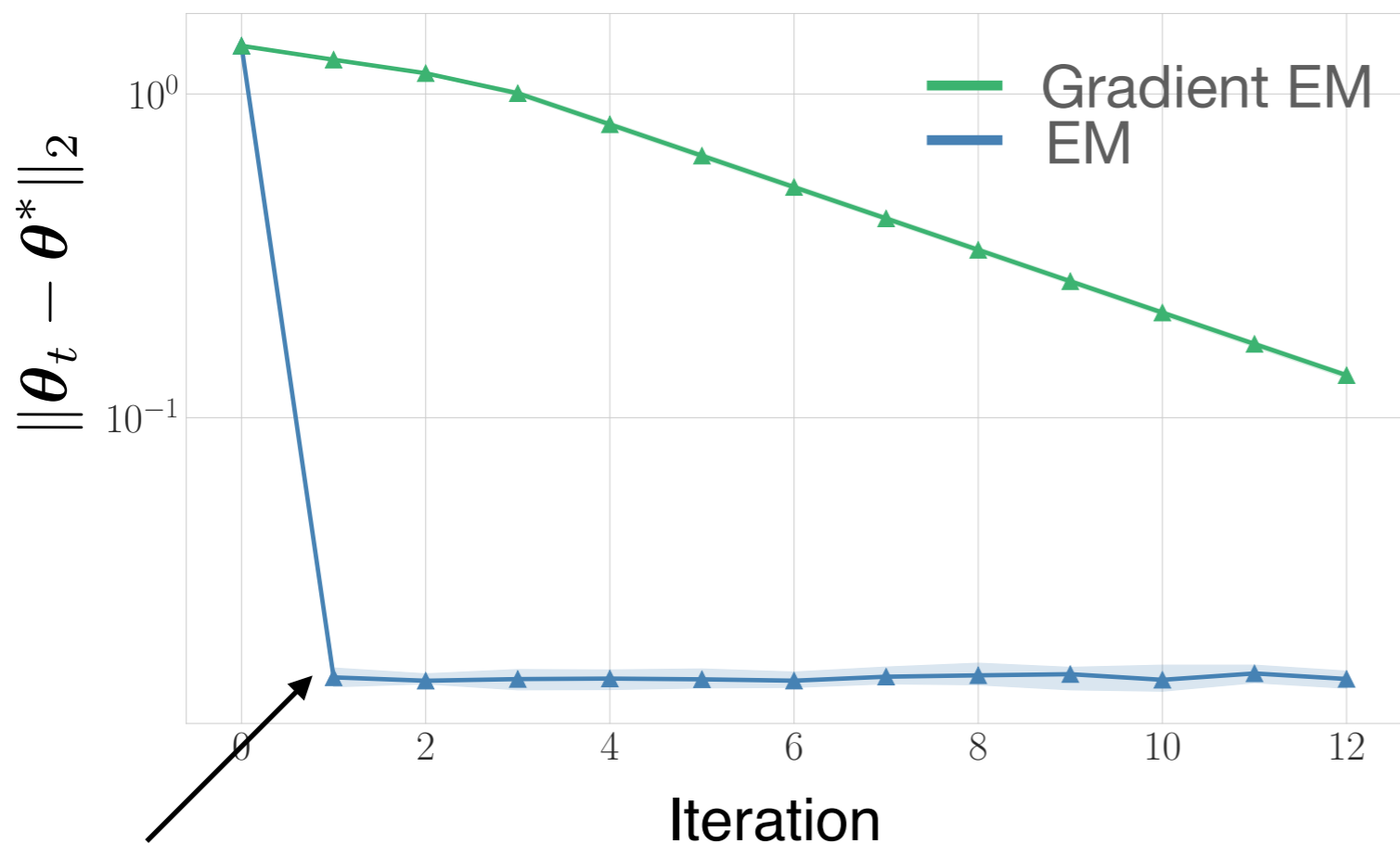
► Different convergence rates depending on problem size: “Larger problems are harder” (but theory does not capture this phenomenon)

# Illustration 2: Expectation maximization algorithms on symmetric Gaussian mixture

$$\mathbf{x}_i \sim \frac{1}{2} \mathcal{N}(\boldsymbol{\theta}^*, \mathbf{I}) + \frac{1}{2} \mathcal{N}(-\boldsymbol{\theta}^*, \mathbf{I})$$

$$\|\boldsymbol{\theta}^*\|_2 = 1$$

$d = 250, \quad n = 10,000$



Noise floor  $\asymp \sigma \sqrt{d/n}$

- ▶ Non-convex optimization problem but parameter estimation possible
- ▶ Gradient EM and EM exhibit different convergence behavior but common analysis tool “cannot capture distinction”



## Desiderata

- Rigorous comparisons between iterative model-fitting methods
- Not just by comparing upper bounds on efficiency!
- Want to answer the following questions in nonconvex problems:
  - Does the algorithm converge (to a statistically useful neighborhood) from a given initialization?
  - Does the algorithm converge *globally*, from a random initialization?
  - How fast does the algorithm converge?

## This talk

- General-purpose, iterate-by-iterate predictions of solution quality if:
  - Each iteration of algorithm is solution to convex program\*
  - Data in optimization problem is suitably Gaussian (i.e. Gaussian conditioned on past iterations)
- Distinguishes convergent behavior from otherwise
- Upper *and* lower bounds on convergence rates and exact error floor

# AMP and first-order algorithms

## Message-passing algorithms for compressed sensing

David L. Donoho<sup>a,1</sup>, Arian Maleki<sup>b</sup>, and Andrea Montanari<sup>a,b,1</sup>

## The dynamics of message passing on dense graphs, with applications to compressed sensing

Mohsen Bayati  
Department of Electrical Engineering  
Stanford University

Andrea Montanari  
Departments of Electrical Engineering and Statistics  
Stanford University

## An Iterative Construction of Solutions of the TAP Equations for the Sherrington–Kirkpatrick Model

Erwin Bolthausen\*

Institute of Mathematics, Universität Zürich, Zürich, Switzerland. E-mail: eb@math.uzh.ch

## The estimation error of general first order methods

Michael Celentano\*      Andrea Montanari\*<sup>†</sup>      Yuchen Wu\*

## General-purpose analysis tool: Population-based

### STATISTICAL GUARANTEES FOR THE EM ALGORITHM: FROM POPULATION TO SAMPLE-BASED ANALYSIS<sup>1</sup>

BY SIVARAMAN BALAKRISHNAN\*<sup>†</sup>,  
MARTIN J. WAINWRIGHT<sup>†</sup> AND BIN YU<sup>†</sup>

University of California, Berkeley\* and Carnegie Mellon University<sup>†</sup>

# Other sharp predictions

## Sharp Time–Data Tradeoffs for Linear Inverse Problems

Samet Oymak\*<sup>†</sup>    Benjamin Recht\*<sup>†</sup>    Mahdi Soltanolkotabi<sup>§</sup>

## Halting Time is Predictable for Large Models: A Universality Property and Average-case Analysis

Courtney Paquette\*<sup>†</sup>    Bart van Merriënboer\*    Elliot Paquette<sup>†</sup>    Fabian Pedregosa\*

## SGD in the Large: Average-case Analysis, Asymptotics, and Stepsize Criticality

Courtney Paquette\*<sup>†</sup>    Kiwon Lee<sup>†</sup>    Fabian Pedregosa\*    Elliot Paquette<sup>†</sup>

▶ Applies under weaker randomness conditions

▶ Specific settings, typically asymptotic

## Global Convergence of the EM Algorithm for Mixtures of Two Component Linear Regression

Jeongyeol Kwon\*  
The University of Texas at Austin  
Wei Qian\*  
Cornell University  
Constantine Caramanis  
The University of Texas at Austin  
Yudong Chen  
Cornell University  
Damek Davis  
Cornell University

KWONCHUNGLI@UTEXAS.EDU  
WQ34@CORNELL.EDU  
CONSTANTINE@UTEXAS.EDU  
YUDONG.CHEN@CORNELL.EDU  
DSD95@CORNELL.EDU

▶ Complementary views: Landscape analysis, properties of loss function verified w.h.p.



# Running example: Phase retrieval with a real signal

Model

$$y_i = |\langle \mathbf{x}_i, \boldsymbol{\theta}^* \rangle| + \epsilon_i, \quad i = 1, 2, \dots$$

$$\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d) \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \|\boldsymbol{\theta}^*\|_2 = 1$$

MLE: Minimizer of nonconvex loss

$$\mathfrak{R}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - |\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle|)^2$$

Algorithms

“Subgradient” method

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \eta \nabla \mathfrak{R}_n(\boldsymbol{\theta}_t) \\ &= \boldsymbol{\theta}_t - \frac{2\eta}{n} \sum_{i=1}^n (|\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle| - y_i) \cdot \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i \end{aligned}$$

Alternating projections

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{n} \cdot \sum_{i=1}^n (y_i - \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle)^2 \\ &= \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i y_i \right) \end{aligned}$$

Sample-splitting: Fresh data  $(\mathbf{X}, \mathbf{y}) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$  in each iteration

$$\Lambda = n/d > 1$$

# Running example: Phase retrieval with a real signal

Model

Algorithms

$$y_i = |\langle \mathbf{x}_i, \boldsymbol{\theta}^* \rangle| + \epsilon_i, \quad i = 1, 2, \dots$$

$$\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d) \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \|\boldsymbol{\theta}^*\|_2 = 1$$

MLE: Minimizer of nonconvex loss

$$\mathfrak{R}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - |\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle|)^2$$

“Subgradient” method

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \eta \nabla \mathfrak{R}_n(\boldsymbol{\theta}_t) \\ &= \boldsymbol{\theta}_t - \frac{2\eta}{n} \sum_{i=1}^n (|\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle| - y_i) \cdot \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i \end{aligned}$$

Alternating projections

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{n} \cdot \sum_{i=1}^n (y_i - \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle)^2 \\ &= \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i y_i \right) \end{aligned}$$

Sample-splitting: Fresh data  $(\mathbf{X}, \mathbf{y}) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$  in each iteration

$$\Lambda = n/d > 1$$

# Running example: Phase retrieval with a real signal

Model

$$y_i = |\langle \mathbf{x}_i, \boldsymbol{\theta}^* \rangle| + \epsilon_i, \quad i = 1, 2, \dots$$

$$\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d) \quad \epsilon_i = \mathcal{N}(0, \sigma^2) \quad \|\boldsymbol{\theta}^*\|_2 = 1$$

MLE: Minimizer of nonconvex loss

$$\mathfrak{R}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - |\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle|)^2$$

Algorithms

“Subgradient” method

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \eta \nabla \mathfrak{R}_n(\boldsymbol{\theta}_t) \\ &= \boldsymbol{\theta}_t - \frac{2\eta}{n} \sum_{i=1}^n (|\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle| - y_i) \cdot \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i \end{aligned}$$

Alternating projections

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{n} \cdot \sum_{i=1}^n (y_i - \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle)^2 \\ &= \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i y_i \right) \end{aligned}$$

Sample-splitting: Fresh data  $(\mathbf{X}, \mathbf{y}) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$  in each iteration

$$\Lambda = n/d > 1$$

# Background on population-based analysis

- ▶ View each iteration of the algorithm as a random operator

$$\mathcal{T}_n : \boldsymbol{\theta}_t \mapsto \boldsymbol{\theta}_{t+1}$$

- ▶ Use infinite sample limit  $\mathcal{T}_\infty$  to guide analysis of the empirical iterates

.....

## Alternating projections

$$\begin{aligned} \mathcal{T}_n(\boldsymbol{\theta}_t) &= \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i y_i \right) \end{aligned}$$

$$\mathcal{T}_\infty(\boldsymbol{\theta}_t) = \mathbb{E}[\text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i y_i]$$

## Subgradient method

$$\begin{aligned} \mathcal{T}_n(\boldsymbol{\theta}_t) &= \boldsymbol{\theta}_t - \frac{2\eta}{n} \sum_{i=1}^n (|\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle| - y_i) \cdot \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i \end{aligned}$$

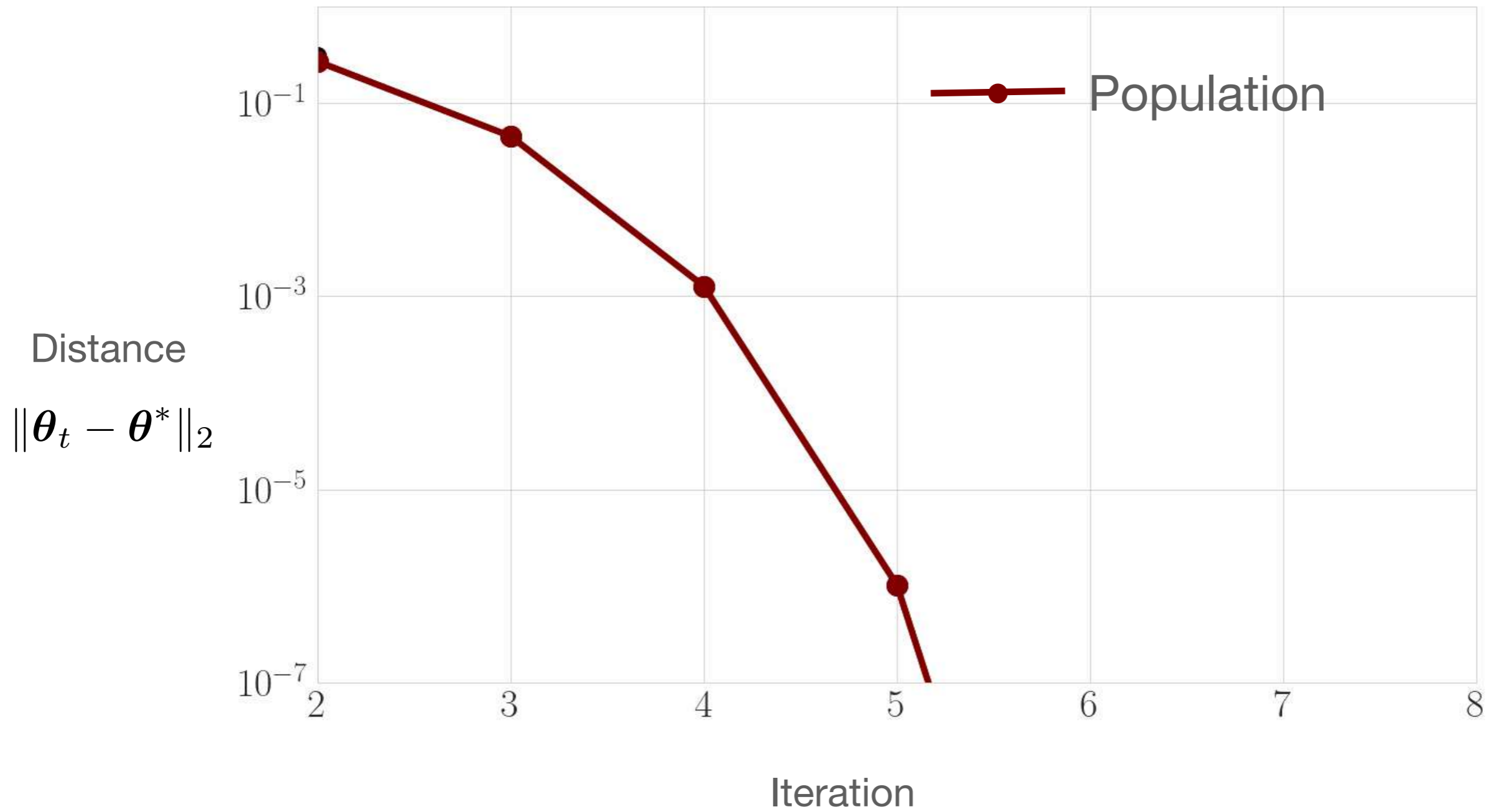
$$\mathcal{T}_\infty(\boldsymbol{\theta}_t) = (1 - 2\eta) \cdot \boldsymbol{\theta}_t + 2\eta \cdot \mathbb{E}[\text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \mathbf{x}_i y_i]$$

**Note**

For  $\eta = 1/2$ , both population updates coincide!

$d = 800, \quad n = 16,000, \quad \sigma = 10^{-6}$

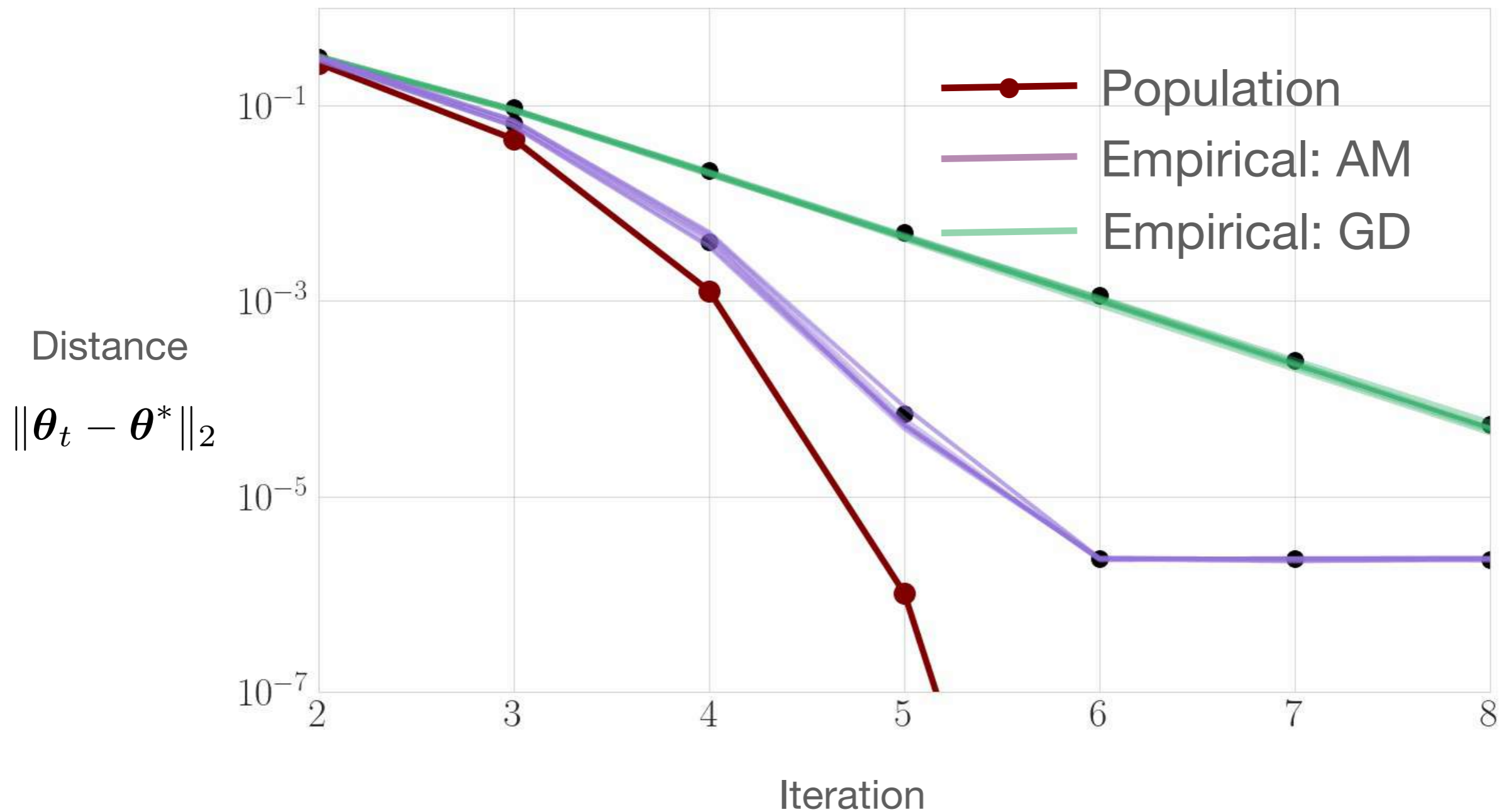
$\eta = 1/2$



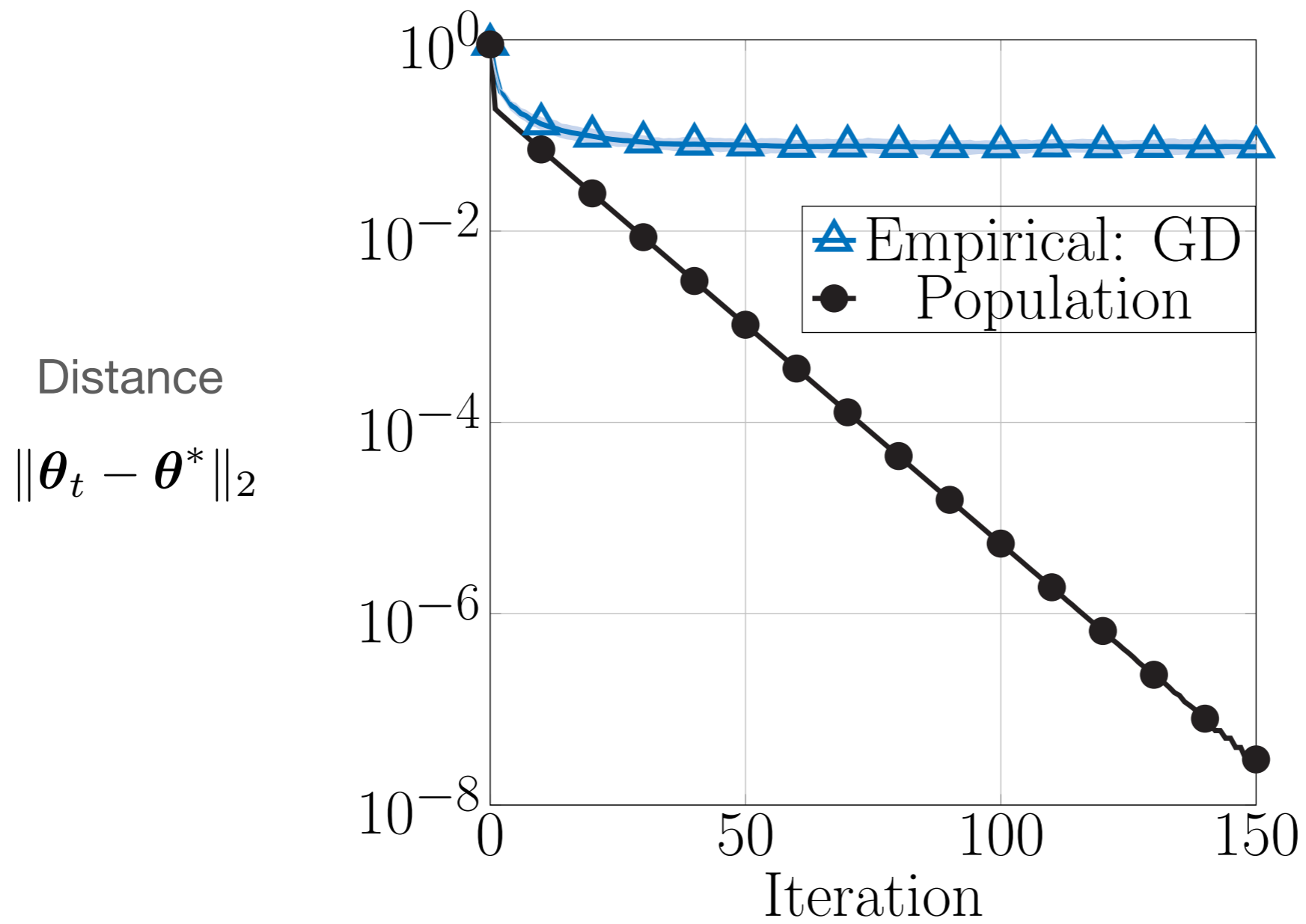


$d = 800, \quad n = 16,000, \quad \sigma = 10^{-6}$

$\eta = 1/2$



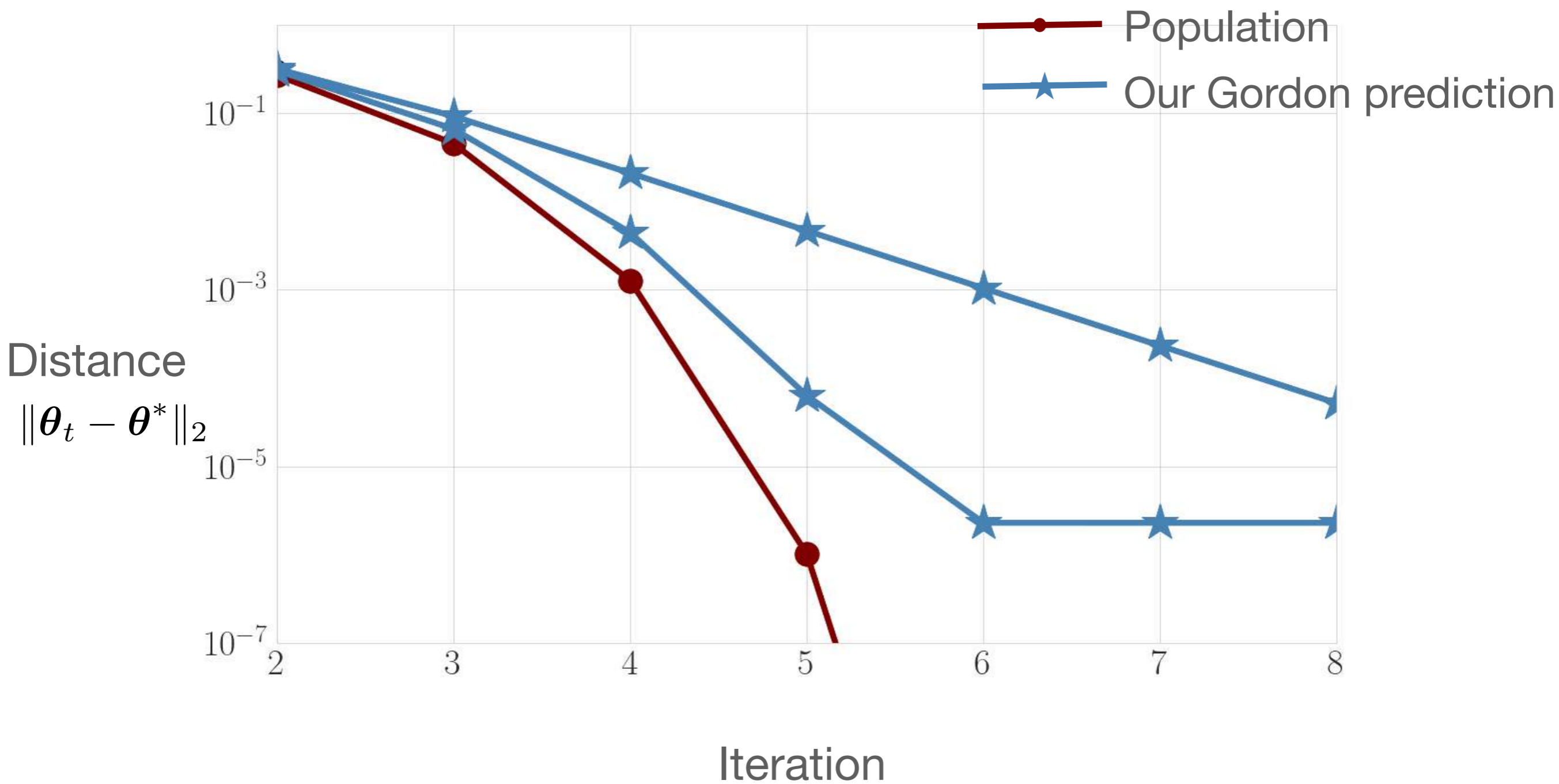
- ▶ Same population update, but different convergence behavior
- ▶ Population update significantly optimistic in both cases



$$d = 600, \quad n = 6,000, \quad \sigma = 0 \quad \eta = 0.95$$

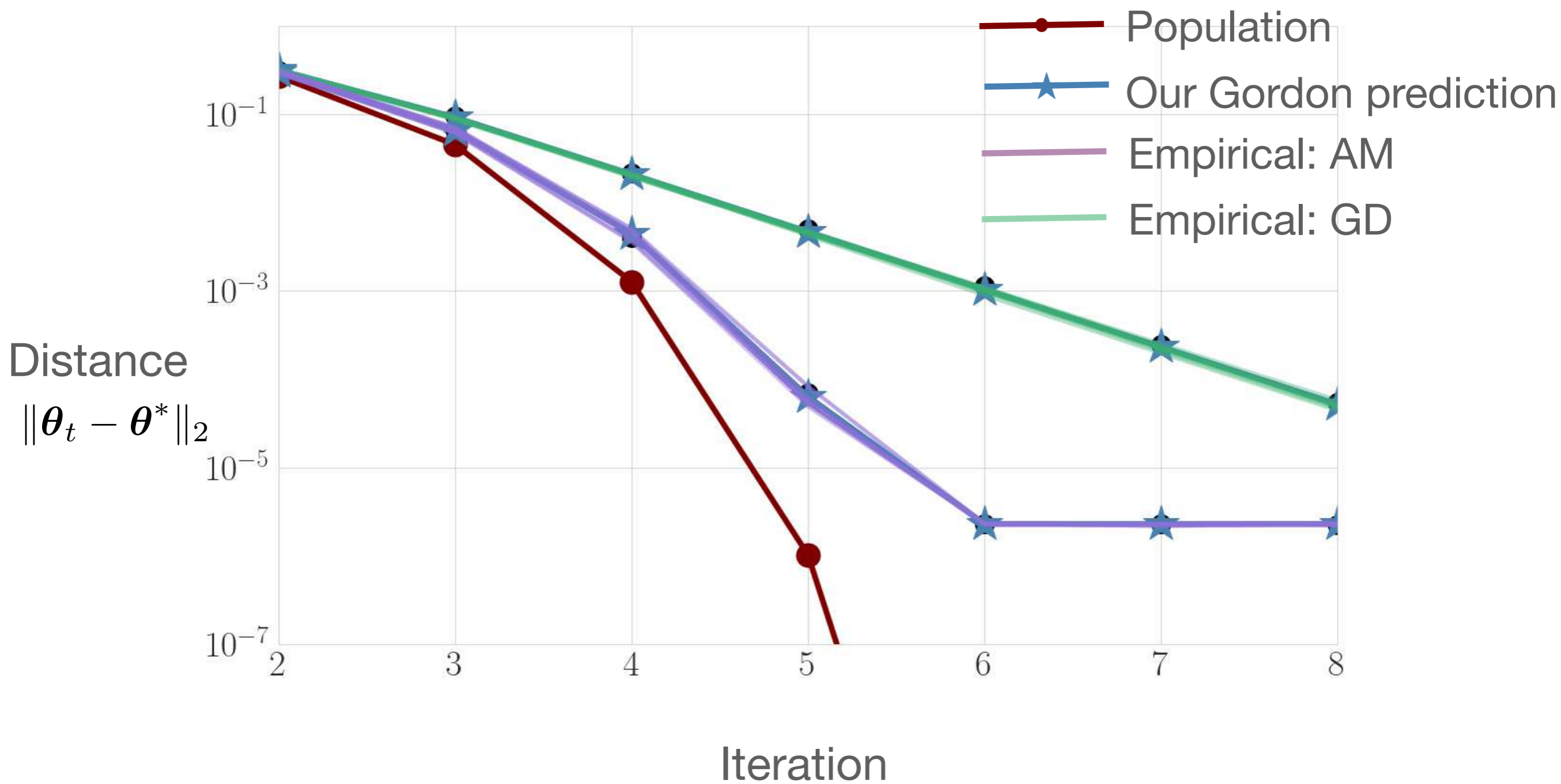
$d = 800, \quad n = 16,000, \quad \sigma = 10^{-6}$

$\eta = 1/2$



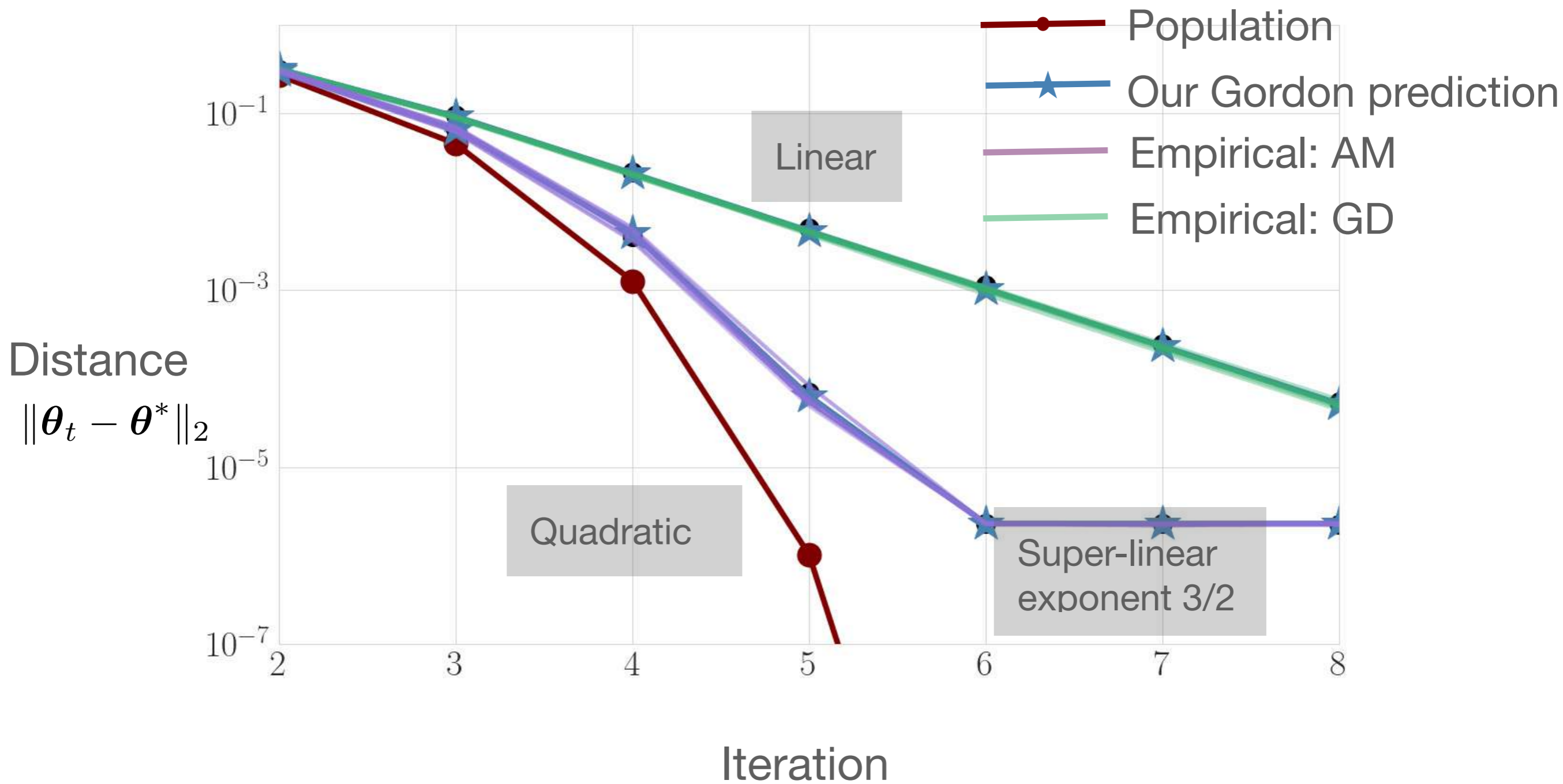
$d = 800, \quad n = 16,000, \quad \sigma = 10^{-6}$

$\eta = 1/2$

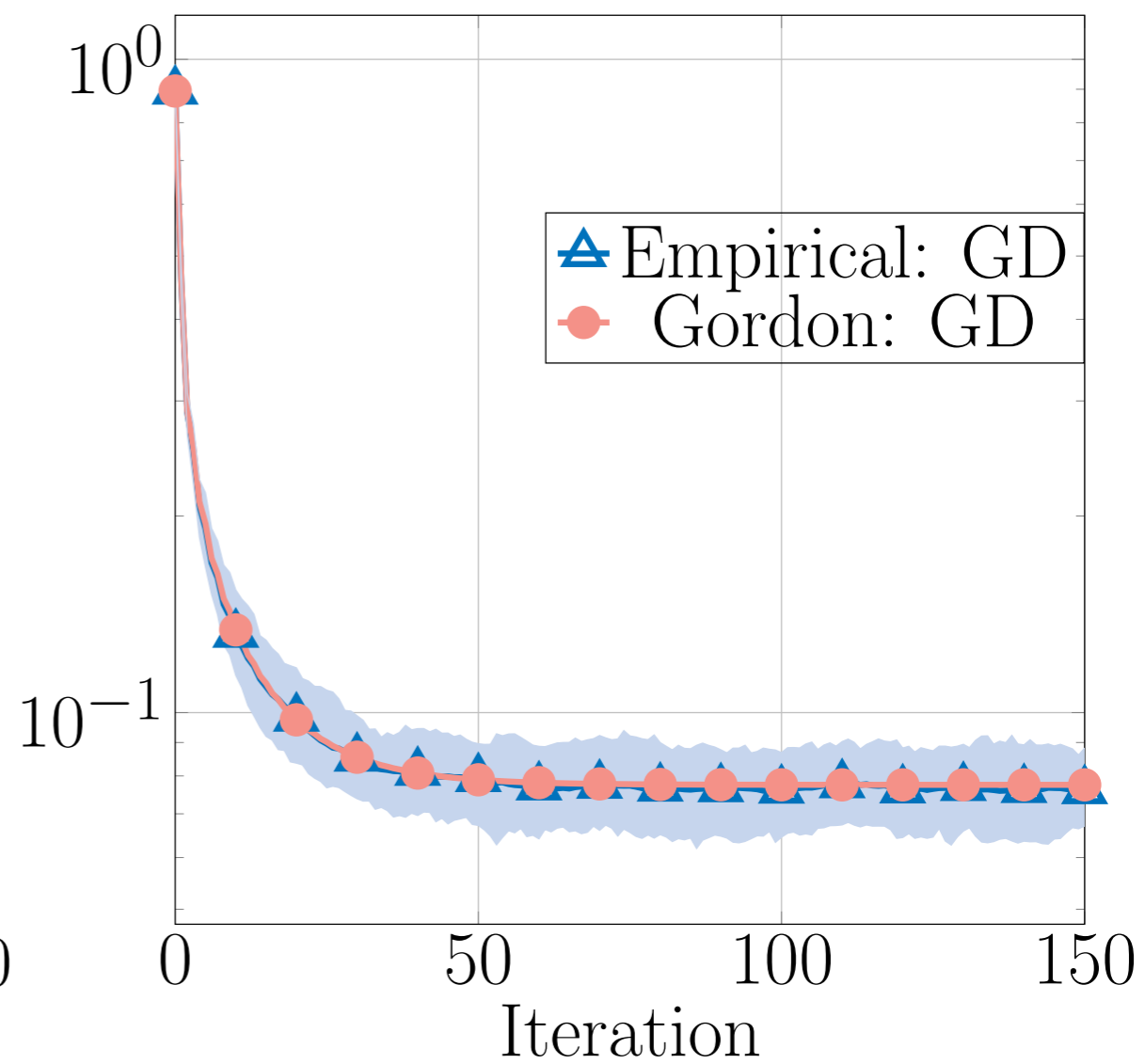
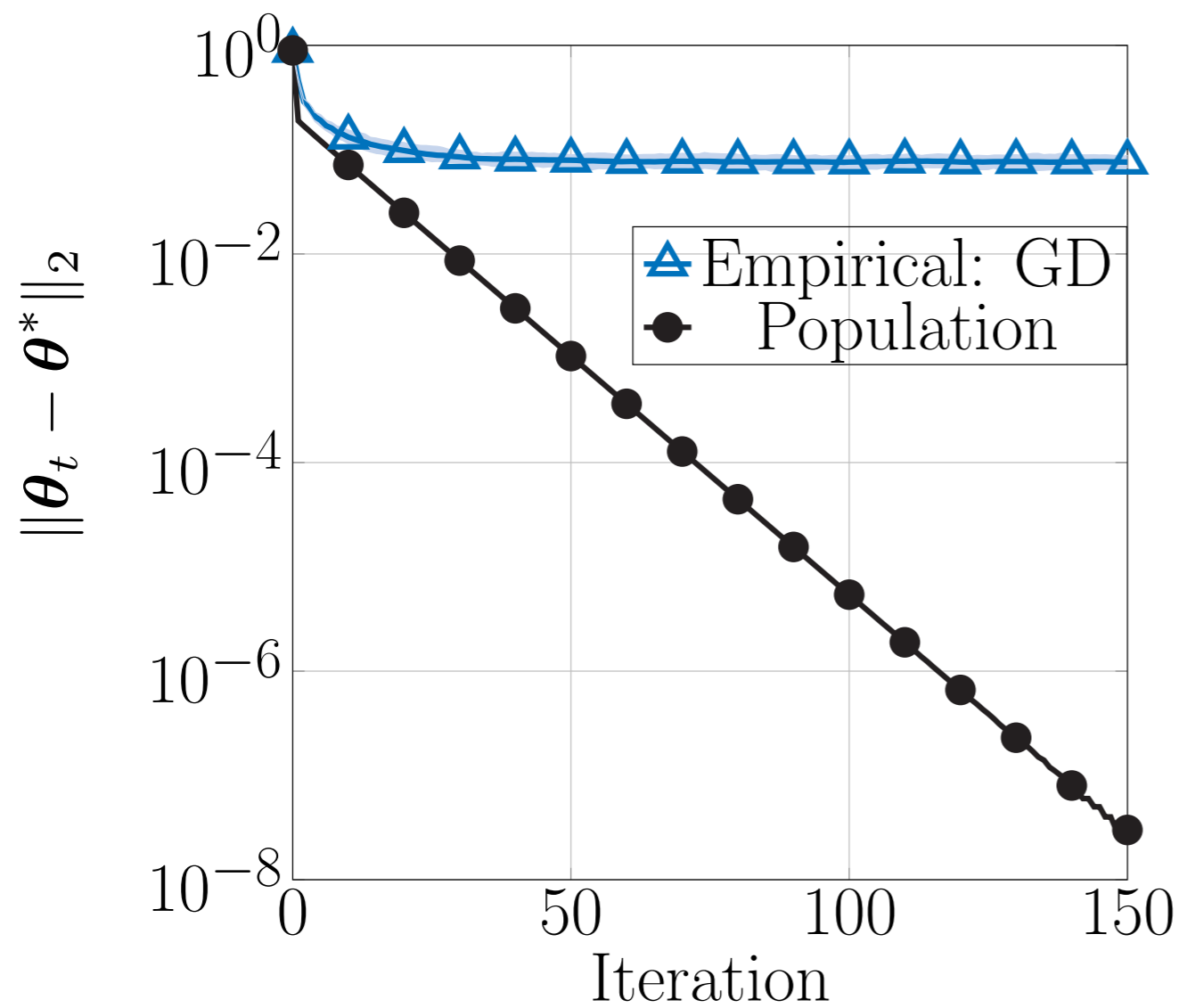


$d = 800, \quad n = 16,000, \quad \sigma = 10^{-6}$

$\eta = 1/2$





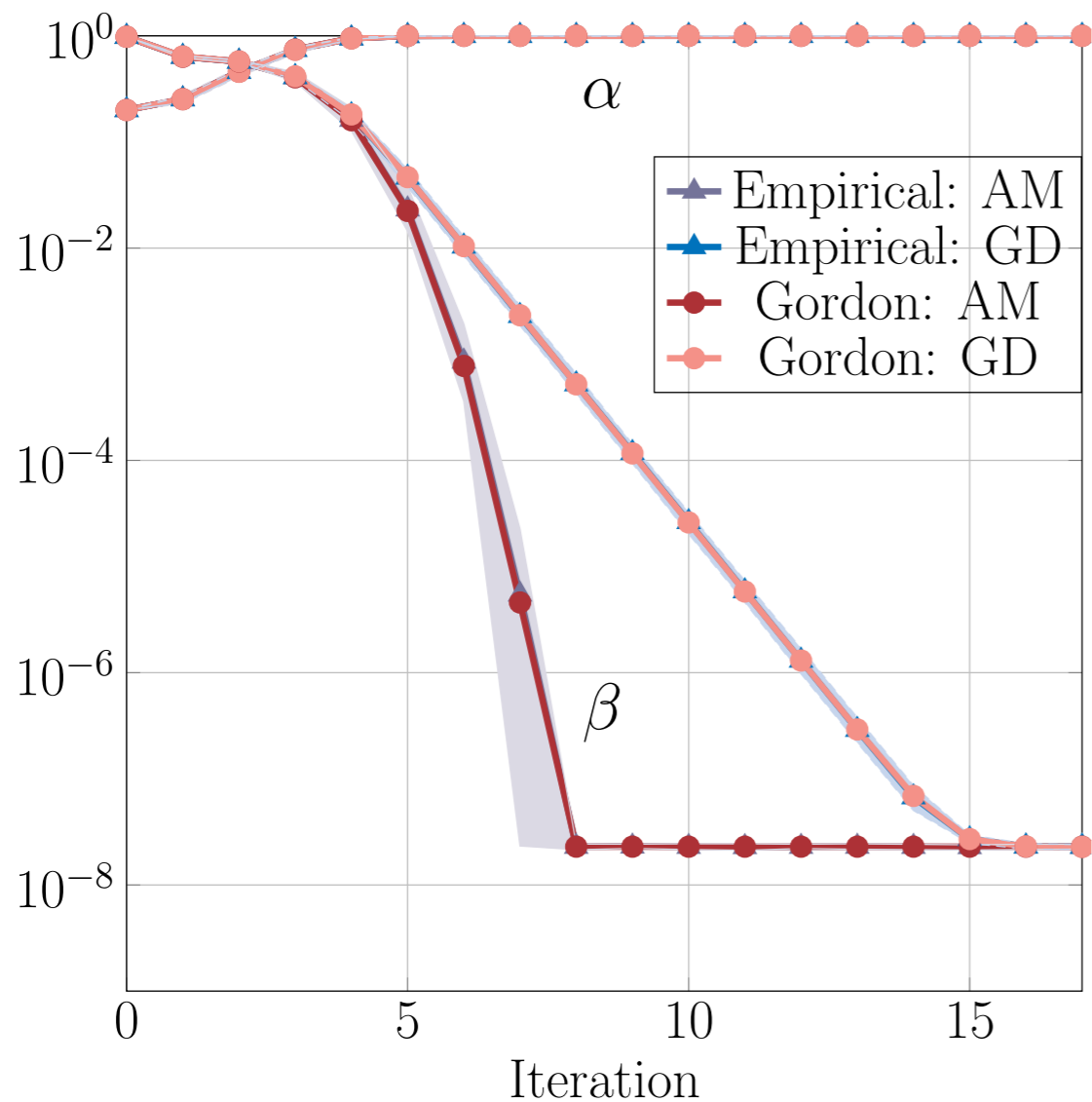


$d = 600, \quad n = 6,000, \quad \sigma = 0 \quad \eta = 0.95$

$n$  : per-iteration sample size  
 $d$  : dimension  
 $\Lambda = n/d$

$$\alpha(\boldsymbol{\theta}) = \langle \boldsymbol{\theta}, \boldsymbol{\theta}^* \rangle$$

$$\beta(\boldsymbol{\theta}) = \|\mathbf{P}_{\boldsymbol{\theta}^*}^\perp \boldsymbol{\theta}\|_2$$



► Gordon prediction is state evolution update

$$(\alpha(\boldsymbol{\theta}_t), \beta(\boldsymbol{\theta}_t)) \mapsto (\alpha_{t+1}^{\text{gor}}, \beta_{t+1}^{\text{gor}})$$

$$\approx (\alpha(\boldsymbol{\theta}_{t+1}), \beta(\boldsymbol{\theta}_{t+1}))$$

► Finite sample “correction” to population prediction:

$$(\beta_{t+1}^{\text{gor}})^2 = (\beta_{t+1}^{\text{pop}})^2 + \mathcal{O}(\Lambda^{-1}) \cdot \Delta(\alpha_t, \beta_t; \sigma)$$

**Part I:**

General iterate-by-iterate recipe  
if each iteration convex, Gaussianity

**Part II:**

Explicit one-step prediction for general class  
of models and methods

**Part III:**

Consequences for nonconvex model-fitting:  
Global convergence prediction

# Workhorse: The Convex Gaussian Minimax theorem

$$\mathbf{G} \in \mathbb{R}^{n \times d} \quad \mathbf{G}_{i,j} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$$

Primary

$$\mathcal{P}(\mathbf{G}) := \min_{\mathbf{u} \in \mathcal{U}} \max_{\mathbf{v} \in \mathcal{V}} \langle \mathbf{v}, \mathbf{G}\mathbf{u} \rangle + Q(\mathbf{u}, \mathbf{v})$$

Auxiliary

$$\mathcal{A}(\gamma_d, \gamma_n) := \min_{\mathbf{u} \in \mathcal{U}} \max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{v}\|_2 \cdot \langle \gamma_d, \mathbf{u} \rangle + \|\mathbf{u}\|_2 \cdot \langle \gamma_n, \mathbf{v} \rangle + Q(\mathbf{u}, \mathbf{v})$$

$$\gamma_d \sim \mathcal{N}(0, \mathbf{I}_d) \quad \gamma_n \sim \mathcal{N}(0, \mathbf{I}_n)$$

## Theorem

Suppose  $Q$  is continuous. Then the following hold for all scalars  $t$  :

(a) We have

$$\mathbb{P}(\mathcal{P}(\mathbf{G}) \leq t) \leq 2\mathbb{P}(\mathcal{A}(\gamma_d, \gamma_n) \leq t)$$

$$\zeta = (\alpha, \beta)$$

# I: The recipe

**Step 1:** Write an iteration as solution to convex program; then write objective in bilinear form

$$\theta_{t+1} = \arg \min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta; \theta_t, \mathbf{X}, \mathbf{y})$$

Variational form/Fenchel conjugate

$$\min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta; \theta_t, \mathbf{X}, \mathbf{y}) \stackrel{(d)}{=} \min_{\mathbf{u} \in \mathcal{U}} \max_{\mathbf{v} \in \mathcal{V}} \langle \mathbf{v}, \mathbf{G}\mathbf{u} \rangle + Q(\mathbf{u}, \mathbf{v})$$

**Step 2:** Invoke CGMT to replace matrix of Gaussian variables with two vectors

$$\begin{aligned} & \min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta; \theta_t, \mathbf{X}, \mathbf{y}) \\ & \approx \min_{\mathbf{u} \in \mathcal{U}} \max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{v}\|_2 \cdot \langle \gamma_d, \mathbf{u} \rangle + \|\mathbf{u}\|_2 \cdot \langle \gamma_n, \mathbf{v} \rangle + Q(\mathbf{u}, \mathbf{v}) \\ & =: \min_{\theta \in \mathbb{R}^d} \mathfrak{L}(\theta; \theta_t, \gamma_d, \gamma_n) \end{aligned}$$

**Step 3:** Scalarize: Obtain equiv. low-dimensional problem, solve

$$\begin{aligned} & \min_{\theta \in \mathbb{R}^d} \mathfrak{L}(\theta; \theta_t, \gamma_d, \gamma_n) \\ & \approx \min_{\zeta \in \mathbb{R}^2} \bar{L}(\zeta; \zeta_t) \end{aligned}$$

Obtain deterministic Gordon state evolution prediction

**Step 4:** Use growth conditions on the losses to make statements about minimizers. Empirical process theory



## II: Prediction for a general class of problems

### Model

i.i.d. observations:  $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$

$$y_i = f(\langle \mathbf{x}_i, \boldsymbol{\theta}^* \rangle; q_i) + \epsilon_i$$

$$\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d)$$

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$q_i \sim \mathcal{Q}$$

$$\|\boldsymbol{\theta}^*\|_2 = 1$$

- ▶ Phase retrieval
- ▶ Mixtures of regressions
- ▶ Mixtures of single-index models

### Algorithms

Higher-order methods

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{n} \|\omega(\mathbf{X}\boldsymbol{\theta}_t, \mathbf{y}) - \mathbf{X}\boldsymbol{\theta}\|_2^2 \\ &= \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \omega(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle, y_i) \cdot \mathbf{x}_i \right) \end{aligned}$$

First-order methods

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left\| \boldsymbol{\theta} - \boldsymbol{\theta}_t + \frac{2\eta}{n} \sum_{i=1}^n \omega(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle, y_i) \cdot \mathbf{x}_i \right\|_2^2 \\ &= \boldsymbol{\theta}_t - \frac{2\eta}{n} \sum_{i=1}^n \omega(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle, y_i) \cdot \mathbf{x}_i \end{aligned}$$

- ▶ Alternating projections, Newton methods
- ▶ Expectation maximization (EM), Newton EM
- ▶ (Sub-)gradient descent, gradient EM

# One-step Gordon update and deviation bound

$$\alpha = \alpha(\boldsymbol{\theta}_t)$$

$$\beta = \beta(\boldsymbol{\theta}_t)$$

$$Z_1, Z_2, Z_3 \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1) \quad Q \sim \mathcal{Q} \quad \Omega(\alpha, \beta) := \omega(\alpha Z_1 + \beta Z_2, f(Z_1; Q) + \sigma Z_3)$$

	First-order	Higher-order
$\alpha^{\text{gor}}$	$\alpha - 2\eta \cdot \mathbb{E}[Z_1 \Omega]$	$\mathbb{E}[Z_1 \Omega]$
$\beta^{\text{gor}}$	$\sqrt{(\beta - 2\eta \cdot \mathbb{E}[Z_2 \Omega])^2 + 4\eta^2 \mathbb{E}[\Omega^2] / \Lambda}$	$\sqrt{\mathbb{E}[Z_2 \Omega]^2 + (\Lambda - 1)^{-1} \cdot (\mathbb{E}[\Omega^2] - \mathbb{E}[Z_1 \Omega]^2 - \mathbb{E}[Z_2 \Omega]^2)}$

## Theorem

If  $\Lambda \geq C$  and  $n \gtrsim \log(1/\delta)$ , then with probability greater than  $1 - \delta$ :

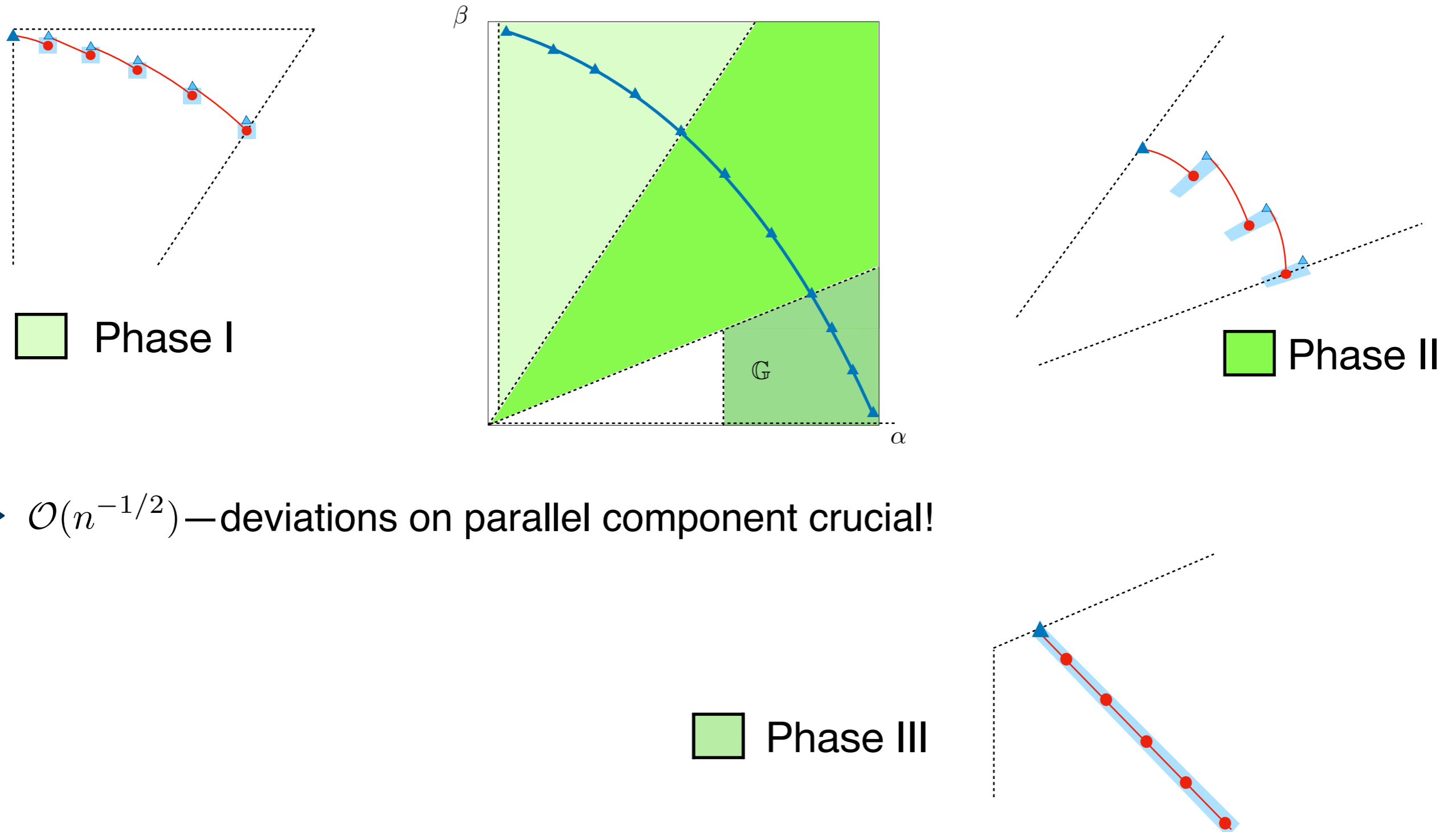
$$|\beta^{\text{gor}} - \beta(\boldsymbol{\theta}_{t+1})| \lesssim \left( \frac{\log(1/\delta)}{n} \right)^{1/4} \quad \text{and} \quad |\alpha^{\text{gor}} - \alpha(\boldsymbol{\theta}_{t+1})| \lesssim \left( \frac{\log^7(1/\delta)}{n} \right)^{1/2}.$$

► Fully non-asymptotic result, parallel component concentration requires additional argument

# III: Convergence guarantees

$$\alpha(\theta) = \langle \theta, \theta^* \rangle$$

$$\beta(\theta) = \|P_{\theta^*}^\perp \theta\|_2$$



►  $\mathcal{O}(n^{-1/2})$ —deviations on parallel component crucial!

### III: Convergence guarantees

- ▶ Natural parameter estimation metrics can be expressed in terms of state only:

$$\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 = \sqrt{(1 - \alpha)^2 + \beta^2} \quad \angle(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \tan^{-1}(\beta/\alpha)$$

- ▶ Gordon state evolution operator:  $\mathcal{S}_{\text{gor}} : (\alpha, \beta) \mapsto (\alpha^{\text{gor}}, \beta^{\text{gor}})$

Linear:  $c \cdot d(\zeta) + \varepsilon/2 \leq d(\mathcal{S}_{\text{gor}}(\zeta)) \leq C \cdot d(\zeta) + \varepsilon$

Superlinear:  $c \cdot [d(\zeta)]^\xi + \varepsilon/2 \leq d(\mathcal{S}_{\text{gor}}(\zeta)) \leq C \cdot [d(\zeta)]^\xi + \varepsilon$

# Example: Global convergence of AM for phase retrieval

## Theorem

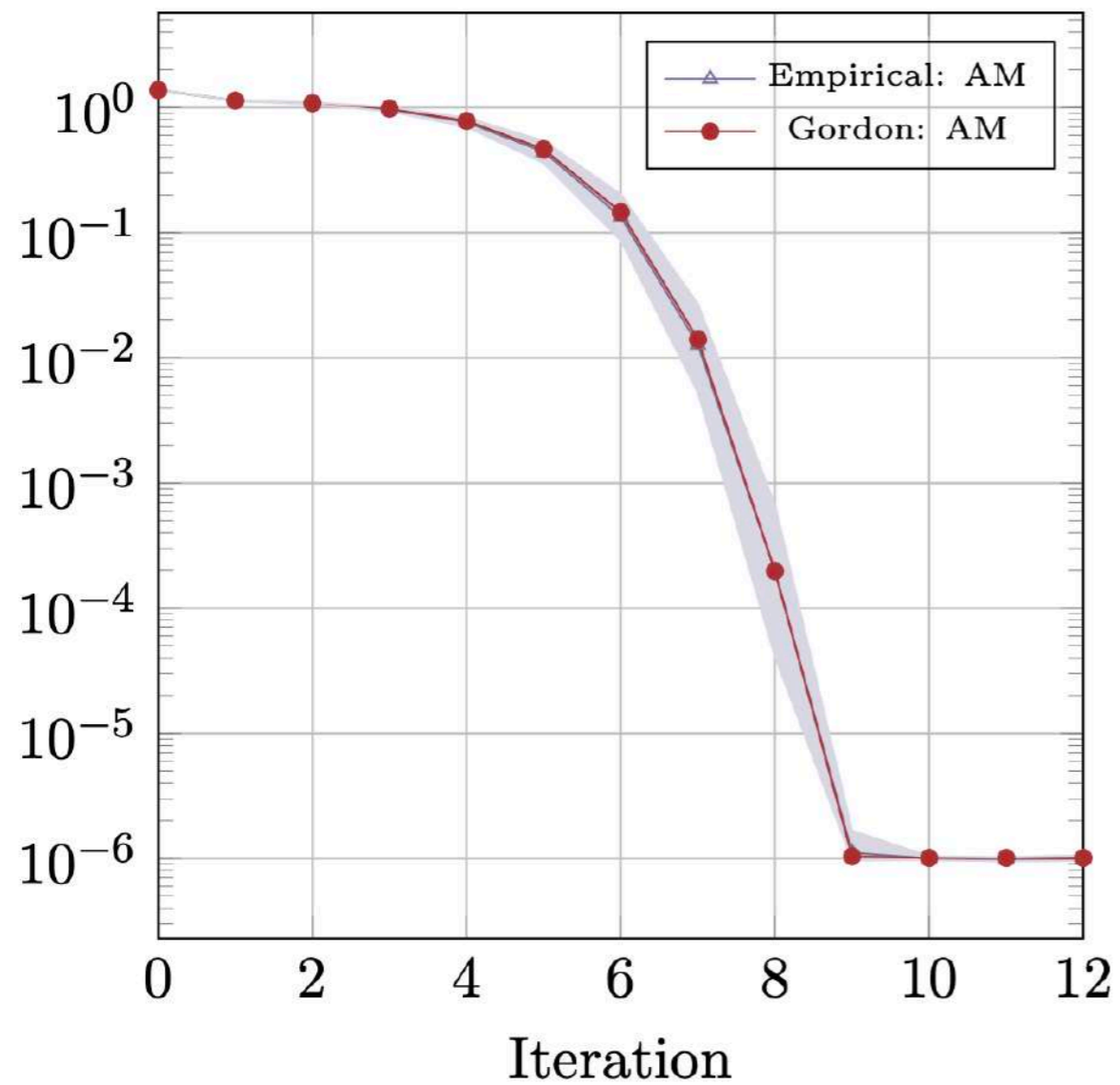
(a) The Gordon state evolution update converges in L2 superlinearly with exponent  $3/2$  within the local region  $\mathbb{G}$  to level  $\varepsilon = \sigma \sqrt{d/n}$ .

(b) If  $\theta \in \mathbb{G}$ , then with probability at least  $1 - 2Tn^{-10}$ :

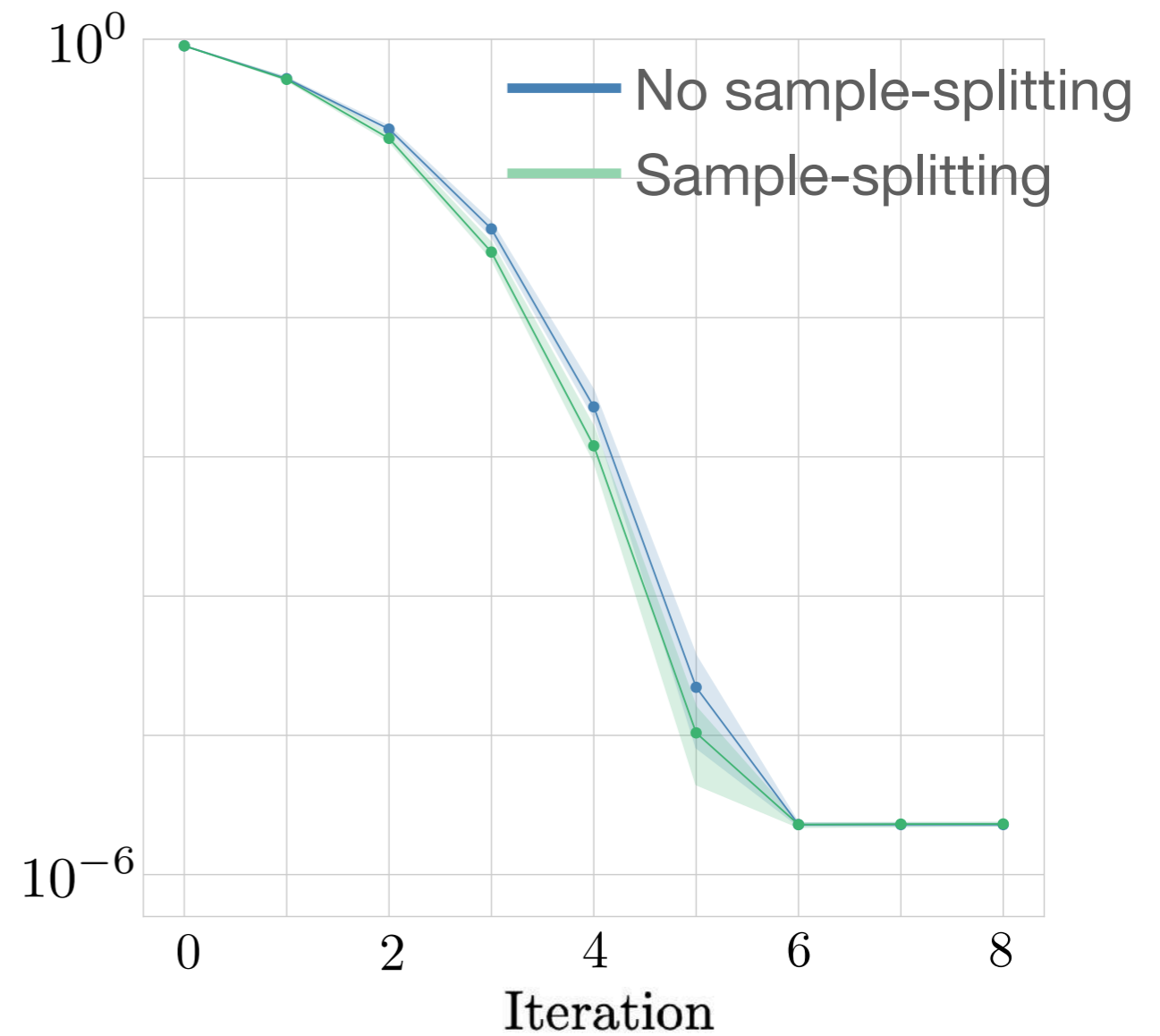
$$\max_{1 \leq t \leq T} |d_{\ell_2}(\mathcal{S}_{\text{gor}}^t(\zeta)) - \|\mathcal{T}_n^t(\theta) - \theta^*\|_2| \lesssim \left(\frac{\log n}{n}\right)^{1/4}.$$

- ▶ Parallel results for mixtures of regressions (angular not L2, AM converges linearly)
- ▶ Sample-splitting results in logarithmic blowup in total sample size

## Global convergence prediction



## Comparison: No sample-splitting



## Zooming out: A vignette

- ▶ Key “meta” observation: Nonconvex optimization can be reduced to iterative convex M-estimation in high dimensions
- ▶ Drawback of Gordon approach: Suboptimal non asymptotic guarantees that necessitate additional work to prove global convergence

### Question

Are there other ways to arrive at deterministic predictions with optimal concentration rates?

$$y_i = \langle \mathbf{x}_i, \boldsymbol{\mu}^* \rangle \cdot \langle \mathbf{z}_i, \boldsymbol{\nu}^* \rangle + \epsilon_i$$

$$\mathbf{x}_i, \mathbf{z}_i \sim \mathcal{N}(0, \mathbf{I})$$

$$\mathfrak{R}_n(\boldsymbol{\mu}, \boldsymbol{\nu}) = \frac{1}{n} \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \boldsymbol{\mu} \rangle \cdot \langle \mathbf{z}_i, \boldsymbol{\nu} \rangle)^2$$

⋮

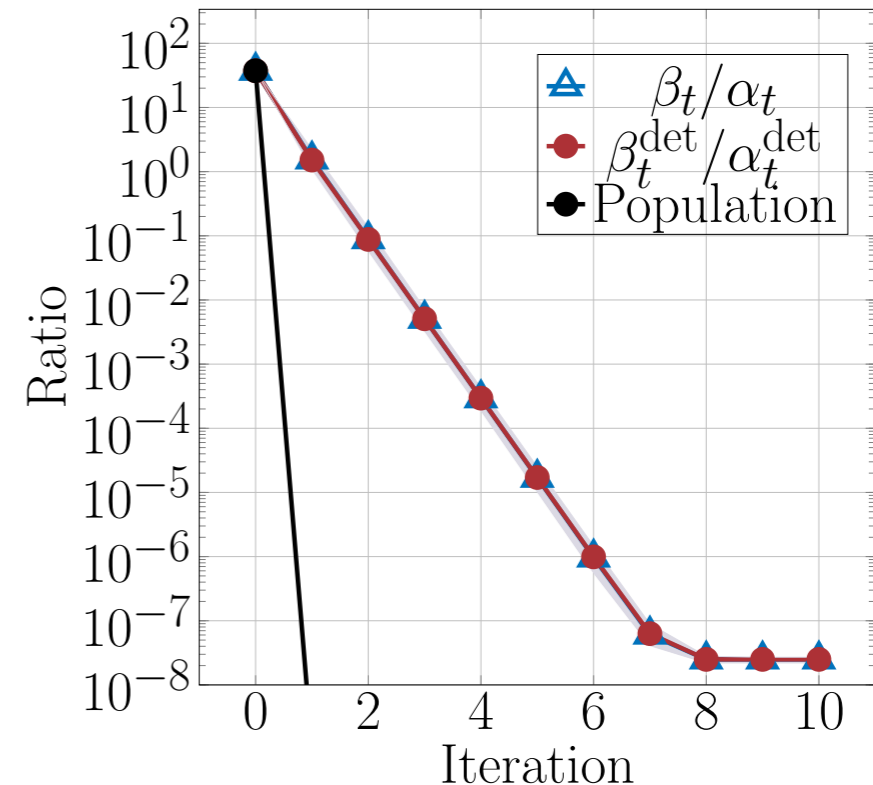
$$\boldsymbol{\nu}_{t+1} = \arg \min_{\boldsymbol{\nu} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \boldsymbol{\mu}_t \rangle \cdot \langle \mathbf{z}_i, \boldsymbol{\nu} \rangle)^2$$

$$\boldsymbol{\mu}_{t+1} = \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \boldsymbol{\mu} \rangle \cdot \langle \mathbf{z}_i, \boldsymbol{\nu}_{t+1} \rangle)^2$$

- ▶ Also other popular higher-order methods, e.g., composite optimization

# RMT-based prediction for AM in rank one bilinear sensing

- ▶ Uninformative population update
- ▶ Ratio (tangent of angle) converges linearly from random init. to noise floor
- ▶ Analysis enabled by “direct”, optimal, non-asymptotic concentration bound



## Theorem

If  $\Lambda \geq C$  and  $n \gtrsim \log(1/\delta)$ , then with probability greater than  $1 - \delta$ :

$$|\alpha^{\text{det}} - \alpha(\mu_{t+1})| \lesssim \left( \frac{\text{polylog}(n/\delta)}{n} \right)^{1/2} \quad \text{and} \quad |\beta^{\text{det}} - \beta(\mu_{t+1})| \lesssim \left( \frac{\text{polylog}(n/\delta)}{n} \right)^{1/2} .$$



- The population method can mis-predict efficiency in model-fitting
- Sharp characterizations of **convergence behavior** for iterative algorithms as well as **statistical accuracy** post-convergence.
- *Key properties:*
  - Each iteration is solution to convex optimization problem
  - Data is Gaussian conditioned on the past

Takeaways

- Removing the sample-splitting assumption
- Weakening the Gaussianity assumption
- Using sharp upper and lower bounds for “algorithmic” model-selection and hyperparameter tuning
- Broadly applicable machinery of reducing to iterative convex M-estimation:  
Can this say anything about your favorite model-fitting algorithm?

Open questions

**Sharp global convergence guarantees for iterative nonconvex optimization with random data**, with Chandrasekher and Thrampoulidis (under revision in Annals of Statistics)

**Higher order methods for rank one bilinear sensing: Random initialization and sharp predictions**, with Chandrasekher and Lou (coming soon)

# Backup

$$S_t = \text{span}(\boldsymbol{\theta}^*, \boldsymbol{\theta}_t)$$

## Example derivation: AM for phase retrieval

**Step 1:** Write an iteration as solution to convex program; then write objective in bilinear form

$$\boldsymbol{\theta}_{t+1} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{\theta}_t, \mathbf{X}, \mathbf{y})$$

Variational form/Fenchel conjugate

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{\theta}_t, \mathbf{X}, \mathbf{y}) \stackrel{(d)}{=} \min_{\mathbf{u} \in \mathcal{U}} \max_{\mathbf{v} \in \mathcal{V}} \langle \mathbf{v}, \mathbf{G}\mathbf{u} \rangle + Q(\mathbf{u}, \mathbf{v})$$

**Step 2:** Invoke **CGMT** to replace matrix of Gaussians with two vectors

$$\begin{aligned} & \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{\theta}_t, \mathbf{X}, \mathbf{y}) \\ & \approx \min_{\mathbf{u} \in \mathcal{U}} \max_{\mathbf{v} \in \mathcal{V}} \|\mathbf{v}\|_2 \cdot \langle \boldsymbol{\gamma}_d, \mathbf{u} \rangle + \|\mathbf{u}\|_2 \cdot \langle \boldsymbol{\gamma}_n, \mathbf{v} \rangle \\ & \quad + Q(\mathbf{u}, \mathbf{v}) \end{aligned}$$

$$=: \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{\theta}_t, \boldsymbol{\gamma}_d, \boldsymbol{\gamma}_n)$$

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{n} \cdot \sum_{i=1}^n (y_i - \text{sign}(\langle \mathbf{x}_i, \boldsymbol{\theta}_t \rangle) \cdot \langle \mathbf{x}_i, \boldsymbol{\theta} \rangle)^2 \\ &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \frac{1}{\sqrt{n}} \|\mathbf{X}\boldsymbol{\theta} - \text{diag}(\text{sign}(\mathbf{X}\boldsymbol{\theta}_t)) \cdot \mathbf{y}\|_2 \end{aligned}$$

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \max_{\|\mathbf{v}\|_2 \leq 1} \langle \mathbf{v}, \mathbf{X}\mathbf{P}_{S_t}^\perp \boldsymbol{\theta} \rangle + \langle \mathbf{v}, \mathbf{X}\mathbf{P}_{S_t} \boldsymbol{\theta} \rangle - \langle \mathbf{v}, \text{diag}(\text{sign}(\mathbf{X}\boldsymbol{\theta}_t)) \cdot \mathbf{y} \rangle$$

$$\begin{aligned} & \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \max_{\|\mathbf{v}\|_2 \leq 1} \|\mathbf{v}\|_2 \cdot \langle \boldsymbol{\gamma}_d, \mathbf{P}_{S_t}^\perp \boldsymbol{\theta} \rangle + \|\mathbf{P}_{S_t}^\perp \boldsymbol{\theta}\|_2 \cdot \langle \mathbf{v}, \boldsymbol{\gamma}_n \rangle \\ & \quad + \langle \mathbf{v}, \mathbf{X}\mathbf{P}_{S_t} \boldsymbol{\theta} \rangle - \langle \mathbf{v}, \text{diag}(\text{sign}(\mathbf{X}\boldsymbol{\theta}_t)) \cdot \mathbf{y} \rangle \end{aligned}$$

$$\alpha = \langle \boldsymbol{\theta}, \boldsymbol{\theta}^* \rangle \quad \mu = \frac{\langle \boldsymbol{\theta}, \mathbf{P}_{\boldsymbol{\theta}^*}^\perp \boldsymbol{\theta}_t \rangle}{\|\mathbf{P}_{\boldsymbol{\theta}^*}^\perp \boldsymbol{\theta}_t\|_2} \quad \nu = \|\mathbf{P}_{S_t}^\perp \boldsymbol{\theta}\|_2 \quad \beta^2 = \mu^2 + \nu^2$$

$$S_t = \text{span}(\boldsymbol{\theta}^*, \boldsymbol{\theta}_t)$$

**Step 3: Scalarize:** Obtain equiv. low-dimensional problem, solve

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{\theta}_t, \gamma_d, \gamma_n) \\ \approx \min_{\boldsymbol{\zeta} \in \mathbb{R}^3} \bar{L}(\boldsymbol{\zeta}; \boldsymbol{\zeta}_t) \end{aligned}$$

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left( -\frac{\|\mathbf{P}_{S_t}^\perp \boldsymbol{\theta}\|_2 \|\mathbf{P}_{S_t}^\perp \boldsymbol{\gamma}_d\|_2}{\sqrt{n}} + \|\text{diag}(\text{sign}(\mathbf{X}\boldsymbol{\theta}_t)) \cdot \mathbf{y} - \mathbf{X}\mathbf{P}_{S_t} \boldsymbol{\theta} - \|\mathbf{P}_{S_t}^\perp \boldsymbol{\theta}\|_2 \cdot \boldsymbol{\gamma}_n\|_2 \right)_+$$

$$\approx \min_{\substack{\alpha, \mu \\ \nu \geq 0}} \left( -\frac{\nu}{\sqrt{\Lambda}} + \sqrt{\mathbb{E}(\Omega_t - \alpha Z_1 - \mu Z_2 - \nu Z')^2} \right)_+$$

$$\text{sign}(\alpha_t Z_1 + \mu_t Z_2) \cdot |Z_1|$$

$$\alpha_{t+1}^{\text{gor}} = 1 - \frac{1}{\pi} (2\phi_t - \sin(2\phi_t)) \quad \phi_t = \tan^{-1}(\beta_t / \alpha_t)$$

$$\beta_{t+1}^{\text{gor}} = \sqrt{\frac{4}{\pi^2} \sin^4(\phi_t) + \frac{1 - (1 - \frac{1}{\pi} (2\phi_t - \sin(2\phi_t)))^2 - \frac{4}{\pi^2} \sin^4(\phi_t)}{\Lambda - 1}}$$

$$\mathcal{O}(\beta_t^4)$$

$$\mathcal{O}(\beta_t^3)$$

**Step 4:** Use growth conditions on the losses to make statements about minimizers.