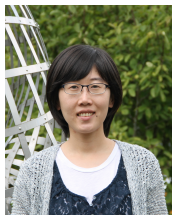


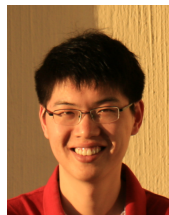
Nonconvex Optimization for High-Dimensional Signal Estimation: Spectral and Iterative Methods – Part I



Yuejie Chi
Carnegie Mellon



Yuxin Chen
Princeton



Cong Ma
UC Berkeley

EUSIPCO Tutorial, December 2020

Acknowledgement

- Our students and collaborators: Emmanuel J. Candès, Jianqing Fan, Yuanxin Li, Yingbin Liang, Yue M. Lu, Laixi Shi, Vincent Monardo, Tian Tong, Kaizheng Wang, Huishuai Zhang.
- This work is supported in part by ARO, AFOSR, ONR and NSF.



Outline

- Part I: Introduction and Warm-Up
Why nonconvex? basic concepts and a warm-up example (PCA)
- Part II: Gradient Descent and Implicit Regularization
phase retrieval, matrix completion, random initialization
- Part III: Spectral Methods
a general recipe, ℓ_2 and ℓ_∞ guarantees, community detection
- Part IV: Robustness to Corruptions and Ill-Conditioning
median truncation, least absolute deviation, scaled gradient descent

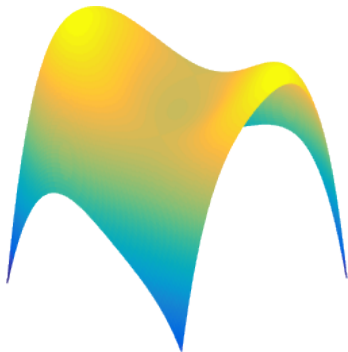
Outline

- Part I: Introduction and Warm-Up
Why nonconvex? basic concepts and a warm-up example (PCA)
- Part II: Gradient Descent and Implicit Regularization
phase retrieval, matrix completion, random initialization
- Part III: Spectral Methods
a general recipe, ℓ_2 and ℓ_∞ guarantees, community detection
- Part IV: Robustness to Corruptions and Ill-Conditioning
median truncation, least absolute deviation, scaled gradient descent

Nonconvex estimation problems are everywhere

Empirical risk minimization is usually nonconvex

$\text{minimize}_{\mathbf{x}} f(\mathbf{x}; \text{data}) \rightarrow$ loss function may be nonconvex

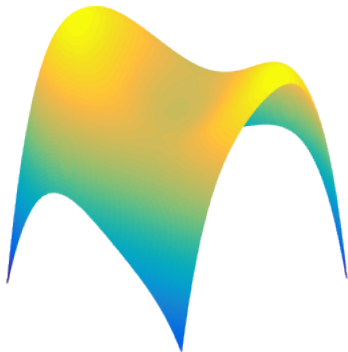


Nonconvex estimation problems are everywhere

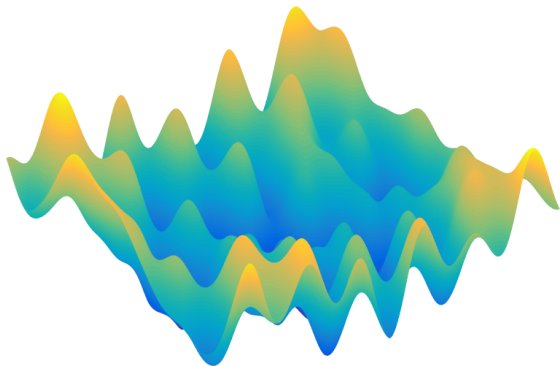
Empirical risk minimization is usually nonconvex

minimize _{\mathbf{x}} $f(\mathbf{x}; \text{data}) \rightarrow$ loss function may be nonconvex

- low-rank matrix completion
- blind deconvolution
- dictionary learning
- mixture models
- deep learning
- ...



Nonconvex optimization may be super scary



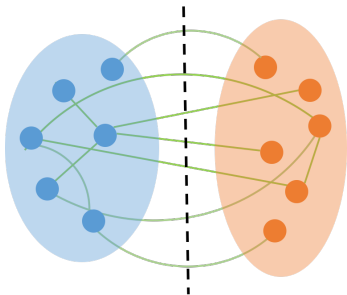
There may be bumps everywhere and exponentially many local optima

e.g. 1-layer neural net (Auer, Herbster, Warmuth '96; Vu '98)

Example: solving quadratic programs is hard

Finding maximum cut in a graph is about solving a quadratic program

$$\begin{array}{ll} \text{maximize}_x & \mathbf{x}^\top \mathbf{W} \mathbf{x} \\ \text{subj. to} & x_i^2 = 1, \quad i = 1, \dots, n \end{array}$$



Example: solving quadratic programs is hard



"I can't find an efficient algorithm, but neither can all these people."

figure credit: coding horror

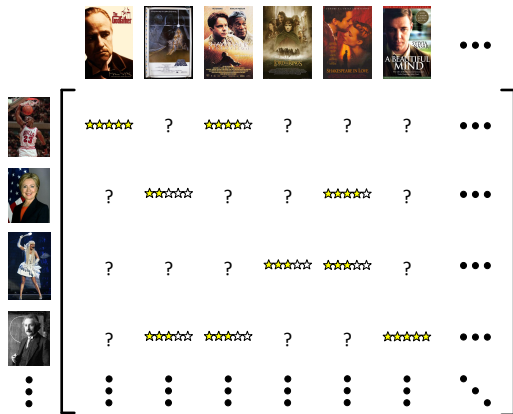
\$1,000,000 question

One strategy: convex relaxation

Can relax into convex problems by

- finding convex surrogates (e.g. matrix completion)
- lifting into higher dimensions (e.g. Max-Cut)

Example of convex surrogate: matrix completion



Netflix challenge

Predict unseen ratings

figure credit: Candès et al.

Low-rank modeling

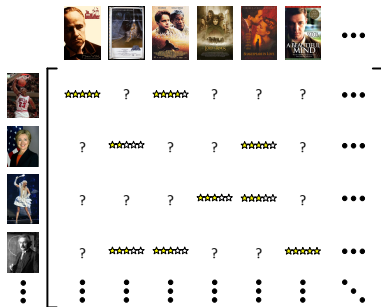
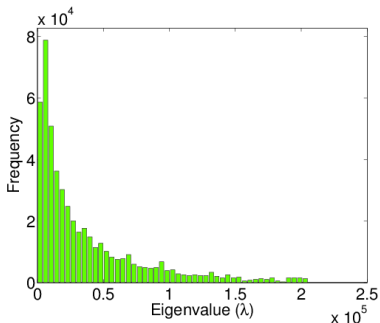


figure credit: E. Candès



A few factors explain most of the data

Low-rank modeling

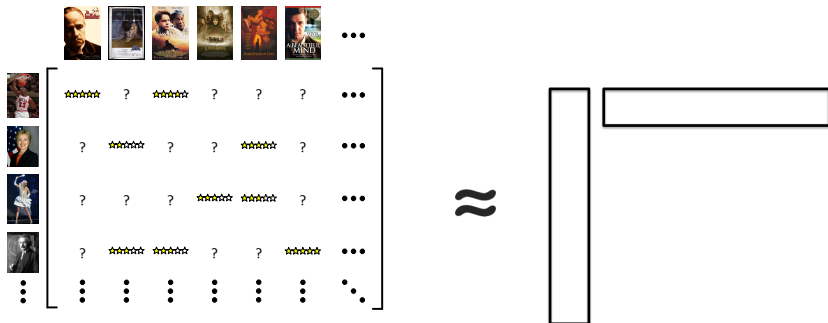


figure credit: E. Candès

A few factors explain most of the data → **low-rank** approximation

How to exploit (approx.) low-rank structure in prediction?

Example of convex surrogate: matrix completion

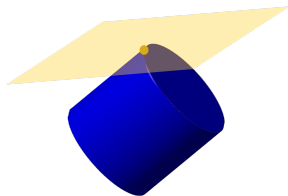
— Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09

minimize _{M} rank(M) subj. to data constraints



cvx surrogate

minimize _{M} nuc-norm(M) subj. to data constraints



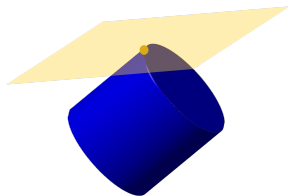
Example of convex surrogate: matrix completion

— Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09

minimize $_M$ rank(M) subj. to data constraints

↓ cvx surrogate

minimize $_M$ nuc-norm(M) subj. to data constraints



robust variation used by Netflix

— Candès, Li, Ma, Wright '10

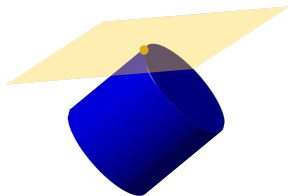
Example of convex surrogate: matrix completion

— Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09

minimize $_M$ rank(M) subj. to data constraints

↓ cvx surrogate

minimize $_M$ nuc-norm(M) subj. to data constraints



robust variation used by Netflix

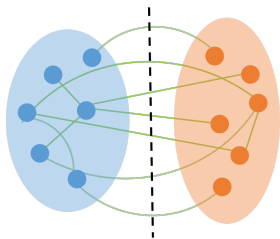
— Candès, Li, Ma, Wright '10

Problem: operate in *full* matrix space even though X is low-rank

Example of lifting: Max-Cut

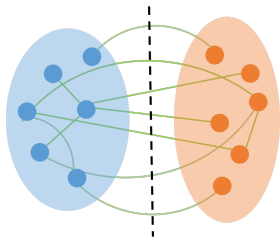
— Goemans, Williamson '95

$$\begin{aligned} \text{maximize}_x \quad & x^\top W x \\ \text{subj. to} \quad & x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$



Example of lifting: Max-Cut

— Goemans, Williamson '95



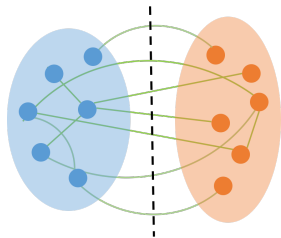
$$\begin{aligned} & \text{maximize}_{\mathbf{x}} && \mathbf{x}^\top \mathbf{W} \mathbf{x} \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

↓ let \mathbf{X} be $\mathbf{x}\mathbf{x}^\top$

$$\begin{aligned} & \text{maximize}_{\mathbf{X}} && \langle \mathbf{X}, \mathbf{W} \rangle \\ & \text{subj. to} && \mathbf{X}_{i,i} = 1, \quad i = 1, \dots, n \\ & && \mathbf{X} \succeq \mathbf{0} \\ & && \text{rank}(\mathbf{X}) = 1 \end{aligned}$$

Example of lifting: Max-Cut

— Goemans, Williamson '95



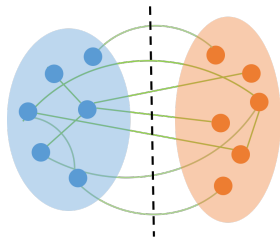
$$\begin{aligned} & \text{maximize}_{\mathbf{x}} && \mathbf{x}^\top \mathbf{W} \mathbf{x} \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

↓ let \mathbf{X} be $\mathbf{x}\mathbf{x}^\top$

$$\begin{aligned} & \text{maximize}_{\mathbf{X}} && \langle \mathbf{X}, \mathbf{W} \rangle \\ & \text{subj. to} && \mathbf{X}_{i,i} = 1, \quad i = 1, \dots, n \\ & && \mathbf{X} \succeq \mathbf{0} \\ & && \text{rank}(\mathbf{X}) = 1 \end{aligned}$$

Example of lifting: Max-Cut

— Goemans, Williamson '95



$$\begin{aligned} & \text{maximize}_x && x^\top W x \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

↓ let X be xx^\top

$$\begin{aligned} & \text{maximize}_X && \langle X, W \rangle \\ & \text{subj. to} && X_{i,i} = 1, \quad i = 1, \dots, n \\ & && X \succeq 0 \\ & && \text{rank}(X) = 1 \end{aligned}$$

Problem: explosion in dimensions ($\mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$)

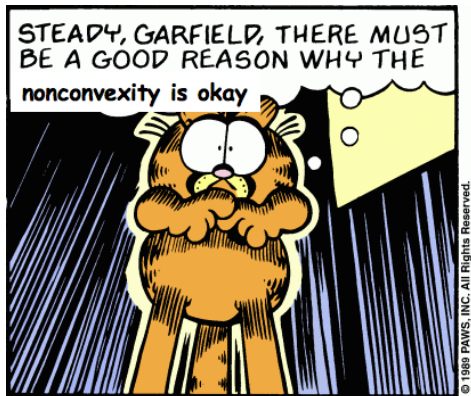
*How about optimizing nonconvex problems directly
without lifting?*

Nonconvex optimization

Complicated nonconvex problems are solved on a daily basis via simple algorithms such as stochastic gradient descent

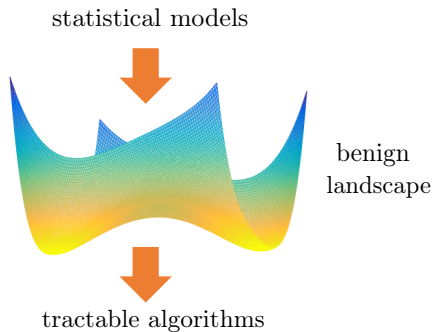
Nonconvex optimization

Complicated nonconvex problems are solved on a daily basis via simple algorithms such as stochastic gradient descent



- How come simple nonconvex algorithms work so well in practice?

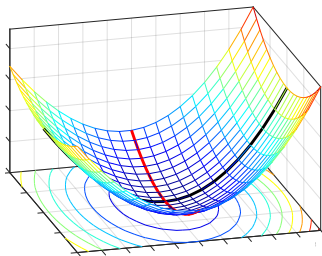
Statistical models come to rescue



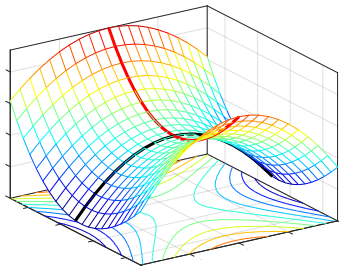
When data are generated by certain statistical models, problems are often much nicer than worst-case instances

Sometimes they are much nicer than we think

Under certain **statistical models**,
we see benign global geometry: **no spurious local optima**



global minimum



saddle point

*Even the simplest possible nonconvex methods
might be remarkably efficient under suitable statistical models*

A bit preliminaries of optimization

Unconstrained optimization

Consider an unconstrained optimization problem

$$\text{minimize}_x \quad f(\mathbf{x})$$

Definition 1 (first-order critical points)

A first-order critical point of f satisfies

$$\nabla f(\mathbf{x}) = \mathbf{0}$$

Unconstrained optimization

Consider an unconstrained optimization problem

$$\text{minimize}_x \quad f(\mathbf{x})$$

Definition 2 (second-order critical points)

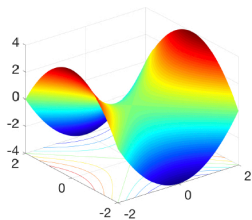
A second-order critical point \mathbf{x} satisfies

$$\nabla f(\mathbf{x}) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$$

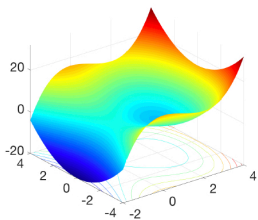
Several types of critical points

For any first-order critical point \mathbf{x} :

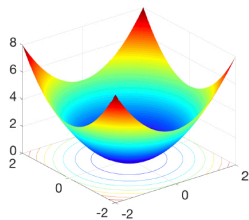
- $\nabla^2 f(\mathbf{x}) \prec \mathbf{0}$ \rightarrow local maximum
- $\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ \rightarrow local minimum
- $\lambda_{\min}(\nabla^2 f(\mathbf{x})) < 0$ \rightarrow *strict saddle point*



(a) strict saddle



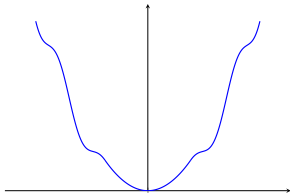
(b) local minimum



(c) global minimum

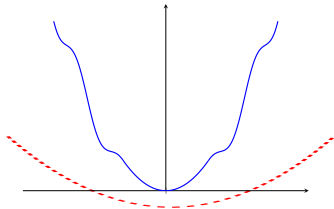
figure credit: Li et al. '16

Gradient descent theory



Two standard conditions that enable geometric convergence of GD

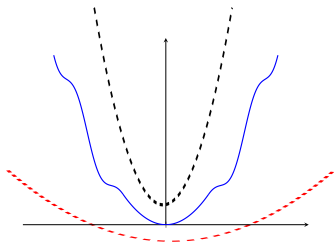
Gradient descent theory



Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity (or regularity condition)

Gradient descent theory



Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity (or regularity condition)
- (local) smoothness

$$\nabla^2 f(\mathbf{x}) \succ \mathbf{0} \quad \text{and} \quad \text{is well-conditioned}$$

Gradient descent theory revisited

f is said to be α -strongly convex and β -smooth if

$$\mathbf{0} \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

ℓ_2 error contraction: GD ($\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t)$) with $\eta = 1/\beta$ obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2$$

Gradient descent theory revisited

f is said to be α -strongly convex and β -smooth if

$$\mathbf{0} \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

ℓ_2 error contraction: GD ($\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t)$) with $\eta = 1/\beta$ obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2$$

- Condition number β/α determines rate of convergence

Gradient descent theory revisited

f is said to be α -strongly convex and β -smooth if

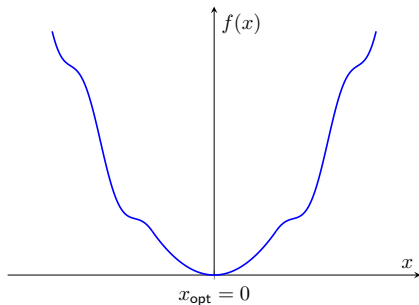
$$\mathbf{0} \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

ℓ_2 error contraction: GD ($\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t)$) with $\eta = 1/\beta$ obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2$$

- Condition number β/α determines rate of convergence
- Attains ε -accuracy within $O\left(\frac{\beta}{\alpha} \log \frac{1}{\varepsilon}\right)$ iterations

Regularity Condition (RC)



Definition 3 (Regularity Condition (RC))

$g(\cdot)$ is said to obey $\text{RC}(\mu, \lambda, \zeta)$ for some $\mu, \lambda, \zeta > 0$ if

$$2\langle g(\mathbf{x}), \mathbf{x} - \mathbf{x}_{\text{opt}} \rangle \geq \mu \|g(\mathbf{x})\|_2^2 + \lambda \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2^2 \quad \forall \mathbf{x}$$

Convergence under RC

ℓ_2 **error contraction:** The update rule ($\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \mathbf{g}(\mathbf{x}^t)$) with $\eta = \mu$ obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq (1 - \mu\lambda) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2$$

- $\mathbf{g}(\cdot)$: more general search directions
 - example: in vanilla GD, $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x})$

Convergence under RC

ℓ_2 **error contraction:** The update rule ($\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \mathbf{g}(\mathbf{x}^t)$) with $\eta = \mu$ obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq (1 - \mu\lambda) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2$$

- $\mathbf{g}(\cdot)$: more general search directions
 - example: in vanilla GD, $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x})$
- The product $\mu\lambda$ determines the rate of convergence

Convergence under RC

ℓ_2 **error contraction:** The update rule ($\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \mathbf{g}(\mathbf{x}^t)$) with $\eta = \mu$ obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq (1 - \mu\lambda) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2$$

- $\mathbf{g}(\cdot)$: more general search directions
 - example: in vanilla GD, $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x})$
- The product $\mu\lambda$ determines the rate of convergence
- Attains ε -accuracy within $O\left(\frac{1}{\mu\lambda} \log \frac{1}{\varepsilon}\right)$ iterations

RC = one-point strong convexity + smoothness

- One-point α -strong convexity:

$$f(\mathbf{x}_{\text{opt}}) - f(\mathbf{x}) \geq \langle \nabla f(\mathbf{x}), \mathbf{x}_{\text{opt}} - \mathbf{x} \rangle + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2^2 \quad (1)$$

- β -smoothness:

$$\begin{aligned} f(\mathbf{x}_{\text{opt}}) - f(\mathbf{x}) &\leq f\left(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})\right) - f(\mathbf{x}) \\ &\leq \left\langle \nabla f(\mathbf{x}), -\frac{1}{\beta} \nabla f(\mathbf{x}) \right\rangle + \frac{\beta}{2} \left\| \frac{1}{\beta} \nabla f(\mathbf{x}) \right\|_2^2 \\ &= -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2 \end{aligned} \quad (2)$$

RC = one-point strong convexity + smoothness

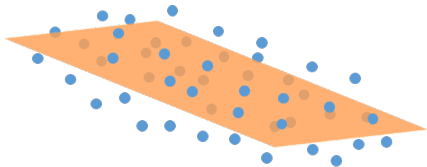
Combining (1) and (2) yields

$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}_{\text{opt}} \rangle \geq \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2^2 + \frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2 \quad (3)$$

— *RC holds with $\mu = 1/\beta$ and $\lambda = \alpha$*

A toy example: rank-1 matrix factorization

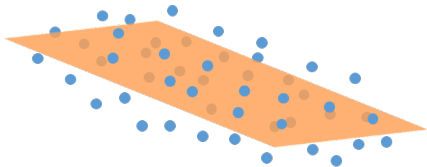
Revisiting PCA



Given $M \succeq \mathbf{0} \in \mathbb{R}^{n \times n}$ (not necessarily low-rank), find its best rank- r approximation:

$$\underbrace{\widehat{M} = \operatorname{argmin}_{\mathbf{Z}} \|\mathbf{Z} - M\|_{\mathbb{F}}^2 \quad \text{s.t.} \quad \operatorname{rank}(\mathbf{Z}) \leq r}_{\text{nonconvex optimization!}}$$

Revisiting PCA



This problem admits a closed-form solution

- let $M = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$ be eigen-decomposition of M ($\lambda_1 \geq \dots \geq \lambda_n$), then

$$\widehat{M} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$$

— *nonconvex, but tractable*

Optimization viewpoint

If we factorize $\mathbf{Z} = \mathbf{X}\mathbf{X}^\top$ with $\mathbf{X} \in \mathbb{R}^{n \times r}$, then it leads to a nonconvex problem:

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \frac{1}{4} \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_{\text{F}}^2$$

To simplify exposition, set $r = 1$:

$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{4} \|\mathbf{x}\mathbf{x}^\top - \mathbf{M}\|_{\text{F}}^2$$

Questions

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4} \|\mathbf{x}\mathbf{x}^\top - \mathbf{M}\|_{\text{F}}^2$$

- Where / what are the critical points?
- What does the curvature behave like, at least locally around the global minimizer?

Critical points of $f(\cdot)$

\mathbf{x} is a critical point, i.e. $\nabla f(\mathbf{x}) = (\mathbf{x}\mathbf{x}^\top - \mathbf{M})\mathbf{x} = \mathbf{0}$

\Leftrightarrow

$$\mathbf{M}\mathbf{x} = \|\mathbf{x}\|_2^2 \mathbf{x}$$

\Leftrightarrow

\mathbf{x} aligns with an eigenvector of \mathbf{M} or $\mathbf{x} = \mathbf{0}$

Since $\mathbf{M}\mathbf{u}_i = \lambda_i \mathbf{u}_i$, the set of critical points is given by

$$\{\mathbf{0}\} \cup \{\pm \sqrt{\lambda_i} \mathbf{u}_i, i = 1, \dots, n\}$$

Categorization of critical points

The critical points can be further categorized based on the **Hessians**:

$$\nabla^2 f(\mathbf{x}) := 2\mathbf{x}\mathbf{x}^\top + \|\mathbf{x}\|_2^2 \mathbf{I} - \mathbf{M}$$

- For any non-zero critical point $\mathbf{x}_k = \pm\sqrt{\lambda_k}\mathbf{u}_k$:

$$\begin{aligned}\nabla^2 f(\mathbf{x}_k) &= 2\lambda_k \mathbf{u}_k \mathbf{u}_k^\top + \lambda_k \mathbf{I} - \mathbf{M} \\ &= 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top + \lambda_k \left(\sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^\top \right) - \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top \\ &= \sum_{i:i \neq k} (\lambda_k - \lambda_i) \mathbf{u}_i \mathbf{u}_i^\top + 2\lambda_k \mathbf{u}_k \mathbf{u}_k^\top\end{aligned}$$

Categorization of critical points

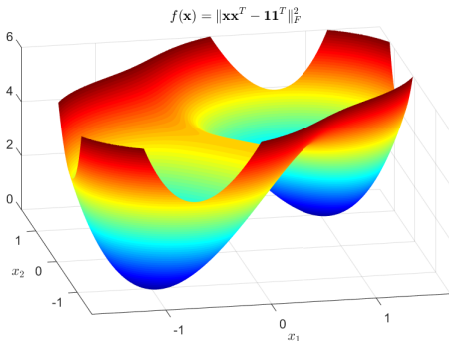
The critical points can be further categorized based on the **Hessians**:

$$\nabla^2 f(\mathbf{x}) := 2\mathbf{x}\mathbf{x}^\top + \|\mathbf{x}\|_2^2 \mathbf{I} - \mathbf{M}$$

- If $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_n \geq 0$, then
 - $\nabla^2 f(\mathbf{x}_1) \succ \mathbf{0}$ → local minima
 - $1 < k \leq n$: $\lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) < 0$, $\lambda_{\max}(\nabla^2 f(\mathbf{x}_k)) > 0$
→ strict saddle
 - $\mathbf{x} = \mathbf{0}$: $\nabla^2 f(\mathbf{0}) \preceq \mathbf{0}$ → local maxima

Good news: benign landscape

For example, for 2-dimensional case $f(\mathbf{x}) = \left\| \mathbf{x}\mathbf{x}^\top - \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \right\|_F^2$



global minima: $\mathbf{x} = \pm \begin{bmatrix} 1 \\ 1 \end{bmatrix}$; strict saddles: $\mathbf{x} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, and $\pm \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

— No “spurious” local minima!

Local strong convexity and local linear convergence

- The global minimizers: $\mathbf{x}_{\text{opt}} = \pm\sqrt{\lambda_1}\mathbf{u}_1$
- For all \mathbf{x} obeying $\underbrace{\|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2 \leq \frac{\lambda_1 - \lambda_2}{15\sqrt{\lambda_1}}}_{\text{basin of attraction}}$, one has

$$0.25(\lambda_1 - \lambda_2)\mathbf{I}_n \preceq \nabla^2 f(\mathbf{x}) \preceq 4.5\lambda_1\mathbf{I}_n$$

Local strong convexity and local linear convergence

- The global minimizers: $\mathbf{x}_{\text{opt}} = \pm\sqrt{\lambda_1}\mathbf{u}_1$
- For all \mathbf{x} obeying $\underbrace{\|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2}_{\text{basin of attraction}} \leq \frac{\lambda_1 - \lambda_2}{15\sqrt{\lambda_1}}$, one has

$$0.25(\lambda_1 - \lambda_2)\mathbf{I}_n \preceq \nabla^2 f(\mathbf{x}) \preceq 4.5\lambda_1\mathbf{I}_n$$

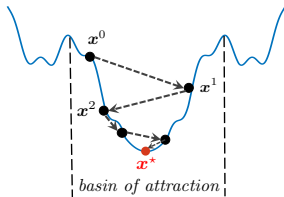
ℓ_2 **error contraction:** The GD iterates obey

$$\|\mathbf{x}^t - \sqrt{\lambda_1}\mathbf{u}_1\|_2 \leq \left(1 - \frac{\lambda_1 - \lambda_2}{18\lambda_1}\right)^t \|\mathbf{x}^0 - \sqrt{\lambda_1}\mathbf{u}_1\|_2, \quad t \geq 0,$$

as long as $\|\mathbf{x}^0 - \sqrt{\lambda_1}\mathbf{u}_1\|_2 \leq \frac{\lambda_1 - \lambda_2}{15\sqrt{\lambda_1}}$

Two vignettes

Two-stage approach:



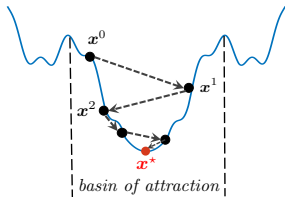
smart initialization

+

local refinement

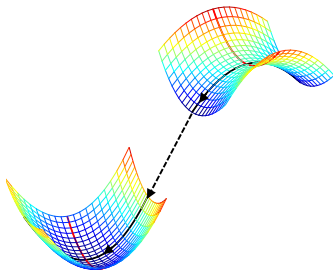
Two vignettes

Two-stage approach:



smart initialization
+
local refinement

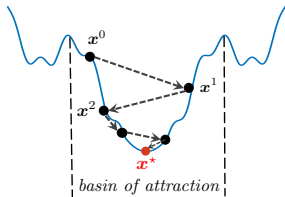
Global landscape:



benign landscape
+
saddle-point escaping

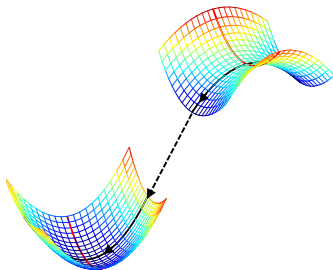
Two vignettes

Two-stage approach:



smart initialization
+
local refinement

Global landscape:



benign landscape
+
saddle-point escaping

This tutorial will mostly focus on the two-stage approach.

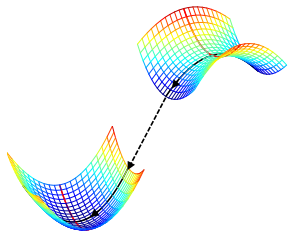
Global landscape

Benign landscape:

- all local minima = global minima
- other critical points = strict saddle points

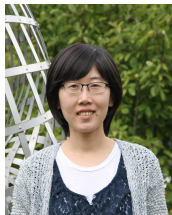
Saddle-point escaping algorithms:

- trust-region methods;
- perturbed gradient descent;
- perturbed SGD;
- etc...



Check the recent overview: *Zhang, Qu, Wright "From Symmetry to Geometry: Tractable Nonconvex Problems"*

Nonconvex Optimization for High-Dimensional Signal Estimation: Spectral and Iterative Methods – Part II



Yuejie Chi
Carnegie Mellon



Yuxin Chen
Princeton



Cong Ma
UC Berkeley

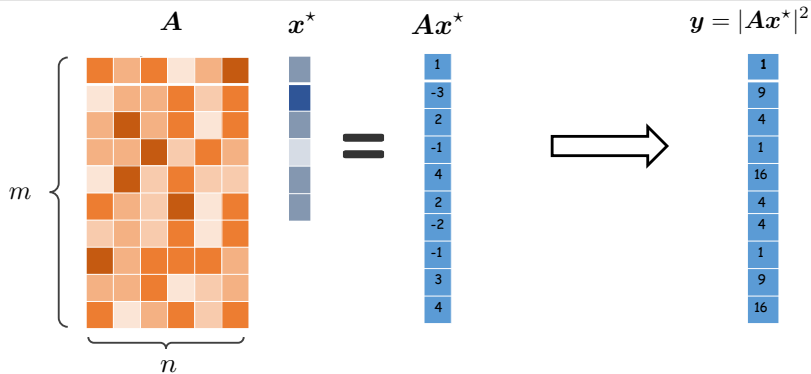
EUSIPCO Tutorial, December 2020

Outline

- Part I: Introduction and Warm-Up
Why nonconvex? basic concepts and a warm-up example (PCA)
- Part II: Gradient Descent and Implicit Regularization
phase retrieval, matrix completion, random initialization
- Part III: Spectral Methods
a general recipe, ℓ_2 and ℓ_∞ guarantees, community detection
- Part IV: Robustness to Corruptions and Ill-Conditioning
median truncation, least absolute deviation, scaled gradient descent

A case study: solving quadratic systems of equations

Solving quadratic systems of equations



Recover $\mathbf{x}^* \in \mathbb{R}^n$ from m random quadratic measurements

$$y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2, \quad k = 1, \dots, m$$

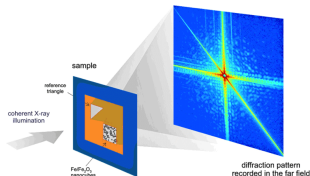
assume w.l.o.g. $\|\mathbf{x}^*\|_2 = 1$

Motivation: phase retrieval

Detectors record **intensities** of diffracted rays

- electric field $x(t_1, t_2) \rightarrow$ Fourier transform $\hat{x}(f_1, f_2)$

figure credit: Stanford SLAC



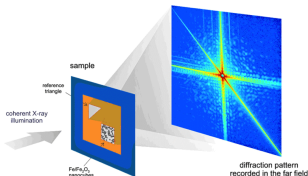
$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

Motivation: phase retrieval

Detectors record **intensities** of diffracted rays

- electric field $x(t_1, t_2) \rightarrow$ Fourier transform $\hat{x}(f_1, f_2)$

figure credit: Stanford SLAC

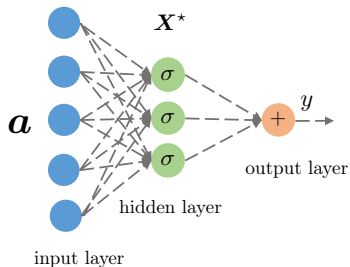


$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

Phase retrieval: recover signal $x(t_1, t_2)$ from intensity $|\hat{x}(f_1, f_2)|^2$

Motivation: learning neural nets with quadratic activation

— *Soltanolkotabi, Javanmard, Lee '17, Li, Ma, Zhang '17*

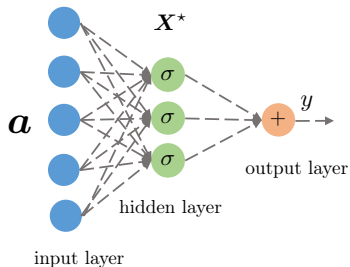


input features: \mathbf{a} ; weights: $\mathbf{X}^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_r^*]$

$$\text{output: } y = \sum_{i=1}^r \sigma(\mathbf{a}^\top \mathbf{x}_i^*)$$

Motivation: learning neural nets with quadratic activation

— *Soltanolkotabi, Javanmard, Lee '17, Li, Ma, Zhang '17*

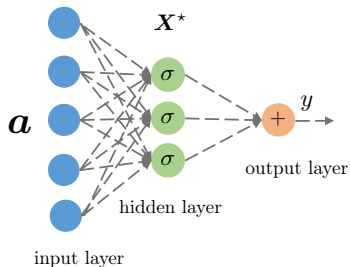


input features: \mathbf{a} ; weights: $\mathbf{X}^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_r^*]$

$$\text{output: } y = \sum_{i=1}^r \sigma(\mathbf{a}^\top \mathbf{x}_i^*) \stackrel{\sigma(z)=z^2}{:=} \sum_{i=1}^r (\mathbf{a}^\top \mathbf{x}_i^*)^2$$

Motivation: learning neural nets with quadratic activation

— Soltanolkotabi, Javanmard, Lee '17, Li, Ma, Zhang '17



input features: \mathbf{a} ; weights: $\mathbf{X}^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_r^*]$

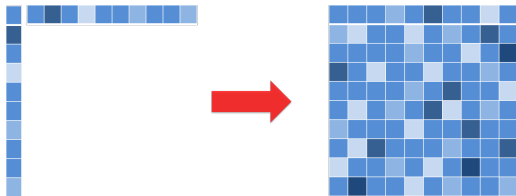
$$\text{output: } y = \sum_{i=1}^r \sigma(\mathbf{a}^\top \mathbf{x}_i^*) \stackrel{\sigma(z)=z^2}{:=} \sum_{i=1}^r (\mathbf{a}^\top \mathbf{x}_i^*)^2$$

We consider simplest model when $r = 1$ (higher r is similar)

An equivalent view: low-rank factorization

Introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$ to linearize constraints

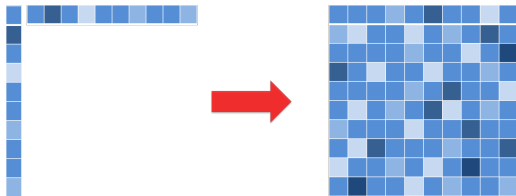
$$y_k = (\mathbf{a}_k^\top \mathbf{x})^2 = \mathbf{a}_k^\top (\mathbf{x}\mathbf{x}^\top) \mathbf{a}_k \quad \implies \quad y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k$$



An equivalent view: low-rank factorization

Introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$ to linearize constraints

$$y_k = (\mathbf{a}_k^\top \mathbf{x})^2 = \mathbf{a}_k^\top (\mathbf{x}\mathbf{x}^\top) \mathbf{a}_k \quad \Longrightarrow \quad y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k$$

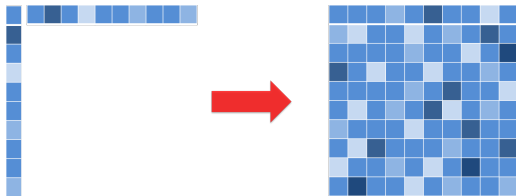


$$\begin{aligned} \text{find} \quad & \mathbf{X} \\ \text{s.t.} \quad & y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k, \quad k = 1, \dots, m \\ & \text{rank}(\mathbf{X}) = 1 \end{aligned}$$

An equivalent view: low-rank factorization

Introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$ to linearize constraints

$$y_k = (\mathbf{a}_k^\top \mathbf{x})^2 = \mathbf{a}_k^\top (\mathbf{x}\mathbf{x}^\top) \mathbf{a}_k \quad \Longrightarrow \quad y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k$$



$$\begin{aligned} \text{find} \quad & \mathbf{X} \\ \text{s.t.} \quad & y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k, \quad k = 1, \dots, m \\ & \text{rank}(\mathbf{X}) = 1 \end{aligned}$$

Solving quadratic systems is essentially **low-rank matrix completion**

A natural least-squares formulation

$$\text{given: } y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2, \quad 1 \leq k \leq m$$

↓

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m [(\mathbf{a}_k^\top \mathbf{x})^2 - y_k]^2$$

A natural least-squares formulation

$$\text{given: } y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2, \quad 1 \leq k \leq m$$

↓

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m [(\mathbf{a}_k^\top \mathbf{x})^2 - y_k]^2$$

- **pros:** often exact as long as sample size is sufficiently large

A natural least-squares formulation

$$\text{given: } y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2, \quad 1 \leq k \leq m$$

↓

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m [(\mathbf{a}_k^\top \mathbf{x})^2 - y_k]^2$$

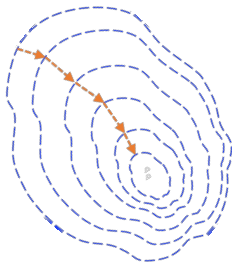
- **pros:** often exact as long as sample size is sufficiently large
- **cons:** $f(\cdot)$ is highly nonconvex
→ *computationally challenging!*

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

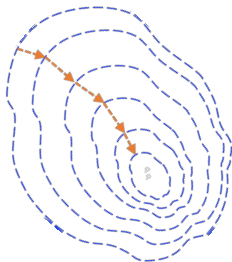
$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m [(\mathbf{a}_k^\top \mathbf{x})^2 - y_k]^2$$



- **spectral initialization:** $\mathbf{x}^0 \leftarrow$ leading eigenvector of certain data matrix

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$



- **spectral initialization:** $\mathbf{x}^0 \leftarrow$ leading eigenvector of certain data matrix
- **gradient descent:**

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t), \quad t = 0, 1, \dots$$

Spectral initialization

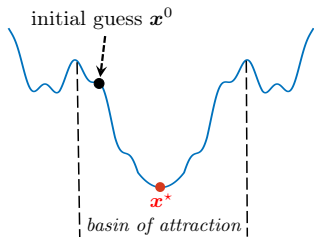
$\mathbf{x}^0 \leftarrow$ leading eigenvector of

$$\mathbf{Y} := \frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^\top$$

Rationale: under random Gaussian design $\mathbf{a}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$,

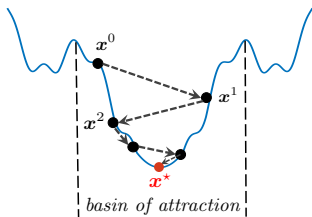
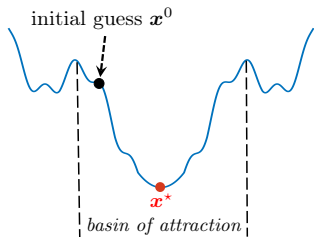
$$\mathbb{E}[\mathbf{Y}] := \mathbb{E} \left[\frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^\top \right] = \underbrace{\|\mathbf{x}^*\|_2^2 \mathbf{I} + 2\mathbf{x}^* \mathbf{x}^{*\top}}_{\text{leading eigenvector: } \pm \mathbf{x}^*}$$

Rationale of two-stage approach



1. initialize within local basin sufficiently close to x^*
(restricted) strongly convex; no saddles / spurious local mins

Rationale of two-stage approach



1. initialize within local basin sufficiently close to x^*
(restricted) strongly convex; no saddles / spurious local mins
2. iterative refinement

A highly incomplete list of two-stage methods

phase retrieval:

- Netrapalli, Jain, Sanghavi '13
- Candès, Li, Soltanolkotabi '14
- Chen, Candès '15
- Cai, Li, Ma '15
- Wang, Giannakis, Eldar '16
- Zhang, Zhou, Liang, Chi '16
- Kolte, Ozgur '16
- Zhang, Chi, Liang '16
- Soltanolkotabi '17
- Vaswani, Nayer, Eldar '16
- Chi, Lu '16
- Wang, Zhang, Giannakis, Akcakaya, Chen '16
- Tan, Vershynin '17
- Ma, Wang, Chi, Chen '17
- Duchi, Ruan '17
- Jeong, Gunturk '17
- Yang, Yang, Fang, Zhao, Wang, Neykov '17
- Qu, Zhang, Wright '17
- Goldstein, Studer '16
- Bahmani, Romberg '16
- Hand, Voroninski '16
- Wang, Giannakis, Saad, Chen '17
- Barmherzig, Sun '17
- ...

other problems:

- Keshavan, Montanari, Oh '09
- Sun, Luo '14
- Chen, Wainwright '15
- Tu, Boczar, Simchowitz, Soltanolkotabi, Recht '15
- Zheng, Lafferty '15
- Balakrishnan, Wainwright, Yu '14
- Chen, Suh '15
- Chen, Candès '16
- Li, Ling, Strohmer, Wei '16
- Yi, Park, Chen, Caramanis '16
- Jin, Kakade, Netrapalli '16
- Huang, Kakade, Kong, Valiant '16
- Ling, Strohmer '17
- Li, Ma, Chen, Chi '18
- Aghasi, Ahmed, Hand '17
- Lee, Tian, Romberg '17
- Li, Chi, Zhang, Liang '17
- Cai, Wang, Wei '17
- Abbe, Bandeira, Hall '14
- Chen, Kamath, Suh, Tse '16
- Zhang, Zhou '17
- Boumal '16
- Zhong, Boumal '17
- ...

First theory of WF

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min\{\|\mathbf{x}^t \pm \mathbf{x}^*\|_2\}$$

Theorem 1 (Candès, Li, Soltanolkotabi '14)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{4}\right)^{t/2} \|\mathbf{x}^*\|_2,$$

*with high prob., provided that step size $\eta \lesssim 1/n$ and sample size:
 $m \gtrsim n \log n$.*

First theory of WF

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min\{\|\mathbf{x}^t \pm \mathbf{x}^*\|_2\}$$

Theorem 1 (Candès, Li, Soltanolkotabi '14)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{4}\right)^{t/2} \|\mathbf{x}^*\|_2,$$

with high prob., provided that step size $\eta \lesssim 1/n$ and sample size: $m \gtrsim n \log n$.

- Iteration complexity: $O(n \log \frac{1}{\epsilon})$

First theory of WF

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min\{\|\mathbf{x}^t \pm \mathbf{x}^*\|_2\}$$

Theorem 1 (Candès, Li, Soltanolkotabi '14)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{4}\right)^{t/2} \|\mathbf{x}^*\|_2,$$

*with high prob., provided that step size $\eta \lesssim 1/n$ and sample size:
 $m \gtrsim n \log n$.*

- Iteration complexity: $O(n \log \frac{1}{\epsilon})$
- Sample complexity: $O(n \log n)$

First theory of WF

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min\{\|\mathbf{x}^t \pm \mathbf{x}^*\|_2\}$$

Theorem 1 (Candès, Li, Soltanolkotabi '14)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{4}\right)^{t/2} \|\mathbf{x}^*\|_2,$$

with high prob., provided that step size and sample size: .

- Iteration complexity: $O(n \log \frac{1}{\epsilon})$
- Sample complexity: $O(n \log n)$
- Derived based on (worst-case) local geometry

Improved theory of WF

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min\{\|\mathbf{x}^t \pm \mathbf{x}^*\|_2\}$$

Theorem 2 (Ma, Wang, Chi, Chen '17)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{2}\right)^t \|\mathbf{x}^*\|_2$$

with high prob., provided that step size $\eta \asymp 1/\log n$ and sample size $m \gtrsim n \log n$.

- Iteration complexity: $O(n \log \frac{1}{\epsilon}) \searrow O(\log n \log \frac{1}{\epsilon})$
- Sample complexity: $O(n \log n)$
- Derived based on finer analysis of GD trajectory

What does optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

What does optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

$$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$$

What does optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ but ill-conditioned (even locally)
condition number $\asymp n$

What does optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ but ill-conditioned (even locally)
condition number $\asymp n$

Consequence (Candès et al '14): WF attains ε -accuracy within $O(n \log \frac{1}{\varepsilon})$ iterations if $m \asymp n \log n$

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations



Step size taken to be $\eta = O(1/n)$

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations



Step size taken to be $\eta = O(1/n)$



This choice is suggested by **worst-case** optimization theory

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations



Step size taken to be $\eta = O(1/n)$

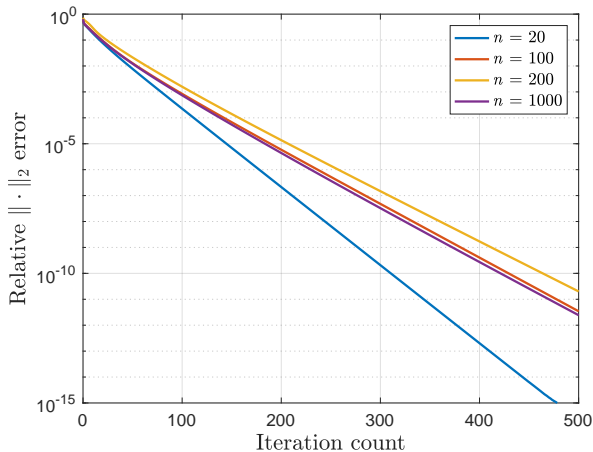


This choice is suggested by **worst-case** optimization theory



Does it capture what really happens?

Numerical efficiency with $\eta_t = 0.1$



Vanilla GD (WF) converges fast for a constant step size!

A second look at gradient descent theory

Which local region enjoys both strong convexity and smoothness?

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \left[3(\mathbf{a}_k^\top \mathbf{x})^2 - (\mathbf{a}_k^\top \mathbf{x}^*)^2 \right] \mathbf{a}_k \mathbf{a}_k^\top$$

A second look at gradient descent theory

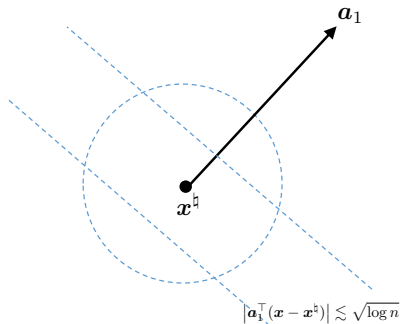
Which local region enjoys both strong convexity and smoothness?

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \left[3(\mathbf{a}_k^\top \mathbf{x})^2 - (\mathbf{a}_k^\top \mathbf{x}^*)^2 \right] \mathbf{a}_k \mathbf{a}_k^\top$$

- Not sufficiently smooth if \mathbf{x} and \mathbf{a}_k are too close (coherent)

A second look at gradient descent theory

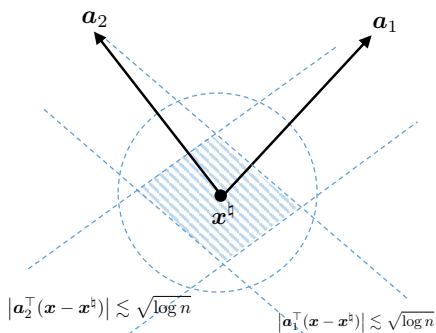
Which local region enjoys both strong convexity and smoothness?



- x is incoherent w.r.t. sampling vectors $\{a_k\}$ (incoherence region)

A second look at gradient descent theory

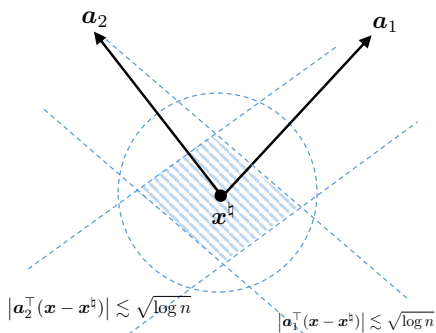
Which local region enjoys both strong convexity and smoothness?



- x is incoherent w.r.t. sampling vectors $\{a_k\}$ (incoherence region)

A second look at gradient descent theory

Which local region enjoys both strong convexity and smoothness?

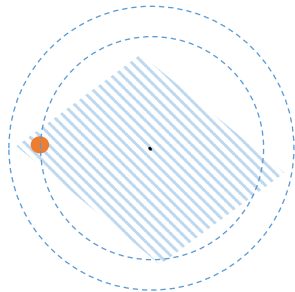


- x is incoherent w.r.t. sampling vectors $\{a_k\}$ (incoherence region)

Prior works suggest enforcing **regularization** (e.g. truncation, projection, regularized loss) to promote incoherence

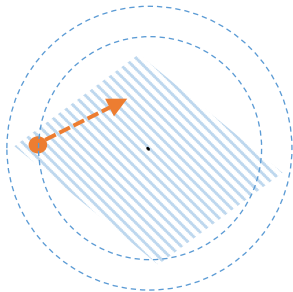
Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



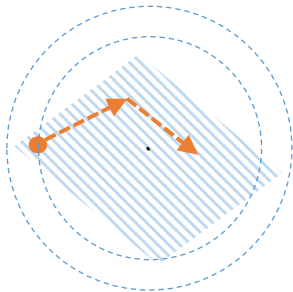
Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



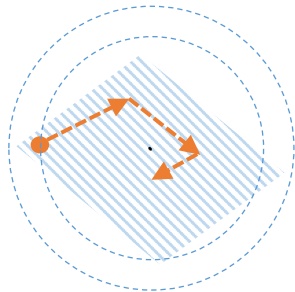
Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



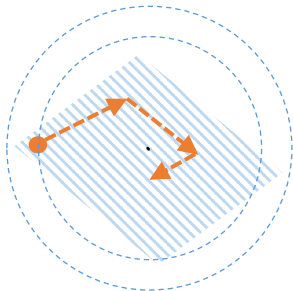
Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



Encouraging message: GD is implicitly regularized

- region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent with** $\{\mathbf{a}_k\}$

$$\max_k |\mathbf{a}_k^\top (\mathbf{x}^t - \mathbf{x}^*)| \lesssim \sqrt{\log n} \|\mathbf{x}^*\|_2, \quad \forall t$$

- cannot be derived from generic optimization theory; relies on finer statistical analysis for entire trajectory of GD

Theoretical guarantees for local refinement stage

Theorem 3 (Ma, Wang, Chi, Chen '17)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

- $\max_k |\mathbf{a}_k^\top \mathbf{x}^t| \lesssim \sqrt{\log n} \|\mathbf{x}^*\|_2$ (incoherence)

Theoretical guarantees for local refinement stage

Theorem 3 (Ma, Wang, Chi, Chen '17)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

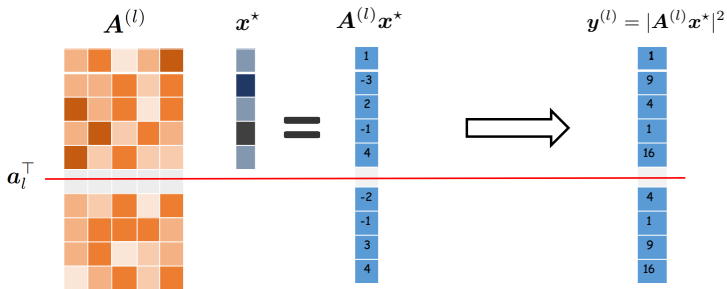
- $\max_k |\mathbf{a}_k^\top \mathbf{x}^t| \lesssim \sqrt{\log n} \|\mathbf{x}^*\|_2$ (incoherence)
- $\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim (1 - \frac{\eta}{2})^t \|\mathbf{x}^*\|_2$ (linear convergence)

provided that step size $\eta \asymp 1/\log n$ and sample size $m \gtrsim n \log n$.

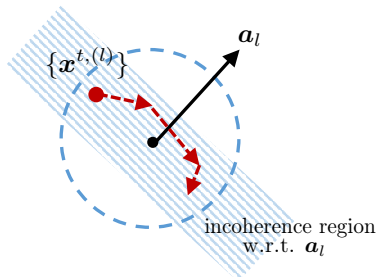
- Attains ε accuracy within $O(\log n \log \frac{1}{\varepsilon})$ iterations

Key proof idea: leave-one-out analysis

For each $1 \leq l \leq m$, introduce leave-one-out iterates $\mathbf{x}^{t,(l)}$ by dropping l th measurement

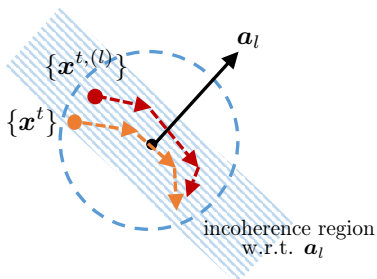


Key proof idea: leave-one-out analysis



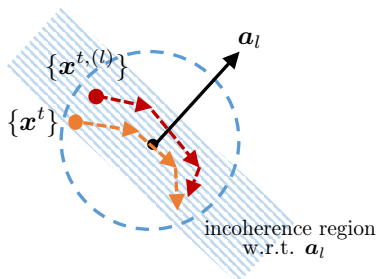
- Leave-one-out iterate $\mathbf{x}^{t,(l)}$ is independent of \mathbf{a}_l

Key proof idea: leave-one-out analysis



- Leave-one-out iterate $x^{t,(l)}$ is independent of \mathbf{a}_l
- Leave-one-out iterate $x^{t,(l)} \approx$ true iterate x^t

Key proof idea: leave-one-out analysis

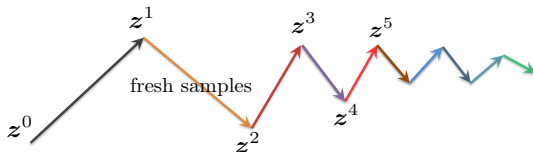


- Leave-one-out iterate $x^{t,(l)}$ is independent of a_l
- Leave-one-out iterate $x^{t,(l)} \approx$ true iterate x^t

$\implies x^t$ is nearly independent of a_l
nearly orthogonal to

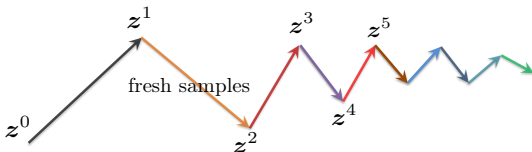
No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but helps analysis

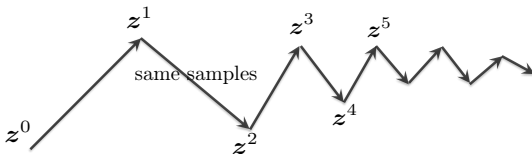


No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but helps analysis



- **This tutorial:** reuses all samples in all iterations



Other examples: low-rank matrix estimation

Low-rank matrix completion

Complete M from partial entries $M_{i,j}$, $(i,j) \in \Omega$

where (i,j) is included in Ω independently with prob. p

$$\text{find low-rank } \widehat{M} \quad \text{s.t.} \quad \mathcal{P}_{\Omega}(\widehat{M}) = \mathcal{P}_{\Omega}(M)$$

In matrix completion, strong convexity and smoothness do not hold in general

→ need to regularize the loss function by promoting **incoherent** solutions

Incoherence for matrix completion

Definition 4 (Incoherence for matrix completion)

A rank- r matrix M with eigendecomposition $M = U\Sigma U^\top$ is said to be μ -incoherent if

$$\|U\|_{2,\infty} \leq \sqrt{\frac{\mu}{n}} \|U\|_F = \sqrt{\frac{\mu r}{n}}$$

e.g.

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}}_{\text{hard } \mu=n} \quad \text{vs.} \quad \underbrace{\begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & & \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}}_{\text{easy } \mu=1}$$

Gradient descent for matrix completion

Let $M = X^* X^{*\top}$. Observe

$$Y_{i,j} = M_{i,j} + E_{i,j}, \quad (i, j) \in \Omega$$

where $(i, j) \in \Omega$ independently with prob. p , and $E_{i,j} \sim \mathcal{N}(0, \sigma^2)$ ¹

$$\text{minimize } \left\| \mathcal{P}_\Omega(\widehat{M} - \mathbf{Y}) \right\|_F^2 \quad \text{s.t.} \quad \text{rank}(\widehat{M}) \leq r$$

¹can be relaxed to sub-Gaussian noise and the asymmetric case

Gradient descent for matrix completion

Let $M = X^* X^{*\top}$. Observe

$$Y_{i,j} = M_{i,j} + E_{i,j}, \quad (i,j) \in \Omega$$

where $(i,j) \in \Omega$ independently with prob. p , and $E_{i,j} \sim \mathcal{N}(0, \sigma^2)$ ¹

$$\text{minimize } \left\| \mathcal{P}_\Omega(\widehat{M} - Y) \right\|_F^2 \quad \text{s.t.} \quad \text{rank}(\widehat{M}) \leq r$$

$$\text{minimize}_{X \in \mathbb{R}^{n \times r}} \underbrace{f(X) = \sum_{(j,k) \in \Omega} (e_j^\top X X^\top e_k - Y_{j,k})^2}_{\text{unregularized least-squares loss}}$$

¹can be relaxed to sub-Gaussian noise and the asymmetric case

Gradient descent for matrix completion

- (1) **Spectral initialization:** let $U^0 \Sigma^0 U^{0\top}$ be rank- r eigendecomposition of

$$\frac{1}{p} \mathcal{P}_\Omega(\mathbf{Y}).$$

and set $\mathbf{X}^0 = U^0 (\Sigma^0)^{1/2}$

- (2) **Gradient descent updates:**

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta_t \nabla f(\mathbf{X}^t), \quad t = 0, 1, \dots$$

Gradient descent for matrix completion

Define the optimal transform from the t th iterate \mathbf{X}^t to \mathbf{X}^* as

$$Q^t := \operatorname{argmin}_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{X}^t \mathbf{R} - \mathbf{X}^*\|_{\text{F}}$$

where $\mathcal{O}^{r \times r}$ is the set of $r \times r$ orthonormal matrices

- orthogonal Procrustes problem

Gradient descent for matrix completion

Theorem 5 (Noiseless MC, Ma, Wang, Chi, Chen '17)

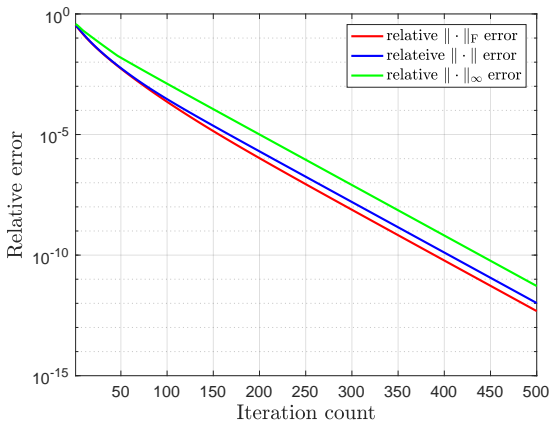
Suppose $M = \mathbf{X}^* \mathbf{X}^{*\top}$ is rank- r , incoherent and well-conditioned. Vanilla GD (with spectral initialization) achieves

- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\|_{\text{F}} \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^*\|_{\text{F}},$
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\| \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^*\|, \quad (\text{spectral})$
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\|_{2,\infty} \lesssim \rho^t \mu r \sqrt{\frac{\log n}{np}} \|\mathbf{X}^*\|_{2,\infty}, \quad (\text{incoherence})$

where $0 < \rho < 1$, if the step size $\eta \asymp 1/\sigma_{\max}$ and the sample complexity $n^2 p \gtrsim \mu^3 n r^3 \log^3 n$

- vanilla gradient descent converges linearly for matrix completion!

Numerical evidence for noiseless data



Relative error of $\mathbf{X}^t \mathbf{X}^{t\top}$ (measured by $\|\cdot\|_F$, $\|\cdot\|$, $\|\cdot\|_\infty$) vs. iteration count for MC, where $n = 1000$, $r = 10$, $p = 0.1$, and $\eta_t = 0.2$

Related theory

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left(\mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - Y_{j,k} \right)^2$$

Related theory promotes incoherence explicitly:

Related theory

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left(\mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - Y_{j,k} \right)^2$$

Related theory promotes incoherence explicitly:

- regularized loss (solve $\min_{\mathbf{X}} f(\mathbf{X}) + Q(\mathbf{X})$ instead)
 - e.g. Keshavan, Montanari, Oh '10, Sun, Luo '14, Ge, Lee, Ma '16

Related theory

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left(\mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - Y_{j,k} \right)^2$$

Related theory promotes incoherence explicitly:

- regularized loss (solve $\min_{\mathbf{X}} f(\mathbf{X}) + Q(\mathbf{X})$ instead)
 - e.g. Keshavan, Montanari, Oh '10, Sun, Luo '14, Ge, Lee, Ma '16
- projection onto set of incoherent matrices
 - e.g. Chen, Wainwright '15, Zheng, Lafferty '16

$$\mathbf{X}^{t+1} = \mathcal{P}_{\mathcal{C}} \left(\mathbf{X}^t - \eta_t \nabla f(\mathbf{X}^t) \right), \quad t = 0, 1, \dots$$

Gradient descent with spectral initialization

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) = \frac{1}{4m} \sum_{k=1}^m \left(\|\mathbf{a}_k^\top \mathbf{X}\|_2^2 - y_k \right)^2$$

Gradient descent with spectral initialization

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) = \frac{1}{4m} \sum_{k=1}^m \left(\|\mathbf{a}_k^\top \mathbf{X}\|_2^2 - y_k \right)^2$$

Theorem 6 (Quadratic sampling)

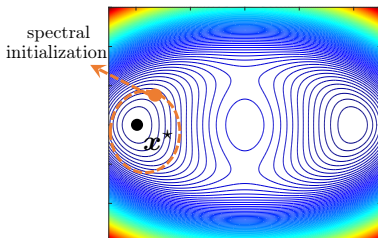
Under i.i.d. Gaussian designs $\mathbf{a}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$, GD (with spectral initialization) achieves

- $\max_l \|\mathbf{a}_l^\top (\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*)\|_2 \lesssim \sqrt{\log n} \frac{\sigma_r^2(\mathbf{X}^*)}{\|\mathbf{X}^*\|_F}$ (incoherence)
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\|_F \lesssim \left(1 - \frac{\sigma_r^2(\mathbf{X}^*)\eta}{2}\right)^t \|\mathbf{X}^*\|_F$ (linear convergence)

provided that $\eta \asymp \frac{1}{(\log n \vee r)^2 \sigma_r^2(\mathbf{X}^*)}$ and $m \gtrsim nr^4 \log n$

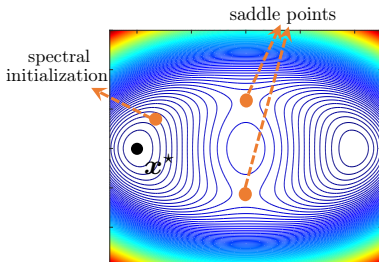
Are carefully-designed initialization or saddle-point escaping schemes necessary for fast convergence?

Initialization



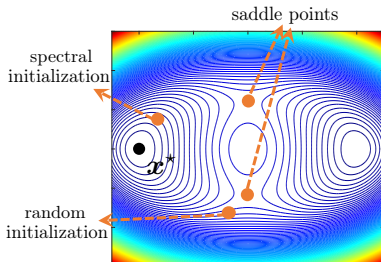
- Spectral initialization gets us reasonably close to truth

Initialization



- Spectral initialization gets us reasonably close to truth
- Cannot initialize GD from anywhere, e.g. it might get stucked at local stationary points (e.g. saddle points)

Initialization

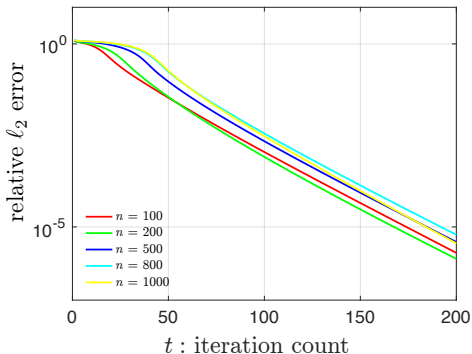


- Spectral initialization gets us reasonably close to truth
- Cannot initialize GD from anywhere, e.g. it might get stucked at local stationary points (e.g. saddle points)

Can we initialize GD randomly, which is **simpler** and **model-agnostic**?

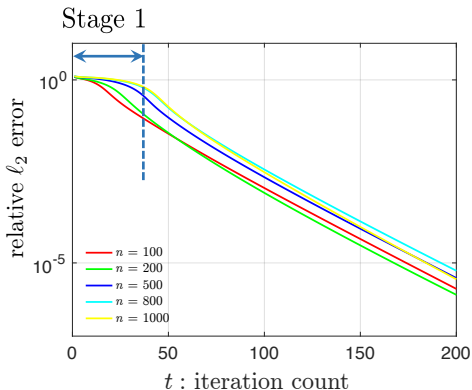
Numerical efficiency of randomly initialized GD

$$\eta_t = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



Numerical efficiency of randomly initialized GD

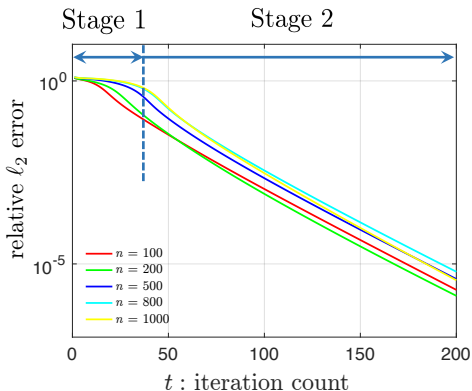
$$\eta_t = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



Randomly initialized GD enters local basin within **a few iterations**

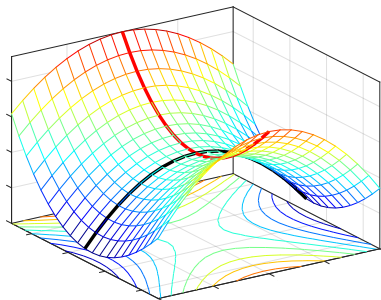
Numerical efficiency of randomly initialized GD

$$\eta_t = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



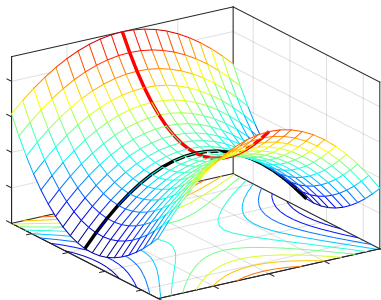
Randomly initialized GD enters local basin within a few iterations

A geometric analysis



- if $m \gtrsim n \log^3 n$, then (Sun et al. '16)
 - there is no spurious local mins
 - all saddle points are strict (i.e. associated Hessian matrices have at least one sufficiently negative eigenvalue)

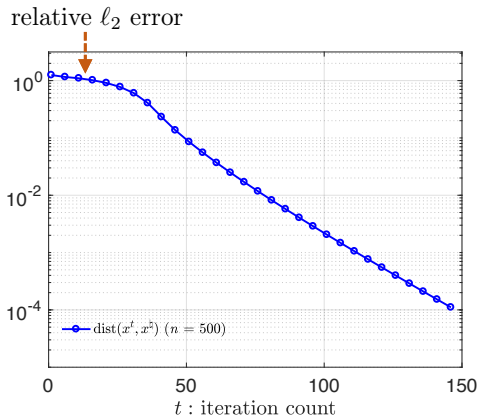
A geometric analysis



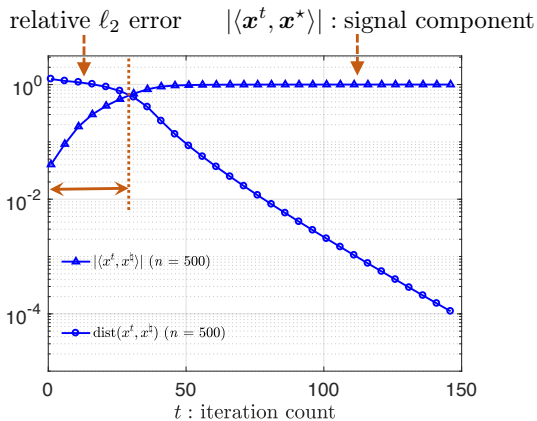
- With such benign landscape, GD with random initialization converges to global min **almost surely** (Lee et al. '16)

No convergence rate guarantees for vanilla GD!

Exponential growth of signal strength in Stage 1

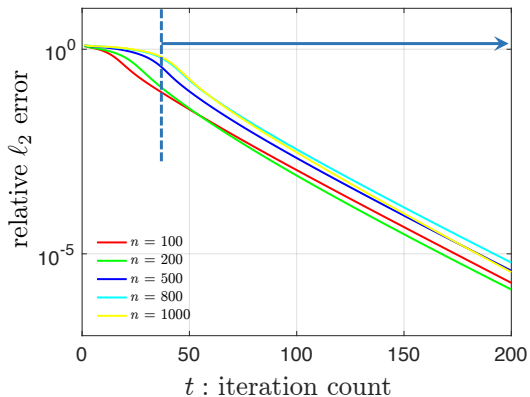


Exponential growth of signal strength in Stage 1

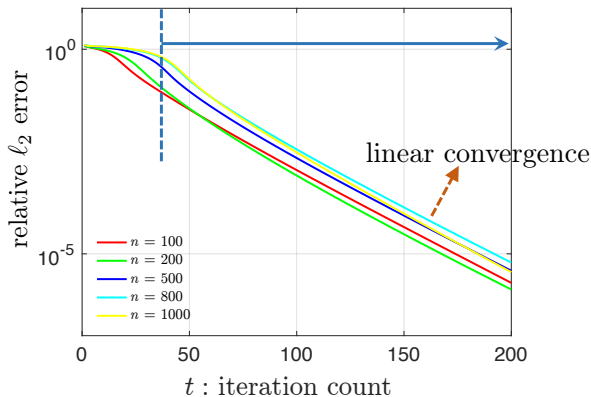


Numerically, $O(\log n)$ iterations are enough to enter local region

Linear / geometric convergence in Stage 2



Linear / geometric convergence in Stage 2



Numerically, GD converges linearly within local region

Theoretical guarantees for randomly initialized GD

These numerical findings can be formalized when $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$:

Theorem 7 (Chen, Chi, Fan, Ma '18)

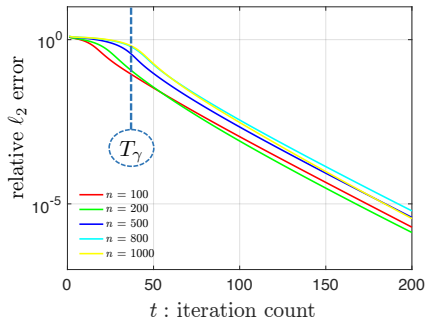
Under i.i.d. Gaussian design, GD with $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1}\mathbf{I}_n)$ achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma$$

for $T_\gamma \lesssim \log n$ and some constants $\gamma, \rho > 0$, provided that step size $\eta \asymp 1$ and sample size $m \gtrsim n \text{ poly} \log m$

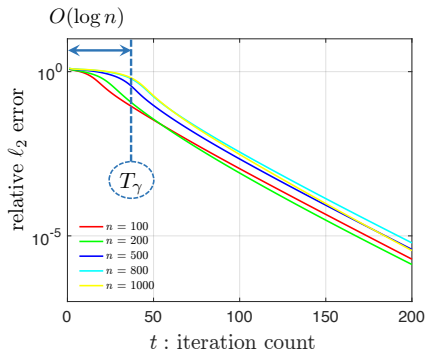
Theoretical guarantees for randomly initialized GD

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



Theoretical guarantees for randomly initialized GD

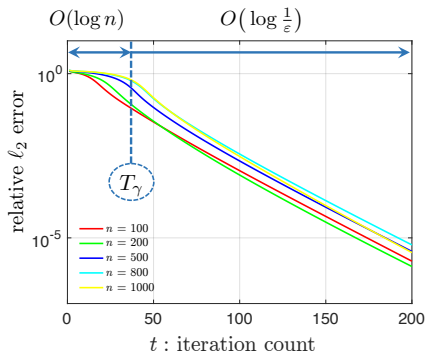
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



- *Stage 1*: takes $O(\log n)$ iterations to reach $\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma$

Theoretical guarantees for randomly initialized GD

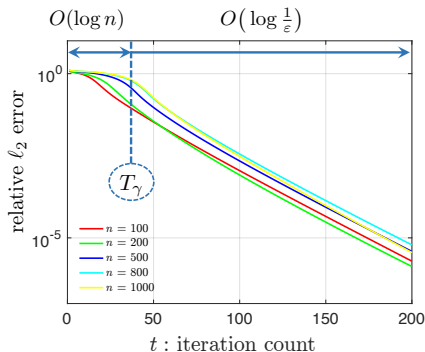
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



- *Stage 1*: takes $O(\log n)$ iterations to reach $\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma$
- *Stage 2*: linear convergence

Theoretical guarantees for randomly initialized GD

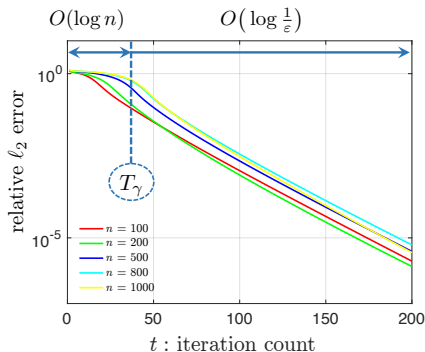
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



- *near-optimal computational cost:*
 - $O(\log n + \log \frac{1}{\epsilon})$ iterations to yield ϵ accuracy

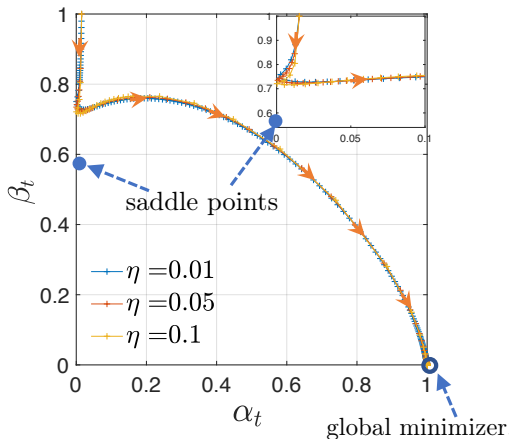
Theoretical guarantees for randomly initialized GD

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



- *near-optimal computational cost:*
 - $O(\log n + \log \frac{1}{\epsilon})$ iterations to yield ϵ accuracy
- *near-optimal sample size:* $m \gtrsim n \text{poly} \log m$

Saddle-escaping schemes?



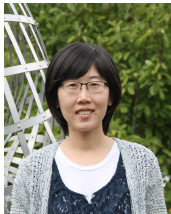
Randomly initialized GD never hits saddle points in phase retrieval!

Other saddle-escaping schemes

	iteration complexity	num of iterations needed to escape saddles	local iteration complexity
Trust-region (Sun et al. '16)	$n^7 + \log \log \frac{1}{\epsilon}$	n^7	$\log \log \frac{1}{\epsilon}$
Perturbed GD (Jin et al. '17)	$n^3 + n \log \frac{1}{\epsilon}$	n^3	$n \log \frac{1}{\epsilon}$
Perturbed accelerated GD (Jin et al. '17)	$n^{2.5} + \sqrt{n} \log \frac{1}{\epsilon}$	$n^{2.5}$	$\sqrt{n} \log \frac{1}{\epsilon}$
GD (Chen et al. '18)	$\log n + \log \frac{1}{\epsilon}$	$\log n$	$\log \frac{1}{\epsilon}$

Generic optimization theory yields highly suboptimal convergence guarantees

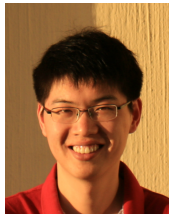
Nonconvex Optimization for High-Dimensional Signal Estimation: Spectral and Iterative Methods – Part III



Yuejie Chi
Carnegie Mellon



Yuxin Chen
Princeton



Cong Ma
UC Berkeley

EUSIPCO Tutorial, December 2020

Outline

- Part I: Introduction and Warm-Up
Why nonconvex? basic concepts and a warm-up example (PCA)
- Part II: Gradient Descent and Implicit Regularization
phase retrieval, matrix completion, random initialization
- Part III: Spectral Methods
a general recipe, ℓ_2 and ℓ_∞ guarantees, community detection
- Part IV: Robustness to Corruptions and Ill-Conditioning
median truncation, least absolute deviation, scaled gradient descent

Outline for Part III

- A motivating application: community detection
- A general recipe for spectral methods (with more applications)
- Classical spectral analysis: ℓ_2 perturbation theory
- Fine-grained analysis: ℓ_∞ perturbation theory
- A bird's-eye view of extensions

A motivating application: community detection

Community detection / graph clustering

Community structures are common in many social networks

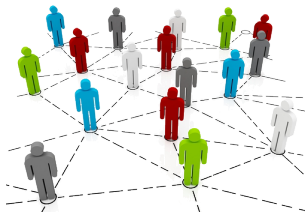


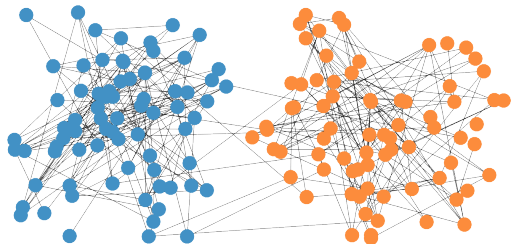
figure credit: The Future Buzz



figure credit: S. Papadopoulos

Goal: partition users into several clusters based on their friendships / similarities

A simple model: stochastic block model (SBM)

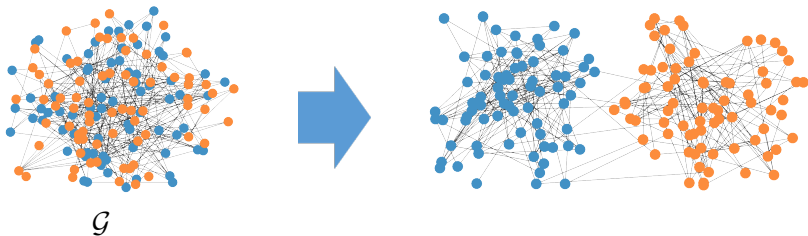


$x_i^* = 1$: 1st community

$x_i^* = -1$: 2nd community

- n nodes $\{1, \dots, n\}$
- 2 communities
- n unknown variables: $x_1^*, \dots, x_n^* \in \{1, -1\}$
 - encode community memberships

A simple model: stochastic block model (SBM)



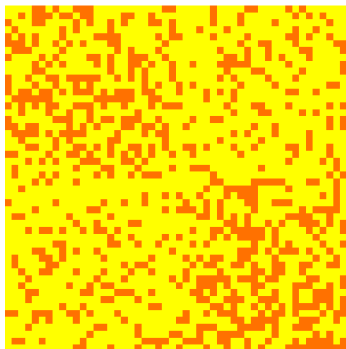
- observe a graph \mathcal{G}

$$(i, j) \in \mathcal{G} \text{ with prob. } \begin{cases} p, & \text{if } i \text{ and } j \text{ are from same community} \\ q, & \text{else} \end{cases}$$

Here, $p > q$

- **Goal:** recover community memberships of all nodes, i.e., $\{x_i^*\}$

Adjacency matrix

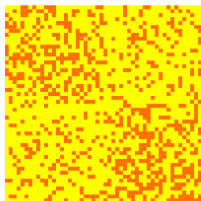


Consider the adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$ of \mathcal{G} :

$$A_{i,j} = \begin{cases} 1, & \text{if } (i, j) \in \mathcal{G} \\ 0, & \text{else} \end{cases}$$

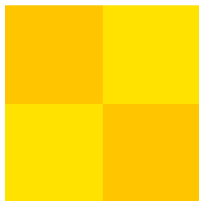
- WLOG, suppose $x_1^* = \dots = x_{n/2}^* = 1$; $x_{n/2+1}^* = \dots = x_n^* = -1$

Adjacency matrix



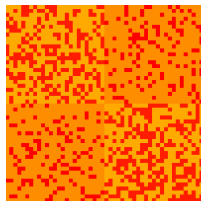
\mathbf{A}

=



$\mathbb{E}[\mathbf{A}]$
rank 2

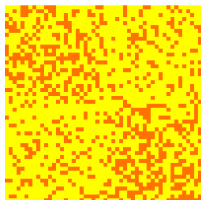
+



$\mathbf{A} - \mathbb{E}[\mathbf{A}]$

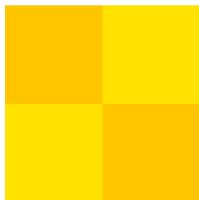
$$\mathbb{E}[\mathbf{A}] = \begin{bmatrix} p\mathbf{1}\mathbf{1}^\top & q\mathbf{1}\mathbf{1}^\top \\ q\mathbf{1}\mathbf{1}^\top & p\mathbf{1}\mathbf{1}^\top \end{bmatrix} = \underbrace{\frac{p+q}{2}\mathbf{1}\mathbf{1}^\top}_{\text{uninformative bias}} + \frac{p-q}{2} \underbrace{\begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix}}_{=\mathbf{x}^*=[x_i]_{1 \leq i \leq n}} [\mathbf{1}^\top, -\mathbf{1}^\top]$$

Spectral clustering



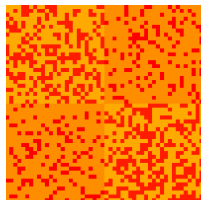
\mathbf{A}

=



$\underbrace{\mathbb{E}[\mathbf{A}]}_{\text{rank 2}}$

+



$\mathbf{A} - \mathbb{E}[\mathbf{A}]$

1. computing the leading eigenvector $\mathbf{u} = [u_i]_{1 \leq i \leq n}$ of $\mathbf{A} - \frac{p+q}{2} \mathbf{1}\mathbf{1}^\top$
2. rounding: output $x_i = \begin{cases} 1, & \text{if } u_i > 0 \\ -1, & \text{if } u_i < 0 \end{cases}$

Rationale behind spectral clustering

Recovery is reliable if $\underbrace{\mathbf{A} - \mathbb{E}[\mathbf{A}]}_{\text{perturbation}}$ is sufficiently small

- if $\mathbf{A} - \mathbb{E}[\mathbf{A}] = \mathbf{0}$, then

$$\mathbf{u} \propto \pm \begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix} \implies \text{perfect clustering}$$

A general recipe for spectral methods

Three key steps:

- identify a key matrix M^* , whose eigenvectors disclose crucial information
- construct a surrogate matrix M of M^* using data
- compute corresponding eigenvectors of M

Low-rank matrix completion



figure credit: Candès

- consider a low-rank matrix $M^* = U^* \Sigma^* V^{*\top}$
- each entry $M_{i,j}^*$ is observed independently with prob. p
- **intermediate goal:** estimate U^*, V^*

Spectral method for matrix completion

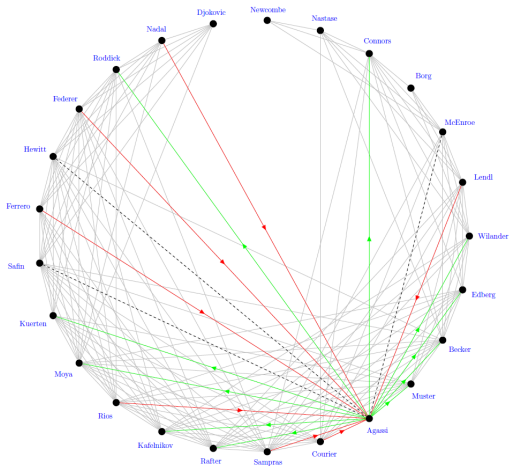
1. identify the key matrix M^*
2. construct surrogate matrix $M \in \mathbb{R}^{n \times n}$ as

$$M_{i,j} = \begin{cases} \frac{1}{p} M_{i,j}^*, & \text{if } M_{i,j}^* \text{ is observed} \\ 0, & \text{else} \end{cases}$$

- **rationale for rescaling:** ensures $\mathbb{E}[M] = M^*$

3. compute the rank- r SVD $U\Sigma V^\top$ of M , and return (U, Σ, V)

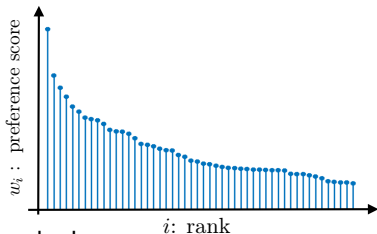
Ranking from pairwise comparisons



pairwise comparisons for ranking tennis players

figure credit: Bozóki, Csató, Temesi

Bradley-Terry-Luce (logistic) model



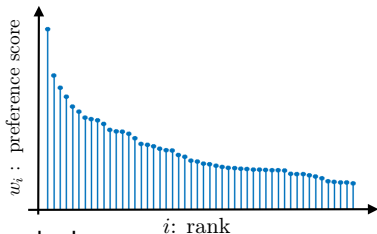
- n items to be ranked
- assign a latent score $\{w_i^*\}_{1 \leq i \leq n}$ to each item, so that

$$\text{item } i \succ \text{item } j \quad \text{if} \quad w_i^* > w_j^*$$

- each pair of items (i, j) is compared independently

$$\mathbb{P} \{ \text{item } j \text{ beats item } i \} = \frac{w_j^*}{w_i^* + w_j^*}$$

Bradley-Terry-Luce (logistic) model



- n items to be ranked
- assign a latent score $\{w_i^*\}_{1 \leq i \leq n}$ to each item, so that

$$\text{item } i \succ \text{item } j \quad \text{if} \quad w_i^* > w_j^*$$

- each pair of items (i, j) is compared independently

$$y_{i,j} \stackrel{\text{ind.}}{=} \begin{cases} 1, & \text{with prob. } \frac{w_j^*}{w_i^* + w_j^*} \\ 0, & \text{else} \end{cases}$$

- **intermediate goal:** estimate score vector w^* (up to scaling)

Spectral ranking

1. identify key matrix P^* —probability transition matrix

$$P_{i,j}^* = \begin{cases} \frac{1}{n} \cdot \frac{w_j^*}{w_i^* + w_j^*}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}^*, & \text{if } i = j \end{cases}$$

Rationale:

- P^* obeys

$$w_i^* P_{i,j}^* = w_j^* P_{j,i}^* \quad (\text{detailed balance})$$

- Thus, the stationary distribution π^* of P^* obeys

$$\pi^* = \frac{1}{\sum_l w_l^*} \mathbf{w}^* \quad (\text{reveals true scores})$$

Spectral ranking

2. construct a surrogate matrix \mathbf{P} obeying

$$P_{i,j} = \begin{cases} \frac{1}{n}y_{i,j}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}, & \text{if } i = j \end{cases}$$

3. return leading left eigenvector $\boldsymbol{\pi}$ of \mathbf{P} as score estimate

— closely related to PageRank

Spectral ranking

2. construct a surrogate matrix P obeying

$$P_{i,j} = \begin{cases} \frac{1}{n}y_{i,j}, & \text{if } i \neq j \\ 1 - \sum_{l:l \neq i} P_{i,l}, & \text{if } i = j \end{cases}$$

3. return leading left eigenvector π of P as score estimate

— closely related to PageRank

Key: stability of eigenspace against perturbation $M - M^*$?

**Classical spectral analysis:
 l_2 perturbation theory**

Setup and notation

Consider two symmetric matrices M^* and its perturbed version

$$M = M^* + E \in \mathbb{R}^{n \times n}$$

with eigendecompositions

$$M^* = \sum_{i=1}^n \lambda_i^* \mathbf{u}_i^* \mathbf{u}_i^{*\top} = \begin{bmatrix} U^* & U_{\perp}^* \end{bmatrix} \begin{bmatrix} \Lambda^* & \mathbf{0} \\ \mathbf{0} & \Lambda_{\perp}^* \end{bmatrix} \begin{bmatrix} U^{*\top} \\ U_{\perp}^{*\top} \end{bmatrix};$$
$$M = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top} = \begin{bmatrix} U & U_{\perp} \end{bmatrix} \begin{bmatrix} \Lambda & \mathbf{0} \\ \mathbf{0} & \Lambda_{\perp} \end{bmatrix} \begin{bmatrix} U^{\top} \\ U_{\perp}^{\top} \end{bmatrix}$$

Eigenspace perturbation theory

Main focus: how does the perturbation E affect the distance between U and U^* ?

Question: how to define distance between two subspaces?

- $\|U - U^*\|_F$ and $\|U - U^*\|$ are not appropriate, since they fall short of accounting for global orthonormal transformation

\forall orthonormal $R \in \mathbb{R}^{r \times r}$, U and UR represent same subspace

Distance between two subspaces

One solution: taking best rotation into consideration

$$\text{dist}(\mathbf{U}, \mathbf{U}^*) := \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{UR} - \mathbf{U}^*\|;$$

$$\text{dist}_F(\mathbf{U}, \mathbf{U}^*) := \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{UR} - \mathbf{U}^*\|_F$$

Davis-Kahan $\sin \Theta$ Theorem: a simple case



Chandler Davis



William Kahan

Theorem 1

Suppose $M^* \succeq 0$ and has rank r . If $\|E\| < \lambda_r(M^*)$, then

$$\text{dist}(U, U^*) \leq \frac{\|EU^*\|}{\lambda_r(M^*) - \|E\|} \leq \frac{\|E\|}{\lambda_r(M^*) - \|E\|}$$

- depends on $\underbrace{\text{smallest non-zero eigenvalue of } M^*}_{\text{eigengap between } \lambda_r(M^*) \text{ and } \lambda_{r+1}(M^*)}$ and perturbation size

Back to stochastic block model

$$\text{Let } \mathbf{M}^* := \mathbb{E}[\mathbf{A}] - \frac{p+q}{2} \mathbf{1}\mathbf{1}^\top = \frac{p-q}{2} \begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{1}^\top & -\mathbf{1}^\top \end{bmatrix},$$
$$\mathbf{M} := \mathbf{A} - \frac{p+q}{2} \mathbf{1}\mathbf{1}^\top, \text{ and } \mathbf{u}^* := \frac{1}{\sqrt{n}} \begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix}$$

Then the Davis-Kahan sin Θ Theorem yields

$$\text{dist}(\mathbf{u}, \mathbf{u}^*) \leq \frac{\|\mathbf{M} - \mathbf{M}^*\|}{\lambda_1(\mathbf{M}^*) - \|\mathbf{M} - \mathbf{M}^*\|} = \frac{\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|}{\frac{(p-q)n}{2} - \|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|} \quad (1)$$

as long as $\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\| < \lambda_1(\mathbf{M}^*) = \frac{(p-q)n}{2}$

Bounding $\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|$

Matrix concentration inequalities tell us that

Lemma 2

Consider SBM with $p > q$ and $p \gtrsim \frac{\log n}{n}$. Then with high prob.

$$\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\| \lesssim \sqrt{np} \quad (2)$$

Statistical accuracy of spectral clustering

Substitute (2) into (1) to reach

$$\text{dist}(\mathbf{u}, \mathbf{u}^*) \leq \frac{\|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|}{\frac{(p-q)n}{2} - \|\mathbf{A} - \mathbb{E}[\mathbf{A}]\|} \lesssim \frac{\sqrt{np}}{(p-q)n}$$

provided that $(p-q)n \gg \sqrt{np}$

Thus, under condition $\frac{p-q}{\sqrt{p}} \gg \sqrt{\frac{1}{n}}$, with high prob. one has

$$\text{dist}(\mathbf{u}, \mathbf{u}^*) \ll 1 \quad \implies \quad \text{nearly perfect clustering}$$

Statistical accuracy of spectral clustering

$$\frac{p - q}{\sqrt{p}} \gg \sqrt{\frac{1}{n}} \implies \text{nearly perfect clustering}$$

- **dense regime:** if $p \asymp q \asymp 1$, then this condition reads

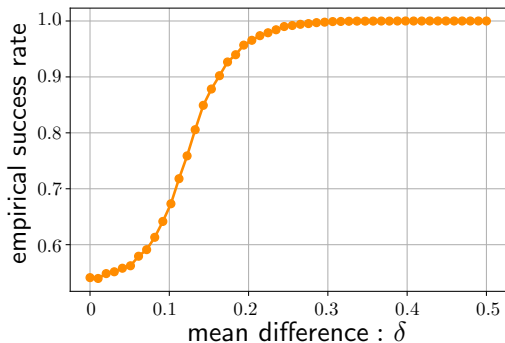
$$p - q \gg \sqrt{\frac{1}{n}}$$

- **“sparse” regime:** if $p = \frac{a \log n}{n}$ and $q = \frac{b \log n}{n}$ for $a, b \asymp 1$, then

$$a - b \gg \sqrt{a \log n}$$

This condition is information-theoretically optimal (up to log factor)
— Mossel, Neeman, Sly '15, Abbe '18

Empirical performance of spectral clustering



ℓ_2 perturbation theory alone cannot explain exact recovery guarantees

— call for fine-grained analysis

Reverse engineering

Spectral clustering uses signs of u to cluster nodes

Reverse engineering

Spectral clustering uses signs of \mathbf{u} to cluster nodes



It achieves exact recovery iff $u_i u_i^* > 0$ for all i

Reverse engineering

Spectral clustering uses signs of \mathbf{u} to cluster nodes



It achieves exact recovery iff $u_i u_i^* > 0$ for all i



A sufficient condition is* $\|\mathbf{u} - \mathbf{u}^*\|_\infty < 1/\sqrt{n}$

Reverse engineering

Spectral clustering uses signs of \mathbf{u} to cluster nodes



It achieves exact recovery iff $u_i u_i^* > 0$ for all i



A sufficient condition is* $\|\mathbf{u} - \mathbf{u}^*\|_\infty < 1/\sqrt{n}$



Need ℓ_∞ perturbation theory

Fine-grained analysis:
 l_∞ **perturbation theory**

Setup and notation (rank-1 case)

Groundtruth: consider a rank-1 psd matrix $M^* = \lambda^* \mathbf{u}^* \mathbf{u}^{*\top} \in \mathbb{R}^{n \times n}$

Incoherence: define

$$\mu := n \|\mathbf{u}^*\|_\infty^2 \quad (1 \leq \mu \leq n)$$

Observations:

$$M = M^* + \mathbf{E} \in \mathbb{R}^{n \times n}$$

with \mathbf{E} a symmetric noise matrix

Noise assumptions

The entries in the lower triangular part of $\mathbf{E} = [E_{i,j}]_{1 \leq i, j \leq n}$ are independently generated obeying

$$\mathbb{E}[E_{i,j}] = 0, \quad \mathbb{E}[E_{i,j}^2] \leq \sigma^2, \quad |E_{i,j}| \leq B, \quad \text{for all } i \geq j$$

Further, assume that

$$c_b := \frac{B}{\sigma \sqrt{n/(\mu \log n)}} = O(1)$$

ℓ_∞ perturbation theory

Theorem 3

With high prob, there exists $z \in \{1, -1\}$ such that

$$\|z\mathbf{u} - \mathbf{u}^*\|_\infty \lesssim \frac{\sigma\sqrt{\mu} + \sigma\sqrt{\log n}}{\lambda^*}, \quad (3a)$$

$$\|z\mathbf{u} - \frac{1}{\lambda^*}\mathbf{M}\mathbf{u}^*\|_\infty \lesssim \frac{\sigma\sqrt{\mu}}{\lambda^*} + \frac{\sigma^2\sqrt{n\log n} + \sigma B\sqrt{\mu\log^3 n}}{(\lambda^*)^2} \quad (3b)$$

provided that $\sigma\sqrt{n\log n} \leq c_\sigma\lambda^*$ for some sufficiently small constant $c_\sigma > 0$.

Key message:

- when $\mu \lesssim \sqrt{\log n}$, (3a) is $\sqrt{n/\log n}$ smaller than ℓ_2 bound

$$\|z\mathbf{u} - \mathbf{u}^*\|_2 \lesssim \sigma\sqrt{n}/\lambda^*$$

Back to stochastic block model

$$\text{Recall } M^* := \mathbb{E}[A] - \frac{p+q}{2}\mathbf{1}\mathbf{1}^\top = \frac{p-q}{2} \begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{1}^\top & -\mathbf{1}^\top \end{bmatrix},$$
$$M := A - \frac{p+q}{2}\mathbf{1}\mathbf{1}^\top, \text{ and } \mathbf{u}^* := \frac{1}{\sqrt{n}} \begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix}$$

These imply

$$\lambda^* = \frac{n(p-q)}{2}$$

$$\mu = 1$$

$$B = 1$$

$$\sigma^2 \leq \max\{p, q\} = p$$

Invoke ℓ_∞ perturbation theory

ℓ_∞ perturbation theory (3b) yields

$$\begin{aligned}\|z\lambda^* \mathbf{u} - \mathbf{M}\mathbf{u}^*\|_\infty &\lesssim \sigma + \frac{\sigma^2 \sqrt{n \log n}}{\lambda^*} + \frac{\sigma B \log^{3/2} n}{\lambda^*} \\ &\leq C \left(\sqrt{p} + \frac{p \sqrt{\log n}}{\sqrt{n}(p-q)} + \frac{\sqrt{p} \log^{3/2} n}{n(p-q)} \right) =: \Delta\end{aligned}$$

for some constant $C > 0$

it boils down to controlling the entrywise behavior of $\mathbf{M}\mathbf{u}^*$

Bounding entries in $M\mathbf{u}^*$

Again concentration inequalities tell us that

Lemma 4

Suppose that

$$(\sqrt{p} - \sqrt{q})^2 \geq (1 + \varepsilon) \frac{2 \log n}{n} \quad (4)$$

for some quantity $\varepsilon > 0$. Let $\varepsilon_0 := \frac{\varepsilon \log n}{\sqrt{n} \log \frac{p(1-q)}{q(1-p)}} - \frac{1}{\sqrt{n}}$. Then with probability exceeding $1 - n^{-\varepsilon/2}$, one has

$$M_{l,\cdot} \mathbf{u}^* \geq \varepsilon_0 \text{ for all } l \leq \frac{n}{2} \quad \text{and} \quad M_{l,\cdot} \mathbf{u}^* \leq -\varepsilon_0 \text{ for all } l > \frac{n}{2}.$$

Key message: entries in $M\mathbf{u}^*$ are bounded away from 0 with correct sign

Completing the picture

On one hand

$$\mathbf{M}_l \cdot \mathbf{u}^* \geq \varepsilon_0 \text{ for all } l \leq \frac{n}{2} \quad \text{and} \quad \mathbf{M}_l \cdot \mathbf{u}^* \leq -\varepsilon_0 \text{ for all } l > \frac{n}{2}.$$

On the other hand

$$\|z\lambda^* \mathbf{u} - \mathbf{M}\mathbf{u}^*\|_\infty \leq \Delta$$

In sum, if one can show

$$\varepsilon_0 > \Delta$$

then it follows that

$$z u_l u_l^* > 0 \quad \text{for all } 1 \leq l \leq n \quad \implies \quad \text{exact recovery}$$

Exact recovery of SBM

Theorem 5

Fix any constant $\varepsilon > 0$. Suppose $p = \frac{\alpha \log n}{n}$ and $q = \frac{\beta \log n}{n}$ for some sufficiently large constants $\alpha > \beta > 0$. In addition, assume that

$$(\sqrt{p} - \sqrt{q})^2 \geq 2(1 + \varepsilon) \frac{\log n}{n}. \quad (5)$$

With probability $1 - o(1)$, spectral clustering achieves exact recovery.

This condition is information-theoretically optimal

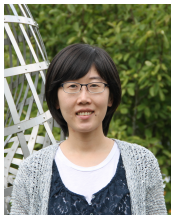
A bird's-eye view of extensions

- Davis-Kahan for general symmetric matrices (not necessarily PSD)
- Wedin's theorem on singular subspace perturbation theory
- Eigenvector perturbation for probability transition matrices
- General $\ell_{2,\infty}$ eigenspace and singular space perturbation

Advertisement: "Spectral Methods for Data Science: A Statistical Perspective",
Y. Chen, Y. Chi, J. Fan and C. Ma, 2020

Foundations and Trends[®] in Machine Learning
Spectral Methods for Data Science
A Statistical Perspective

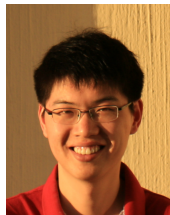
Nonconvex Optimization for High-Dimensional Signal Estimation: Spectral and Iterative Methods – Part IV



Yuejie Chi
Carnegie Mellon



Yuxin Chen
Princeton



Cong Ma
UC Berkeley

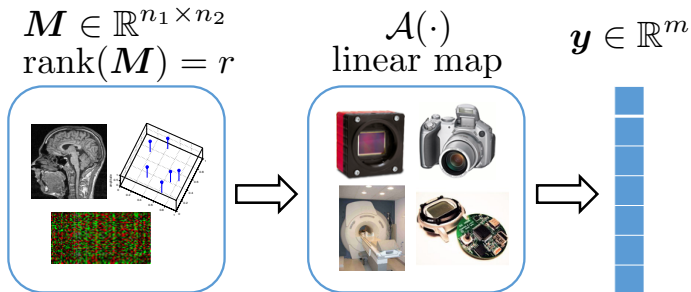
EUSIPCO Tutorial, December 2020

Outline

- Part I: Introduction and Warm-Up
Why nonconvex? basic concepts and a warm-up example (PCA)
- Part II: Gradient Descent and Implicit Regularization
phase retrieval, matrix completion, random initialization
- Part III: Spectral Methods
a general recipe, ℓ_2 and ℓ_∞ guarantees, community detection
- Part IV: Robustness to Corruptions and Ill-Conditioning
median truncation, least absolute deviation, scaled gradient descent

Robustness to ill-conditioning?

A factorization approach to low-rank matrix sensing



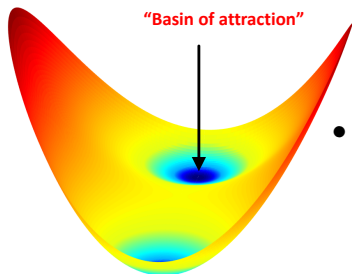
$$y = \mathcal{A}(M) + \text{noise}$$

find $X \in \mathbb{R}^{n_1 \times r}$, $Y = \mathbb{R}^{n_2 \times r}$, such that $y \approx \mathcal{A}(XY^T)$

Prior art: GD with balancing regularization

$$\min_{\mathbf{X}, \mathbf{Y}} f_{\text{reg}}(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_2^2 + \frac{1}{8} \left\| \mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y} \right\|_F^2$$

- **Spectral initialization:** find an initial point in the “basin of attraction”.



$$(\mathbf{X}_0, \mathbf{Y}_0) \leftarrow \text{SVD}_r(\mathcal{A}^*(\mathbf{y}))$$

- **Gradient iterations:**

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f_{\text{reg}}(\mathbf{X}_t, \mathbf{Y}_t)$$

$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta \nabla_{\mathbf{Y}} f_{\text{reg}}(\mathbf{X}_t, \mathbf{Y}_t)$$

for $t = 0, 1, \dots$

Prior theory for vanilla GD

Theorem 1 (Tu et al., ICML 2016)

Suppose $M = X_* Y_*^\top$ is rank- r and has a condition number $\kappa = \sigma_{\max}(M)/\sigma_{\min}(M)$. For low-rank matrix sensing with i.i.d. Gaussian design, vanilla GD (with spectral initialization) achieves

$$\|X_t Y_t^\top - M\|_F \leq \varepsilon \cdot \sigma_{\min}(M)$$

- **Computational:** within $O(\kappa \log \frac{1}{\varepsilon})$ iterations;
- **Statistical:** as long as the sample complexity satisfies

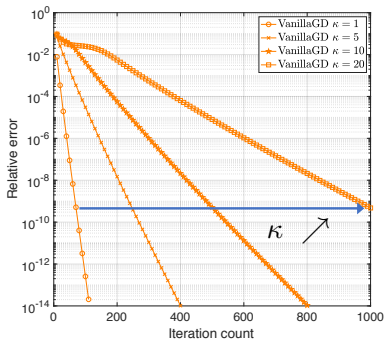
$$m \gtrsim (n_1 + n_2)r^2\kappa^2.$$

Similar results hold for many low-rank problems.

(Netrapalli et al. '13, Candès, Li, Soltanolkotabi '14, Sun and Luo '15, Chen and Wainwright '15, Zheng and Lafferty '15, Ma et al. '17,)

Convergence slows down for ill-conditioned matrices

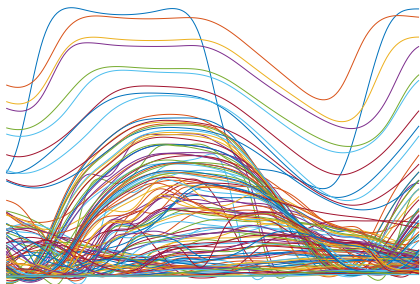
$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathcal{P}_{\Omega}(\mathbf{X}\mathbf{Y}^{\top} - \mathbf{M}) \right\|_{\text{F}}^2$$



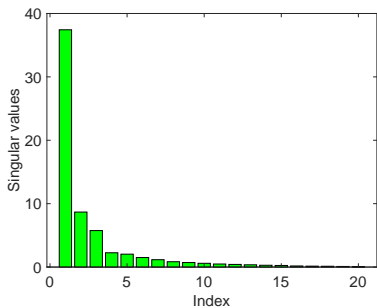
Vanilla GD converges in $O(\kappa \log \frac{1}{\epsilon})$ iterations.

— *Can we provably accelerate the convergence to $O(\log \frac{1}{\epsilon})$?*

Condition number can be large

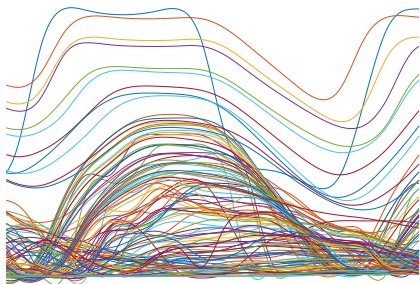


chlorine concentration levels
120 junctions, 180 time slots

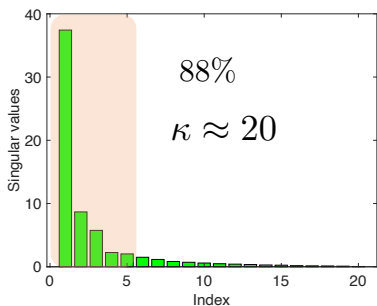


power-law spectrum

Condition number can be large

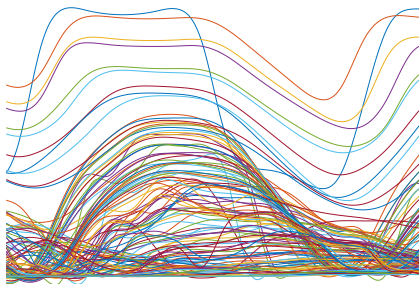


chlorine concentration levels
120 junctions, 180 time slots

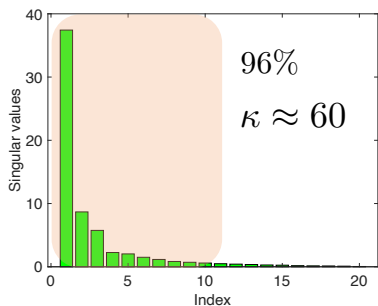


rank-5 approximation

Condition number can be large



chlorine concentration levels
120 junctions, 180 time slots



rank-10 approximation

A new algorithm: scaled gradient descent

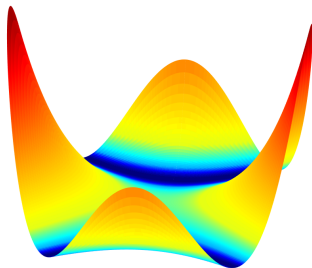
$$f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top)\|_2^2$$

- **Spectral initialization:** find an initial point in the “basin of attraction”.
- **Scaled gradient iterations:**

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

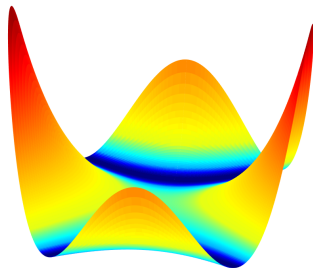
$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta \nabla_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{X}_t^\top \mathbf{X}_t)^{-1}}_{\text{preconditioner}}$$

for $t = 0, 1, \dots$



A new algorithm: scaled gradient descent

$$f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top)\|_2^2$$



- **Spectral initialization:** find an initial point in the “basin of attraction”.
- **Scaled gradient iterations:**

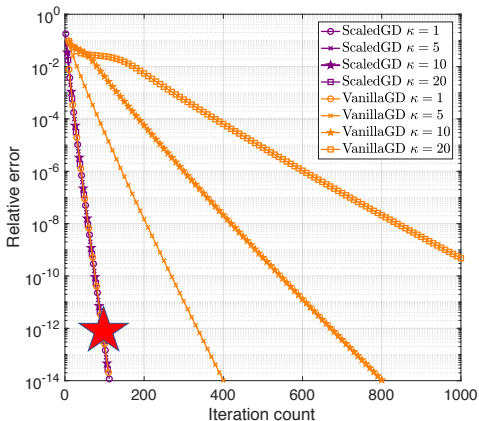
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta \nabla_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{X}_t^\top \mathbf{X}_t)^{-1}}_{\text{preconditioner}}$$

for $t = 0, 1, \dots$

ScaledGD is a *preconditioned* gradient method
without balancing regularization!

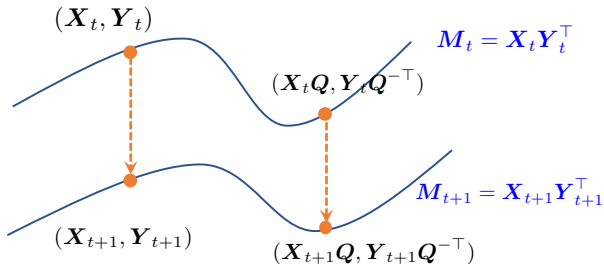
ScaledGD for low-rank matrix completion



Huge computational saving: ScaledGD converges in an κ -independent manner with a minimal overhead!

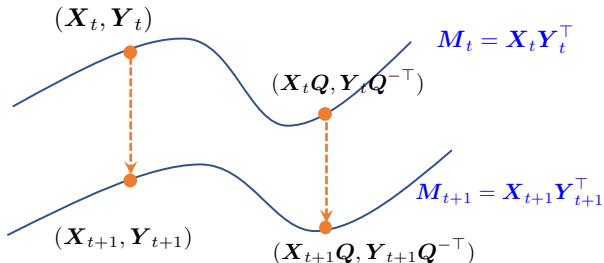
A closer look at ScaledGD

Invariance to invertible transforms: (Tanner and Wei, '16; Mishra '16)



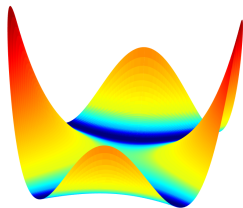
A closer look at ScaledGD

Invariance to invertible transforms: (Tanner and Wei, '16; Mishra '16)



New distance metric as Lyapunov function:

$$\text{dist}^2 \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}, \begin{bmatrix} \mathbf{X}_* \\ \mathbf{Y}_* \end{bmatrix} \right) = \inf_{\mathbf{Q} \in \text{GL}(r)} \left\| (\mathbf{X} \mathbf{Q} - \mathbf{X}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{Y} \mathbf{Q}^{-\top} - \mathbf{Y}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{F}}^2.$$



Theoretical guarantees of ScaledGD

Theorem 2 (Tong, Ma and Chi, 2020)

For low-rank matrix sensing with i.i.d. Gaussian design, ScaledGD with spectral initialization achieves

$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

- **Computational:** within $O(\log \frac{1}{\varepsilon})$ iterations;
- **Statistical:** the sample complexity satisfies

$$m \gtrsim nr^2 \kappa^2.$$

Theoretical guarantees of ScaledGD

Theorem 2 (Tong, Ma and Chi, 2020)

For low-rank matrix sensing with i.i.d. Gaussian design, ScaledGD with spectral initialization achieves

$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

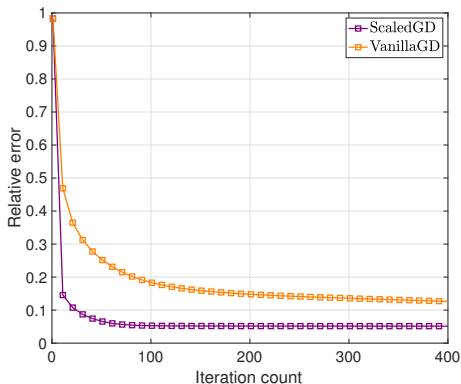
- **Computational:** within $O(\log \frac{1}{\varepsilon})$ iterations;
- **Statistical:** the sample complexity satisfies

$$m \gtrsim nr^2 \kappa^2.$$

Compared with Tu et. al.: ScaledGD provably accelerates vanilla GD at the same sample complexity!

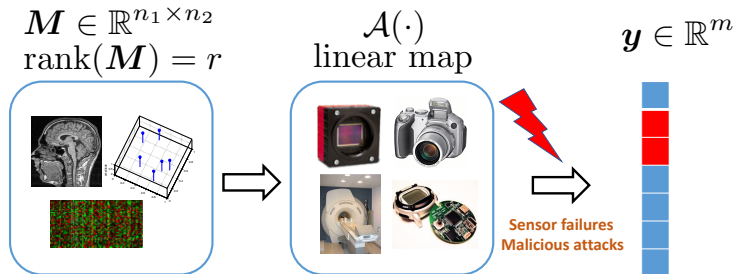
Stability: a numerical result

For the chlorine concentration levels dataset, ScaledGD converges faster than vanilla GD in a small number of iterations.



Robustness to outliers and corruptions?

Outlier-corrupted low-rank matrix sensing



$$y = \mathcal{A}(M) + \underbrace{s}_{\text{outliers}}, \quad \mathcal{A}(M) = \{\langle A_i, M \rangle\}_{i=1}^m$$

Arbitrary but sparse outliers: $\|s\|_0 \leq \alpha \cdot m$, where $0 \leq \alpha < 1$ is fraction of outliers.

Existing approaches fail

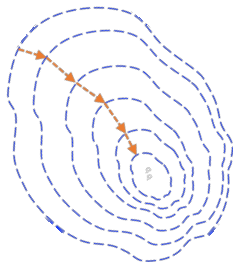
- **Spectral initialization would fail:** $\mathbf{X}_0 \leftarrow$ top- r SVD of

$$\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{A}_i$$

- **Gradient iterations would fail:**

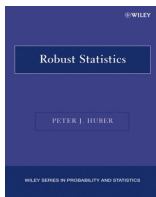
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \frac{\eta}{m} \sum_{i=1}^m \nabla l_i(y_i; \mathbf{X}_t)$$

for $t = 0, 1, \dots$



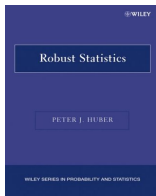
Even a single outlier can fail the algorithm!

Median-truncated gradient descent



Key idea: “median-truncation” —
discard samples *adaptively* based on how
large sample gradients / values deviate
from median

Median-truncated gradient descent

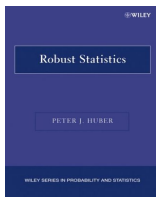


Key idea: “median-truncation” — discard samples *adaptively* based on how large sample gradients / values deviate from median

- **Robustify spectral initialization:** $X_0 \leftarrow$ top- r SVD of

$$Y = \frac{1}{m} \sum_{i: |y_i| \lesssim \text{median}(|y_i|)} y_i A_i$$

Median-truncated gradient descent



Key idea: “median-truncation” — discard samples *adaptively* based on how large sample gradients / values deviate from median

- **Robustify spectral initialization:** $\mathbf{X}_0 \leftarrow$ top- r SVD of

$$\mathbf{Y} = \frac{1}{m} \sum_{i: |y_i| \lesssim \text{median}(|y_i|)} y_i \mathbf{A}_i$$

- **Robustify gradient descent:**

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \frac{\eta}{m} \sum_{i: |r_t^i| \lesssim \text{median}(|r_t^i|)} \nabla \ell_i(y_i; \mathbf{X}_t), \quad t = 0, 1, \dots$$

where $r_t^i := \left| y_i - \langle \mathbf{A}_i, \mathbf{X}_t \mathbf{X}_t^\top \rangle \right|$ is the size of the gradient.

Theoretical guarantees

Theorem 3 (Li, Chi, Zhang, and Liang, IMIAI 2020)

For low-rank matrix sensing with i.i.d. Gaussian design, median-truncated GD (with robust spectral initialization) achieves

$$\|\mathbf{X}_t \mathbf{X}_t^\top - \mathbf{M}\|_F \leq \varepsilon \cdot \sigma_{\min}(\mathbf{M}),$$

- **Computational:** within $O(\kappa \log \frac{1}{\varepsilon})$ iterations;
- **Statistical:** the sample complexity satisfies

$$m \gtrsim nr^2 \text{poly}(\kappa, \log n);$$

- **Robustness:** and the fraction of outliers

$$\alpha \lesssim 1/\sqrt{r}.$$

Theoretical guarantees

Theorem 3 (Li, Chi, Zhang, and Liang, IMIAI 2020)

For low-rank matrix sensing with i.i.d. Gaussian design, median-truncated GD (with robust spectral initialization) achieves

$$\|\mathbf{X}_t \mathbf{X}_t^\top - \mathbf{M}\|_F \leq \varepsilon \cdot \sigma_{\min}(\mathbf{M}),$$

- **Computational:** within $O(\kappa \log \frac{1}{\varepsilon})$ iterations;
- **Statistical:** the sample complexity satisfies

$$m \gtrsim nr^2 \text{poly}(\kappa, \log n);$$

- **Robustness:** and the fraction of outliers

$$\alpha \lesssim 1/\sqrt{r}.$$

Median-truncated GD adds robustness to GD *obliviously*.

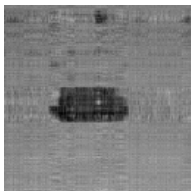
Numerical example

Low-rank matrix sensing:

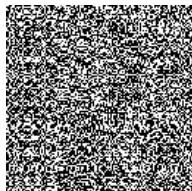
$$y_i = \langle \mathbf{A}_i, \mathbf{M} \rangle + s_i, \quad i = 1, \dots, m$$



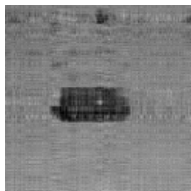
Ground truth



GD
no outliers



GD
1% outliers



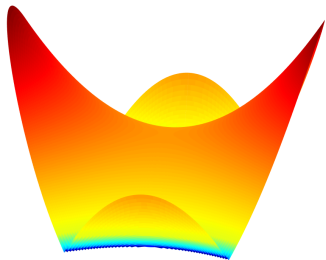
median-TGD
1% outliers

Median-truncated GD achieves similar performance as if performing GD on the clean data.

Dealing with outliers: subgradient methods

Least absolute deviation (LAD): (Charisopoulos et.al.'19; Li et al'18)

$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_1$$



Subgradient iterations:

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta_t \partial_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t)$$

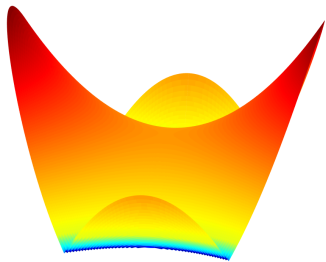
$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta_t \partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t)$$

where η_t is set as Polyak's or geometric decaying stepsize.

Dealing with outliers: scaled subgradient methods

Least absolute deviation (LAD):

$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_1$$



Scaled subgradient iterations:

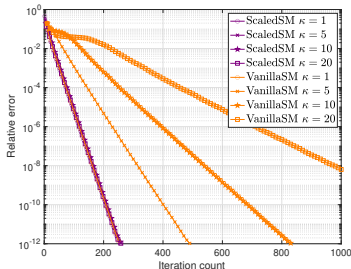
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta_t \partial_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta_t \partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{X}_t^\top \mathbf{X}_t)^{-1}}_{\text{preconditioner}}$$

where η_t is set as Polyak's or geometric decaying stepsize.

Performance guarantees

	matrix sensing	quadratic sensing
Subgradient Method (Charisopoulos et al, '19)	$\frac{\kappa}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$	$\frac{r\kappa}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$
ScaledSM (Tong, Ma, Chi, '20)	$\frac{1}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$	$\frac{r}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$



Robustness to both ill-conditioning and adversarial corruptions!

Demixing sparse and low-rank matrices

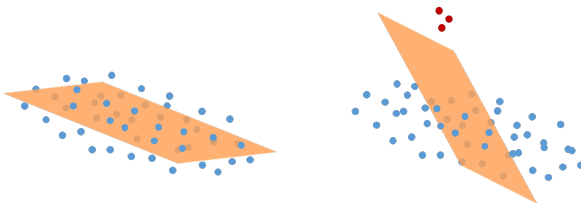
Suppose we are given a matrix

$$M = \underbrace{L}_{\text{low-rank}} + \underbrace{S}_{\text{sparse}} \in \mathbb{R}^{n \times n}$$

Question: can we hope to recover both L and S from M ?

Applications

- Robust PCA



- Video surveillance: separation of background and foreground



Nonconvex approach

- rank(L) $\leq r$; if we write the SVD of $L = U\Sigma V^\top$, set

$$X^* = U_L \Sigma^{1/2}; \quad Y^* = V \Sigma^{1/2}$$

- non-zero entries of S are “spread out” (no more than s fraction of non-zeros per row/column), but otherwise arbitrary

$$\mathcal{S}_s = \{S \in \mathbb{R}^{n \times n} : \|S_{i,:}\|_0 \leq s \cdot n; \|S_{:,j}\|_0 \leq s \cdot n\}$$

$$\underset{X, Y, S \in \mathcal{S}_s}{\text{minimize}} \quad F(X, Y, S) := \underbrace{\|M - XY^\top - S\|_F^2}_{\text{least-squares loss}}$$

where $X, Y \in \mathbb{R}^{n \times r}$.

Gradient descent and hard thresholding

$$\text{minimize}_{\mathbf{X}, \mathbf{Y}, \mathbf{S} \in \mathcal{S}_s} F(\mathbf{X}, \mathbf{Y}, \mathbf{S})$$

- **Spectral initialization:** Set $\mathbf{S}^0 = \mathcal{H}_{\gamma_s}(\mathbf{M})$. Let $U^0 \Sigma^0 V^{0\top}$ be rank- r SVD of $\mathbf{M}^0 := \mathcal{P}_{\Omega}(\mathbf{M} - \mathbf{S}^0)$; set $\mathbf{X}^0 = U^0 (\Sigma^0)^{1/2}$ and $\mathbf{Y}^0 = V^0 (\Sigma^0)^{1/2}$

Gradient descent and hard thresholding

$$\text{minimize}_{\mathbf{X}, \mathbf{Y}, \mathbf{S} \in \mathcal{S}_s} F(\mathbf{X}, \mathbf{Y}, \mathbf{S})$$

- **Spectral initialization:** Set $\mathbf{S}^0 = \mathcal{H}_{\gamma_s}(\mathbf{M})$. Let $U^0 \Sigma^0 V^{0\top}$ be rank- r SVD of $\mathbf{M}^0 := \mathcal{P}_{\Omega}(\mathbf{M} - \mathbf{S}^0)$; set $\mathbf{X}^0 = U^0 (\Sigma^0)^{1/2}$ and $\mathbf{Y}^0 = V^0 (\Sigma^0)^{1/2}$
- for $t = 0, 1, 2, \dots$
 - **Hard thresholding:** $\mathbf{S}^{t+1} = \mathcal{H}_{\gamma_s}(\mathbf{M} - \mathbf{X}^t \mathbf{Y}^{t\top})$
 - **Scaled gradient updates:**
$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta \nabla_{\mathbf{X}} F(\mathbf{X}^t, \mathbf{Y}^t, \mathbf{S}^{t+1}) (\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}$$
$$\mathbf{Y}^{t+1} = \mathbf{Y}^t - \eta \nabla_{\mathbf{Y}} F(\mathbf{X}^t, \mathbf{Y}^t, \mathbf{S}^{t+1}) (\mathbf{X}_t^\top \mathbf{X}_t)^{-1}$$

Efficient nonconvex recovery

Theorem 4 (Nonconvex RPCA, Tian, Ma, Chi '20)

Set $\gamma = 2$ and $0.1 \leq \eta \leq 2/3$. Suppose that

$$s \lesssim \frac{1}{\mu r^{3/2} \kappa}$$

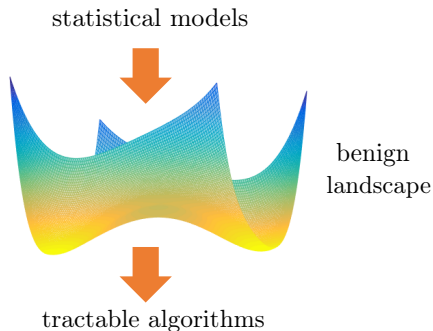
Then GD+HT satisfies

$$\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{L}\|_F \lesssim (1 - 0.6\eta)^t \sigma_{\min}(\mathbf{L})$$

- $O(\log \frac{1}{\varepsilon})$ iterations to reach ε accuracy
- for adversarial outliers, optimal fraction is $s = O(1/\mu r)$;
Theorem 4 is suboptimal by a factor of $\kappa\sqrt{r}$
- Improves over GD (Yi et al '16) which requires
 $s \lesssim \frac{1}{\max\{\mu r^{3/2} \kappa^{3/2}, \mu r \kappa^2\}} \lesssim \frac{1}{\mu r^{3/2} \kappa}$ and $O(\kappa \log \frac{1}{\varepsilon})$ iterations;

Concluding remarks

Statistical thinking + Optimization efficiency



When data are generated by certain statistical models, problems are often much nicer than worst-case instances

A growing list of “benign” nonconvex problems

- phase retrieval
- matrix sensing
- matrix completion
- blind deconvolution / self-calibration
- dictionary learning
- tensor decomposition / completion
- robust PCA
- mixed linear regression
- learning one-layer neural networks
- ...

Open problems

- characterize generic landscape properties that enable fast convergence of gradient methods from random initialization
- relax the stringent assumptions on the statistical models underlying the data
- develop robust and scalable nonconvex methods that can handle distributed data with strong statistical guarantees
- identify new classes of nonconvex problems that admit efficient optimization procedures
- ...

Advertisement: overview and monographs

- “Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview”, Y. Chi, Y. M. Lu and Y. Chen.

IEEE TRANSACTIONS ON SIGNAL PROCESSING, VOL. 67, NO. 20, OCTOBER 15, 2019

5239

Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview

Yuejie Chi[✉], Yue M. Lu[✉], and Yuxin Chen[✉]

(Overview Article)

- “Spectral Methods for Data Science: A Statistical Perspective”, Y. Chen, Y. Chi, J. Fan and C. Ma.

Foundations and Trends[®] in Machine Learning Spectral Methods for Data Science A Statistical Perspective

An incomplete list of reference

- [1] Dr. Ju Sun's webpage: "<http://sunju.org/research/nonconvex/>".
- [2] "*Harnessing Structures in Big Data via Guaranteed Low-Rank Matrix Estimation*," Y. Chen, and Y. Chi, *IEEE Signal Processing Magazine*, 2018.
- [3] "*Phase retrieval via Wirtinger flow: Theory and algorithms*," E. Candès, X. Li, M. Soltanolkotabi, *IEEE Transactions on Information Theory*, 2015.
- [4] "*Solving random quadratic systems of equations is nearly as easy as solving linear systems*," Y. Chen, E. Candes, *Communications on Pure and Applied Mathematics*, 2017.
- [5] "*Provable non-convex phase retrieval with outliers: Median truncated Wirtinger flow*," H. Zhang, Y. Chi, and Y. Liang, ICML 2016.
- [6] "*Implicit Regularization in Nonconvex Statistical Estimation: Gradient Descent Converges Linearly for Phase Retrieval, Matrix Completion and Blind Deconvolution*," C. Ma, K. Wang, Y. Chi and Y. Chen, *Foundations of Computational Mathematics*, 2019.

An incomplete list of reference

- [7] "*Gradient Descent with Random Initialization: Fast Global Convergence for Nonconvex Phase Retrieval*," Y. Chen, Y. Chi, J. Fan, C. Ma, *Mathematical Programming*, 2018.
- [8] "*Solving systems of random quadratic equations via truncated amplitude flow*," G. Wang, G. Giannakis, and Y. Eldar, *IEEE Transactions on Information Theory*, 2017.
- [9] "*Matrix completion from a few entries*," R. Keshavan, A. Montanari, and S. Oh, *IEEE Transactions on Information Theory*, 2010.
- [10] "*Guaranteed matrix completion via non-convex factorization*," R. Sun, T. Luo, *IEEE Transactions on Information Theory*, 2016.
- [11] "*Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees*," Y. Chen and M. Wainwright, *arXiv preprint arXiv:1509.03025*, 2015.
- [12] "*Fast Algorithms for Robust PCA via Gradient Descent*," X. Yi, D. Park, Y. Chen, and C. Caramanis, *NIPS*, 2016.

An incomplete list of reference

- [13] “*Matrix completion has no spurious local minimum*,” R. Ge, J. Lee, and T. Ma, *NIPS*, 2016.
- [14] “*No Spurious Local Minima in Nonconvex Low Rank Problems: A Unified Geometric Analysis*,” R. Ge, C. Jin, and Y. Zheng, *ICML*, 2017.
- [15] “*Symmetry, Saddle Points, and Global Optimization Landscape of Nonconvex Matrix Factorization*,” X. Li et al., *IEEE Transactions on Information Theory*, 2019.
- [16] “*Kaczmarz Method for Solving Quadratic Equations*,” Y. Chi and Y. M. Lu, *IEEE Signal Processing Letters*, vol. 23, no. 9, pp. 1183-1187, 2016.
- [17] “*A Geometric Analysis of Phase Retrieval*,” S. Ju, Q. Qu, and J. Wright, to appear, *Foundations of Computational Mathematics*, 2016.
- [18] “*Gradient descent converges to minimizers*,” J. Lee, M. Simchowitz, M. Jordan, B. Recht, *Conference on Learning Theory*, 2016.
- [19] “*Phase retrieval using alternating minimization*,” P. Netrapalli, P. Jain, and S. Sanghavi, *NIPS*, 2013.

An incomplete list of reference

- [20] "*How to escape saddle points efficiently*," C. Jin, R. Ge, P. Netrapalli, S. Kakade, M. Jordan, arXiv:1703.00887, 2017.
- [21] "*Complete dictionary recovery over the sphere*," J. Sun, Q. Qu, J. Wright, *IEEE Transactions on Information Theory*, 2017.
- [22] "*A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization*," S. Burer, and R. Monteiro, *Mathematical Programming*, 2003.
- [23] "*Memory-efficient Kernel PCA via Partial Matrix Sampling and Nonconvex Optimization: a Model-free Analysis of Local Minima*," J. Chen, X. Li, *Journal of Machine Learning Research*, 2019.
- [24] "*Rapid, robust, and reliable blind deconvolution via nonconvex optimization*," X. Li, S. Ling, T. Strohmer, K. Wei, *Applied and computational harmonic analysis*, 2019.
- [25] "*Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming*," E. Candes, T. Strohmer, V. Voroninski, *Communications on Pure and Applied Mathematics*, 2012.

An incomplete list of reference

- [26] “*Exact matrix completion via convex optimization*,” E. Candès, B. Recht, *Foundations of Computational mathematics*, 2009.
- [27] “*Low-rank solutions of linear matrix equations via Procrustes flow*,” S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, B. Recht, arXiv:1507.03566, 2015.
- [28] “*Global optimality of local search for low rank matrix recovery*,” S. Bhojanapalli, B. Neyshabur, and N. Srebro, NIPS, 2016.
- [29] “*Optimal rates of convergence for noisy sparse phase retrieval via thresholded Wirtinger flow*,” T. Cai, X. Li, Z. Ma, *The Annals of Statistics*, 2016.
- [30] “*The landscape of empirical risk for non-convex losses*,” S. Mei, Y. Bai, and A. Montanari, arXiv:1607.06534, 2016.
- [31] “*Non-convex robust PCA*,” P. Netrapalli, U. Niranjan, S. Sanghavi, A. Anandkumar, and P. Jain, NIPS, 2014.
- [32] “*Non-square matrix sensing without spurious local minima via the Burer-Monteiro approach*,” D. Park, A. Kyrillidis, C. Caramanis, and S. Sanghavi, arXiv:1609.03240, 2016.

An incomplete list of reference

- [33] “*From symmetry to geometry: Tractable nonconvex problems*”, Y. Zhang, Q. Qu, J. Wright, arXiv preprint arXiv:2007.06753, 2020.
- [34] “*A Nonconvex Approach for Phase Retrieval: Reshaped Wirtinger Flow and Incremental Algorithms*,” H. Zhang, Y. Zhou, Y. Liang and Y. Chi, *Journal of Machine Learning Research*, 2017.
- [35] “*Nonconvex Matrix Factorization from Rank-One Measurements*,” Y. Li, C. Ma, Y. Chen and Y. Chi, arXiv:1802.06286, 2018.
- [36] “*Manifold Gradient Descent Solves Multi-Channel Sparse Blind Deconvolution Provably and Efficiently*”, L. Shi and Y. Chi, arXiv preprint arXiv:1911.11167, 2019.
- [37] “*Median-Truncated Gradient Descent: A Robust and Scalable Nonconvex Approach for Signal Estimation*”, Y. Chi, Y. Li, H. Zhang, and Y. Liang, *Compressed Sensing and Its Applications*, Springer, Birkhauser, 2019.
- [38] “*Nonconvex Low-Rank Matrix Recovery with Arbitrary Outliers via Median-Truncated Gradient Descent*”, Y. Li, Y. Chi, H. Zhang, and Y. Liang, *Information and Inference*, 2020.

An incomplete list of reference

- [39] “*Nonconvex robust low-rank matrix recovery*”, Li et al., SIAM Journal on Optimization, 2020.
- [40] “*Median-Truncated Nonconvex Approach for Phase Retrieval with Outliers*”, H. Zhang, Y. Chi and Y. Liang, IEEE Trans. on Information Theory, 2018.
- [41] “*Low-rank matrix recovery with composite optimization: good conditioning and rapid convergence*”, V. Charisopoulos et al., arXiv:1904.10020, 2019.
- [42] “*Accelerating Ill-Conditioned Low-Rank Matrix Estimation via Scaled Gradient Descent*”, T. Tong, C. Ma, and Y. Chi, arXiv:2005.08898, 2020.
- [43] “*Low-Rank Matrix Recovery with Scaled Subgradient Methods: Fast and Robust Convergence Without the Condition Number*”, T. Tong, C. Ma, and Y. Chi, arXiv:2010.13364, 2020.