

## Introduction



Cong Ma

University of Chicago, Autumn 2021

# A Message from Harvard Data Science Review

---

## 4. Balance of Statistical and Computational Efficiencies

When we have limited data, the emphasis on statistical efficiency to make the best use of the available data has naturally become an important focus of statistics research. We do not think statistical efficiency will become irrelevant in the big data era; often inference is made locally and the relevant data that are available to infer around a specific subpopulation remain limited. On the other hand, useful statistical modeling and data analysis must take into account constraints on data storage, communication across sites, and the quality of numerical approximations in the computation. An 'optimally efficient' statistical approach is far from optimal in practice if it relies on optimization of a highly nonconvex and nonsmooth objective function, for instance. The need to work with streaming data for real-time actions also calls for a balanced approach. This is where statisticians and computer scientists, as well as experts from related domains (e.g., operation research, mathematics, and subject-matter science) can work together to address efficiency in a holistic way.

*Challenges and Opportunities in Statistics and Data Science: Ten Research Areas*

*by Xuming He and Xihong Lin*

# Key points

---

- Statistical efficiency is still relevant in big data era

# Key points

---

- Statistical efficiency is still relevant in big data era
  - *big data vs big parameters (high-dimensional statistics)*

# Key points

---

- Statistical efficiency is still relevant in big data era
  - *big data vs big parameters (high-dimensional statistics)*
- Computational efficiency cannot be ignored

# Key points

---

- Statistical efficiency is still relevant in big data era
  - *big data vs big parameters (high-dimensional statistics)*
- Computational efficiency cannot be ignored
  - *due to limited computation/memory constraints*

# Key points

---

- Statistical efficiency is still relevant in big data era
  - *big data vs big parameters (high-dimensional statistics)*
- Computational efficiency cannot be ignored
  - *due to limited computation/memory constraints*
- “A ... procedure is far from optimal in practice if it relies on optimization of a highly nonconvex and nonsmooth objective function”

# Key points

---

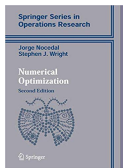
- Statistical efficiency is still relevant in big data era
  - *big data vs big parameters (high-dimensional statistics)*
- Computational efficiency cannot be ignored
  - *due to limited computation/memory constraints*
- “A ... procedure is far from optimal in practice if it relies on optimization of a highly nonconvex and nonsmooth objective function”
  - *hmm...maybe...nonconvexity maybe our friend*



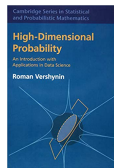
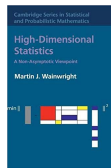
# Main theme of this course

---

By blending **statistical** and **computational** theory, we can extract useful information from big data more efficiently



(nonconvex) optimization



(high-dimensional) statistics

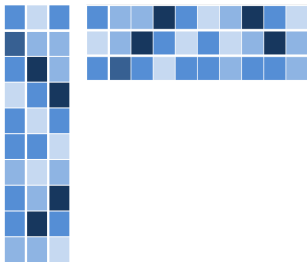
# Outline

---

- A motivating example: low-rank matrix completion
- Topics covered in this course
- Course logistics

**A motivating example:  
low-rank matrix completion**

# Noisy low-rank matrix completion

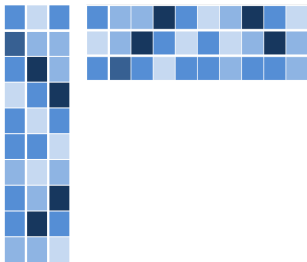


unknown rank- $r$  matrix  $\Theta^* \in \mathbb{R}^{d \times d}$



sampling set  $\Omega$

# Noisy low-rank matrix completion



unknown rank- $r$  matrix  $\Theta^* \in \mathbb{R}^{d \times d}$













sampling set  $\Omega$

observations:  $Y_{i,j} = \Theta_{i,j}^* + \text{noise}, \quad (i,j) \in \Omega$

goal: estimate  $\Theta^*$





# Motivation 1: recommendation systems

							...
	★★★★☆	?	★★★★☆	?	?	?	...
	?	★★★★☆	?	?	★★★★☆	?	...
	?	?	?	★★★★☆	★★★★☆	?	...
	?	★★★★☆	★★★★☆	?	?	★★★★☆	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

- Netflix challenge: Netflix provides highly incomplete ratings from nearly 0.5 million users & 20k movies
- How to predict unseen user ratings for movies?

# In general, we cannot infer missing ratings

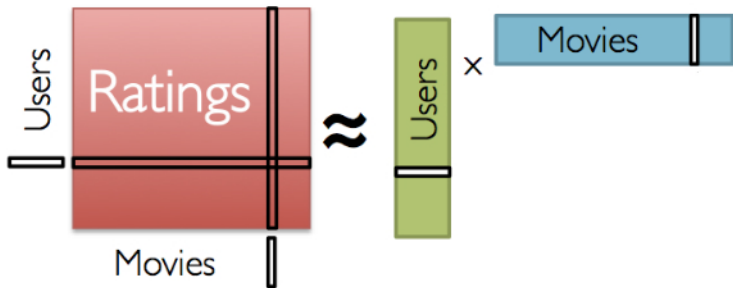
$$\begin{bmatrix}
 \checkmark & ? & ? & ? & \checkmark & ? \\
 ? & ? & \checkmark & \checkmark & ? & ? \\
 \checkmark & ? & ? & \checkmark & ? & ? \\
 ? & ? & \checkmark & ? & ? & \checkmark \\
 \checkmark & ? & ? & ? & ? & ? \\
 ? & \checkmark & ? & ? & \checkmark & ? \\
 ? & ? & \checkmark & \checkmark & ? & ?
 \end{bmatrix}$$

							...
	★★★★★	?	★★★★★	?	?	?	...
	?	★★★★★	?	?	★★★★★	?	...
	?	?	?	★★★★★	★★★★★	?	...
	?	★★★★★	★★★★★	?	?	★★★★★	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Underdetermined system (more unknowns than equations)

... unless rating matrix has some structure

---

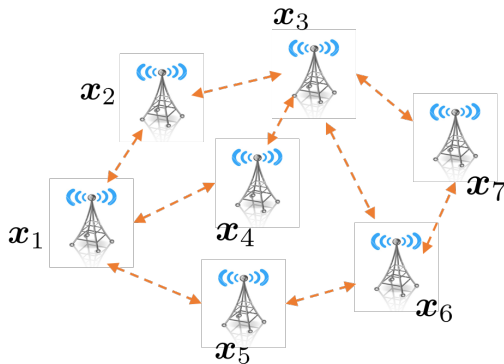


low-rank approximation  $\rightarrow$  a few factors explain most of data



## Motivation 2: sensor localization

---



- Observe partial pairwise distances
- Goal: infer distance between every pair of nodes

## Motivation 2: sensor localization

Introduce location matrix

$$\mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_d^\top & - \end{bmatrix} \in \mathbb{R}^{d \times 3}$$

then distance matrix  $\mathbf{D} = [D_{i,j}]_{1 \leq i,j \leq d}$  can be written as

$$\mathbf{D} = \underbrace{\begin{bmatrix} \|\mathbf{x}_1\|_2^2 \\ \vdots \\ \|\mathbf{x}_d\|_2^2 \end{bmatrix}}_{\text{rank 1}} \mathbf{1}^\top + \underbrace{\mathbf{1} \cdot [\|\mathbf{x}_1\|_2^2, \dots, \|\mathbf{x}_d\|_2^2]}_{\text{rank 1}} - \underbrace{2\mathbf{X}\mathbf{X}^\top}_{\text{rank 3}}$$

low rank

$\text{rank}(\mathbf{D}) \ll d \longrightarrow$  low-rank matrix completion

# Least-squares estimator

---

$$\begin{array}{ll} \underset{\Theta \in \mathbb{R}^{d \times d}}{\text{minimize}} & f(\Theta) = \sum_{(i,j) \in \Omega} (\Theta_{i,j} - Y_{i,j})^2 \\ \text{subject to} & \text{rank}(\Theta) = r \end{array}$$

# Least-squares estimator

---

$$\begin{array}{ll} \underset{\Theta \in \mathbb{R}^{d \times d}}{\text{minimize}} & f(\Theta) = \sum_{(i,j) \in \Omega} (\Theta_{i,j} - Y_{i,j})^2 \\ \text{subject to} & \text{rank}(\Theta) = r \end{array}$$

— *This is also MLE when noise follows Gaussian*

# Least-squares estimator

---

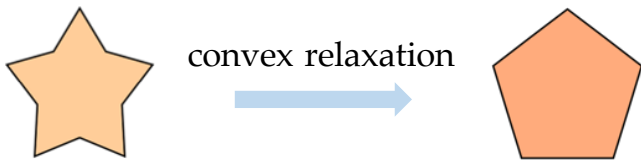
$$\begin{array}{ll} \underset{\Theta \in \mathbb{R}^{d \times d}}{\text{minimize}} & f(\Theta) = \sum_{(i,j) \in \Omega} (\Theta_{i,j} - Y_{i,j})^2 \\ \text{subject to} & \text{rank}(\Theta) = r \end{array}$$

— *This is also MLE when noise follows Gaussian*

**Challenge:** nonconvexity  $\implies$  computational hardness

# Popular workaround: convex relaxation

---



Relax nonconvex problems into convex ones by finding convex surrogates

# Convex relaxation for matrix completion

---

Replace rank constraint by nuclear norm constraint

$$\begin{aligned} & \underset{\Theta \in \mathbb{R}^{d \times d}}{\text{minimize}} && f(\Theta) = \sum_{(i,j) \in \Omega} (\Theta_{i,j} - Y_{i,j})^2 \\ & \text{subject to} && \cancel{\text{rank}(\Theta) \leq r} \quad \|\Theta\|_* \leq t \\ & && \text{--- } \|\Theta\|_* = \sum_{i=1}^d \sigma_i(\Theta) \end{aligned}$$

# Convex relaxation for matrix completion

Replace rank constraint by nuclear norm constraint

$$\begin{aligned} & \underset{\Theta \in \mathbb{R}^{d \times d}}{\text{minimize}} && f(\Theta) = \sum_{(i,j) \in \Omega} (\Theta_{i,j} - Y_{i,j})^2 \\ & \text{subject to} && \cancel{\text{rank}(\Theta) \leq r} \quad \|\Theta\|_* \leq t \\ & && \quad \quad \quad - \quad \|\Theta\|_* = \sum_{i=1}^d \sigma_i(\Theta) \end{aligned}$$

**convex relaxation (regularized version):**

$$\underset{\Theta \in \mathbb{R}^{d \times d}}{\text{minimize}} \quad \sum_{(i,j) \in \Omega} (\Theta_{i,j} - Y_{i,j})^2 + \lambda \|\Theta\|_*$$



# Convex relaxation: pros and cons

---

convex relaxation (regularized version):

$$\underset{\Theta \in \mathbb{R}^{d \times d}}{\text{minimize}} \quad \sum_{(i,j) \in \Omega} (\Theta_{i,j} - Y_{i,j})^2 + \lambda \|\Theta\|_*$$

**Pro:** often achieve statistical optimality

# Convex relaxation: pros and cons

---

convex relaxation (regularized version):

$$\underset{\Theta \in \mathbb{R}^{d \times d}}{\text{minimize}} \quad \sum_{(i,j) \in \Omega} (\Theta_{i,j} - Y_{i,j})^2 + \lambda \|\Theta\|_*$$

**Pro:** often achieve statistical optimality

**Issue:** expensive in computation/memory

*Can we solve matrix completion with lower  
computational cost?*

# Spectral methods

---

- Assumption: each entry is observed indep. with probability  $p$

# Spectral methods

---

- Assumption: each entry is observed indep. with probability  $p$
- Key observation: let

$$\hat{Y}_{i,j} = \begin{cases} \frac{1}{p}Y_{i,j}, & \text{if } (i,j) \text{ is observed,} \\ 0, & \text{otherwise} \end{cases}$$

we have  $\mathbb{E}[\hat{\mathbf{Y}}] = \mathbf{\Theta}^*$

# Spectral methods

---

- Assumption: each entry is observed indep. with probability  $p$
- Key observation: let

$$\hat{Y}_{i,j} = \begin{cases} \frac{1}{p} Y_{i,j}, & \text{if } (i,j) \text{ is observed,} \\ 0, & \text{otherwise} \end{cases}$$

we have  $\mathbb{E}[\hat{\mathbf{Y}}] = \Theta^*$

**spectral method:**

deploy best rank- $r$  approximation to  $\hat{\mathbf{Y}}$  as estimator of  $\Theta^*$

# Spectral methods

---

- Assumption: each entry is observed indep. with probability  $p$
- Key observation: let

$$\hat{Y}_{i,j} = \begin{cases} \frac{1}{p} Y_{i,j}, & \text{if } (i,j) \text{ is observed,} \\ 0, & \text{otherwise} \end{cases}$$

we have  $\mathbb{E}[\hat{\mathbf{Y}}] = \Theta^*$

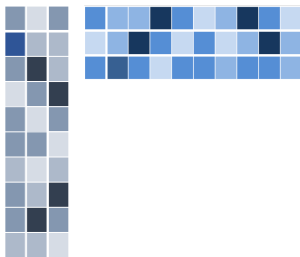
## spectral method:

deploy best rank- $r$  approximation to  $\hat{\mathbf{Y}}$  as estimator of  $\Theta^*$

— simple, but sometimes statistically inefficient

# Nonconvex optimization

Represent low-rank matrix by  $LR^\top$  with  $\underbrace{L, R}_{\text{low-rank factors}} \in \mathbb{R}^{d \times r}$



$$\underset{L, R \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(L, R) = \sum_{(i,j) \in \Omega} \left[ (LR^\top)_{i,j} - Y_{i,j} \right]^2$$



## Two-stage algorithm

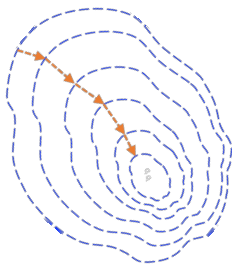
---

$$\underset{\mathbf{L}, \mathbf{R} \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(\mathbf{L}, \mathbf{R}) = \sum_{(i,j) \in \Omega} \left[ (\mathbf{L} \mathbf{R}^\top)_{i,j} - Y_{i,j} \right]^2$$

# Two-stage algorithm

---

$$\underset{\mathbf{L}, \mathbf{R} \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(\mathbf{L}, \mathbf{R}) = \sum_{(i,j) \in \Omega} \left[ (\mathbf{L}\mathbf{R}^\top)_{i,j} - Y_{i,j} \right]^2$$



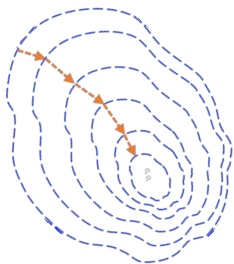
- **spectral initialization:**  $(\mathbf{L}^0, \mathbf{R}^0)$   
— top singular vectors of  $\hat{\mathbf{Y}}$
- **gradient descent:** for  $t = 0, 1, \dots$

$$\mathbf{L}^{t+1} = \mathbf{L}^t - \eta_t \nabla_{\mathbf{L}} f(\mathbf{L}^t, \mathbf{R}^t)$$

$$\mathbf{R}^{t+1} = \mathbf{R}^t - \eta_t \nabla_{\mathbf{R}} f(\mathbf{L}^t, \mathbf{R}^t)$$

# Two-stage algorithm

$$\underset{\mathbf{L}, \mathbf{R} \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(\mathbf{L}, \mathbf{R}) = \sum_{(i,j) \in \Omega} \left[ (\mathbf{L}\mathbf{R}^\top)_{i,j} - Y_{i,j} \right]^2$$



- **spectral initialization:**  $(\mathbf{L}^0, \mathbf{R}^0)$   
— top singular vectors of  $\hat{\mathbf{Y}}$
- **gradient descent:** for  $t = 0, 1, \dots$

$$\mathbf{L}^{t+1} = \mathbf{L}^t - \eta_t \nabla_{\mathbf{L}} f(\mathbf{L}^t, \mathbf{R}^t)$$

$$\mathbf{R}^{t+1} = \mathbf{R}^t - \eta_t \nabla_{\mathbf{R}} f(\mathbf{L}^t, \mathbf{R}^t)$$

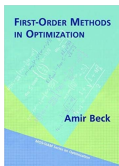
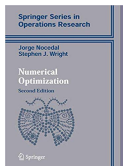
nonconvex estimator achieves optimal estimation error

— Ma, Wang, Chi, Chen '17

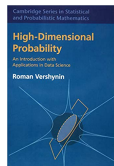
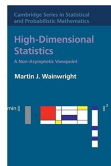
# Main theme of this course

---

By blending **statistical** and **computational** theory, we can extract useful information from big data more efficiently



(nonconvex) optimization



(high-dimensional) statistics

# Tentative topics

---

- Spectral methods
  - Classic  $\ell_2$  matrix perturbation theory
  - Matrix concentration inequalities
  - Applications of spectral methods ( $\ell_2$  theory)
  - $\ell_\infty$  matrix perturbation theory
  - Applications of spectral methods ( $\ell_\infty$  theory)
- Nonconvex optimization
  - Basic optimization theory
  - Generic local analysis for regularized gradient descent (GD)
  - Refined local analysis for vanilla GD
  - Global landscape analysis
  - Gradient descent with random initialization
- Convex relaxation
  - Compressed sensing and sparse recovery
  - Phase transition and convex geometry
  - Low-rank matrix recovery
  - Robust principal component analysis

# Logistics

# Why you **should not** take this course

---

# Why you **should not** take this course

---

- There will be quite a few THEOREMS and PROOFS ...



# Why you **should not** take this course

---

- There will be quite a few THEOREMS and PROOFS ...
- Nonrigorous/heuristic arguments from time to time

# Why you **should** consider taking this course

---

# Why you **should** consider taking this course

---

- There will be quite a few THEOREMS and PROOFS ...

# Why you **should** consider taking this course

---

- There will be quite a few THEOREMS and PROOFS ...
  - promote deeper understanding of scientific results

# Why you **should** consider taking this course

---

- There will be quite a few THEOREMS and PROOFS ...
  - promote deeper understanding of scientific results
- Nonrigorous/heuristic arguments from time to time

# Why you **should** consider taking this course

---

- There will be quite a few THEOREMS and PROOFS ...
  - promote deeper understanding of scientific results
- Nonrigorous/heuristic arguments from time to time
  - “nonrigorous” but grounded in rigorous theory
  - help develop intuition

# Prerequisites

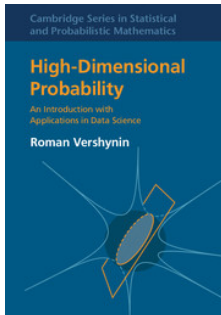
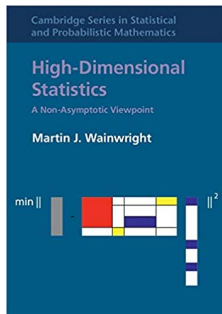
---

- linear algebra
- probability theory
- a programming language (e.g., Matlab, Python, Julia, ...)
- *knowledge in convex optimization*

# Textbooks

---

We recommend these books, but will not follow them closely





# Useful references

---

- *Spectral Methods for Data Science: A Statistical Perspective*, Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma
- *Nonconvex optimization meets low-rank matrix factorization: An overview*, Yuejie Chi, Yue M. Lu, and Yuxin Chen
- *Convex optimization*, Stephen Boyd, and Lieven Vandenberghe

# Grading

---

- Homework:  $\sim 3$  problem sets involving proofs and simulations
  - Due at Thursday lectures
- Course project
  - Either individually or in groups of two
- Your grade:  $\max\{0.4 \times \text{HW} + 0.6 \times \text{project}, \text{project}\}$

# Course project

---

## Two forms

- literature review
- original research
  - *You are strongly encouraged to combine it with your own research*

# Course project

---

## Two forms

- literature review
- original research
  - *You are strongly encouraged to combine it with your own research*

## Three milestones

- proposal (due Oct. 21st): up to 1 page
- in-class presentation: last week of class
- report (due TBA): up to 4 pages with unlimited appendix