Computational Principles for High-dim Data Analysis (Lecture One)

Yi Ma

EECS Department, UC Berkeley

August 26, 2021



EECS208, Fall 2021

1 Administrative Matters

2 Introduction (Chapter 1)

A Universal Task: Pursuit of Low-Dim Structures A Brief History The Modern Era

3 Assignments

- 4 ∃ ▶

Instructors of EECS 208, Fall 2021

- 1 Prof. Yi Ma (main instructor): yima@eecs.berkeley.edu
- **2** Prof. Jiantao Jiao (co-instructor): jiantao@eecs.berkeley.edu
- 3 Simon (Yuexiang) Zhai (head GSI): simonzhai@berkeley.edu
- **4 Xiangyu Yue** (GSI): xyyue@berkeley.edu

Office hours are all posted on the course websites. Office hours are all to be held virtually **via Zoom** (for now).

Main Textbook

High-Dimensional Data Analysis with Low-Dimensional Models Principles, Computation, and Applications

John Wright and Yi Ma Cambridge University Press, 2022.

Pre-production Copy from Website: https://book-wright-ma.github.io

• Github:

https://book-wright-ma.github.io/Book-WM-20210422.pdf

• Dropbox:

https://www.dropbox.com/s/0jm9qhm0t5v21gq/Book-WM-20210422.
pdf?dl=0

イロト イポト イヨト イヨト

Course Websites

• Course Website:

https://pages.github.berkeley.edu/UCB-EECS208/course_site/ Course information, detailed schedules, and resources etc.

• Piazza:

https://piazza.com/berkeley/fall2021/eecs208/

Interactive functions, announcements, Q&A, discussions, and team work etc.

Grading Policy

- Participation: 10%.
- Homeworks: 50%.
 - 4-5 Written Homeworks (principles).
 - 4-5 Programming Homeworks (practices).
- Final Project: 40%.
 - Midterm: 5min pitch of ideas; 2-3 pages of proposal.
 - Final: 15min presentation; 8-10 pages of final report (conference paper style).

A New Paradigm for Modern Data Science Principles (Computation) Applications



Pursue Low-dim Structures in High-dim Data

Introduction: A Universal Task



"Entities should not be multiplied without necessity." – William of Ockham, Law of Parsimony

System Identification: Linear Systems or RNNs



Problem: determine the system (A, B, C, D) from the input and output sequences

$$\{u_0, u_1, u_2, \ldots\}, \{y_0, y_1, y_2, \ldots\}.$$

Fact: If the dimension of $oldsymbol{x} \in \mathbb{R}^n$, then

$$\operatorname{rank}(\mathbf{Y}\mathbf{U}^{\perp}) \le n.$$
 (1)

Variants: Recursive Neural Networks (RNNs):

$$\begin{cases} \boldsymbol{x}(t+1) = \sigma_{\boldsymbol{x}} (\boldsymbol{A} \boldsymbol{x}(t) + \boldsymbol{B} \boldsymbol{u}(t) + \boldsymbol{b}), \\ \boldsymbol{y}(t) = \sigma_{\boldsymbol{y}} (\boldsymbol{C} \boldsymbol{x}(t) + \boldsymbol{d}), \end{cases}$$
(2)

Visual Patterns and Correlations





$$\operatorname{rank}(M) = d \ll n.$$
 (3)



Let M be the multiview matrix associated with corresponding features (points, lines, planes, symmetric structures), we have

$$\operatorname{rank}(M) \le 1 \text{ or } 2.$$
 (4)

An Invitation to 3D Vision, Ma, Soatto, Kosecka, and Sastry, Springer, 2004.

 a_1 a_2 a_3 a_4 a_5 a_6 a_6

Signal Acquisition and Processing

Fact: Sample band-limited signals with Nyquist frequency: $f = 2 \cdot \frac{\Omega}{2\pi}$.



Figure: Comparing Classical Signal Processing and Compressive Sensing Pipelines



Ma (EECS Department, UC Berkeley)

Graphical Models in Machine Learning

 $x \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$ with a covariance matrix $\mathbf{\Sigma} \in \mathbb{R}^{n imes n}$. Let $\mathbf{\Theta} \equiv \mathbf{\Sigma}^{-1}$. Then:

Fact:
$$\theta_{ij} = 0$$
 iff $x_i \perp \perp x_j \mid \boldsymbol{x}_{-i,-j}$. (5)



Figure: A graphical model for dependency among random variables.

Graphical Model Identification

Let x be partitioned into **observed** and **hidden** $x = (x_o, x_h)$. Its covariance matrix:

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_o & \boldsymbol{\Sigma}_{o,h} \\ \boldsymbol{\Sigma}_{o,h}^* & \boldsymbol{\Sigma}_h \end{bmatrix} \equiv \begin{bmatrix} \boldsymbol{\Theta}_o & \boldsymbol{\Theta}_{o,h} \\ \boldsymbol{\Theta}_{o,h}^* & \boldsymbol{\Theta}_h \end{bmatrix}^{-1} \in \mathbb{R}^{n \times n}.$$
(6)

From linear algebra (Schur complement):

Fact:
$$\Sigma_o^{-1} = \Theta_o - \Theta_{o,h} \Theta_h^{-1} \Theta_{o,h}^* \in \mathbb{R}^{n_o \times n_o}.$$
 (7)

Hence to infer Θ from the observable Σ_o , we need to solve a problem of **sparse plus low-rank** decomposition:

$$\sum_{\substack{o \text{bbserved}}}^{-1} = \frac{S}{\text{sparse}} + \frac{L}{\text{low rank}} \in \mathbb{R}^{n_o \times n_o}.$$
(8)

< □ > < 同 > < 回 > < 回 > < 回 >

History: Nature and Neuroscience

Dogma for natural vision [Barlow 1972]: "... to represent the input as completely as possible by activity in as few neurons as possible."



Find sparse $\{x_i\}$ such that

$$oldsymbol{y} = \sum_{i=1}^n x_i oldsymbol{a}_i + oldsymbol{\epsilon} \quad \in \mathbb{R}^m, \quad ext{(9)}$$

< □ > < □ > < □ > < □ > < □ > < □ >

[Nature, Olshausen and Field 1996.]



Ma (EECS Department, UC Berkeley)

EECS208, Fall 2021

History: Signal Processing

Model y as a **linear function** of variables a_1, \ldots, a_n :

$$y = f(\mathbf{a}) = \mathbf{a}^* \mathbf{x} = a_1 x_1 + a_2 x_2 + \dots + a_n x_n,$$
 (10)

from measurements

$$y_i = \boldsymbol{a}_i^* \boldsymbol{x} + \epsilon_i, \quad i = 1, 2, \dots, m,$$
(11)

where ϵ_i is possible measurement noise or error.



History: Error Correction and Denoising (m > n)

Least Absolute Deviations [Roger Joseph Boscovich, 1750]:

$$\min_{\boldsymbol{x}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_{1} = \sum_{i=1}^{m} |y_{i} - \boldsymbol{a}_{i}^{*}\boldsymbol{x}|, \quad \epsilon \sim \frac{1}{2b} \exp\left(-\frac{|\epsilon|}{b}\right).$$
(12)

Least Squares [Legendre in 1805 and Gauss in 1809]:

$$\min_{\boldsymbol{x}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_{2}^{2} = \sum_{i=1}^{m} (y_{i} - \boldsymbol{a}_{i}^{*}\boldsymbol{x})^{2}, \quad \epsilon \sim \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{\epsilon^{2}}{2\sigma^{2}}\right).$$
(13)

Error Correction [Ben Logan 1960]:

$$y(t) = x(t) + e(t): \min ||x - y||_1 \text{ subject to } x \in \mathcal{B}_1(\Omega).$$
(14)
Logan's Phenomenon: $|T| \times \Omega < \frac{\pi}{2}.$

More Recent History: Linear Regression (m < n)

Best Subset Selection [Hocking, Leslie, and Beale 1967]:

$$\min_{\boldsymbol{x}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2 \quad \text{subject to} \quad \|\boldsymbol{x}\|_0 \le k, \tag{15}$$

Stepwise Regression (\mathcal{P}_I projection on A_I) [Efroymson 1966]:

$$i_k = \arg\min_{i \notin I_k} \| \boldsymbol{y} - \mathcal{P}_{I_k \cup \{i\}}(\boldsymbol{y}) \|_2^2, \quad I_{k+1} = I_k \cup \{i_k\}.$$
 (16)

Lasso Regression [Tibshirani 1996]:

$$\min_{\boldsymbol{x}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2 \quad \text{subject to} \quad \|\boldsymbol{x}\|_1 \leq k. \tag{17}$$

Basis Pursuit [Chen, Donoho, and Saunders 1998] :

$$\min_{\boldsymbol{x}} \|\boldsymbol{x}\|_1$$
 subject to $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}.$ (18)

History: Principal Component Analysis (PCA)



Figure: Left: regression; Right: principal component analysis.

A high-dim random vector \boldsymbol{y} is approximated by the d < m components as:

$$oldsymbol{y} = oldsymbol{u}_1 w_1 + oldsymbol{u}_2 w_2 + \dots + oldsymbol{u}_d w_d + oldsymbol{\epsilon} \doteq oldsymbol{U}oldsymbol{w} + oldsymbol{\epsilon} \quad \in \mathbb{R}^m,$$
 (19)

where $U = [u_1, u_2, \dots, u_d] \in \mathbb{R}^{m \times d}$, $w = [w_1, w_2, \dots, w_d]^* \in \mathbb{R}^d$, and the variance of the residual $\epsilon \in \mathbb{R}^m$ is minimized:

$$\min \mathbb{E} \left[\| \boldsymbol{y} - \boldsymbol{U} \boldsymbol{w} \|_2^2 \right].$$
 (20)

History: Low-Rank Matrix Approximation [Eckart and Young 1936]

A Matrix of Samples: $Y = [y_1, y_2, \dots, y_n] \in \mathbb{R}^{m \times n}$. (21)

Matrix approximation by rank-1 factors (Beltrami and Jordan 1870's):

$$\boldsymbol{Y} = \sigma_1 \boldsymbol{u}_1 \boldsymbol{v}_1^* + \sigma_2 \boldsymbol{u}_2 \boldsymbol{v}_2^* + \dots + \sigma_d \boldsymbol{u}_d \boldsymbol{v}_d^* + \boldsymbol{E}, \qquad (22)$$

Low-rank matrix approximation:

$$X_{\star} = \arg\min_{\mathbf{X}} \|\mathbf{Y} - \mathbf{X}\|_2^2$$
 subject to $\operatorname{rank}(\mathbf{X}) \le d.$ (23)

Solution via Singular Value Decomposition:

$$\boldsymbol{X}_{\star} = \boldsymbol{U}_d \boldsymbol{\Sigma}_d \boldsymbol{V}_d^*, \qquad (24)$$

< □ > < 同 > < 三 > < 三 >

where $\boldsymbol{Y} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^*$ be the SVD of the matrix $\boldsymbol{Y} \in \mathbb{R}^{m \times n}$.

A Long and Rich History...

A **long and rich history** of robust estimation with error correction and missing data imputation:



- R. J. Boscovich. *De calculo probailitatum que respondent diversis valoribus summe errorum post plures observationes* ..., before 1756
- A. Legendre. Nouvelles methodes pour la determination des orbites des cometes, 1806



over-determined + dense, Gaussian





C. Gauss. Theory of motion of heavenly bodies, 1809

A. Beurling. Sur les integrales de Fourier absolument convergentes et leur application a une transformation functionelle, 1938

B. Logan. Properties of High-Pass Signals, 1965



underdetermined + sparse, Laplacian

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))



Why a Shift of Paradigm?





A D F A B F A B F A B

Real application data often contain **missing observations**, **corruptions**, or subject to unknown **deformation or misalignment**.

Classical methods (e.g., PCA, least square regression) break down...

From Curses to Blessings of High-Dimensionality

For problems of identifying low-dimensional (e.g. sparse or low-rank) structures of massive data in high-dimensional spaces, we like to answer **two fundamental questions:**

- Why many seemingly intractable high-dimensional problems can be solved efficiently without suffering the curses of dimensionality? (seemingly NP-hard, combinatorial, exponential, astronomical scale...)
- What is the precise characterization of the required data complexity and computational complexity for certain guaranteed accuracy or probability of success? (number of samples needed, number of oracles computed...)

イロト 不得 トイヨト イヨト 二日

Our goal: develop a principled mathematical foundation to answer above two questions for the following problems:

• Compressive Sensing (Parsimony):

$$\min_{\boldsymbol{x}} \|\boldsymbol{x}\|_1 \quad \text{subject to} \quad \boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}, \tag{25}$$

• Error Correction (Robustness):

$$\min_{\boldsymbol{x}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_1 \quad \text{with} \quad \boldsymbol{y} = \boldsymbol{A}\boldsymbol{x} + \boldsymbol{e}. \tag{26}$$

• **Deep Learning** (Nonlinearlity):

$$\begin{cases} \boldsymbol{z}_{\ell+1} = \phi(\boldsymbol{A}^{\ell}\boldsymbol{z}_{\ell}), \quad \boldsymbol{z}_0 = \boldsymbol{x}, \quad \ell = 0, 1, \dots, L-1, \\ \boldsymbol{y} = \phi(\boldsymbol{C}\boldsymbol{z}_L), \end{cases}$$
(27)

where $\phi(\cdot)$ is typically *sparsity-promoting* nonlinear activation.

Guarantees: High-Dim Geometry and Statistics

Minimal data that ensure tractable method for a correct solution?

Phenomena against intuition from low-dim spaces:

• Measure Concentration ($\epsilon \sim O(n^{-1/2})$)

$$\operatorname{Area}\{x \in \mathbb{S}^{n-1} : -\epsilon \le x_n \le \epsilon\} = 0.99 \cdot \operatorname{Area}(\mathbb{S}^{n-1}), \qquad (28)$$

• Neighborly Polytopes (vertices from a Gaussian matrix):

$$oldsymbol{A} = [oldsymbol{a}_1, oldsymbol{a}_2, \dots, oldsymbol{a}_n] \in \mathbb{R}^{m imes n}$$



Computational Cost: Scalable Optimization

Minimal computational cost that ensures an accurate solution?

Solutions defying conventional wisdom:

- **Convex Optimization**: accelerated first order methods, augmented Lagrangian method, alternating minimization, etc.
- **Nonconvex Optimization**: symmetry, stochastic gradient descent, generalized power iteration, etc.



Figure: Left: conventional view. Right: actual landscape.

< □ > < □ > < □ > < □ > < □ > < □ >

The first 20 years of the century

Figure: A **perfect storm** for unprecedented confluence and advancement in mathematics, computation, technology, and science.



Homework 0

- Reading Assignment I: Preface
- Reading Assignment II: Introduction (Chapter 1)
- Reading Assignment III: Linear Algebra (Appendix A)

э

< 1 k