

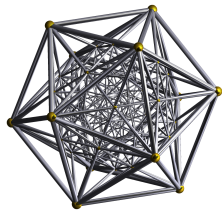
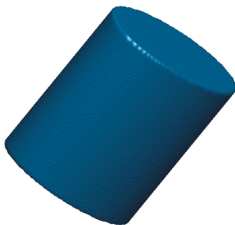
# Computational Principles for High-dim Data Analysis

(Lecture Ten)

**Yi Ma and Jiantao Jiao**

EECS Department, UC Berkeley

September 30, 2021



# Convex Methods for Low-Rank Matrix Recovery (Matrix Completion)

- 1 Motivating Example
- 2 Nuclear Norm Minimization
- 3 Algorithm: Augmented Lagrange Multiplier
- 4 Conditions for Success
- 5 Stable Matrix Completion

*“Mathematics is the art of giving the same name to different things.”*  
– Henri Poincaré

# Example of Low-rank Matrix Completion

**Recommendation Systems (how internet companies make money):**

$$\begin{array}{c} \text{Users} \end{array} \begin{bmatrix} 5 & 3 & \dots & ? \\ ? & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & ? & \dots & ? \end{bmatrix} = \mathcal{P}_{\Omega} \begin{pmatrix} \begin{bmatrix} 5 & 3 & \dots & 5 \\ 4 & 2 & \dots & 4 \\ \vdots & \vdots & \ddots & \vdots \\ 5 & 5 & \dots & 3 \end{bmatrix} \\ \text{Complete Ratings } \mathbf{X} \end{pmatrix}$$

Items

Observed (Incomplete) Ratings  $\mathbf{Y}$

We observe:

$$\begin{array}{c} \mathbf{Y} \\ \text{Observed ratings} \end{array} = \mathcal{P}_{\Omega} \begin{bmatrix} \mathbf{X} \\ \text{Complete ratings} \end{bmatrix},$$

where  $\Omega \doteq \{(i, j) \mid \text{user } i \text{ has rated product } j\}$ .

# Nuclear Norm Minimization

## Problem (Matrix Completion)

Let  $\mathbf{X}_o \in \mathbb{R}^{n \times n}$  be a low-rank matrix. Suppose we are given  $\mathbf{Y} = \mathcal{P}_\Omega[\mathbf{X}_o]$ , where  $\Omega \subseteq [n] \times [n]$ . Fill in the missing entries of  $\mathbf{X}_o$ .

**Notice:** If  $(i, j) \notin \Omega$ ,  $\mathcal{P}_\Omega[\mathbf{E}_{ij}] = \mathbf{0}$ . So  $\mathcal{P}_\Omega$  has matrices of rank one in its null space! So,  $\mathcal{P}_\Omega$  cannot be rank-RIP for any rank  $r > 0$  with  $\delta < 1$ .

**Question:** can we still find  $\mathbf{X}_o$  by solving the nuclear norm minimization:

$$\min \|\mathbf{X}\|_* \quad \text{subject to} \quad \mathcal{P}_\Omega[\mathbf{X}] = \mathbf{Y} \quad (1)$$

**Simulations lead the way of investigation – need an algorithm...**

# Algorithm via Augmented Lagrange Multiplier

Nuclear norm minimization for matrix completion:

$$\min \underbrace{\|\mathbf{X}\|_*}_{f(\mathbf{x})} \quad \text{subject to} \quad \underbrace{\mathcal{P}_\Omega[\mathbf{X}] = \mathbf{Y}}_{g(\mathbf{x})=0} \quad (2)$$

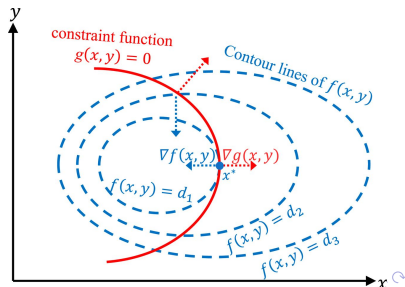
The Lagrangian method:

$$\mathcal{L}(\mathbf{X}, \boldsymbol{\Lambda}) = \|\mathbf{X}\|_* + \langle \boldsymbol{\Lambda}, \mathbf{Y} - \mathcal{P}_\Omega[\mathbf{X}] \rangle. \quad (3)$$

Optimality conditions:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}} = 0, \quad \frac{\partial \mathcal{L}}{\partial \boldsymbol{\Lambda}} = 0. \quad (4)$$

However, it only holds  
at the point of the optimal solution  $\mathbf{x}^*$ .



# Algorithm via Augmented Lagrange Multiplier

The *augmented* Lagrangian is to regularize the landscape around the optimal solution  $\mathbf{x}^*$ :

$$\mathcal{L}_\mu(\mathbf{X}, \mathbf{\Lambda}) = \|\mathbf{X}\|_* + \langle \mathbf{\Lambda}, \mathbf{Y} - \mathcal{P}_\Omega[\mathbf{X}] \rangle + \frac{\mu}{2} \|\mathbf{Y} - \mathcal{P}_\Omega[\mathbf{X}]\|_F^2. \quad (5)$$

Amenable for alternating optimization to converge to the optimal solution  $\mathbf{x}^*$  more easily and efficiently:

$$\text{Primal: } \mathbf{X}_{k+1} \in \arg \min_{\mathbf{X}} \mathcal{L}_\mu(\mathbf{X}, \mathbf{\Lambda}_k), \quad (6)$$

$$\text{Dual: } \mathbf{\Lambda}_{k+1} = \mathbf{\Lambda}_k + \mu \mathcal{P}_\Omega[\mathbf{Y} - \mathbf{X}_{k+1}]. \quad (7)$$

# Algorithm: Proximal Gradient Descent

How to minimize the augmented Lagrangian  $\mathcal{L}_\mu$ :

$$\min_{\mathbf{X}} F(\mathbf{X}) \doteq \underbrace{\|\mathbf{X}\|_*}_{g(\mathbf{X}) \text{ convex}} + \underbrace{\langle \mathbf{\Lambda}, \mathbf{Y} - \mathcal{P}_\Omega[\mathbf{X}] \rangle + \frac{\mu}{2} \|\mathbf{Y} - \mathcal{P}_\Omega[\mathbf{X}]\|_F^2}_{f(\mathbf{X}) \text{ smooth, convex, } \mu\text{-Lipschitz}}. \quad (8)$$

At each iterate  $\mathbf{X}_k$ , construct a local (quadratic) upper bound for  $F$ :

$$\hat{F}(\mathbf{X}, \mathbf{X}_k) = g(\mathbf{X}) + f(\mathbf{X}_k) + \langle \nabla f(\mathbf{X}_k), \mathbf{X} - \mathbf{X}_k \rangle + \frac{\mu}{2} \|\mathbf{X} - \mathbf{X}_k\|_2^2. \quad (9)$$

**Proximal gradient descent:** the next iterate  $\mathbf{X}_{k+1}$  is computed as

$$\mathbf{X}_{k+1} = \arg \min_{\mathbf{X}} \left\{ g(\mathbf{X}) + \frac{\mu}{2} \left\| \mathbf{X} - \underbrace{\left( \mathbf{X}_k - \frac{1}{\mu} \nabla f(\mathbf{X}_k) \right)}_M \right\|_F^2 \right\} \quad (10)$$

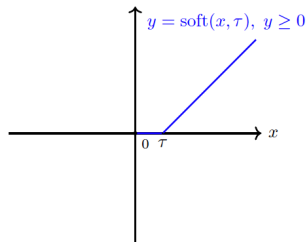
$$= \text{prox}_{g/\mu}(M) \quad (\text{see details in Chapter 8}). \quad (11)$$

## Algorithm: Proximal Operator for Nuclear Norm

For a matrix  $M$  with SVD  $M = U\Sigma V^*$ , its singular value thresholding operator is:

$$\mathcal{D}_\tau[M] = U\mathcal{S}_\tau[\Sigma]V^*,$$

where  $\mathcal{S}_\tau[X] = \text{sign}(X) \circ (|X| - \tau)_+$  is the entry-wise soft thresholding operator.



### Theorem

*The unique solution  $X_\star$  to the program:*

$$\min_{\mathbf{X}} \{ \|\mathbf{X}\|_* + \frac{\mu}{2} \|\mathbf{X} - \mathbf{M}\|_F^2 \}, \quad (12)$$

*is given by*

$$\mathbf{X}_\star = \mathcal{D}_{\mu^{-1}}[\mathbf{M}]. \quad (13)$$



# Algorithm via Augmented Lagrange Multiplier

## Outer Loop: Matrix Completion by ALM

- 1: **initialize:**  $\mathbf{X}_0 = \mathbf{\Lambda}_0 = 0, \mu > 0$ .
- 2: **while** not converged **do**
- 3:   compute  $\mathbf{X}_{k+1} \in \arg \min_{\mathbf{X}} \mathcal{L}_{\mu}(\mathbf{X}, \mathbf{\Lambda}_k)$  (say by PG);
- 4:   compute  $\mathbf{\Lambda}_{k+1} = \mathbf{\Lambda}_k + \mu(\mathbf{Y} - \mathcal{P}_{\Omega}[\mathbf{X}_{k+1}])$ .
- 5: **end while**

## Inner Loop: Proximal Gradient

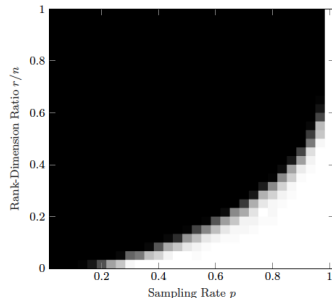
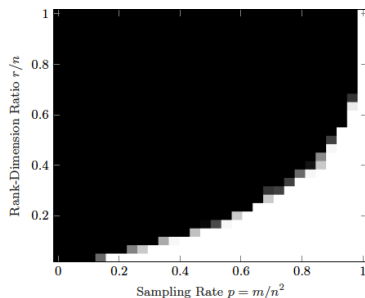
- 1: **initialize:**  $\mathbf{X}_0$  starts with the  $\mathbf{X}_k$  from the outer loop.
- 2: **while** not converged **do**
- 3:   compute

$$\begin{aligned} \mathbf{X}_{\ell+1} &= \text{prox}_{g/\mu}(\mathbf{X}_{\ell} - \mu^{-1} \nabla f(\mathbf{X}_{\ell})) \\ &= \mathcal{D}_{\mu^{-1}} \left[ \underbrace{\mathcal{P}_{\Omega^c}[\mathbf{X}_{\ell}] + \mathbf{Y} + \mu^{-1} \mathcal{P}_{\Omega}[\mathbf{\Lambda}_k]}_{\text{exercise}} \right]. \end{aligned}$$

- 4: **end while**

# Similar Phenomena of Success

**Comparison:** low-rank matrix recovery from random linear measurements versus matrix completion from random sampled entries.



**Figure:** Left: phase transition for matrix recovery; Right: phase transition for matrix completion.

# When Nuclear Norm Minimization Succeeds?

## When it fails?

- ① if  $\mathbf{X}_o$  is itself *sparse* (as in the example of  $\mathbf{E}_{ij}$ )
- ② if  $\Omega$  is chosen adversarially (e.g., an entire row or column of  $\mathbf{X}_o$ ).

Notice for any rank- $r$  orthogonal matrix  $\mathbf{U}$ :

$$\sum_i \|\mathbf{e}_i^* \mathbf{U}\|_2^2 = \|\mathbf{U}\|_F^2 = r \quad \implies \quad \max_i \|\mathbf{e}_i^* \mathbf{U}\|_2^2 \geq r/n.$$

## Definition

We say that  $\mathbf{X}_o = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  is  $\nu$ -*incoherent* if the following hold:

$$\forall i \in [n], \quad \|\mathbf{e}_i^* \mathbf{U}\|_2^2 \leq \nu r/n, \quad (14)$$

$$\forall j \in [n], \quad \|\mathbf{e}_j^* \mathbf{V}\|_2^2 \leq \nu r/n. \quad (15)$$

## When Nuclear Norm Minimization Succeeds?

**Bernoulli  $\text{Ber}(p)$  sampling model:** each entry  $(i, j)$  belongs to the observed set  $\Omega$  independently with probability  $p \in [0, 1]$ . Hence, the expected number of observed entries is:

$$m = \mathbb{E}[|\Omega|] = pn^2. \quad (16)$$

### Theorem (Matrix Completion via Nuclear Norm Minimization)

Let  $\mathbf{X}_o \in \mathbb{R}^{n \times n}$  be a rank- $r$  matrix with incoherence parameter  $\nu$ . Suppose that we observe  $\mathbf{Y} = \mathcal{P}_\Omega[\mathbf{X}_o]$ , with  $\Omega$  sampled according to the Bernoulli model with probability

$$p \geq C_1 \frac{\nu r \log^2(n)}{n}. \quad (17)$$

Then with probability at least  $1 - C_2 n^{-c_3}$ ,  $\mathbf{X}_o$  is the unique optimal solution to

$$\text{minimize } \|\mathbf{X}\|_* \quad \text{subject to } \mathcal{P}_\Omega[\mathbf{X}] = \mathbf{Y}. \quad (18)$$

# When Nuclear Norm Minimization Succeeds?

## Lemma (Subdifferential of nuclear norm)

*Let  $\mathbf{X} \in \mathbb{R}^{n \times n}$  have compact singular value decomposition  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ . The subdifferential of the nuclear norm at  $\mathbf{X}$  is given by*

$$\partial \|\cdot\|_* (\mathbf{X}) = \{ \mathbf{Z} \mid \mathcal{P}_{\mathbf{T}}[\mathbf{Z}] = \mathbf{U}\mathbf{V}^*, \|\mathcal{P}_{\mathbf{T}^\perp}[\mathbf{Z}]\| \leq 1 \}. \quad (19)$$

# When Nuclear Norm Minimization Succeeds?

## Key ideas for the Theorem:

For the program:

$$\min \| \mathbf{X} \|_* \quad \text{subject to} \quad \mathcal{P}_\Omega[\mathbf{X}] = \mathcal{P}_\Omega[\mathbf{X}_o]. \quad (20)$$

Similar to the  $\ell^1$  case, find a dual certificate  $\mathbf{\Lambda}$  that satisfies (the KKT condition):

- (i)  $\mathbf{\Lambda}$  is supported on  $\Omega$ :  $\mathcal{P}_\Omega[\mathbf{\Lambda}] = \mathbf{\Lambda}$  and
- (ii)  $\mathbf{\Lambda} \in \partial \| \cdot \|_* (\mathbf{X}_o)$  – i.e.,  $\mathcal{P}_\mathbf{T}[\mathbf{\Lambda}] = \mathbf{UV}^*$  and  $\| \mathcal{P}_{\mathbf{T}^\perp}[\mathbf{\Lambda}] \| \leq 1$ ,

**Strategy:** look for a matrix  $\mathbf{\Lambda}$  of smallest 2-norm that satisfies the equality constraints

$$\mathcal{P}_{\Omega^c}[\mathbf{\Lambda}] = \mathbf{0}, \quad \mathcal{P}_\mathbf{T}[\mathbf{\Lambda}] = \mathbf{UV}^*, \quad (21)$$

and then hope to check that it satisfies the inequality constraints

$$\| \mathcal{P}_{\mathbf{T}^\perp}[\mathbf{\Lambda}] \| \leq 1.$$

## When Nuclear Norm Minimization Succeeds?

Unfortunately, this straightforward strategy does not work out directly as solution to the equalities is not so easy to analyze...

**An alternative strategy:** an set of (relaxed) conditions for optimality:

### Proposition (KKT Conditions – Approximate Version)

*The matrix  $\mathbf{X}_o$  is the unique optimal solution to the nuclear minimization problem (18) if the following set of conditions hold*

- ① *The operator norm of the operator  $p^{-1}\mathcal{P}_\mathsf{T}\mathcal{P}_\Omega\mathcal{P}_\mathsf{T} - \mathcal{P}_\mathsf{T}$  is small:*

$$\|p^{-1}\mathcal{P}_\mathsf{T}\mathcal{P}_\Omega\mathcal{P}_\mathsf{T} - \mathcal{P}_\mathsf{T}\| \leq \frac{1}{2}.$$

- ② *There exists a dual certificate  $\mathbf{\Lambda}$  that satisfies  $\mathcal{P}_\Omega[\mathbf{\Lambda}] = \mathbf{\Lambda}$  and*
  - *(a)  $\|\mathcal{P}_{\mathsf{T}^\perp}[\mathbf{\Lambda}]\| \leq \frac{1}{2}$ ;*
  - *(b)  $\|\mathcal{P}_\mathsf{T}[\mathbf{\Lambda}] - \mathbf{UV}^*\|_F \leq \frac{1}{4n}$ .*

## Matrix Completion with Noise

**Problem:** the observed entries are often corrupted with some noise:

$$Y_{ij} = [\mathbf{X}_o]_{ij} + Z_{ij}, \quad (i, j) \in \Omega; \quad \text{or} \quad \mathcal{P}_\Omega[\mathbf{Y}] = \mathcal{P}_\Omega[\mathbf{X}_o] + \mathcal{P}_\Omega[\mathbf{Z}], \quad (22)$$

where  $Z_{ij}$  can be some small noise, say  $\|\mathcal{P}_\Omega[\mathbf{Z}]\|_F < \epsilon$ .

$$\min \|\mathbf{X}\|_* \quad \text{subject to} \quad \|\mathcal{P}_\Omega[\mathbf{X}] - \mathcal{P}_\Omega[\mathbf{Y}]\|_F < \epsilon. \quad (23)$$

### Theorem (Stable Matrix Completion)

Let  $\mathbf{X}_o \in \mathbb{R}^{n \times n}$  be a rank- $r$ ,  $\nu$ -incoherent matrix. Suppose that we observe  $\mathcal{P}_\Omega[\mathbf{Y}] = \mathcal{P}_\Omega[\mathbf{X}_o] + \mathcal{P}_\Omega[\mathbf{Z}]$ , where  $\Omega$  is uniformly sampled from subsets of size

$$m \geq C_1 \nu n r \log^2(n), \quad (24)$$

then with high probability, the optimal solution  $\hat{\mathbf{X}}$  to the convex program (23) satisfies

$$\|\hat{\mathbf{X}} - \mathbf{X}_o\|_F \leq c \frac{n\sqrt{n} \log(n)}{\sqrt{m}} \epsilon \leq c' \frac{n}{\sqrt{r}} \epsilon, \quad \text{for some } c > 0. \quad (25)$$



# Summary

Nuclear norm minimization can recover w.h.p. a low-rank matrix  $\mathbf{X}_o$  from

- ①  $m = O(nr)$  random linear measurements:  $\mathbf{y} = \mathcal{A}[\mathbf{X}]$ ;
- ②  $m = O(nr \log^2 n)$  randomly sampled entries:  $\mathbf{Y} = \mathcal{P}_\Omega[\mathbf{X}]$ ;
- ③ the estimate  $\hat{\mathbf{X}}$  is stable to small noise.

# Assignments

- Reading: Section 4.4-4.6 of Chapter 4.
- Programming Homework # 2.