

ECE 8201: Low-dimensional Signal Models for High-dimensional Data Analysis

Lecture 7: Matrix completion

Yuejie Chi
The Ohio State University



THE OHIO STATE UNIVERSITY

Reference

- “Guaranteed Minimum-Rank Solutions of Linear Matrix Equations via Nuclear Norm Minimization”, Recht, Fazel, and Parrilo, 2007.
- “The power of convex relaxation: Near-optimal matrix completion”, E. J. Candès and T. Tao, 2007.

Outline

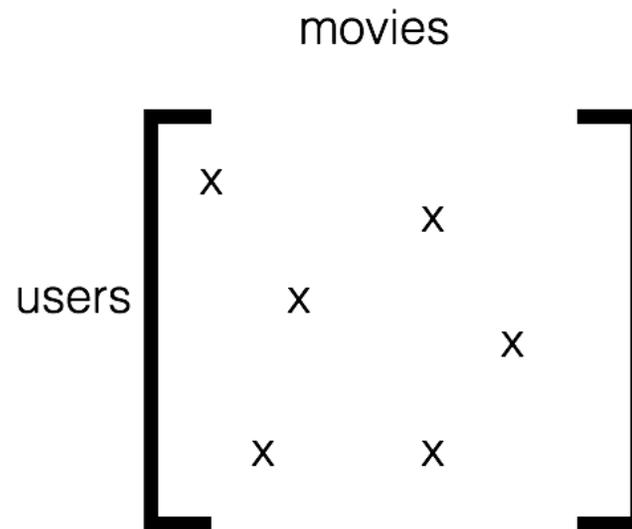
Matrix completion

- Motivation
- Theoretical aspects:
 - nuclear norm
 - low-rank matrix sensing
 - low-rank matrix completion
- efficient algorithm

The Netflix problem

The Netflix problem, or collaborative filtering

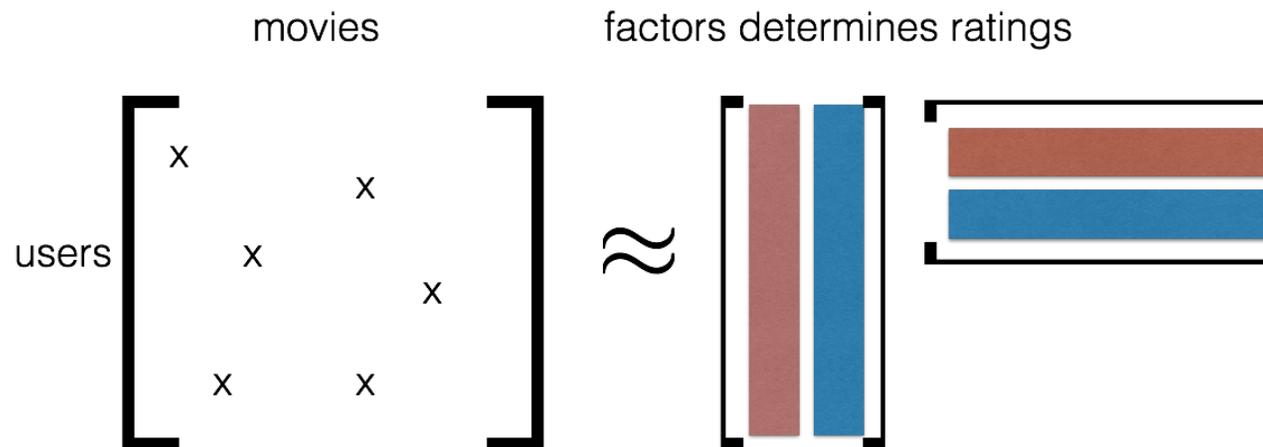
- How to estimate the missing ratings?



- About a million users, and 25,000 movies, with sparsely sampled ratings

Solution: low-rank matrix completion

- Matrix completion problem: consider $M \in \mathbb{R}^{n_1 \times n_2}$ to represent the Netflix data set, we may model it through factorization:



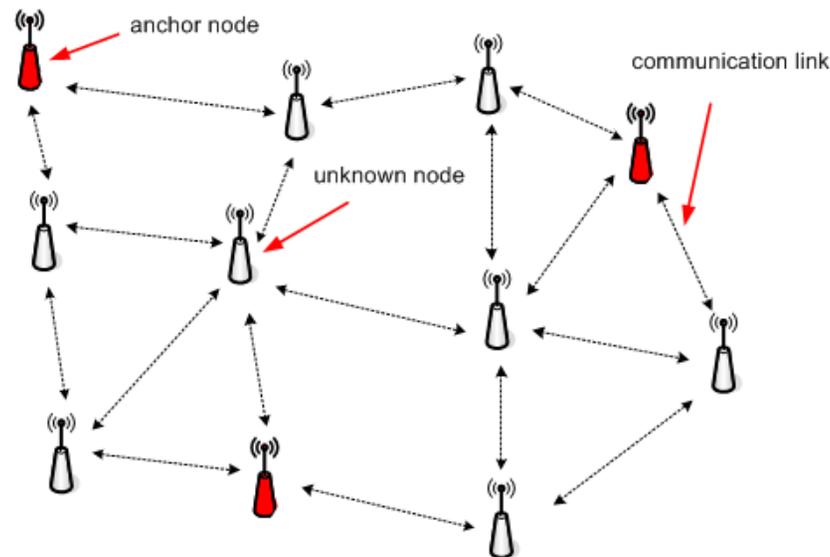
- The rank r of M is much smaller than its dimension $r \ll \min\{n_1, n_2\}$.

Sensor localization

- Given n points $\{\mathbf{x}_j\}_{j=1}^n \in \mathbb{R}^3$
- Observe partial information about distances:

$$M_{i,j} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$$

e.g. in wireless sensor network, each sensor can measure the distance to its neighbors, would like to globally locate all sensors.



Solution: low-rank matrix completion

- Write the matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times 3}$$

then

$$M_{i,j} = \mathbf{x}_i^T \mathbf{x}_i + \mathbf{x}_j^T \mathbf{x}_j - 2\mathbf{x}_i^T \mathbf{x}_j$$

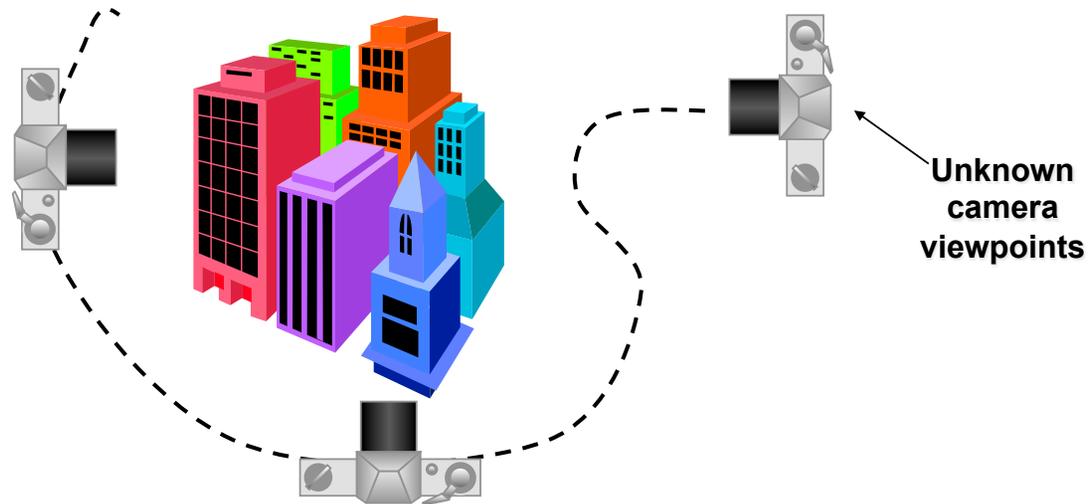
- Matrix completion problem: Let $\mathbf{Y} = \mathbf{X}\mathbf{X}^T$. The distance matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$ between points can be written as

$$\mathbf{M} = \text{diag}(\mathbf{Y})\mathbf{e}^T + \mathbf{e}\text{diag}(\mathbf{Y})^T - 2\mathbf{Y}$$

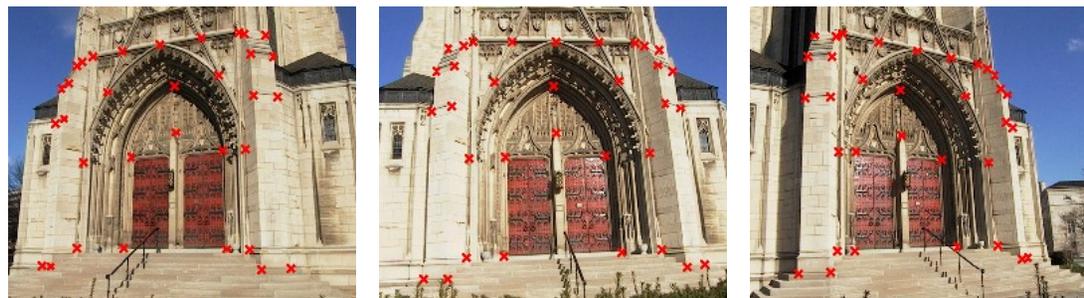
- The rank r of \mathbf{M} is much smaller than its dimension $r \ll n$.

Structure from motion

Structure from motion: reconstruct scene geometry and camera motion from multiple images.



In the pipeline of performing SfM, assume we've found a set of good feature points with their corresponding 2D locations in the images.



Structure from motion - continued

Tomasi and Kanade's factorization: Given n points $\mathbf{x}_{i,j}^T \in \mathbb{R}^2$ corresponding to the location of the i th point in the j th frame, define the matrix

$$\mathbf{M} = \begin{bmatrix} \mathbf{x}_{1,1} & \cdots & \mathbf{x}_{1,m} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{n,1} & \cdots & \mathbf{x}_{n,m} \end{bmatrix} \in \mathbb{R}^{n \times 2m}$$

In the absence of noise, this matrix admits a low-rank factorization:

$$\mathbf{M} = \underbrace{\begin{bmatrix} \mathbf{s}_1^T \\ \vdots \\ \mathbf{s}_n^T \end{bmatrix}}_{\text{3D structure matrix}} \underbrace{\begin{bmatrix} \mathbf{w}_1 & \mathbf{w}_2 & \cdots & \mathbf{w}_{2m} \end{bmatrix}}_{\text{camera motion matrix}}$$

where $\mathbf{s}_i \in \mathbb{R}^3$, which gives $\text{rank}(\mathbf{M}) = 3$.

Due to occlusions, there are many missing entries in the matrix \mathbf{M} . Can we complete the missing entries?

Many more...

Many more applications:

- spatial-temporal data: low-rank due to correlations, e.g. MRI video, network traffic, etc..
- quantum space tomography
- linear system identification

Problem of interest: Can we recover the matrices of interest from “incomplete” observations, using efficient algorithms?

- the problem is ill-posed without additional constraints

Low-rank matrices

- Let $\mathbf{M} \in \mathbb{R}^{n \times n}$ (square case for simplicity) be a matrix of rank $r \ll n$.
- The Singular Value Decomposition (SVD) of \mathbf{M} is given as

$$\mathbf{M} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

where $\{\sigma_i\}_{i=1}^r$ are the singular values; and $\{\mathbf{u}_i\}_{i=1}^r$, $\{\mathbf{v}_i\}_{i=1}^r$ are the singular vectors.

- \mathbf{M} has $(2n - r)r$ degrees of freedom.

Linear measurements of low-rank matrices

We make linear measurements of \mathbf{M} :

$$y_i = \langle \mathbf{A}_i, \mathbf{M} \rangle = \text{Tr}(\mathbf{A}_i^T \mathbf{M}), \quad i = 1, \dots, m,$$

which can be written more concisely in an operator form:

$$\mathbf{y} = \mathcal{A}(\mathbf{M})$$

where $\mathcal{A} : \mathbb{R}^{n \times n} \mapsto \mathbb{R}^m$ denotes the measurement process. Its adjoint operator $\mathcal{A}^* : \mathbb{R}^m \mapsto \mathbb{R}^{n \times n}$ is defined as

$$\mathcal{A}^*(\mathbf{y}) = \sum_{i=1}^m y_i \mathbf{A}_i.$$

The problem of rank minimization:

$$\hat{\mathbf{M}} = \underset{\mathbf{X}}{\text{argmin}} \text{rank}(\mathbf{X}) \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}(\mathbf{X}).$$

Nuclear norm

Just as ℓ_1 norm provides a convex relaxation to cardinality minimization, we use the nuclear norm which give a convex relaxation to rank minimization.

Definition 1. *The nuclear norm of \mathbf{X} is defined as*

$$\|\mathbf{X}\|_* = \sum_{i=1}^n \sigma_i(\mathbf{X})$$

where $\sigma_i(\mathbf{X})$ is the i th largest singular value of \mathbf{X} .

- Since the rank is $\sum_{i=1}^n 1(\sigma_i(\mathbf{X}) \neq 0)$, the nuclear norm can be thought as an ℓ_1 norm relaxation of the vector of singular values.
- This is a norm. Relationships between different norms:

$$\|\mathbf{X}\| \leq \|\mathbf{X}\|_F \leq \|\mathbf{X}\|_* \leq \sqrt{r}\|\mathbf{X}\|_F \leq r\|\mathbf{X}\|.$$

- Tightest convex relaxation: $\{\mathbf{X} : \|\mathbf{X}\|_* \leq 1\}$ is the convex hull of rank-1 matrices obeying $\|\mathbf{x}\mathbf{y}^T\| \leq 1$.

Additivity of the nuclear norm

Lemma 1. *Let A and B be matrices of the same dimensions. If $AB^T = 0$ and $A^T B = 0$, then $\|A + B\|_* = \|A\|_* + \|B\|_*$.*

Remark: this implies that, if the row and column spaces of A and B are orthogonal, then $\|A + B\|_* = \|A\|_* + \|B\|_*$.

This is similar to the ℓ_1 norm when x and y have disjoint support:

$$\|x + y\|_1 = \|x\|_1 + \|y\|_1$$

which is essentially all we need to get the proof of ℓ_1 minimization with RIP...

Compute the nuclear norm via SDP

Lemma 2.

$$\|\mathbf{X}\|_* = \min_{\mathbf{W}_1, \mathbf{W}_2} \left\{ \frac{1}{2} \text{Tr}(\mathbf{W}_1) + \frac{1}{2} \text{Tr}(\mathbf{W}_2) \mid \begin{bmatrix} \mathbf{W}_1 & \mathbf{X} \\ \mathbf{X}^T & \mathbf{W}_2 \end{bmatrix} \succeq 0 \right\}.$$

This means we can compute the nuclear norm efficiently via semidefinite programming (SDP).

Proof: on the blackboard.

Dual norm

Definition 2. For a given norm $\|\cdot\|_{\mathcal{A}}$ in an inner product space $\langle \cdot, \cdot \rangle$, the dual norm is defined as

$$\|\mathbf{X}\|_{\mathcal{A}}^* := \max\{\langle \mathbf{X}, \mathbf{Y} \rangle : \|\mathbf{Y}\|_{\mathcal{A}} \leq 1\}.$$

By definition, this gives a general version of Cauchy-Schwarz inequality:

$$\langle \mathbf{X}, \mathbf{Y} \rangle \leq \|\mathbf{X}\|_{\mathcal{A}} \|\mathbf{Y}\|_{\mathcal{A}}^*.$$

Examples:

- The dual norm of $\|\cdot\|_F$ is $\|\cdot\|_F$;
- The dual norm of $\|\cdot\|_1$ is $\|\cdot\|_{\infty}$;
- The dual norm of $\|\cdot\|_*$ is $\|\cdot\|$;

Summary

rank minimization vs cardinality minimization:

parsimony concept	cardinality	rank
Hilbert Space norm	Euclidean	Frobenius
sparsity inducing norm	ℓ_1	nuclear
dual norm	ℓ_∞	operator
norm additivity	disjoint support	orthogonal row and column spaces
convex optimization	linear programming	semidefinite programming

Table 1: A dictionary relating the concepts of cardinality and rank minimization.

Low-rank matrix recovery via nuclear norm minimization

- The rank minimization problem:

$$\hat{\mathbf{M}} = \underset{\mathbf{X}}{\operatorname{argmin}} \operatorname{rank}(\mathbf{X}) \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}(\mathbf{X}).$$

- We pose the following nuclear norm minimization algorithm:

$$\hat{\mathbf{M}} = \underset{\mathbf{X}}{\operatorname{argmin}} \|\mathbf{X}\|_* \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}(\mathbf{X}),$$

which can be solved efficiently via SDP:

$$\begin{aligned} \hat{\mathbf{M}} = \underset{\mathbf{X}, \mathbf{W}_1, \mathbf{W}_2}{\operatorname{argmin}} & \frac{1}{2} \operatorname{Tr}(\mathbf{W}_1) + \frac{1}{2} \operatorname{Tr}(\mathbf{W}_2) \\ \text{s.t.} \quad \mathbf{y} = \mathcal{A}(\mathbf{X}), & \quad \begin{bmatrix} \mathbf{W}_1 & \mathbf{X} \\ \mathbf{X}^T & \mathbf{W}_2 \end{bmatrix} \succeq 0. \end{aligned}$$

Low-rank matrix sensing

- If \mathcal{A} satisfies the restricted isometry property for low-rank matrices:

Definition 3. *The operator \mathcal{A} satisfies the RIP of rank- r , if for any rank- r matrix, we have*

$$(1 - \delta_r) \|\mathbf{X}\|_F^2 \leq \|\mathcal{A}(\mathbf{X})\|_F^2 \leq (1 + \delta_r) \|\mathbf{X}\|_F^2$$

for $0 \leq \delta_r \leq 1$.

- If $\{\mathbf{A}_i\}_{i=1}^m$ are composed of i.i.d. Gaussian entries, then it satisfies the matrix RIP of order r with high probability, as soon as $m \gtrsim nr$.
- This allows us to develop almost parallel results to compressed sensing.

Theoretical guarantees

Theorem 3. *If \mathbf{A} satisfies the RIP of rank $4r$ with $\delta_{4r} \leq \sqrt{2} - 1$, then for all rank- r matrices, the nuclear norm minimization algorithm recovers \mathbf{M} exactly.*

Exact recovery from $O(nr)$ measurements!!

- For the noisy case,

$$\mathbf{y} = \mathcal{A}(\mathbf{M}) + \mathbf{w}$$

where \mathbf{w} is composed of i.i.d. $\mathcal{N}(0, \sigma^2)$ entries. We could similarly propose the matrix LASSO algorithm:

$$\hat{\mathbf{M}} = \operatorname{argmin}_{\mathbf{X}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2 + \lambda \|\mathbf{X}\|_*,$$

where λ is a regularization parameter.

- If $\|\mathcal{A}^*(\mathbf{w})\| \leq \lambda/2$ and $\delta_{4r} < (3\sqrt{2} - 1)/17$, then

$$\|\hat{\mathbf{M}} - \mathbf{M}\|_F \leq C\sqrt{r}\lambda$$

for some constant C . For the Gaussian case,

$$\|\mathcal{A}^*(\mathbf{w})\| \leq c_1\sqrt{n}\sigma := \lambda$$

for some large enough constant c_1 with probability at least $1 - 2e^{-cn}$.

- If \mathbf{M} is an approximately low-rank matrix, we further have

$$\|\hat{\mathbf{M}} - \mathbf{M}\|_F \leq C_1 \frac{\|\mathbf{M} - \mathbf{M}_r\|_*}{\sqrt{r}} + C_2\sqrt{nr}\sigma$$

with probability at least $1 - 2e^{-cn}$, in the Gaussian sampling case.

Low-rank matrix completion

In the matrix completion setting, we are given partial observations of the entries of \mathbf{M} , and wish to recover the missing entries.

- Denote $\Omega = \{(i, j) \in [n] \times [n]\}$ as the index set of observed entries.
- The observation can be written as

$$\mathbf{Y} = \mathcal{P}_\Omega(\mathbf{M})$$

where $Y_{ij} = M_{ij}$ if $(i, j) \in \Omega$ and $Y_{ij} = 0$ otherwise.

- Consider the following algorithm:

$$\min \|\mathbf{X}\|_* \quad \text{s.t.} \quad \mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X})$$

- The observation operator doesn't satisfy RIP!

Which sampling pattern?

Consider a rank-one matrix $M = \mathbf{x}\mathbf{y}^T$ with the following sampling pattern:

$$\begin{bmatrix} \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times \end{bmatrix}$$

If single row (or column) is not sampled, recovery is not possible.

Fix the number of observed entries $m = |\Omega|$, would like to get performance bound that holds for almost all sampling patterns.

\implies We'll consider subset of m entries selected uniformly at random.

Which low-rank matrices can we recover?

Compare the following two rank-one matrices:

$$\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix} \quad \text{vs} \quad \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & & \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix} \quad \text{vs} \quad \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}$$

The middle one would be “easier” to complete.

Column and row spaces cannot be aligned with basis vectors.

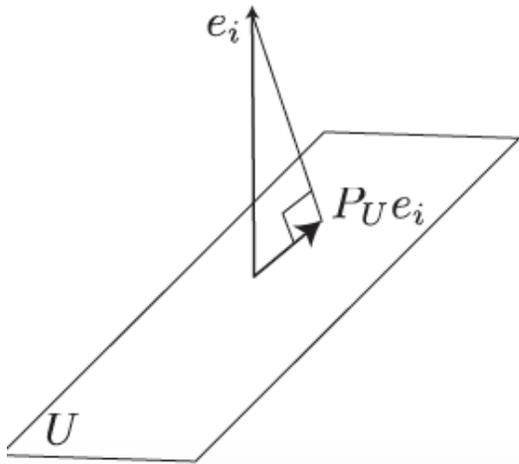
Coherence

Let M be a rank- r matrix with the SVD $M = U\Sigma V^T$, where $U, V \in \mathbb{R}^{n \times r}$.

Definition 4. [Coherence] *Smallest scalar μ obeying*

$$\max_{1 \leq i \leq n} \|U^T e_i\|_2^2 \leq \mu \frac{r}{n}, \quad \max_{1 \leq i \leq n} \|V^T e_i\|_2^2 \leq \mu \frac{r}{n},$$

where e_i is the i th standard basis vector.



We would like $\mu = O(1)$.

- Geometric condition: $U = \text{colspan}(M)$
- Since $\sum_{i=1}^n \|U^T e_i\|_2^2 = r$, $\mu \geq 1$.
- If $e_i \in U$, $\mu = n/r$;
- If $\frac{1}{\sqrt{n}}\mathbf{1} = U$, $\mu = 1$.

Information-theoretic lower bound

Theorem 4. [Candes and Tao, 2009] *No method can succeed with*

$$m \lesssim \mu \times nr \times \log n \approx \text{dof} \times \mu \log n$$

Remarks:

- When $\mu = O(1)$, we need $m \lesssim nr \log n$.
- Need at least one observation /row and column – related to the coupon collector's problem: Suppose that there is an urn of n different coupons, from which coupons are being collected, equally likely, with replacement. How many trials do we need to collect all n coupons?
- The adjacency graph needs to be fully connected

Performance Guarantee

Theorem 5. [Chen, Gross, Recht, Candes and Tao, etc..] *There exists universal constant $c_0, c_1, c_2 > 0$ such that if*

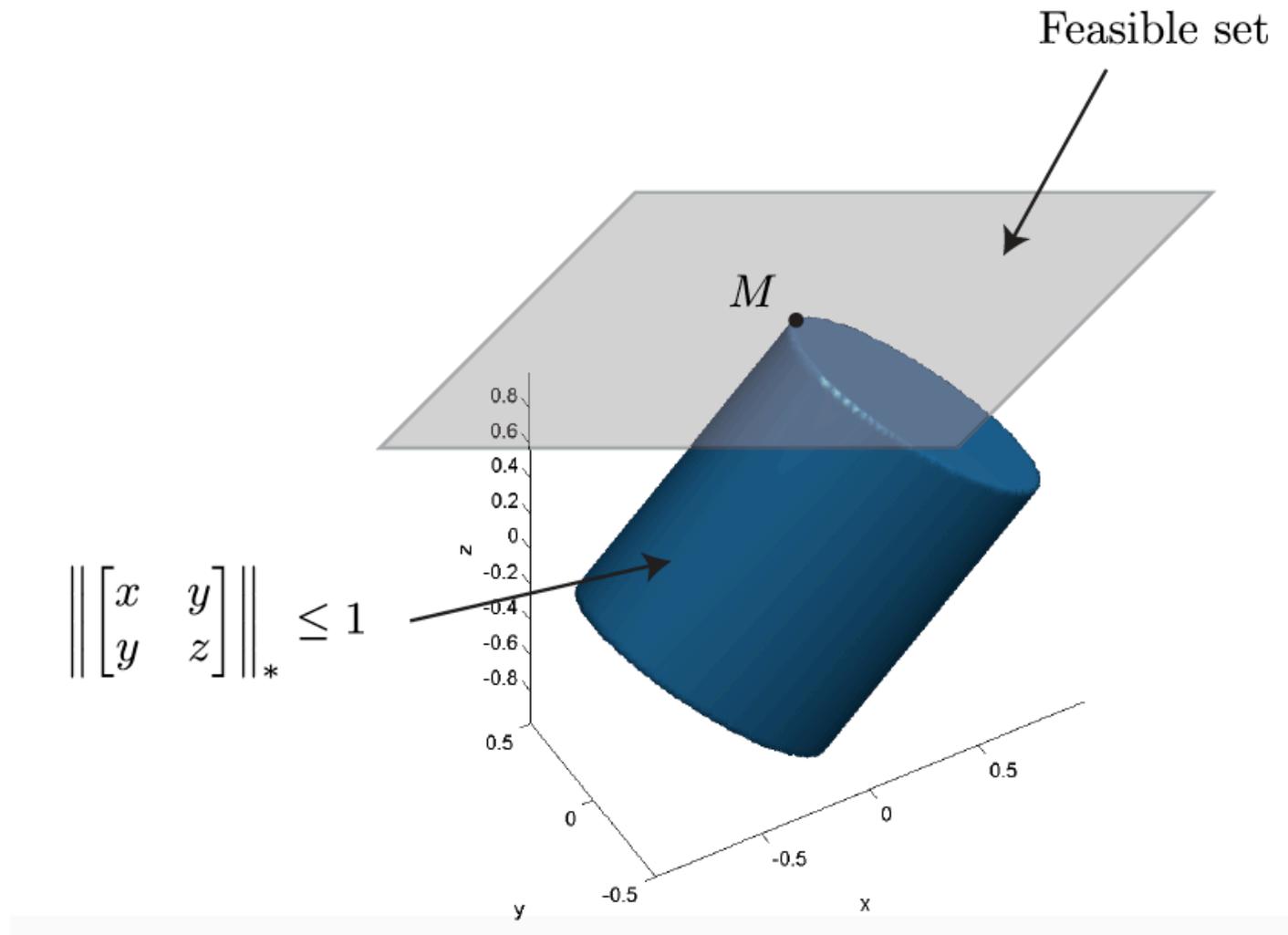
$$m \geq c_0 \mu n r \log^2 n,$$

then M is the unique optimal solution to the nuclear norm minimization problem with probability at least $1 - c_1 n^{-c_2}$.

Remark:

- This result is optimal up to a logarithmic factor in n . See [Chen, Incoherence-Optimal Matrix Completion].

Geometry



A few notations

- P_U is the orthogonal projection to the column space of M ;
- P_V is the orthogonal projection to the row space of M ;
- Let T be the span of matrices of the form:

$$T = \{UX^T + YV^T : X, Y \in \mathbb{R}^{n \times r}\}$$

- Let \mathcal{P}_T be the orthogonal projection onto T :

$$\mathcal{P}_T(\mathbf{X}) = P_U \mathbf{X} + \mathbf{X} P_V - P_U \mathbf{X} P_V$$

- The complement projection $\mathcal{P}_{T^\perp} = \mathcal{I} - \mathcal{P}_T$:

$$\mathcal{P}_{T^\perp}(\mathbf{X}) = (\mathbf{I} - P_U) \mathbf{X} (\mathbf{I} - P_V)$$

Subgradient of the nuclear norm

The subgradient of $\|\cdot\|_*$ at \mathbf{M} can be written as

$$\partial\|\mathbf{M}\|_* = \left\{ \mathbf{UV}^T + \mathbf{W} : \mathcal{P}_T(\mathbf{W}) = 0, \|\mathbf{W}\| \leq 1 \right\}$$

$\mathbf{Z} \in \partial\|\mathbf{M}\|_*$ if and only if

$$\mathcal{P}_T(\mathbf{Z}) = \mathbf{UV}^T, \quad \|\mathcal{P}_{T^\perp}(\mathbf{Z})\| \leq 1.$$

The subgradient doesn't depend on the singular values of \mathbf{M} .

Basic consequence of incoherence

For any $(i, j) \in [n] \times [n]$:

$$\|\mathcal{P}_T(\mathbf{e}_i \mathbf{e}_j^T)\|_F^2 \leq \frac{2\mu r}{n}.$$

The sampling basis is incoherent to the tangent space T .

Sampling with replacement

It turns out it is easier to use a sampling with replacement model, where we assume each observed entry (i_k, j_k) , $k = 1, \dots, m$ is **i.i.d.** observed uniformly at random from $[n] \times [n]$.

This is much easier to analyze, however it is different from the sampling without replacement model stated earlier because we may sample the same entry several times.

Proposition 1. *The probability that the nuclear norm heuristic fails when the set of observed entries is sampled uniformly from the collection of sets of size m is less than or equal to the probability that the heuristic fails when m entries are sampled independently with replacement.*

Bound the number of repetitions

Proposition 2. *With probability at least $1 - n^{2-2\beta}$, the maximum number of repetitions of any entry in Ω is less than $\frac{8}{3}\beta \log(n)$ for $n \geq 9$ and $\beta > 1$.*

Define the operator

$$\mathcal{R}_\Omega(\mathbf{X}) = \sum_{k=1}^m \langle \mathbf{X}, \mathbf{e}_{i_k} \mathbf{e}_{j_k}^T \rangle \mathbf{e}_{i_k} \mathbf{e}_{j_k}^T = \sum_{k=1}^m X_{i_k, j_k} \mathbf{e}_{i_k} \mathbf{e}_{j_k}^T$$

where (i_k, j_k) is uniformly drawn from $[n] \times [n]$. From the above proposition, we have

$$\|\mathcal{R}_\Omega\| \leq \frac{8}{3}\beta \log(n)$$

with probability at least $1 - n^{2-2\beta}$.

Optimality condition

Proposition 3. [Exact Dual Certificate] *M is the unique minimizer of the nuclear norm minimization problem if the following holds:*

- *the sampling operator \mathcal{P}_Ω restricted to elements in T is injective;*
- *there exists \mathbf{Y} supported on Ω such that $\mathbf{Y} \in \partial\|\mathbf{M}\|_*$, i.e.*

$$\mathcal{P}_T(\mathbf{Y}) = \mathbf{U}\mathbf{V}^T, \quad \|\mathcal{P}_{T^\perp}(\mathbf{Y})\| \leq 1.$$

The first equality constraint is not easy to satisfy, see [Candes and Tao, 2009].

Dual certificate

Under a stronger injectivity requirement, we can relax the second requirement a bit, which much simplifies the analysis.

Proposition 4. [Inexact Dual Certificate] *Suppose that*

$$\frac{n^2}{m} \left\| \mathcal{P}_T \mathcal{R}_\Omega \mathcal{P}_T - \frac{m}{n^2} \mathcal{P}_T \right\| \leq \frac{1}{2},$$

and there exists \mathbf{Y} supported on Ω such that

$$\left\| \mathcal{P}_T(\mathbf{Y}) - \mathbf{UV}^T \right\|_F \leq \sqrt{\frac{r}{2n}}, \quad \left\| \mathcal{P}_{T^\perp}(\mathbf{Y}) \right\| < \frac{1}{2},$$

then \mathbf{M} is the unique minimizer of the nuclear norm minimization problem if the following holds:

Injectivity of \mathcal{R}_Ω on T

Proposition 5. For all $\beta > 1$,

$$\frac{n^2}{m} \left\| \mathcal{P}_T \mathcal{R}_\Omega \mathcal{P}_T - \frac{m}{n^2} \mathcal{P}_T \right\| \leq \sqrt{\frac{32\beta\mu_0 nr \log n}{3m}}$$

with probability at least $1 - 2n^{2-2\beta}$ provided $m \geq \frac{32}{3}\beta\mu_0 nr \log n$.

Remark: Provided $\sqrt{\frac{32\beta\mu_0 nr \log n}{3m}} \leq \frac{1}{2}$, i.e.

$$m \geq \frac{128\beta\mu_0 nr \log n}{3}$$

we have with probability at least $1 - 2n^{2-2\beta}$,

$$\frac{n^2}{m} \left\| \mathcal{P}_T \mathcal{R}_\Omega \mathcal{P}_T - \frac{m}{n^2} \mathcal{P}_T \right\| \leq \frac{1}{2}.$$

Constructing the dual certificate

We introduce the clever golfing scheme proposed by David Gross.

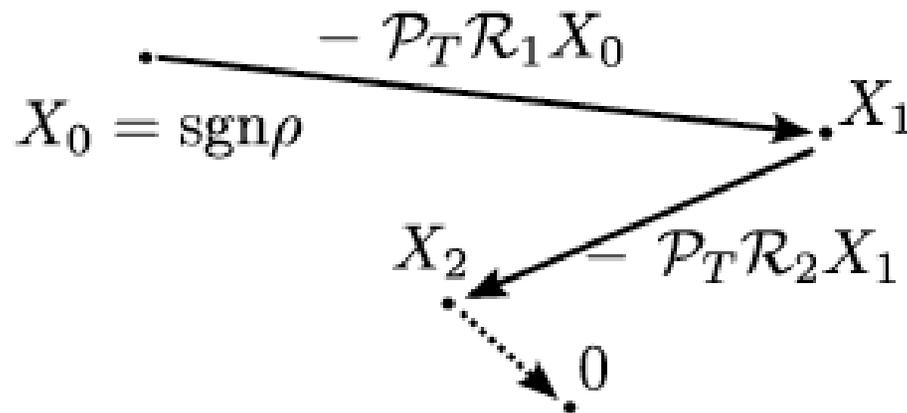


Fig. 3. Caricature of the “golfing scheme” used to construct the certificate. In the i th step, X_{i-1} designates the vector we aim to represent. The approximation of X_{i-1} actually obtained is $\mathcal{P}_T \mathcal{R}_i X_{i-1}$. The distance of the new goal $X_i = X_{i-1} - \mathcal{P}_T \mathcal{R}_i X_{i-1}$ to the origin is guaranteed to be only half the previous one. The sequence X_i thus converges exponentially fast to the origin.

Computational aspect using FISTA

Recall the FISTA algorithm we discussed to solve

$$\hat{\mathbf{M}} = \underset{\mathbf{X}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2 + \lambda \|\mathbf{X}\|_*$$

- Initialization: $\mathbf{x}_0 = \mathbf{x}_{-1} \in \mathbb{R}^n$, $\theta_0 = 1$,
- For $k = 1, 2, \dots$,

$$\theta_k = \frac{1 + \sqrt{1 + 4\theta_{k-1}^2}}{2}$$

$$\mathbf{Y}_k = \mathbf{X}_{k-1} + \left(\frac{\theta_{k-1} - 1}{\theta_k} \right) (\mathbf{X}_{k-1} - \mathbf{X}_{k-2})$$

$$\mathbf{X}_k = \operatorname{prox}_{t_k \lambda \|\cdot\|_*} (\mathbf{Y}_k - t_k \mathcal{A}^*(\mathcal{A}(\mathbf{Y}_k) - \mathbf{y}))$$

- What is the proximal operator for $\|\cdot\|_*$?

Proximal operator for $\|\cdot\|_*$

Proposition 6.

$$\text{prox}_{t_k \lambda \|\cdot\|_*}(\mathbf{X}) = \underset{\mathbf{Z}}{\text{argmin}} \left\{ \frac{1}{2} \|\mathbf{Z} - \mathbf{X}\|_2^2 + t_k \lambda \|\mathbf{Z}\|_* \right\} = \mathcal{T}_{t_k \lambda}(\mathbf{X})$$

where

$$\mathcal{T}_\tau(\mathbf{X}) = \mathbf{U} \mathcal{T}_\tau(\mathbf{\Sigma}) \mathbf{V}^T,$$

where the SVD of \mathbf{X} is given as $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$, $\mathbf{\Sigma} = \text{diag}(\{\sigma_i\})$, and

$$\mathcal{T}_\tau(\mathbf{\Sigma}) = \text{diag}(\{(\sigma_i - \tau)_+\}).$$