

ECE 18-898G: Special Topics in Signal Processing: Sparsity, Structure, and Inference

Low-rank matrix recovery via convex relaxations

Yuejie Chi

Department of Electrical and Computer Engineering

Carnegie Mellon University

Spring 2018

Outline

- Low-rank matrix completion and recovery
- Nuclear norm minimization (this lecture)
 - RIP and low-rank matrix recovery
 - Matrix completion
 - Algorithms for nuclear norm minimization
- Non-convex methods (next lecture)
 - Spectral methods
 - (Projected) gradient descent

Low-rank matrix completion and recovery: motivation

Motivation 1: recommendation systems

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

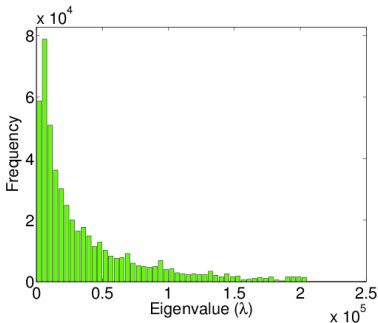
- Netflix challenge: Netflix provides highly incomplete ratings from 0.5 million users for & 17,770 movies
- How to predict unseen user ratings for movies?

In general, we cannot infer missing ratings

✓	?	?	?	✓	?
?	?	✓	✓	?	?
✓	?	?	✓	?	?
?	?	✓	?	?	✓
✓	?	?	?	?	?
?	✓	?	?	✓	?
?	?	✓	✓	?	?

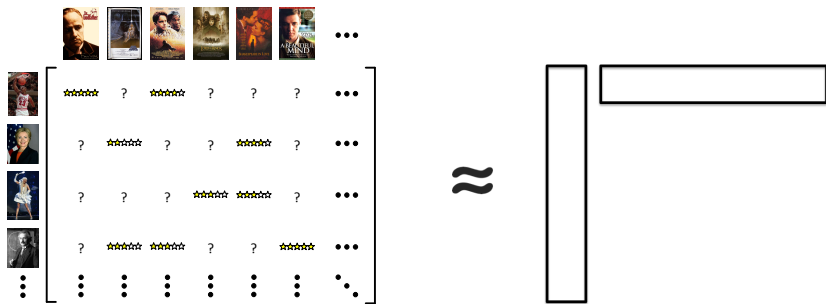
Underdetermined system (more unknowns than observations)

... unless rating matrix has other structure



A few factors explain most of the data

... unless rating matrix has other structure



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

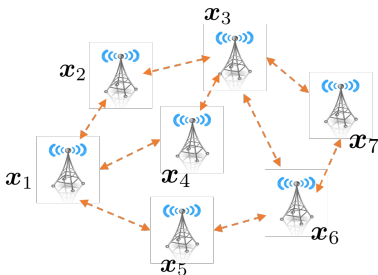
How to exploit (approx.) low-rank structure in prediction?

Motivation 2: sensor localization

- n sensors / points $\mathbf{x}_j \in \mathbb{R}^3$, $j = 1, \dots, n$
- Observe partial information about pairwise distances

$$D_{i,j} = \|\mathbf{x}_i - \mathbf{x}_j\|^2 = \|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2 - 2\mathbf{x}_i^\top \mathbf{x}_j$$

- Want to infer distance between every pair of nodes



Motivation 2: sensor localization

Introduce

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

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

$$\mathbf{D} = \underbrace{\mathbf{d}_2 \mathbf{e}^\top + \mathbf{e} \mathbf{d}_2^\top - 2\mathbf{X}\mathbf{X}^\top}_{\text{low rank}}$$

where $\mathbf{d}_2 := [\|\mathbf{x}_1\|^2, \dots, \|\mathbf{x}_n\|^2]^\top$

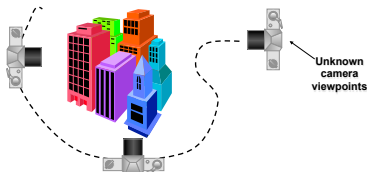
$\text{rank}(\mathbf{D}) \ll n \quad \longrightarrow \quad \text{low-rank matrix completion}$

Motivation 3: structure from motion

Structure from motion: reconstruct 3D scene geometry and camera poses from multiple images

camera poses **motion**

3D scene geometry **structure**



Given multiple images and a few correspondences between image features, how to estimate locations of 3D points?



Motivation 3: structure from motion

Tomasi and Kanade's factorization:

- Consider n 3D points in m different 2D frames
- $\mathbf{x}_{i,j} \in \mathbb{R}^{2 \times 1}$: locations of j^{th} point in i^{th} frame

$$\mathbf{x}_{i,j} = \underbrace{\mathbf{P}_i}_{\text{projection matrix } \in \mathbb{R}^{2 \times 3}} \underbrace{\mathbf{s}_j}_{\text{3D position } \in \mathbb{R}^3}$$

- Matrix of all 2D locations $\text{rank}(\mathbf{X}) = 3$.

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_{1,1} & \cdots & \mathbf{x}_{1,n} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{m,1} & \cdots & \mathbf{x}_{m,n} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{P}_1 \\ \vdots \\ \mathbf{P}_m \end{bmatrix} \begin{bmatrix} \mathbf{s}_1 & \cdots & \mathbf{s}_n \end{bmatrix}}_{\text{low-rank factorization}} \in \mathbb{R}^{2m \times n}$$

Due to occlusions, there are many missing entries in the matrix \mathbf{X} .

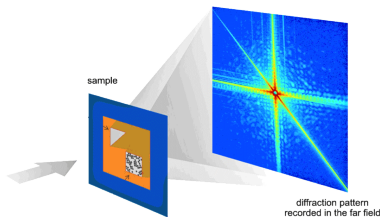
Goal: Can we complete the missing entries?

Motivation 4: missing phase problem

Detectors record **intensities** of diffracted rays

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

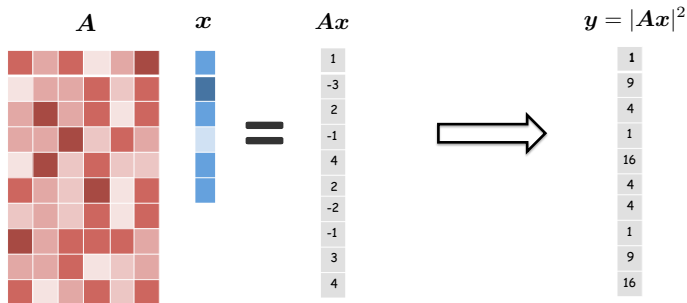
Fig credit: Stanford SLAC



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

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

A discrete-time model: solving quadratic systems



Solve for $x \in \mathbb{C}^n$ in m quadratic equations

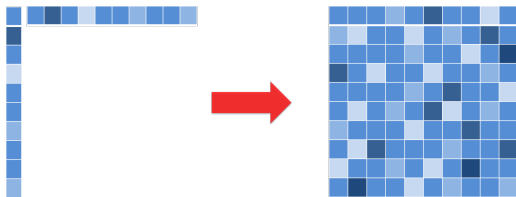
$$y_k = |\langle a_k, x \rangle|^2, \quad k = 1, \dots, m$$

or $y = |Ax|^2$ where $|z|^2 := \{|z_1|^2, \dots, |z_m|^2\}$

An equivalent view: low-rank factorization

Lifting: introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^*$ to linearize constraints

$$y_k = |\mathbf{a}_k^* \mathbf{x}|^2 = \mathbf{a}_k^* (\mathbf{x}\mathbf{x}^*) \mathbf{a}_k \quad \implies \quad y_k = \mathbf{a}_k^* \mathbf{X} \mathbf{a}_k \quad (6.1)$$



$$\begin{aligned} \text{find} \quad & \mathbf{X} \succeq \mathbf{0} \\ \text{s.t.} \quad & y_k = \langle \mathbf{a}_k \mathbf{a}_k^*, \mathbf{X} \rangle, \quad k = 1, \dots, m \\ & \text{rank}(\mathbf{X}) = 1 \end{aligned}$$

The list continues

- system identification and time series analysis;
- spatial-temporal data: low-rank due to correlations, e.g. MRI video, network traffic, ..
- face recognition;
- quantum state tomography;
- community detection;
-

Low-rank matrix completion and recovery: setup and algorithms

Setup

- Consider $M \in \mathbb{R}^{n \times n}$ (square case for simplicity)
- $\text{rank}(M) = r \ll n$
- The **thin** Singular value decomposition (SVD) of M :

$$M = \underbrace{U \Sigma V^T}_{(2n-r)r \text{ degrees of freedom}} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

where $\Sigma = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix}$ contain all singular values $\{\sigma_i\}$;

$U := [\mathbf{u}_1, \dots, \mathbf{u}_r]$, $V := [\mathbf{v}_1, \dots, \mathbf{v}_r]$ consist of singular vectors

Low-rank matrix completion

Observed entries

$$M_{i,j}, \quad (i,j) \in \underbrace{\Omega}_{\text{sampling set}}$$

Completion via rank minimization

$$\text{minimize}_{\mathbf{X}} \text{rank}(\mathbf{X}) \quad \text{s.t.} \quad X_{i,j} = M_{i,j}, \quad (i,j) \in \Omega$$

Low-rank matrix completion

Observed entries

$$M_{i,j}, \quad (i,j) \in \underbrace{\Omega}_{\text{sampling set}}$$

- An operator \mathcal{P}_Ω : orthogonal projection onto subspace of matrices supported on Ω

$$[\mathcal{P}_\Omega(\mathbf{X})]_{i,j} = \begin{cases} X_{i,j} & \text{if } (i,j) \in \Omega \\ 0 & \text{otherwise} \end{cases}$$

Completion via rank minimization

$$\text{minimize}_{\mathbf{X}} \text{rank}(\mathbf{X}) \quad \text{s.t.} \quad \mathcal{P}_\Omega(\mathbf{X}) = \mathcal{P}_\Omega(\mathbf{M})$$

More general: low-rank matrix recovery

Linear measurements

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

- An operator form, with $\mathcal{A} : \mathbb{R}^{n \times n} \mapsto \mathbb{R}^m$:

$$\mathbf{y} = \mathcal{A}(\mathbf{M}) := \begin{bmatrix} \langle \mathbf{A}_1, \mathbf{M} \rangle \\ \vdots \\ \langle \mathbf{A}_m, \mathbf{M} \rangle \end{bmatrix}$$

Recovery via rank minimization

$$\text{minimize}_{\mathbf{X}} \text{rank}(\mathbf{X}) \quad \text{s.t.} \quad \mathbf{y} = \mathcal{A}(\mathbf{X})$$

Nuclear norm minimization

Convex relaxation

Low-rank matrix completion:

$$\begin{array}{ll} \text{minimize}_{\mathbf{X}} & \underbrace{\text{rank}(\mathbf{X})}_{\text{nonconvex}} \\ \text{s.t.} & \mathcal{P}_{\Omega}(\mathbf{X}) = \mathcal{P}_{\Omega}(\mathbf{M}) \end{array}$$

Low-rank matrix recovery:

$$\begin{array}{ll} \text{minimize}_{\mathbf{X}} & \underbrace{\text{rank}(\mathbf{X})}_{\text{nonconvex}} \\ \text{s.t.} & \mathcal{A}(\mathbf{X}) = \mathcal{A}(\mathbf{M}) \end{array}$$

Question: what is convex surrogate for $\text{rank}(\cdot)$?

Analogy with sparse recovery

For a given matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$, take the full SVD of \mathbf{X} :

$$\begin{aligned}\mathbf{X} &= \widehat{\mathbf{U}}\widehat{\mathbf{\Sigma}}\widehat{\mathbf{V}}^\top \\ &= \begin{bmatrix} \mathbf{u}_1, \dots, \mathbf{u}_n \end{bmatrix} \underbrace{\begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix}}_{\text{diag}(\boldsymbol{\sigma})} \begin{bmatrix} \mathbf{v}_1, \dots, \mathbf{v}_n \end{bmatrix}^\top.\end{aligned}$$

Then

$$\text{rank}(\mathbf{X}) = \sum_{i=1}^n \mathbf{1}\{\sigma_i \neq 0\} = \|\boldsymbol{\sigma}\|_0$$

Convex relaxation: from $\|\boldsymbol{\sigma}\|_0$ to $\|\boldsymbol{\sigma}\|_1$?

Nuclear norm

Definition 6.1

The nuclear norm of \mathbf{X} is

$$\|\mathbf{X}\|_* := \sum_{i=1}^n \underbrace{\sigma_i(\mathbf{X})}_{i^{\text{th}} \text{ largest singular value}}$$

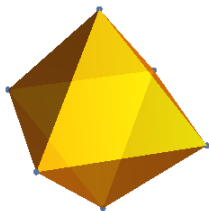
- Nuclear norm is a counterpart of ℓ_1 norm for rank function
- Equivalence among different norms ($r = \text{rank}(\mathbf{X})$)

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

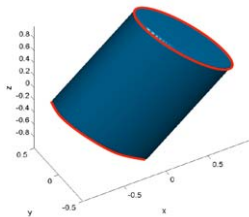
where $\|\mathbf{X}\| = \sigma_1(\mathbf{X})$; $\|\mathbf{X}\|_F = (\sum_{i=1}^n \sigma_i^2(\mathbf{X}))^{1/2}$.

Tightness of relaxation

Recall: the ℓ_1 norm ball is convex hull of 1-sparse, unit-norm vectors.



(a) ℓ_1 norm ball



(b) nuclear norm ball

Fact 6.2

The nuclear norm ball $\{\mathbf{X} : \|\mathbf{X}\|_ \leq 1\}$ is the convex hull of rank-1 matrices, unit-norm matrices obeying $\|\mathbf{u}\mathbf{v}^\top\| = 1$.*

Additivity of nuclear norm

Fact 6.3

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

- If row (resp. column) spaces of A and B are orthogonal, then $\|A + B\|_* = \|A\|_* + \|B\|_*$
- Similar to ℓ_1 norm: when x and y have disjoint support,

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

which is a key to study ℓ_1 -min under RIP.

Proof of Fact 6.3

Suppose $\mathbf{A} = \mathbf{U}_A \boldsymbol{\Sigma}_A \mathbf{V}_A^\top$ and $\mathbf{B} = \mathbf{U}_B \boldsymbol{\Sigma}_B \mathbf{V}_B^\top$, which gives

$$\begin{aligned} \mathbf{A}\mathbf{B}^\top &= \mathbf{0} \\ \mathbf{A}^\top \mathbf{B} &= \mathbf{0} \end{aligned} \iff \begin{aligned} \mathbf{V}_A^\top \mathbf{V}_B &= \mathbf{0} \\ \mathbf{U}_A^\top \mathbf{U}_B &= \mathbf{0} \end{aligned}$$

Thus, one can write

$$\begin{aligned} \mathbf{A} &= [\mathbf{U}_A, \mathbf{U}_B, \mathbf{U}_C] \begin{bmatrix} \boldsymbol{\Sigma}_A & & \\ & \mathbf{0} & \\ & & \mathbf{0} \end{bmatrix} [\mathbf{V}_A, \mathbf{V}_B, \mathbf{V}_C]^\top \\ \mathbf{B} &= [\mathbf{U}_A, \mathbf{U}_B, \mathbf{U}_C] \begin{bmatrix} & & \\ \mathbf{0} & \boldsymbol{\Sigma}_B & \\ & & \mathbf{0} \end{bmatrix} [\mathbf{V}_A, \mathbf{V}_B, \mathbf{V}_C]^\top \end{aligned}$$

and hence

$$\|\mathbf{A} + \mathbf{B}\|_* = \left\| [\mathbf{U}_A, \mathbf{U}_B] \begin{bmatrix} \boldsymbol{\Sigma}_A & \\ & \boldsymbol{\Sigma}_B \end{bmatrix} [\mathbf{V}_A, \mathbf{V}_B]^\top \right\|_* = \|\mathbf{A}\|_* + \|\mathbf{B}\|_*$$

Dual norm

Definition 6.4 (Dual norm)

For a given norm $\|\cdot\|_{\mathcal{A}}$, the dual norm is defined as

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

- ℓ_1 norm $\xleftrightarrow{\text{dual}}$ ℓ_∞ norm
- ℓ_2 norm $\xleftrightarrow{\text{dual}}$ ℓ_2 norm
- Frobenius norm $\xleftrightarrow{\text{dual}}$ Frobenius norm
- nuclear norm $\xleftrightarrow{\text{dual}}$ spectral norm

Schur complement

Given a block matrix

$$D = \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix} \in \mathbb{R}^{(p+q) \times (p+q)},$$

where $A \in \mathbb{R}^{p \times p}$, $B \in \mathbb{R}^{p \times q}$ and $C \in \mathbb{R}^{q \times q}$.

The **Schur complement** of the block A (assume it's invertible) in D is given as

$$C - B^\top A^{-1} B.$$

Fact 6.5

$$D \succeq 0 \iff A \succeq 0, C - B^\top A^{-1} B \succeq 0$$

Representing nuclear norm via SDP

Since spectral norm is dual norm of nuclear norm,

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

The constraint is equivalent to

$$\|\mathbf{Y}\| \leq 1 \iff \mathbf{Y}\mathbf{Y}^\top \preceq \mathbf{I} \quad \text{Schur complement} \iff \begin{bmatrix} \mathbf{I} & \mathbf{Y} \\ \mathbf{Y}^\top & \mathbf{I} \end{bmatrix} \succeq \mathbf{0}$$

Fact 6.6

$$\|\mathbf{X}\|_* = \max_{\mathbf{Y}} \left\{ \langle \mathbf{X}, \mathbf{Y} \rangle \mid \begin{bmatrix} \mathbf{I} & \mathbf{Y} \\ \mathbf{Y}^\top & \mathbf{I} \end{bmatrix} \succeq \mathbf{0} \right\}$$

Representing nuclear norm via SDP

Since spectral norm is dual norm of nuclear norm,

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

The constraint is equivalent to

$$\|\mathbf{Y}\| \leq 1 \iff \mathbf{Y}\mathbf{Y}^\top \preceq \mathbf{I} \quad \text{Schur complement} \iff \begin{bmatrix} \mathbf{I} & \mathbf{Y} \\ \mathbf{Y}^\top & \mathbf{I} \end{bmatrix} \succeq \mathbf{0}$$

Fact 6.7 (Dual characterization)

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

- Optimal point: $\mathbf{W}_1 = \mathbf{U}\Sigma\mathbf{U}^\top$, $\mathbf{W}_2 = \mathbf{V}\Sigma\mathbf{V}^\top$ (where $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$)

Aside: dual of semidefinite program

$$\begin{aligned} \text{(primal)} \quad & \text{minimize}_{\mathbf{X}} && \langle \mathbf{C}, \mathbf{X} \rangle \\ & \text{s.t.} && \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i, \quad 1 \leq i \leq m \\ & && \mathbf{X} \succeq \mathbf{0} \end{aligned}$$

\Leftrightarrow

$$\begin{aligned} \text{(dual)} \quad & \text{maximize}_{\mathbf{y}} && \mathbf{b}^\top \mathbf{y} \\ & \text{s.t.} && \sum_{i=1}^m y_i \mathbf{A}_i + \mathbf{S} = \mathbf{C} \\ & && \mathbf{S} \succeq \mathbf{0} \end{aligned}$$

Exercise: use this to verify Fact 6.7

Nuclear norm minimization via SDP

Nuclear norm minimization

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

This is solvable via SDP

$$\begin{aligned} & \operatorname{minimize}_{\mathbf{X}, \mathbf{W}_1, \mathbf{W}_2} && \frac{1}{2} \operatorname{Tr}(\mathbf{W}_1) + \frac{1}{2} \operatorname{Tr}(\mathbf{W}_2) \\ & \text{s.t.} && \mathbf{y} = \mathcal{A}(\mathbf{X}), \\ & && \begin{bmatrix} \mathbf{W}_1 & \mathbf{X} \\ \mathbf{X}^\top & \mathbf{W}_2 \end{bmatrix} \succeq \mathbf{0} \end{aligned}$$

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.

Fig. credit: Fazel et.al. '10

Proximal algorithm

In the presence of noise, one needs to solve

$$\text{minimize}_{\mathbf{X}} \quad \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_{\text{F}}^2 + \lambda \|\mathbf{X}\|_*$$

which can be solved via proximal methods.

Algorithm 6.1 Proximal gradient methods

for $t = 0, 1, \dots$:

$$\mathbf{X}^{t+1} = \text{prox}_{\mu_t \lambda \|\cdot\|_*} \left(\mathbf{X}^t - \mu_t \mathcal{A}^*(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}) \right)$$

where μ_t : step size / learning rate

Proximal operator for nuclear norm

Proximal operator:

$$\begin{aligned}\text{prox}_{\lambda\|\cdot\|_*}(\mathbf{X}) &= \arg \min_{\mathbf{Z}} \left\{ \frac{1}{2} \|\mathbf{Z} - \mathbf{X}\|_{\text{F}}^2 + \lambda \|\mathbf{Z}\|_* \right\} \\ &= \mathbf{U} \mathcal{T}_{\lambda}(\mathbf{\Sigma}) \mathbf{V}^{\top}\end{aligned}$$

where SVD of \mathbf{X} is $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\top}$ with $\mathbf{\Sigma} = \text{diag}(\{\sigma_i\})$, and

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

Accelerated proximal gradient

Algorithm 6.2 Accelerated proximal gradient methods

for $t = 0, 1, \dots$:

$$\mathbf{X}^{t+1} = \text{prox}_{\mu_t \lambda \|\cdot\|_*} \left(\mathbf{Z}^t - \mu_t \mathcal{A}^*(\mathcal{A}(\mathbf{Z}^t) - \mathbf{y}) \right)$$

$$\mathbf{Z}^{t+1} = \mathbf{X}^{t+1} + \underbrace{\alpha_t (\mathbf{X}^{t+1} - \mathbf{X}^t)}_{\text{momentum term}}$$

where μ_t : step size / learning rate, α_t is the momentum.

- Convergence rate: $O\left(\frac{1}{\sqrt{\epsilon}}\right)$ iterations to reach ϵ -accuracy.
- Per-iteration cost is a (partial-)SVD.

Frank-Wolfe for nuclear norm minimization

Consider the constrained problem:

$$\text{minimize}_{\mathbf{X}} \quad \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_{\text{F}}^2 \quad \text{s.t.} \quad \|\mathbf{X}\|_* \leq \tau.$$

which can be solved via conditional gradient method (Frank-Wolfe 1956).

- More generally, consider the problem:

$$\text{minimize}_{\beta} \quad f(\beta) \quad \text{s.t.} \quad \beta \in \mathcal{C},$$

where $f(x)$ is smooth and convex, and \mathcal{C} is a convex set.

- Recall projected gradient descent:

$$\beta^{t+1} = \mathcal{P}_{\mathcal{C}} \left(\beta^t - \mu_t \nabla f(\beta^t) \right)$$

Conditional Gradient Method

Algorithm 6.3 Frank-Wolfe

for $t = 0, 1, \dots$:

$$\mathbf{s}^t = \operatorname{argmin}_{\mathbf{s} \in \mathcal{C}} \nabla f(\boldsymbol{\beta}^t)^\top \mathbf{s}$$

$$\boldsymbol{\beta}^{t+1} = (1 - \gamma_t)\boldsymbol{\beta}^t + \gamma_t \mathbf{s}^t$$

where $\gamma_t := 2/(t + 1)$ is the (default) step size.

- The first step is a constrained optimization of a linear approximation at $f(\boldsymbol{\beta}^t)$;
- The second step controls how much we move towards \mathbf{s}^t .

Figure illustration

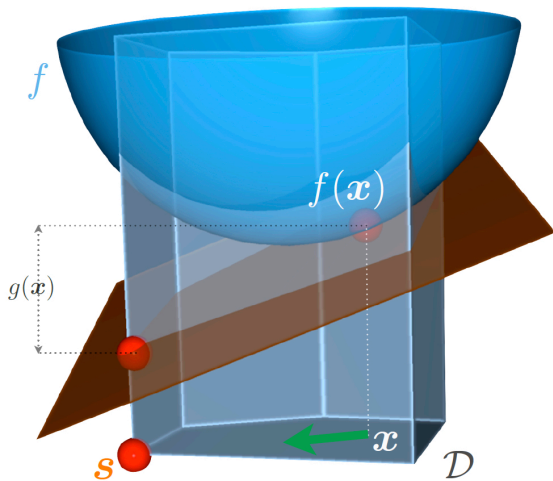


Figure credit: Jaggi 2011

Frank-Wolfe for nuclear norm minimization

Algorithm 6.4 Frank-Wolfe for nuclear norm minimization

for $t = 0, 1, \dots$:

$$\mathbf{S}^t = \operatorname{argmin}_{\|\mathbf{S}\|_* \leq \tau} \langle \nabla f(\mathbf{X}^t), \mathbf{S} \rangle,$$

$$\mathbf{X}^{t+1} = (1 - \gamma_t) \mathbf{X}^t + \gamma_t \mathbf{S}^t;$$

where $\gamma_t := 2/(t + 1)$ is the (default) step size.

- (Homework) Note that $\nabla f(\mathbf{X}^t) = \mathcal{A}^*(\mathcal{A}(\mathbf{X}^t) - \mathbf{y})$, and

$$\begin{aligned} \mathbf{S}^t &= \tau \cdot \operatorname{argmin}_{\|\mathbf{S}\|_* \leq 1} \langle \nabla f(\mathbf{X}^t), \mathbf{S} \rangle \\ &= \tau \mathbf{u} \mathbf{v}^T, \end{aligned}$$

where \mathbf{u} , \mathbf{v} are the left and right top singular vector of $-\nabla f(\mathbf{X}^t)$.

Further comments on Frank-Wolfe

- Extremely low per-iteration cost (only top singular vectors are needed);
- Every iteration is a rank-1 update;
- Convergence rate: $O(\frac{1}{\epsilon})$ to reach ϵ -accuracy, which can be very slow.
- Various ways to speed up; active research area.

RIP and low-rank matrix recovery

RIP for low-rank matrices

Almost parallel results to compressed sensing ...¹

Definition 6.8

The r -restricted isometry constants $\delta_r^{\text{ub}}(\mathcal{A})$ and $\delta_r^{\text{lb}}(\mathcal{A})$ are smallest quantities s.t.

$$(1 - \delta_r^{\text{lb}}) \|\mathbf{X}\|_F \leq \|\mathcal{A}(\mathbf{X})\|_F \leq (1 + \delta_r^{\text{ub}}) \|\mathbf{X}\|_F, \quad \forall \mathbf{X} : \text{rank}(\mathbf{X}) \leq r$$

¹One can also define RIP w.r.t. $\|\cdot\|_F^2$ rather than $\|\cdot\|_F$.

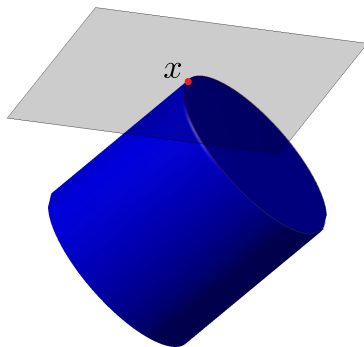
RIP and low-rank matrix recovery

Theorem 6.9 (Recht, Fazel, Parrilo '10, Candes, Plan '11)

Suppose $\text{rank}(\mathbf{M}) = r$. For any fixed integer $K > 0$, if $\frac{1 + \delta_{Kr}^{\text{ub}}}{1 - \delta_{(2+K)r}^{\text{lb}}} < \sqrt{\frac{K}{2}}$, then nuclear norm minimization is exact

- Can be easily extended to account for noise and imperfect structural assumption

Geometry of nuclear norm ball



Level set of nuclear norm ball: $\left\| \begin{bmatrix} x & y \\ y & z \end{bmatrix} \right\|_* \leq 1$

Fig. credit: Candes '14

Some notation

Recall $M = U\Sigma V^\top$

- Let T be span of matrices of the form (called *tangent space*)

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

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

$$\mathcal{P}_T(\mathbf{X}) = UU^\top \mathbf{X} + \mathbf{X}VV^\top - UU^\top \mathbf{X}VV^\top$$

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

$$\mathcal{P}_{T^\perp}(\mathbf{X}) = (\mathcal{I} - UU^\top)\mathbf{X}(\mathcal{I} - VV^\top)$$

- $M\mathcal{P}_{T^\perp}^\top(\mathbf{X}) = \mathbf{0}$ and $M^\top\mathcal{P}_{T^\perp}(\mathbf{X}) = \mathbf{0}$

Proof of Theorem 6.9

Suppose $\mathbf{X} = \mathbf{M} + \mathbf{H}$ is feasible and obeys $\|\mathbf{M} + \mathbf{H}\|_* \leq \|\mathbf{M}\|_*$.
The goal is to show that $\mathbf{H} = \mathbf{0}$ under RIP.

The key is to decompose \mathbf{H} into $\mathbf{H}_0 + \underbrace{\mathbf{H}_1 + \mathbf{H}_2 + \dots}_{\mathbf{H}_c}$.

- $\mathbf{H}_0 = \mathcal{P}_T(\mathbf{H})$ (rank $2r$)
- $\mathbf{H}_c = \mathcal{P}_T^\perp(\mathbf{H})$ (obeying $\mathbf{M}\mathbf{H}_c^\top = \mathbf{0}$ and $\mathbf{M}^\top\mathbf{H}_c = \mathbf{0}$)
- \mathbf{H}_1 : best rank- (Kr) approximation of \mathbf{H}_c (K is const)
- \mathbf{H}_2 : best rank- (Kr) approximation of $\mathbf{H}_c - \mathbf{H}_1$
- ...

Proof of Theorem 6.9

Informally, the proof proceeds by showing that

1. \mathbf{H}_0 dominates $\sum_{i \geq 2} \mathbf{H}_i$ (by objective function)
2. (converse) $\sum_{i \geq 2} \mathbf{H}_i$ dominates $\mathbf{H}_0 + \mathbf{H}_1$ (by RIP + feasibility)

These can happen simultaneously only when $\mathbf{H} = \mathbf{0}$

Proof of Theorem 6.9

Step 1 (which does not rely on RIP). Show that

$$\sum_{j \geq 2} \|\mathbf{H}_j\|_{\text{F}} \leq \|\mathbf{H}_0\|_* / \sqrt{Kr}. \quad (6.2)$$

This follows immediately by combining the following 2 observations:

(i) Since $\mathbf{M} + \mathbf{H}$ is assumed to be a better estimate:

$$\begin{aligned} \|\mathbf{M}\|_* &\geq \|\mathbf{M} + \mathbf{H}\|_* \geq \|\mathbf{M} + \mathbf{H}_c\|_* - \|\mathbf{H}_0\|_* \\ &= \underbrace{\|\mathbf{M}\|_* + \|\mathbf{H}_c\|_*}_{\text{Fact 6.3 } (\mathbf{M}\mathbf{H}_c^{\top}=\mathbf{0} \text{ and } \mathbf{M}^{\top}\mathbf{H}_c=\mathbf{0})} - \|\mathbf{H}_0\|_* \end{aligned} \quad (6.3)$$

$$\implies \|\mathbf{H}_c\|_* \leq \|\mathbf{H}_0\|_* \quad (6.4)$$

(ii) Since nonzero singular values of \mathbf{H}_{j-1} dominate those of \mathbf{H}_j ($j \geq 2$):

$$\begin{aligned} \|\mathbf{H}_j\|_{\text{F}} &\leq \sqrt{Kr} \|\mathbf{H}_j\| \leq \sqrt{Kr} [\|\mathbf{H}_{j-1}\|_* / (Kr)] \leq \|\mathbf{H}_{j-1}\|_* / \sqrt{Kr} \\ \implies \sum_{j \geq 2} \|\mathbf{H}_j\|_{\text{F}} &\leq \frac{1}{\sqrt{Kr}} \sum_{j \geq 2} \|\mathbf{H}_{j-1}\|_* = \frac{1}{\sqrt{Kr}} \|\mathbf{H}_c\|_* \end{aligned} \quad (6.5)$$

Proof of Theorem 6.9

Step 2 (using feasibility + RIP). Show that $\exists \rho < \sqrt{K/2}$ s.t.

$$\|\mathbf{H}_0 + \mathbf{H}_1\|_F \leq \rho \sum_{j \geq 2} \|\mathbf{H}_j\|_F \quad (6.6)$$

If this claim holds, then

$$\begin{aligned} \|\mathbf{H}_0 + \mathbf{H}_1\|_F &\leq \rho \sum_{j \geq 2} \|\mathbf{H}_j\|_F \stackrel{(6.2)}{\leq} \rho \frac{1}{\sqrt{Kr}} \|\mathbf{H}_0\|_* \\ &\leq \rho \frac{1}{\sqrt{Kr}} \left(\sqrt{2r} \|\mathbf{H}_0\|_F \right) = \rho \sqrt{\frac{2}{K}} \|\mathbf{H}_0\|_F \\ &\leq \rho \sqrt{\frac{2}{K}} \|\mathbf{H}_0 + \mathbf{H}_1\|_F. \end{aligned} \quad (6.7)$$

This cannot hold with $\rho < \sqrt{K/2}$ unless $\underbrace{\mathbf{H}_0 + \mathbf{H}_1}_{\text{equivalently, } \mathbf{H}_0 = \mathbf{H}_1 = \mathbf{0}} = \mathbf{0}$

Proof of Theorem 6.9

We now prove (6.6). To connect $\mathbf{H}_0 + \mathbf{H}_1$ with $\sum_{j \geq 2} \mathbf{H}_j$, we use feasibility:

$$\mathcal{A}(\mathbf{H}) = \mathbf{0} \iff \mathcal{A}(\mathbf{H}_0 + \mathbf{H}_1) = - \sum_{j \geq 2} \mathcal{A}(\mathbf{H}_j),$$

which taken collectively with RIP yields

$$\begin{aligned} (1 - \delta_{(2+K)r}^{\text{lb}}) \|\mathbf{H}_0 + \mathbf{H}_1\|_{\mathbf{F}} &\leq \|\mathcal{A}(\mathbf{H}_0 + \mathbf{H}_1)\|_{\mathbf{F}} = \left\| \sum_{j \geq 2} \mathcal{A}(\mathbf{H}_j) \right\|_{\mathbf{F}} \\ &\leq \sum_{j \geq 2} \|\mathcal{A}(\mathbf{H}_j)\|_{\mathbf{F}} \\ &\leq \sum_{j \geq 2} (1 + \delta_{Kr}^{\text{ub}}) \|\mathbf{H}_j\|_{\mathbf{F}} \end{aligned}$$

This establishes (6.6) as long as $\rho := \frac{1 + \delta_{Kr}^{\text{ub}}}{1 - \delta_{(2+K)r}^{\text{lb}}} < \sqrt{\frac{K}{2}}$.

Gaussian sampling operators satisfy RIP

If entries of $\{\mathbf{A}_i\}_{i=1}^m$ are i.i.d. $\mathcal{N}(0, 1/m)$, then

$$\delta_{5r}(\mathcal{A}) < \frac{\sqrt{3} - \sqrt{2}}{\sqrt{3} + \sqrt{2}}$$

with high prob., provided that

$$m \gtrsim nr \quad (\text{near-optimal sample size})$$

This satisfies assumption of Theorem 6.9 with $K = 3$

Precise phase transition

Using statistical dimension machinery, we can locate precise phase transition (Amelunxen, Lotz, McCoy & Tropp '13)

$$\text{nuclear norm min} \begin{cases} \text{works if} & m > \text{stat-dim}(\mathcal{D}(\|\cdot\|_*, \mathbf{X})) \\ \text{fails if} & m < \text{stat-dim}(\mathcal{D}(\|\cdot\|_*, \mathbf{X})) \end{cases}$$

where

$$\text{stat-dim}(\mathcal{D}(\|\cdot\|_*, \mathbf{X})) \approx n^2 \psi\left(\frac{r}{n}\right)$$

and

$$\psi(\rho) = \inf_{\tau \geq 0} \left\{ \rho + (1 - \rho) \left[\rho(1 + \tau^2) + (1 - \rho) \int_{\tau}^2 (u - \tau)^2 \frac{\sqrt{4 - u^2}}{\pi} du \right] \right\}$$

Numerical phase transition ($n = 30$)

Low-rank matrix recovery via Schatten 1-norm minimization

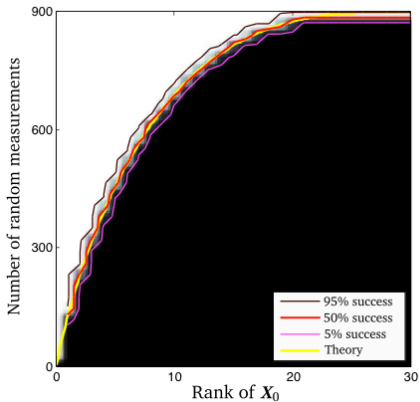


Figure credit: Amelunxen, Lotz, McCoy, & Tropp '13

Sampling operators that do NOT satisfy RIP

Unfortunately, many sampling operators fail to satisfy RIP
(e.g. none of the 4 motivating examples in this lecture satisfies RIP)

Matrix completion

Sampling operators for matrix completion

Observation operator (projection onto matrices supported on Ω)

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

where $(i, j) \in \Omega$ with prob. p (random sampling)

- \mathcal{P}_Ω does NOT satisfy RIP when $p \ll 1$!
- For example,

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}}_M \quad \underbrace{\begin{bmatrix} ? & \checkmark & ? & \checkmark & \checkmark \\ \checkmark & ? & \checkmark & ? & \checkmark \\ ? & \checkmark & \checkmark & ? & ? \\ \checkmark & ? & ? & \checkmark & ? \\ \checkmark & ? & \checkmark & ? & \checkmark \end{bmatrix}}_\Omega$$

$$\|\mathcal{P}_\Omega(M)\|_F = 0, \text{ or equivalently, } \frac{1+\delta_K^{\text{ub}}}{1-\delta_{2+K}^{\text{lb}}} = \infty$$

Which sampling pattern?

Consider the following sampling pattern

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

- If some rows/columns are not sampled, recovery is impossible.

Which low-rank matrices can we recover?

Compare following rank-1 matrices:

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

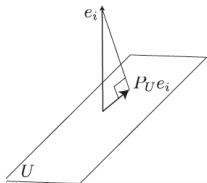
Column / row spaces cannot be aligned with canonical basis vectors

Coherence

Definition 6.10

Coherence parameter μ of $M = U\Sigma V^\top$ is smallest quantity s.t.

$$\max_i \|U^\top e_i\|^2 \leq \frac{\mu r}{n} \quad \text{and} \quad \max_i \|V^\top e_i\|^2 \leq \frac{\mu r}{n}$$



- $\mu \geq 1$ (since $\sum_{i=1}^n \|U^\top e_i\|^2 = \|U\|_F^2 = r$)
- $\mu = 1$ if $\frac{1}{\sqrt{n}}\mathbf{1} = U = V$ (most incoherent)
- $\mu = \frac{n}{r}$ if $e_i \in U$ (most coherent)

Performance guarantee

Theorem 6.11 (Candes & Recht '09, Candes & Tao '10, Gross '11, ...)

Nuclear norm minimization is exact and unique with high probability, provided that

$$m \gtrsim \mu nr \log^2 n$$

- This result is optimal up to a logarithmic factor
- Established via a RIPless theory

Numerical performance of nuclear-norm minimization

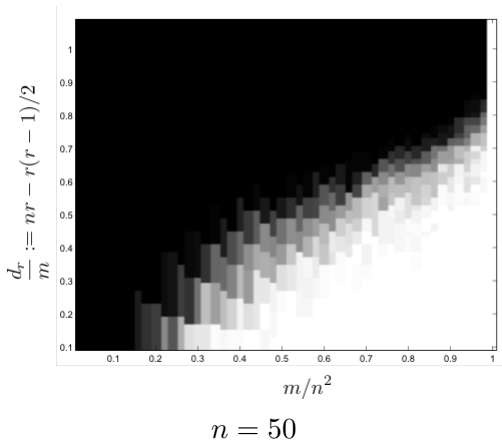


Fig. credit: Candes, Recht '09

Subgradient of nuclear norm

Subdifferential (set of subgradients) of $\|\cdot\|_*$ at M is

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

- Does not depend on singular values of M
- $Z \in \partial\|M\|_*$ iff

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

KKT condition

Lagrangian:

$$\mathcal{L}(\mathbf{X}, \mathbf{\Lambda}) = \|\mathbf{X}\|_* + \langle \mathbf{\Lambda}, \mathcal{P}_\Omega(\mathbf{X}) - \mathcal{P}_\Omega(\mathbf{M}) \rangle = \|\mathbf{X}\|_* + \langle \mathcal{P}_\Omega(\mathbf{\Lambda}), \mathbf{X} - \mathbf{M} \rangle$$

When \mathbf{M} is minimizer, KKT condition reads

$$\mathbf{0} \in \partial_{\mathbf{X}} \mathcal{L}(\mathbf{X}, \mathbf{\Lambda}) \big|_{\mathbf{X}=\mathbf{M}} \iff \exists \mathbf{\Lambda} \text{ s.t. } -\mathcal{P}_\Omega(\mathbf{\Lambda}) \in \partial \|\mathbf{M}\|_*$$

$$\iff \exists \mathbf{W} \text{ s.t. } \quad \mathbf{U}\mathbf{V}^\top + \mathbf{W} \text{ is supported on } \Omega, \\ \mathcal{P}_T(\mathbf{W}) = \mathbf{0}, \text{ and } \|\mathbf{W}\| \leq 1$$

Optimality condition via dual certificate

Slightly stronger condition than KKT guarantees uniqueness:

Lemma 6.12

M is unique minimizer of nuclear norm minimization if

- sampling operator \mathcal{P}_Ω restricted to T is injective, i.e.

$$\mathcal{P}_\Omega(\mathbf{H}) \neq \mathbf{0} \quad \forall \text{ nonzero } \mathbf{H} \in T$$

- $\exists \mathbf{W}$ s.t.

$UV^\top + \mathbf{W}$ is supported on Ω ,

$$\mathcal{P}_T(\mathbf{W}) = \mathbf{0}, \text{ and } \|\mathbf{W}\| < 1$$

Proof of Lemma 6.12

For any \mathbf{W}_0 obeying $\|\mathbf{W}_0\| \leq 1$ and $\mathcal{P}_T(\mathbf{W}) = \mathbf{0}$, one has

$$\begin{aligned}\|\mathbf{M} + \mathbf{H}\|_* &\geq \|\mathbf{M}\|_* + \langle \mathbf{UV}^\top + \mathbf{W}_0, \mathbf{H} \rangle \\ &= \|\mathbf{M}\|_* + \langle \mathbf{UV}^\top + \mathbf{W}, \mathbf{H} \rangle + \langle \mathbf{W}_0 - \mathbf{W}, \mathbf{H} \rangle \\ &= \|\mathbf{M}\|_* + \langle \mathcal{P}_\Omega(\mathbf{UV}^\top + \mathbf{W}), \mathbf{H} \rangle + \langle \mathcal{P}_{T^\perp}(\mathbf{W}_0 - \mathbf{W}), \mathbf{H} \rangle \\ &= \|\mathbf{M}\|_* + \langle \mathbf{UV}^\top + \mathbf{W}, \mathcal{P}_\Omega(\mathbf{H}) \rangle + \langle \mathbf{W}_0 - \mathbf{W}, \mathcal{P}_{T^\perp}(\mathbf{H}) \rangle \\ &\quad \text{if we take } \mathbf{W}_0 \text{ s.t. } \langle \mathbf{W}_0, \mathcal{P}_{T^\perp}(\mathbf{H}) \rangle = \|\mathcal{P}_{T^\perp}(\mathbf{H})\|_* \\ &\geq \|\mathbf{M}\|_* + \|\mathcal{P}_{T^\perp}(\mathbf{H})\|_* - \|\mathbf{W}\| \cdot \|\mathcal{P}_{T^\perp}(\mathbf{H})\|_* \\ &= \|\mathbf{M}\|_* + (1 - \|\mathbf{W}\|) \|\mathcal{P}_{T^\perp}(\mathbf{H})\|_* > \|\mathbf{M}\|_*\end{aligned}$$

unless $\mathcal{P}_{T^\perp}(\mathbf{H}) = \mathbf{0}$.

But if $\mathcal{P}_{T^\perp}(\mathbf{H}) = \mathbf{0}$, then $\mathbf{H} = \mathbf{0}$ by injectivity. Thus, $\|\mathbf{M} + \mathbf{H}\|_* > \|\mathbf{M}\|_*$ for any $\mathbf{H} \neq \mathbf{0}$, concluding the proof.

Constructing dual certificates

Use “golfing scheme” to produce approximate dual certificate (Gross '11)

- Think of it as an iterative algorithm (with sample splitting) to solve KKT

Reference

- [1] "Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization," B. Recht, M. Fazel, P. Parrilo, *SIAM Review*, 2010.
- [2] "Exact matrix completion via convex optimization," E. Candes, and B. Recht, *Foundations of Computational Mathematics*, 2009
- [3] "Matrix rank minimization with applications," M. Fazel, *Ph.D. Thesis*, 2002.
- [4] "The power of convex relaxation: Near-optimal matrix completion," E. Candes, and T. Tao, 2010.
- [5] "Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements," E. Candes, and Y. Plan, *IEEE Transactions on Information Theory*, 2011.
- [6] "Shape and motion from image streams under orthography: a factorization method," C. Tomasi and T. Kanade, *International Journal of Computer Vision*, 1992.

Reference

- [7] "An accelerated proximal gradient algorithm for nuclear norm regularized linear least squares problems," K. C. Toh and S. Yun, *Pacific Journal of optimization*, 2010.
- [8] "Topics in random matrix theory," T. Tao, *American mathematical society*, 2012.
- [9] "A singular value thresholding algorithm for matrix completion," J. Cai, E. Candes, Z. Shen, *SIAM Journal on Optimization*, 2010.
- [10] "Recovering low-rank matrices from few coefficients in any basis," D. Gross, *IEEE Transactions on Information Theory*, 2011.
- [11] "Incoherence-optimal matrix completion," Y. Chen, *IEEE Transactions on Information Theory*, 2015.
- [12] "Revisiting Frank-Wolfe: Projection-Free Sparse Convex Optimization," M. Jaggi, *ICML*, 2013.