

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

Low-rank matrix recovery via nonconvex optimization

Yuejie Chi

Department of Electrical and Computer Engineering

Carnegie Mellon University

Spring 2018

Outline

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

Why nonconvex?

- Consider completing an $n \times n$ matrix, with rank r :

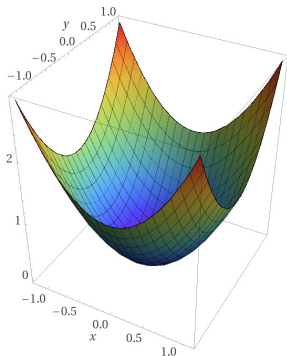
$$\text{minimize}_{\mathbf{X}} \quad \|\mathcal{P}_{\Omega}(\mathbf{X} - \mathbf{M})\|_{\text{F}}^2 \quad \text{s.t.} \quad \text{rank}(\mathbf{X}) \leq r,$$

where $r \ll n$.

- The size of observation $|\Omega|$ is about $nr \text{polylog}n$;
- The degrees of freedom in \mathbf{X} is about nr ;
- **Question:** Can we develop algorithms that work with computational and memory complexity that nearly linear in n ?
- This means that we don't even want to store the matrix \mathbf{X} which takes n^2 storage.
- A nonconvex approach will store and update a "low-dimensional" representation of \mathbf{X} throughout the execution of the algorithm.

Convex vs. nonconvex

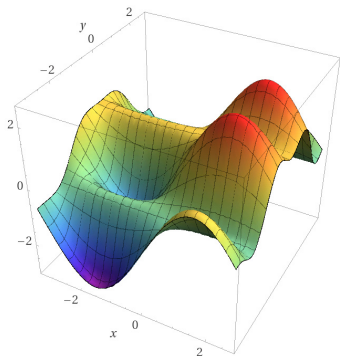
$$\text{minimize}_x f(x)$$



Computed by WolframAlpha

convex

vs.



Computed by WolframAlpha

nonconvex

Prelude: low-rank matrix approximation — an optimization perspective

Low-rank matrix approximation / PCA

Given $M \in \mathbb{R}^{n \times n}$ (not necessarily low-rank), solve the *low-rank approximation problem* (best rank- r approximation):

$$\widehat{M} = \operatorname{argmin}_{\mathbf{X}} \|\mathbf{X} - M\|_F^2 \quad \text{s.t.} \quad \operatorname{rank}(\mathbf{X}) \leq r.$$

this is a nonconvex optimization problem.

The solution is known as the **Eckart-Young theorem**:

- denote the SVD of $M = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$, where σ_i 's are in a descending order; then

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

nonconvex, but tractable.

Optimization viewpoint

Let us factorize $\mathbf{X} = \mathbf{UV}^\top$, where $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times r}$. Our problem is equivalent to

$$\text{minimize}_{\mathbf{U}, \mathbf{V}} f(\mathbf{U}, \mathbf{V}) := \|\mathbf{UV}^\top - \mathbf{M}\|_{\text{F}}^2.$$

- The size of \mathbf{U}, \mathbf{V} are of $O(nr)$, which is much smaller than \mathbf{X} ;
- Identifiability issues: for any orthonormal $\mathbf{R} \in \mathbb{R}^{r \times r}$, we have

$$\mathbf{UV}^\top = (\alpha \mathbf{UR})(\alpha^{-1} \mathbf{VR})^\top.$$

If (\mathbf{U}, \mathbf{V}) is a global minimizer (\cdot), so does $(\alpha \mathbf{UR}, \alpha^{-1} \mathbf{VR})$.

Question: what does $f(\mathbf{U}, \mathbf{V})$ look like (landscape)? (we already found its global minima.)

The PSD case

For simplicity, consider the PSD case.

- Let M be PSD, so that $M = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{u}_i^\top$.
- Let $X = UU^\top$, where $U \in \mathbb{R}^{n \times r}$.

We're interested in the landscape of

$$f(U) := \frac{1}{4} \|UU^\top - M\|_F^2.$$

Identifiability: for any orthonormal $R \in \mathbb{R}^{r \times r}$, we have

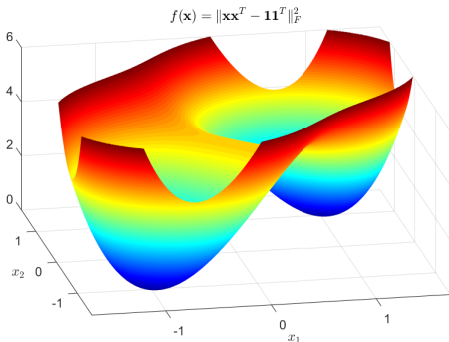
$$UU^\top = (UR)(UR)^\top.$$

make the exposition even simpler: set $r = 1$.

$$f(\mathbf{u}) = \frac{1}{4} \|\mathbf{u}\mathbf{u}^\top - M\|_F^2.$$

Good news: benign landscape

$$\text{Take } f(\mathbf{u}) = \left\| \mathbf{u}\mathbf{u}^\top - \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \right\|_F^2.$$



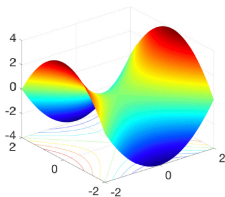
Global optima: $\mathbf{x} = \pm \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, strict saddle $\mathbf{x} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$. No “spurious” local minima.

Critical points

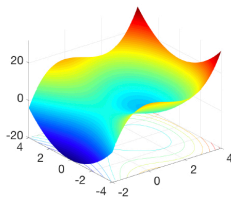
Definition 7.1

A first-order critical point (stationary point) satisfies

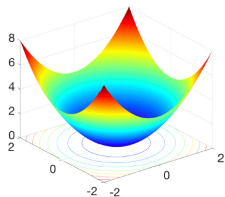
$$\nabla f(\mathbf{u}) = 0.$$



(a) strict saddle



(b) local minimum



(c) global minimum

Figure credit: Li et al., 2016

Critical points of $f(\mathbf{u})$

Any $\mathbf{u} \in \mathbb{R}^n$ satisfies

$$\nabla f(\mathbf{u}) = (\mathbf{u}\mathbf{u}^\top - \mathbf{M})\mathbf{u} = \mathbf{0}.$$

$$\Leftrightarrow$$

$$\mathbf{M}\mathbf{u} = \|\mathbf{u}\|_2^2 \mathbf{u}$$

$$\Leftrightarrow$$

\mathbf{u} aligns with eigenvectors of \mathbf{M} .

or

$$\mathbf{u} = \mathbf{0}.$$

Critical points of $f(\mathbf{u})$

Any $\mathbf{u} \in \mathbb{R}^n$ satisfies

$$\nabla f(\mathbf{u}) = (\mathbf{u}\mathbf{u}^\top - \mathbf{M})\mathbf{u} = \mathbf{0}.$$

$$\Leftrightarrow$$

$$\mathbf{M}\mathbf{u} = \|\mathbf{u}\|_2^2 \mathbf{u}$$

$$\Leftrightarrow$$

\mathbf{u} aligns with eigenvectors of \mathbf{M} .

or

$$\mathbf{u} = \mathbf{0}.$$

Since $\mathbf{M}\mathbf{u}_i = \sigma_i \mathbf{u}_i$, the set of critical points are given as

$$\{\sqrt{\sigma_i} \mathbf{u}_i, i = 1, \dots, n\}.$$

Categorization of critical points

Need to examine the Hessian:

$$\nabla^2 f(\mathbf{u}) := 2\mathbf{u}\mathbf{u}^\top + \|\mathbf{u}\|_2^2 \mathbf{I} - \mathbf{M}.$$

- Plug in the non-zero critical points: $\tilde{\mathbf{u}}_k := \sqrt{\sigma_k} \mathbf{u}_k$,

$$\begin{aligned}\nabla^2 f(\tilde{\mathbf{u}}_k) &= 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top + \sigma_k \mathbf{I} - \mathbf{M} \\ &= 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top + \sigma_k \left(\sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^\top \right) - \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{u}_i^\top \\ &= \sum_{i \neq k} (\sigma_k - \sigma_i) \mathbf{u}_i \mathbf{u}_i^\top + 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top\end{aligned}$$

- Assume $\sigma_1 > \sigma_2 > \dots > \sigma_n > 0$:
 - $k = 1$: $\nabla^2 f(\tilde{\mathbf{u}}_1) \succ 0 \rightarrow$ **local minima**
 - $1 < k \leq n$: $\lambda_{\min}(\nabla^2 f(\tilde{\mathbf{u}}_k)) < 0$, $\lambda_{\max}(\nabla^2 f(\tilde{\mathbf{u}}_k)) > 0$, \rightarrow **strict saddle**
 - $\mathbf{u} = 0$: $\nabla^2 f(0) \preceq 0 \rightarrow$ **local maxima**

Categorization of critical points

Need to examine the Hessian:

$$\nabla^2 f(\mathbf{u}) := 2\mathbf{u}\mathbf{u}^\top + \|\mathbf{u}\|_2^2 \mathbf{I} - \mathbf{M}.$$

- Plug in the non-zero critical points: $\tilde{\mathbf{u}}_k := \sqrt{\sigma_k} \mathbf{u}_k$,

$$\begin{aligned} \nabla^2 f(\tilde{\mathbf{u}}_k) &= 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top + \sigma_k \mathbf{I} - \mathbf{M} \\ &= 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top + \sigma_k \left(\sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^\top \right) - \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{u}_i^\top \\ &= \sum_{i \neq k} (\sigma_k - \sigma_i) \mathbf{u}_i \mathbf{u}_i^\top + 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top \end{aligned}$$

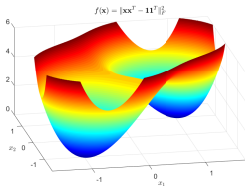
- Assume $\sigma_1 > \sigma_2 \geq \dots \geq \sigma_n \geq 0$:
 - $k = 1$: $\nabla^2 f(\tilde{\mathbf{u}}_1) \succ 0 \rightarrow$ **local minima**
 - $1 < k \leq n$: $\lambda_{\min}(\nabla^2 f(\tilde{\mathbf{u}}_k)) < 0$, $\lambda_{\max}(\nabla^2 f(\tilde{\mathbf{u}}_k)) > 0$, \rightarrow **strict saddle**
 - $\mathbf{u} = 0$: $\nabla^2 f(0) \prec 0 \rightarrow$ **strict saddle**

Summary

$$f(\mathbf{U}) := \frac{1}{4} \|\mathbf{U}\mathbf{U}^\top - \mathbf{M}\|_{\text{F}}^2, \quad \mathbf{U} \in \mathbb{R}^{n \times r},$$

If $\sigma_r > \sigma_{r+1}$,

- **all local minima are global:** \mathbf{U} contains the top- r eigenvectors up to an orthonormal transformation;
- **strict saddle points:** all stationary points are saddle points except the global optimum.



Undersampled regime

Consider linear measurements:

$$\mathbf{y} = \mathcal{A}(\mathbf{M}), \quad \mathbf{y} \in \mathbb{R}^m, \quad m \ll n^2$$

where $\mathbf{M} = \mathbf{U}_0 \mathbf{U}_0^\top \in \mathbb{R}^{n \times n}$ is **rank- r** , $r \ll n$, and PSD (for simplicity).

- The loss function we consider:

$$f(\mathbf{U}) := \frac{1}{4} \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top - \mathbf{M})\|_{\mathbb{F}}^2.$$

- If $\mathbb{E}[\mathcal{A}^* \mathcal{A}] = \mathcal{I}$, then

$$\mathbb{E}[f(\mathbf{U})] = \frac{1}{4} \|\mathbf{U}\mathbf{U}^\top - \mathbf{M}\|_{\mathbb{F}}^2.$$

- Does $f(\mathbf{U})$ inherit the benign landscape?

Landscape preserving under RIP

Recall the definition of RIP:

Definition 7.2

The rank- r restricted isometry constants δ_r is the smallest quantity

$$(1 - \delta_r)\|\mathbf{X}\|_F^2 \leq \|\mathcal{A}(\mathbf{X})\|_F^2 \leq (1 + \delta_r)\|\mathbf{X}\|_F^2, \quad \forall \mathbf{X} : \text{rank}(\mathbf{X}) \leq r$$

Landscape preserving under RIP

Recall the definition of RIP:

Definition 7.2

The rank- r restricted isometry constants δ_r is the smallest quantity

$$(1 - \delta_r) \|\mathbf{X}\|_F^2 \leq \|\mathcal{A}(\mathbf{X})\|_F^2 \leq (1 + \delta_r) \|\mathbf{X}\|_F^2, \quad \forall \mathbf{X} : \text{rank}(\mathbf{X}) \leq r$$

Theorem 7.3 (Bhojanapalli et al.' 2016, Ge et al.' 2017)

If \mathcal{A} satisfies the RIP with $\delta_{2r} < \frac{1}{10}$, then $f(\mathbf{U})$ satisfies

- all local min are global: for any local minimum \mathbf{U} of $f(\mathbf{U})$, it satisfies $\mathbf{U}\mathbf{U}^\top = \mathbf{M}$;
- strict saddle points: for non-local min critical point \mathbf{U} , it satisfies $\lambda_{\min}[\nabla^2 f(\mathbf{U})] \leq -\frac{2}{5}\sigma_r$.

Proof of Theorem 7.3 when $r = 1$

Without loss of generality, assume $M = \mathbf{u}_0 \mathbf{u}_0^\top$, and $\sigma_1 = 1$.

- Step 1: check all the critical points:

$$\nabla f(\mathbf{u}) = \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top - \underbrace{\mathbf{u}_0 \mathbf{u}_0^\top}_M \rangle \mathbf{A}_i \mathbf{u} = 0$$

- Step 2: verify the Hessian at all the critical points:

$$\nabla^2 f(\mathbf{u}) = \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top - \mathbf{u}_0 \mathbf{u}_0^\top \rangle \mathbf{A}_i + 2 \mathbf{A}_i \mathbf{u} \mathbf{u}^\top \mathbf{A}_i^\top$$

Proof of Theorem 7.3 when $r = 1$

Proof: Assume \mathbf{u} is first-order optimal. Consider the descent direction: $\Delta = \mathbf{u} - \mathbf{u}_0$:

$$\begin{aligned}\Delta^\top \nabla^2 f(\mathbf{u}) \Delta &= \sum_{i=1}^m \left[\langle \mathbf{A}_i, \mathbf{u}\mathbf{u}^\top - \mathbf{u}_0\mathbf{u}_0^\top \rangle \langle \mathbf{A}_i, \Delta\Delta^\top \rangle + 2\langle \mathbf{A}_i, \mathbf{u}\Delta^\top \rangle^2 \right] \\ &= \sum_{i=1}^m \left[\langle \mathbf{A}_i, \Delta\Delta^\top \rangle^2 - 3\langle \mathbf{A}_i, \mathbf{u}\mathbf{u}^\top - \mathbf{u}_0\mathbf{u}_0^\top \rangle^2 \right].\end{aligned}$$

where we have used the first order optimality condition.

Proof of Theorem 7.3 when $r = 1$

By the RIP property:

$$\begin{aligned}\Delta^\top \nabla^2 f(\mathbf{u}) \Delta &= \sum_{i=1}^m \left[\langle \mathbf{A}_i, (\mathbf{u} - \mathbf{u}_0)(\mathbf{u} - \mathbf{u}_0)^\top \rangle^2 - 3 \langle \mathbf{A}_i, \mathbf{u}\mathbf{u}^\top - \mathbf{u}_0\mathbf{u}_0^\top \rangle^2 \right] \\ &\leq (1 + \delta) \|(\mathbf{u} - \mathbf{u}_0)(\mathbf{u} - \mathbf{u}_0)^\top\|_F^2 - 3(1 - \delta) \|\mathbf{u}\mathbf{u}^\top - \mathbf{u}_0\mathbf{u}_0^\top\|_F^2 \\ &\leq [2(1 + \delta) - 3(1 - \delta)] \|\mathbf{u}\mathbf{u}^\top - \mathbf{u}_0\mathbf{u}_0^\top\|_F^2 \\ &\leq -(1 - 5\delta) \|\mathbf{u}\mathbf{u}^\top - \mathbf{u}_0\mathbf{u}_0^\top\|_F^2\end{aligned}$$

where we use

$$\|(\mathbf{u} - \mathbf{u}_0)(\mathbf{u} - \mathbf{u}_0)^\top\|_F^2 \leq 2 \|\mathbf{u}\mathbf{u}^\top - \mathbf{u}_0\mathbf{u}_0^\top\|_F^2.$$

Landscape without RIP

In matrix completion, we need to regularize the loss function by promoting **incoherent** solutions: set

$$Q(\mathbf{U}) = \sum_{i=1}^m (\|\mathbf{e}_i^\top \mathbf{U}\|_2 - \alpha)_+^4$$

where α is some regularization parameter, and $z_+ = \max\{z, 0\}$.

Landscape without RIP

In matrix completion, we need to regularize the loss function by promoting **incoherent** solutions: set

$$Q(\mathbf{U}) = \sum_{i=1}^m (\|\mathbf{e}_i^\top \mathbf{U}\|_2 - \alpha)_+^4$$

where α is some regularization parameter, and $z_+ = \max\{z, 0\}$.

Consider the loss function

$$f(\mathbf{U}) = \frac{1}{p} \|\mathcal{P}_\Omega(\mathbf{U}\mathbf{U}^\top - \mathbf{M})\|_F^2 + \lambda Q(\mathbf{U})$$

where λ is a regularization parameter.

- adding $Q(\mathbf{U})$ doesn't affect the global optimizer if α is set properly.

MC doesn't have spurious local minima

Theorem 7.4 (Ge et al, 2016)

If $p \gtrsim \frac{\mu^4 r^6 \log n}{n}$, $\alpha^2 = \Theta(\frac{\mu r \sigma_1}{n})$ and $\lambda = \Theta(\frac{n}{\mu r})$, then with probability at least $1 - n^{-1}$,

- *all local min are global: for any local minimum U of $f(U)$, it satisfies $UU^\top = M$;*
 - *saddle points that are not local minima are strict saddle points.*
-
- saddle-point escaping algorithms can be used to guarantee convergence to local minima, which in our problem are global minima.
 - active research area for constructing saddle-point escaping algorithms: (perturbed) gradient descent, trust-region methods, etc...

Spectral methods: a one-shot approach

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

Signal + noise

$(i, j) \in \Omega$ independently with prob. p

One can write observation $\mathcal{P}_\Omega(\mathbf{M})$ as

$$\frac{1}{p}\mathcal{P}_\Omega(\mathbf{M}) = \underbrace{\mathbf{M}}_{\text{signal}} + \underbrace{\frac{1}{p}\mathcal{P}_\Omega(\mathbf{M}) - \mathbf{M}}_{\text{noise}}$$

- Noise has mean zero: $\mathbb{E}\left[\frac{1}{p}\mathcal{P}_\Omega(\mathbf{M})\right] = \mathbf{M}$

Low-rank denoising

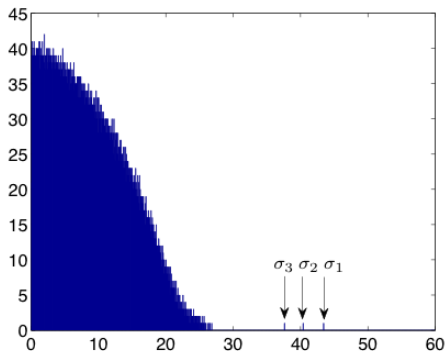
$$\frac{1}{p}\mathcal{P}_\Omega(\mathbf{M}) = \underbrace{\mathbf{M}}_{\text{low-rank signal}} + \underbrace{\frac{1}{p}\mathcal{P}_\Omega(\mathbf{M}) - \mathbf{M}}_{:= \mathbf{E} \text{ (zero-mean noise)}}$$

Algorithm 7.1 Spectral method

$$\hat{\mathbf{M}} \leftarrow \text{best rank-}r \text{ approximation of } \frac{1}{p}\mathcal{P}_\Omega(\mathbf{M})$$

The spectral method can be solved via power methods or Lanczos methods, and we don't need to realize the matrix $\frac{1}{p}\mathcal{P}_\Omega(\mathbf{M})$.

Histograms of singular values of $\mathcal{P}_\Omega(M)$



A $10^4 \times 10^4$ random rank-3 matrix M with $p = 0.003$

Fig. credit: Keshavan, Montanari, Oh '10

Performance of spectral methods

Theorem 7.5 (Keshavan, Montanari, Oh '10)

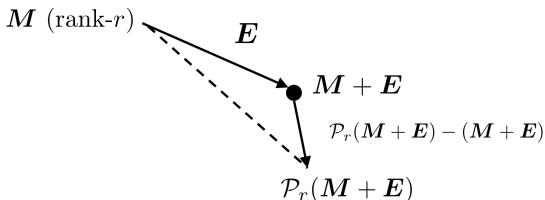
Suppose number of observed entries m obeys $m \gtrsim n \log n$. Then

$$\frac{\|\widehat{\mathbf{M}} - \mathbf{M}\|_F}{\|\mathbf{M}\|_F} \lesssim \underbrace{\frac{\max_{i,j} |M_{i,j}|}{\frac{1}{n} \|\mathbf{M}\|_F}}_{:=\nu} \cdot \sqrt{\frac{nr \log^2 n}{m}},$$

- ν reflects whether energy of \mathbf{M} is spread out, $|M_{i,j}| \lesssim \mu r/n$;
- When $m \gg \nu^2 n \log^2 n$, estimate $\widehat{\mathbf{M}}$ is very close to truth¹
- Degrees of freedom $\asymp nr$
→ nearly-optimal sample complexity for **incoherent** matrices

¹The logarithmic factor can be improved.

Perturbation bounds



To ensure \widehat{M} is good estimate, it suffices to control noise E

Lemma 7.6

Suppose $\text{rank}(M) = r$. For any perturbation E ,

$$\|\mathcal{P}_r(M + E) - M\| \leq 2\|E\|$$

$$\|\mathcal{P}_r(M + E) - M\|_F \leq 2\sqrt{2r}\|E\|$$

where $\mathcal{P}_r(X)$ is best rank- r approximation of X .

Prior on matrix perturbation theory

Lemma 7.7 (Weyl's inequality, 1912)

Let M, E be $n \times n$ matrices. Then

$$|\sigma_k(M + E) - \sigma_k(M)| \leq \|E\|, \quad k = 1, \dots, n.$$

Proof: Invoke the Courant-Fisher Minimax Characterization:

$$\sigma_k(A) = \max_{\dim(S)=k} \min_{0 \neq v \in S} \frac{\|Av\|_2}{\|v\|_2}.$$

Proof of Lemma 7.6

By matrix perturbation theory,

$$\begin{aligned} & \|\mathcal{P}_r(\mathbf{M} + \mathbf{E}) - \mathbf{M}\| \\ & \stackrel{\text{triangle inequality}}{\leq} \|\mathcal{P}_r(\mathbf{M} + \mathbf{E}) - (\mathbf{M} + \mathbf{E})\| + \|(\mathbf{M} + \mathbf{E}) - \mathbf{M}\| \\ & \leq \sigma_{r+1}(\mathbf{M} + \mathbf{E}) + \|\mathbf{E}\| \\ & \stackrel{\text{Weyl's inequality}}{\leq} \underbrace{\sigma_{r+1}(\mathbf{M})}_{=0} + \|\mathbf{E}\| + \|\mathbf{E}\| = 2\|\mathbf{E}\|. \end{aligned}$$

The 2nd inequality of Lemma 7.6 follows since both $\mathcal{P}_r(\mathbf{M} + \mathbf{E})$ and \mathbf{M} are rank- r , and hence

$$\|\mathcal{P}_r(\mathbf{M} + \mathbf{E}) - \mathbf{M}\|_F \leq \sqrt{2r} \|\mathcal{P}_r(\mathbf{M} + \mathbf{E}) - \mathbf{M}\|.$$

Controlling the noise

Recall that entries of $\mathbf{E} = \frac{1}{p}\mathcal{P}_\Omega(\mathbf{M}) - \mathbf{M}$ are zero-mean and independent.

A bit of random matrix theory ...

Lemma 7.8 (Chapter 2.3, Tao '12)

Suppose $\mathbf{X} \in \mathbb{R}^{n \times n}$ is a random symmetric matrix obeying

- $\{X_{i,j} : i < j\}$ are independent
- $\mathbb{E}[X_{i,j}] = 0$ and $\text{Var}[X_{i,j}] \lesssim 1$
- $\max_{i,j} |X_{i,j}| \lesssim \sqrt{n}$

Then $\|\mathbf{X}\| \lesssim \sqrt{n} \log n$.

Proof of Theorem 7.5

If we look at the zero-mean matrix $\tilde{\mathbf{E}} = \frac{\sqrt{p}}{\mu/n} \mathbf{E}$, then

$$\begin{aligned}\text{Var} [\tilde{E}_{i,j}] &= p(1-p) \cdot \left(\frac{\sqrt{p}}{\mu/n} \cdot \frac{1}{p} M_{i,j} \right)^2 \leq \left(\frac{M_{i,j}}{\mu/n} \right)^2 \lesssim 1, \\ |\tilde{E}_{i,j}| &\leq \frac{|M_{i,j}|}{\sqrt{p}\mu/n} \lesssim \frac{1}{\sqrt{p}},\end{aligned}$$

where we have used the fact

$$|M_{i,j}| = |\mathbf{e}_i^\top \mathbf{U} \Sigma \mathbf{V}^\top \mathbf{e}_j| \leq \|\mathbf{U}^\top \mathbf{e}_i\| \cdot \sigma_1 \cdot \|\mathbf{V}^\top \mathbf{e}_j\| \stackrel{\text{(by our assumptions)}}{\lesssim} \frac{\mu r}{n} \asymp \frac{\mu}{n}$$

Lemma 7.8 tells us that if $p \gtrsim \frac{\log n}{n}$, then

$$\|\tilde{\mathbf{E}}\| \lesssim \sqrt{n} \log n \iff \|\mathbf{E}\| \lesssim \frac{\mu}{\sqrt{pn}} \log n$$

This together with Lemma 7.6 and the fact $m \asymp pn^2$ establishes Theorem 7.5.

Gradient methods: iterative refinements

Iterative methods: an overview

$$\text{minimize}_{\mathbf{U}, \mathbf{V}} f(\mathbf{U}, \mathbf{V}) := \|\mathcal{P}_{\Omega}(\mathbf{U}\mathbf{V}^{\top} - \mathbf{M})\|_{\text{F}}^2.$$

- **Gradient descent:** (our focus)

$$\mathbf{U}_{t+1} = \mathcal{P}_{\mathbf{U}} \left[\mathbf{U}_t - \eta_t \nabla_{\mathbf{U}} f(\mathbf{U}_t, \mathbf{V}_t) \right],$$

$$\mathbf{V}_{t+1} = \mathcal{P}_{\mathbf{V}} \left[\mathbf{V}_t - \eta_t \nabla_{\mathbf{V}} f(\mathbf{U}_t, \mathbf{V}_t) \right].$$

where η_t is the step size and $\mathcal{P}_{\mathbf{U}}, \mathcal{P}_{\mathbf{V}}$ denote the Euclidean projection onto some constraint sets;

- **Alternating minimization:** One optimizes \mathbf{U}, \mathbf{V} alternatively while fixing the other, which is a convex problem.

$$\mathbf{U}_{t+1} = \operatorname{argmin}_{\mathbf{U}} f(\mathbf{U}, \mathbf{V}_t),$$

$$\mathbf{V}_{t+1} = \operatorname{argmin}_{\mathbf{V}} f(\mathbf{U}_{t+1}, \mathbf{V}).$$

Gradient descent for matrix completion

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left(\mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - M_{j,k} \right)^2$$

Algorithm 7.2 Gradient descent for MC

Input: $\mathbf{Y} = [Y_{j,k}]_{1 \leq j,k \leq n}$, r , p .

Spectral initialization: Let $\mathbf{U}^0 \boldsymbol{\Sigma}^0 \mathbf{U}^{0\top}$ be the rank- r eigendecomposition of

$$\mathbf{M}^0 := \frac{1}{p} \mathcal{P}_\Omega(\mathbf{Y}),$$

and set $\mathbf{X}^0 = \mathbf{U}^0 (\boldsymbol{\Sigma}^0)^{1/2}$.

Gradient updates: for $t = 0, 1, 2, \dots, T - 1$ do

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta_t \nabla f(\mathbf{X}^t).$$

Gradient descent for matrix completion

Define the optimal transform from the t th iterate \mathbf{X}^t to \mathbf{X}^\natural as

$$\mathbf{Q}^t := \operatorname{argmin}_{\mathbf{R} \in \mathcal{O}^{r \times r}} \left\| \mathbf{X}^t \mathbf{R} - \mathbf{X}^\natural \right\|_{\mathbb{F}}.$$

Theorem 7.9 (Ma et al., 2017)

Suppose $\mathbf{M} = \mathbf{X}^\natural \mathbf{X}^\natural \top$ is rank- r , incoherent and well-conditioned. Vanilla GD (with spectral initialization) achieves

- $\left\| \mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural \right\|_{\mathbb{F}} \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \left\| \mathbf{X}^\natural \right\|_{\mathbb{F}},$
- $\left\| \mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural \right\| \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \left\| \mathbf{X}^\natural \right\|, \quad (\text{spectral})$
- $\left\| \mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural \right\|_{2, \infty} \lesssim \rho^t \mu r \sqrt{\frac{\log n}{np}} \left\| \mathbf{X}^\natural \right\|_{2, \infty}, \quad (\text{incoherence})$

where $\rho = 1 - \frac{\sigma_{\min} \eta}{5} < 1$, if step size $\eta \asymp 1/\sigma_{\max}$ and sample complexity $n^2 p \gtrsim \mu^3 n r^3 \log^3 n$.

Gradient descent for matrix completion

Define the optimal transform from the t th iterate \mathbf{X}^t to \mathbf{X}^\natural as

$$\mathbf{Q}^t := \operatorname{argmin}_{\mathbf{R} \in \mathcal{O}^{r \times r}} \left\| \mathbf{X}^t \mathbf{R} - \mathbf{X}^\natural \right\|_{\mathbb{F}}.$$

Theorem 7.9 (Ma et al., 2017)

Suppose $\mathbf{M} = \mathbf{X}^\natural \mathbf{X}^{\natural \top}$ is rank- r , incoherent and well-conditioned. Vanilla GD (with spectral initialization) achieves

- $\left\| \mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural \right\|_{\mathbb{F}} \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \left\| \mathbf{X}^\natural \right\|_{\mathbb{F}},$
- $\left\| \mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural \right\| \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \left\| \mathbf{X}^\natural \right\|, \quad (\text{spectral})$
- $\left\| \mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural \right\|_{2, \infty} \lesssim \rho^t \mu r \sqrt{\frac{\log n}{np}} \left\| \mathbf{X}^\natural \right\|_{2, \infty}, \quad (\text{incoherence})$

where $\rho = 1 - \frac{\sigma_{\min} \eta}{5} < 1$, if step size $\eta \asymp 1/\sigma_{\max}$ and sample complexity $n^2 p \gtrsim \mu^3 n r^3 \log^3 n$.

- **linear convergence** of $\left\| \mathbf{X}^t \mathbf{X}^{t \top} - \mathbf{M}^\natural \right\|$ in Frobenius, spectral and infinity norms.

Numerical evidence for noiseless data

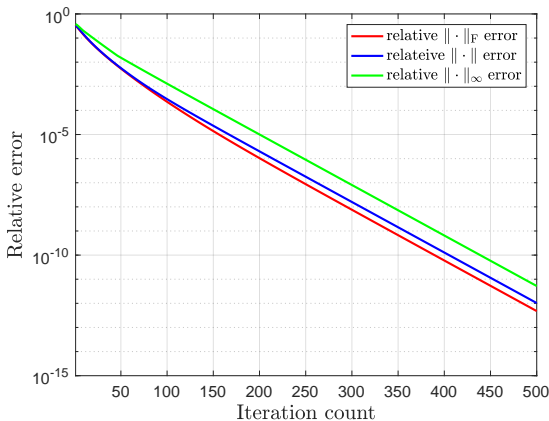


Figure 7.1: Relative error of $\mathbf{X}^t \mathbf{X}^{t\top}$ (measured by $\|\cdot\|_F$, $\|\cdot\|$, $\|\cdot\|_\infty$) vs. iteration count for matrix completion, where $n = 1000$, $r = 10$, $p = 0.1$, and $\eta_t = 0.2$.

Numerical evidence for noisy data

Set $\text{SNR} := \frac{\|M^\dagger\|_F^2}{n^2\sigma^2}$.

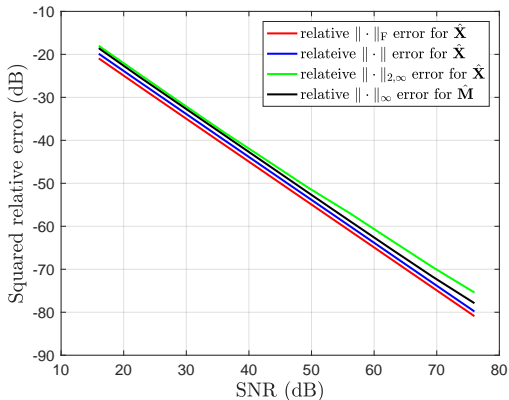


Figure 7.2: Squared relative error of the estimate \hat{X} (measured by $\|\cdot\|_F$, $\|\cdot\|$, $\|\cdot\|_{2,\infty}$) and $\hat{M} = \hat{X}\hat{X}^\top$ (measured by $\|\cdot\|_\infty$) vs. SNR, where $n = 500$, $r = 10$, $p = 0.1$, and $\eta_t = 0.2$.

Restricted strong convexity and smoothness

Lemma 7.10 (Restricted strong convexity and smoothness)

Suppose that $n^2 p \geq C \kappa^2 \mu r n \log n$ for some $C > 0$. Then with high probability, the Hessian $\nabla^2 f(\mathbf{X})$ obeys

$$\text{vec}(\mathbf{V})^\top \nabla^2 f(\mathbf{X}) \text{vec}(\mathbf{V}) \geq \frac{\sigma_{\min}}{2} \|\mathbf{V}\|_F^2 \quad (\text{restricted strong convexity})$$

$$\left\| \nabla^2 f(\mathbf{X}) \right\| \leq \frac{5}{2} \sigma_{\max} \quad (\text{smoothness})$$

for all \mathbf{X} and $\mathbf{V} = \mathbf{Y} \mathbf{H}_Y - \mathbf{Z}$, $\mathbf{H}_Y := \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{Y} \mathbf{R} - \mathbf{Z}\|_F$ satisfying

- $\left\| \mathbf{X} - \mathbf{X}^{\natural} \right\|_{2, \infty} \leq \epsilon \left\| \mathbf{X}^{\natural} \right\|_{2, \infty}$ (incoherence region),
- $\left\| \mathbf{Z} - \mathbf{X}^{\natural} \right\| \leq \delta \left\| \mathbf{X}^{\natural} \right\|$,

where $\epsilon \ll 1/\sqrt{\kappa^3 \mu r \log^2 n}$ and $\delta \ll 1/\kappa$.

Linear convergence induction I

Given the definition of \mathbf{Q}^{t+1} , we have

$$\begin{aligned}\|\mathbf{X}^{t+1}\mathbf{Q}^{t+1} - \mathbf{X}^{\natural}\|_{\mathbb{F}} &\leq \|\mathbf{X}^{t+1}\mathbf{Q}^t - \mathbf{X}^{\natural}\|_{\mathbb{F}} \\ &\stackrel{(i)}{=} \|\left[\mathbf{X}^t - \eta\nabla f(\mathbf{X}^t)\right]\mathbf{Q}^t - \mathbf{X}^{\natural}\|_{\mathbb{F}} \\ &\stackrel{(ii)}{=} \|\mathbf{X}^t\mathbf{Q}^t - \eta\nabla f(\mathbf{X}^t\mathbf{Q}^t) - \mathbf{X}^{\natural}\|_{\mathbb{F}} \\ &\stackrel{(iii)}{=} \underbrace{\|\mathbf{X}^t\mathbf{Q}^t - \eta\nabla f(\mathbf{X}^t\mathbf{Q}^t) - (\mathbf{X}^{\natural} - \eta\nabla f(\mathbf{X}^{\natural}))\|_{\mathbb{F}}}_{:=\alpha},\end{aligned}$$

where (i) follows from the GD rule, (ii) follows from the identity $\nabla f(\mathbf{X}^t\mathbf{R}) = \nabla f(\mathbf{X}^t)\mathbf{R}$ for any $\mathbf{R} \in \mathcal{O}^{r \times r}$, and (iii) follows from $\nabla f(\mathbf{X}^{\natural}) = \mathbf{0}$.

Linear convergence induction II

The fundamental theorem of calculus reveals

$$\begin{aligned}
 & \text{vec} \left[\mathbf{X}^t \mathbf{Q}^t - \eta \nabla f(\mathbf{X}^t \mathbf{Q}^t) - \left(\mathbf{X}^{\natural} - \eta \nabla f(\mathbf{X}^{\natural}) \right) \right] \\
 &= \text{vec} \left[\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural} \right] - \eta \cdot \text{vec} \left[\nabla f(\mathbf{X}^t \mathbf{Q}^t) - \nabla f(\mathbf{X}^{\natural}) \right] \\
 &= \left(\mathbf{I}_{nr} - \underbrace{\eta \int_0^1 \nabla^2 f(\mathbf{X}(\tau)) d\tau}_{:=\mathbf{A}} \right) \text{vec} \left(\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural} \right), \quad (7.1)
 \end{aligned}$$

where we denote $\mathbf{X}(\tau) := \mathbf{X}^{\natural} + \tau(\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural})$. Taking the squared Euclidean norm of both sides of the equality (7.1) leads to

$$\begin{aligned}
 \alpha^2 &= \text{vec}(\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural})^\top (\mathbf{I}_{nr} - \eta \mathbf{A})^2 \text{vec}(\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}) \\
 &\leq \left\| \mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural} \right\|_{\text{F}}^2 + \eta^2 \|\mathbf{A}\|^2 \left\| \mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural} \right\|_{\text{F}}^2 \\
 &\quad - 2\eta \text{vec}(\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural})^\top \mathbf{A} \text{vec}(\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}), \quad (7.2)
 \end{aligned}$$

Linear convergence induction III

Based on the incoherence of \mathbf{X}^{\natural} and \mathbf{X}^t , $\forall \tau \in [0, 1]$,

$$\|\mathbf{X}(\tau) - \mathbf{X}^{\natural}\|_{2,\infty} \leq \underbrace{\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}\|_{2,\infty}}_{\text{incoherence hypothesis}} \leq C\mu r \sqrt{\frac{\log n}{np}} \|\mathbf{X}^{\natural}\|_{2,\infty}.$$

Taking $\mathbf{X} = \mathbf{X}(\tau)$, $\mathbf{Y} = \mathbf{X}^t$ and $\mathbf{Z} = \mathbf{X}^{\natural}$ in Lemma 7.10, one can easily verify the assumptions therein given $n^2 p \gg \kappa^3 \mu^3 r^3 n \log^3 n$.

Hence,

$$\text{vec}(\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural})^\top \mathbf{A} \text{vec}(\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}) \geq \frac{\sigma_{\min}}{2} \|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}\|_{\text{F}}^2$$

and

$$\|\mathbf{A}\| \leq \frac{5}{2} \sigma_{\max}.$$

Linear convergence induction IV

Substituting these two inequalities into (7.2) yields

$$\begin{aligned}\alpha^2 &\leq \left(1 + \frac{25}{4}\eta^2\sigma_{\max}^2 - \sigma_{\min}\eta\right) \|\mathbf{X}^t \hat{\mathbf{H}}^t - \mathbf{X}^\natural\|_{\text{F}}^2 \\ &\leq \left(1 - \frac{\sigma_{\min}}{2}\eta\right) \|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural\|_{\text{F}}^2\end{aligned}$$

as long as $0 < \eta \leq (2\sigma_{\min})/(25\sigma_{\max}^2)$, which further implies that

$$\alpha \leq \left(1 - \frac{\sigma_{\min}}{4}\eta\right) \|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural\|_{\text{F}}.$$

The incoherence hypothesis is important for fast convergence: the fact that \mathbf{X}^t stays incoherent throughout the execution is called “implicit regularization” and can be established by a leave-one-out analysis trick [Ma et al., 2017].

Reference

- [1] "Guaranteed matrix completion via non-convex factorization," R. Sun, T. Luo, *IEEE Transactions on Information Theory*, 2016.
- [2] "The rotation of eigenvectors by a perturbation," C. Davis, and W. Kahan, *SIAM Journal on Numerical Analysis*, 1970.
- [3] "Matrix completion from a few entries," R. Keshavan, A. Montanari, and S. Oh, *IEEE Transactions on Information Theory*, 2010.
- [4] "Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees," Y. Chen and M. Wainwright, *arXiv preprint arXiv:1509.03025*, 2015.
- [5] "Implicit Regularization in Nonconvex Statistical Estimation: Gradient Descent Converges Linearly for Phase Retrieval, Matrix Completion and Blind Deconvolution," C. Ma, K. Wang, Y. Chi and Y. Chen, *arXiv preprint arXiv:1711.10467*, 2017.

Reference

- [6] "*No Spurious Local Minima in Nonconvex Low Rank Problems: A Unified Geometric Analysis*," R. Ge, C. Jin, and Y. Zheng, *ICML*, 2017.
- [7] "*Symmetry, Saddle Points, and Global Optimization Landscape of Nonconvex Matrix Factorization*," X. Li et al., *arXiv preprint arxiv:1612.09296*, 2016.
- [8] "*Topics in random matrix theory*," T. Tao, *American mathematical society*, 2012.
- [9] "*Harnessing Structures in Big Data via Guaranteed Low-Rank Matrix Estimation*," Y. Chen, and Y. Chi, *arXiv preprint arXiv:1802.08397*, 2018.
- [10] "*How to escape saddle points efficiently*," Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I, *arXiv preprint arXiv:1703.00887*, 2017.