

## Lecture 6: Spectral Algorithms I

Lecturer: Yudong Chen

Scribe: Nayoung Lee

In this lecture,<sup>1</sup> we will introduce the singular perturbation theory. Previously, we mainly focused on the operator or Frobenius norm of the difference in the estimated matrix and the true matrix. This lecture will focus on the singular values/vectors of the estimator and the true matrix. We first begin with a motivating example where a small perturbation in the matrix can lead to a big difference in the singular vector of the matrix and show that a gap in the singular values are needed for stability. Then we introduce a new metric, the subspace distance, and two theorems — the Wedin’s  $\sin \Theta$  theorem and the Weyl Inequality. We apply the two theorem in a simple case of the matrix completion problem.

## 1 Notation

A quick summary of the notation.

- For a matrix  $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$ , we use  $\|\mathbf{A}\|_F$  to denote its Frobenius norm and  $\|\mathbf{A}\|_{op}$  its operator/spectral norm (i.e., the largest singular value of  $\mathbf{A}$ ).
- $\sigma_i(\mathbf{A})$  is the  $i$ -th largest singular value of matrix  $\mathbf{A}$ .
- $\text{col}(\mathbf{A})$  denotes the column space of the given matrix  $\mathbf{A}$ .

## 2 Recap

In previous lectures, we covered spectral algorithms for two problems: (i) low-rank matrix completion and (ii) non-parametric Bradley-Terry model. We followed the following setup where given a (noisy/partial) observation  $\mathbf{Y}$  of some unknown matrix  $\mathbf{Y}^*$ , our goal is to find an good estimator  $\hat{\mathbf{Y}}$  of  $\mathbf{Y}^*$ .

$$\begin{array}{ccccc} \mathbf{Y}^* & \longrightarrow & \mathbf{Y} & \longrightarrow & \hat{\mathbf{Y}} \\ \text{unknown} & & \text{observation} & & \text{estimator} \end{array}$$

We derived bounds on the error of estimating  $\mathbf{Y}^*$ , in terms of  $\|\hat{\mathbf{Y}} - \mathbf{Y}^*\|_{op}$  and/or  $\|\hat{\mathbf{Y}} - \mathbf{Y}^*\|_F$ . However, sometimes the row/column spaces and the eigen/singular vectors of  $\mathbf{Y}^*$  are of interest rather than the matrix  $\mathbf{Y}^*$  itself.

## 3 Motivating Examples

**Example 1** (Bradley-Terry Model).

An estimate of true ranking can be extracted from eigenvector of  $\hat{\mathbf{Y}}$  (will be covered in next lecture)

---

<sup>1</sup>Reading:

- Rank centrality: Ranking from pairwise comparisons, S. Negahban, S. Oh, D. Shah, Operations Research, 2016. <https://arxiv.org/abs/1209.1688>
- Spectral method and regularized MLE are both optimal for top-K ranking, Yuxin Chen, Jianqing Fan, Cong Ma, Kaizheng Wang, Annals of Statistics, 2019. <https://arxiv.org/abs/1707.09971>
- Lecture Notes for ELE 520, Yuxin Chen, Princeton University.

**Example 2** (Spectral Initialization for Non-convex Matrix Completion<sup>2</sup>).

Suppose that the spectral estimator and true matrix have rank- $r$  singular value decomposition  $\mathbf{Y}^* = \mathbf{U}^* \boldsymbol{\Sigma}^* \mathbf{V}^{*\top}$  and  $\hat{\mathbf{Y}} = \hat{\mathbf{U}} \hat{\boldsymbol{\Sigma}} \hat{\mathbf{V}}^\top$ . Let

$$\begin{aligned} \mathbf{F}_0 &= \hat{\mathbf{U}} \hat{\boldsymbol{\Sigma}}^{1/2}, & \mathbf{F}^* &= \mathbf{U}^* \boldsymbol{\Sigma}^{*1/2}, \\ \mathbf{G}_0 &= \hat{\mathbf{V}} \hat{\boldsymbol{\Sigma}}^{1/2}, & \mathbf{G}^* &= \mathbf{V}^* \boldsymbol{\Sigma}^{*1/2}, \end{aligned}$$

where  $\hat{\mathbf{F}}, \hat{\mathbf{G}}, \mathbf{F}^*, \mathbf{G}^* \in \mathbb{R}^{n \times r}$ . Then the estimator and true matrix can be written in the factorized form  $\hat{\mathbf{Y}} = \hat{\mathbf{F}} \hat{\mathbf{G}}^\top$  and  $\mathbf{Y}^* = \mathbf{F}^* \mathbf{G}^{*\top}$ .

The matrices  $\hat{\mathbf{F}}, \hat{\mathbf{G}}$ , given by the singular values/vectors of the spectral estimator, are often used as an initial solution for gradient descent method applied to the non-convex formulation:

$$\min_{\mathbf{F}, \mathbf{G} \in \mathbb{R}^{d \times r}} \sum_{(i,j) \text{ observed}} \left( (\mathbf{F} \mathbf{G}^\top)_{i,j} - \mathbf{Y}_{ij}^* \right)^2.$$

We want to show that  $\hat{\mathbf{F}}, \hat{\mathbf{G}}$  is close to  $\mathbf{F}^*, \mathbf{G}^*$ , so they are in a local convex region of the above (globally nonconvex) objective function. In this region, one may use standard arguments from convex optimization to show that converges to the true  $\mathbf{F}^*, \mathbf{G}^*$ .

The question is:

$$\text{If } \hat{\mathbf{Y}} \approx \mathbf{Y}^*, \text{ then do we have } \hat{\mathbf{F}} \approx \mathbf{F}^* \text{ and } \hat{\mathbf{G}} \approx \mathbf{G}^*?$$

We will return to this example at the end of the lecture.

**Remark** (Comparison with spectral algorithm)

Recall our previous error bound for spectral algorithm:

$$\frac{1}{n^2} \|\hat{\mathbf{Y}} - \mathbf{Y}^*\|_F^2 \leq \frac{r \log n}{pn}.$$

Even when  $p$  is close to 1, in which case most entries are observed, the bound above gives a non-zero error. The non-convex formulation can be viewed as a refinement of the spectral estimator; when  $p$  is sufficiently large, this formulation can achieve *exact* recovery of  $\mathbf{Y}^*$  (under appropriate conditions).

## 4 Singular Perturbation Theory

Let  $\mathbf{M} \in \mathbb{R}^{n_1 \times n_2}$  be the *original matrix*, and  $\hat{\mathbf{M}} = \mathbf{M} + \mathbf{H}$  be the *perturbed matrix*, where  $\mathbf{H}$  is the *perturbation* or noise. Assume  $n_1 \geq n_2$ . The singular value decomposition (SVD) of  $\mathbf{M}$  and  $\hat{\mathbf{M}}$  is given by,

$$\begin{aligned} \mathbf{M} &= [\mathbf{U}_0 \quad \mathbf{U}_1] \begin{bmatrix} \boldsymbol{\Sigma}_0 & 0 \\ 0 & \boldsymbol{\Sigma}_1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{V}_0^\top \\ \mathbf{V}_1^\top \end{bmatrix}, \\ \hat{\mathbf{M}} &= [\hat{\mathbf{U}}_0 \quad \hat{\mathbf{U}}_1] \begin{bmatrix} \hat{\boldsymbol{\Sigma}}_0 & 0 \\ 0 & \hat{\boldsymbol{\Sigma}}_1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \hat{\mathbf{V}}_0^\top \\ \hat{\mathbf{V}}_1^\top \end{bmatrix}, \end{aligned} \tag{1}$$

where

2

- Harnessing Structures in Big Data via Guaranteed Low-Rank Matrix Estimation, Yudong Chen and Yuejie Chi. IEEE Signal Processing Magazine, vol. 35, no. 4, pp. 14-31, 2018. <https://arxiv.org/abs/1802.08397>
- Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview, Yuejie Chi, Yue M. Lu, Yuxin Chen. IEEE Transactions on Signal Processing, vol. 67, no. 20, pp. 5239-5269, 2019. <https://arxiv.org/abs/1809.09573>

- $\mathbf{U}_0 \in \mathbb{R}^{n_1 \times r}$ ,  $\mathbf{\Sigma}_0 \in \mathbb{R}^{r \times r}$ ,  $\mathbf{V}_0 \in \mathbb{R}^{n_2 \times r}$  correspond to the top- $r$  singular values/vectors of  $\mathbf{M}$ ;
- $\hat{\mathbf{U}}_0 \in \mathbb{R}^{n_1 \times r}$ ,  $\hat{\mathbf{\Sigma}}_0 \in \mathbb{R}^{r \times r}$ ,  $\hat{\mathbf{V}}_0 \in \mathbb{R}^{n_2 \times r}$  correspond to the top- $r$  singular values/vectors of  $\hat{\mathbf{M}}$ ;
- similarly,  $\mathbf{U}_1 \in \mathbb{R}^{n_1 \times (n_1 - r)}$ ,  $\mathbf{\Sigma}_1 \in \mathbb{R}^{(n_2 - r) \times (n_2 - r)}$ ,  $\mathbf{V}_1 \in \mathbb{R}^{n_2 \times (n_2 - r)}$  correspond to the bottom singular values/vectors of  $\mathbf{M}$ ;
- $\hat{\mathbf{U}}_1 \in \mathbb{R}^{n_1 \times (n_1 - r)}$ ,  $\hat{\mathbf{\Sigma}}_1 \in \mathbb{R}^{(n_2 - r) \times (n_2 - r)}$ ,  $\hat{\mathbf{V}}_1 \in \mathbb{R}^{n_2 \times (n_2 - r)}$  correspond to the bottom singular values/vectors of  $\hat{\mathbf{M}}$ .

**Example 3.** Consider two matrices  $\mathbf{M}$  and  $\hat{\mathbf{M}}$ , where  $\mathbf{M} - \hat{\mathbf{M}}$  is small:

$$\begin{aligned} \mathbf{M} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 + \epsilon \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 + \epsilon & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \\ \hat{\mathbf{M}} &= \begin{bmatrix} 1 & \epsilon \\ \epsilon & 1 \end{bmatrix} = \left( \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \right) \begin{bmatrix} 1 + \epsilon & 0 \\ 0 & 1 - \epsilon \end{bmatrix} \left( \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \right). \end{aligned}$$

We consider two cases:

- Setting I:  $r = 1$

In this case,  $\mathbf{U}_0$  and  $\hat{\mathbf{U}}_0 \in \mathbb{R}^{2 \times 1}$  correspond to the top-1 left singular vectors of  $\mathbf{M}$  and  $\hat{\mathbf{M}}$ , respectively. Although  $\mathbf{M} - \hat{\mathbf{M}}$  is small, their top left singular vectors are very different:

$$\mathbf{U}_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \text{and} \quad \hat{\mathbf{U}}_0 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

This is an example where a small perturbation can significantly change the singular vectors of a matrix. Reason:  $\sigma_1(\mathbf{M}) = \sigma_2(\mathbf{M})$ . For stability, a gap in singular values  $\sigma_1(\mathbf{M}) - \sigma_2(\mathbf{M}) \gtrsim \epsilon$  is needed.

- Setting II:  $r = 2$

In this case,  $\mathbf{U}_0$  and  $\hat{\mathbf{U}}_0 \in \mathbb{R}^{2 \times 2}$  correspond to the top-2 singular vectors of  $\mathbf{M}$  and  $\hat{\mathbf{M}}$ , respectively. We see that  $\mathbf{U}_0$  and  $\hat{\mathbf{U}}_0$  are very different:

$$\mathbf{U}_0 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{and} \quad \hat{\mathbf{U}}_0 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

However,  $\text{col}(\mathbf{U}_0) = \text{col}(\hat{\mathbf{U}}_0)$ ; that is, the columns of  $\mathbf{U}_0$  and  $\hat{\mathbf{U}}_0$  span the the same vector space  $\mathbb{R}^2$ , although they have different basis representation.

Setting II above shows that  $\|\hat{\mathbf{U}}_0 - \mathbf{U}_0\|_F$  or  $\|\hat{\mathbf{U}}_0 - \mathbf{U}_0\|_{op}$  is not the right choice of metric. Instead, we use the following subspace distance metric.

**Definition 1** (Subspace Distance). *The subspace distance between  $\mathbf{U}_0$  and  $\hat{\mathbf{U}}_0$  given in (1) is defined as*

$$\begin{aligned} \text{dist}(\hat{\mathbf{U}}_0, \mathbf{U}_0) &\triangleq \|\hat{\mathbf{U}}_0 \hat{\mathbf{U}}_0^\top - \mathbf{U}_0 \mathbf{U}_0^\top\|_{op} \\ &= \|\hat{\mathbf{U}}_0^\top \mathbf{U}_1\|_{op} = \|\mathbf{U}_0^\top \hat{\mathbf{U}}_1\|_{op} \\ &= \max\{|\sin \theta_1|, \dots, |\sin \theta_r|\}, \end{aligned}$$

where

$$\theta_i \triangleq i\text{-th principal angle between } \text{col}(\mathbf{U}_0) \text{ and } \text{col}(\hat{\mathbf{U}}_0).$$

Some remarks:

- $\text{dist}(\hat{\mathbf{U}}_0, \mathbf{U}_0) \in [0, 1]$ .
- $\hat{\mathbf{U}}_0 \hat{\mathbf{U}}_0^\top$  and  $\mathbf{U}_0 \mathbf{U}_0^\top$  are the projection matrices onto  $\text{col}(\hat{\mathbf{U}}_0)$  and  $\text{col}(\mathbf{U}_0)$ , respectively.

- When  $r = 1$  (so  $\hat{\mathbf{U}}_0, \mathbf{U}_0 \in \mathbb{R}^n$ ), we have

$$\begin{aligned}\theta_1 &= \arccos(\hat{\mathbf{U}}_0^\top \mathbf{U}_0) = \angle(\hat{\mathbf{U}}_0, \mathbf{U}_0), \\ (\sin \theta_1)^2 &= 1 - (\cos \theta_1)^2 = 1 - (\hat{\mathbf{U}}_0^\top \mathbf{U}_0)^2 = (\hat{\mathbf{U}}_0^\top \mathbf{U}_1)^2.\end{aligned}$$

- For the proof of the equalities in Definition 1, see, e.g., Part 2 of Theorem I.5.5 in G. W. Stewart and Ji guang Sun (1990), *Matrix Perturbation Theory*, Academic Press, Boston.
- Technicality: one needs to be a bit careful with the relationship between  $r$  and  $n - r$ , since at most  $k := \min\{r, n - r\}$  of the principal angles are nonzero (equivalently, one may define only the first  $k$  principal angles, as done in Stewart and Sun's book.)

#### 4.1 Wedin's $\sin \Theta$ Theorem and Weyl Inequality

Recall that  $\hat{\mathbf{M}} = \mathbf{M} + \mathbf{H}$ . For the singular vectors, we have Wedin's  $\sin \Theta$  Theorem.

**Theorem 4** (Wedin's  $\sin \Theta$  Theorem).

Suppose  $\sigma_r(\mathbf{M}) - \sigma_{r+1}(\hat{\mathbf{M}}) \geq \Delta > 0$ . Then the following inequality holds:

$$\begin{aligned}\max\{\text{dist}(\hat{\mathbf{U}}_0, \mathbf{U}_0), \text{dist}(\hat{\mathbf{V}}_0, \mathbf{V}_0)\} &\leq \frac{\max\{\|\mathbf{H}\mathbf{V}_0\|_{op}, \|\mathbf{H}^\top \mathbf{U}_0\|_{op}\}}{\Delta} \\ &\leq \frac{\|\mathbf{H}\|_{op}}{\Delta}.\end{aligned}$$

The second inequality above follows from  $\|\mathbf{H}\mathbf{V}_0\|_{op} \leq \|\mathbf{H}\|_{op} \|\mathbf{V}_0\|_{op} = \|\mathbf{H}\|_{op}$ . Note that the above bound is useful only when the right hand side is less than 1.

For the singular values we have Weyl's Inequality.

**Theorem 5** (Weyl's Inequality).

We have

$$|\sigma_i(\mathbf{M}) - \sigma_i(\hat{\mathbf{M}})| \leq \|\mathbf{H}\|_{op}, \quad \forall i = 1, 2, \dots, n.$$

Consequently

$$\sigma_r(\mathbf{M}) - \sigma_{r+1}(\hat{\mathbf{M}}) \geq \sigma_r(\mathbf{M}) - \sigma_{r+1}(\mathbf{M}) - \|\mathbf{H}\|_{op}.$$

Combining the above two theorems, we have:

**Corollary 1.**

The following inequality holds, given that the denominator is positive:

$$\max\{\text{dist}(\hat{\mathbf{U}}_0, \mathbf{M}_0), \text{dist}(\hat{\mathbf{V}}_0, \mathbf{V}_0)\} \leq \frac{\|\mathbf{H}\|_{op}}{\sigma_r(\mathbf{M}) - \sigma_{r+1}(\hat{\mathbf{M}}) - \|\mathbf{H}\|_{op}}.$$

Note that the inequality is valid even if  $\sigma_1(\mathbf{M}) = \sigma_2(\mathbf{M}) = \dots = \sigma_r(\mathbf{M})$  — only the gap between the  $r$ -th and  $(r + 1)$ -th singular value is required. The perturbation may change the ordering of these  $r$  singular values, but it does not matter since we are using a distance metric in the column space.

**Remark**

Similar bounds for eigenvectors of symmetric matrices are given in the **Davis-Kahan  $\sin \Theta$  Theorem**.

## 4.2 Application: Matrix Completion

Suppose  $\mathbf{Y}^* \in \mathbb{R}^{n \times n}$  is rank  $r$  with SVD  $\mathbf{Y}^* = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ , where  $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_r)$ ,  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times r}$ . From Lecture 1, the spectral algorithm gives a rank- $r$  estimator  $\hat{\mathbf{Y}}$  bound on  $\|\hat{\mathbf{Y}} - \mathbf{Y}^*\|_{op}$ .

Now, suppose  $\hat{\mathbf{Y}}$  has SVD  $\hat{\mathbf{Y}} = \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\hat{\mathbf{V}}^\top$ , where  $\hat{\mathbf{U}}, \hat{\mathbf{V}} \in \mathbb{R}^{n \times r}$ . Then the following can be derived.

- Apply Theorem 4 to bound  $\|\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}\|_{op}$ .
- Apply Corollary 1 to get bounds on  $\text{dist}(\hat{\mathbf{U}}, \mathbf{U}), \text{dist}(\hat{\mathbf{V}}, \mathbf{V})$ .
- We can then use  $\hat{\mathbf{F}} \triangleq \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}^{1/2}$  and  $\hat{\mathbf{G}}_0 \triangleq \hat{\mathbf{V}}\hat{\mathbf{\Sigma}}^{1/2}$  at the initial solution for gradient descent for the non-convex formulation in Example 2.

**A simple case (rank  $r = 1$ ):**

Suppose  $\mathbf{Y} = \mathbf{u}\mathbf{v}^\top$  where  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ ,  $\mathbf{u} = \frac{1}{\|\tilde{\mathbf{u}}\|_2} \tilde{\mathbf{u}}$ ,  $\mathbf{v} = \frac{1}{\|\tilde{\mathbf{v}}\|_2} \tilde{\mathbf{v}}$ ,  $\tilde{\mathbf{u}}, \tilde{\mathbf{v}} \sim \mathcal{N}(0, \mathbf{I}_{n \times n})$ . Let the rank-1 estimator  $\hat{\mathbf{Y}} = \hat{\sigma}\hat{\mathbf{u}}\hat{\mathbf{v}}^\top$  be obtained by the spectral algorithm.

**Lemma 1.**

*If the sampling probability satisfies  $p \geq C \frac{\log^3 n}{n}$  for a sufficiently large constant  $C$ , then w.h.p.,*

$$\text{dist}(\hat{\mathbf{u}}, \mathbf{u}) \leq \frac{1}{100}, \quad \text{dist}(\hat{\mathbf{v}}, \mathbf{v}) \leq \frac{1}{100}, \quad |\hat{\sigma} - 1| \leq \frac{1}{100}.$$

**Sketch of Proof**

1. Show  $\|\tilde{\mathbf{u}}\|_2^2 \asymp \|\tilde{\mathbf{v}}\|_2^2 \asymp n$  w.h.p. by scalar Bernstein's inequality
2. Show  $\max_i |\tilde{u}_i| \lesssim \sqrt{\log n}$ ,  $\max_j |\tilde{v}_j| \lesssim \sqrt{\log n}$  w.h.p by bounding each  $|\tilde{u}_i|$  and applying a union bound.
3. Bound  $\max_{i,j} |Y_{ij}^*|$  by combining 1 and 2.
4. Bound  $\|\hat{\mathbf{Y}} - \mathbf{Y}^*\|_{op}$  (similar to Lecture 1, using matrix Bernstein - Hint: Matrix Bernstein require boundedness, ensured by 3.)
5. Bound  $\text{dist}(\hat{\mathbf{u}}, \mathbf{u}), \text{dist}(\hat{\mathbf{v}}, \mathbf{v})$  by Corollary 1.

□