

Lecture 2: Non-parametric Bradley-Terry Model

Lecturer: Yudong Chen

Scribe: Augustine(Runshi) Tang

In this lecture, we will introduce another application of Spectral Algorithm and Matrix Bernstein Inequality, which is to prove an upper bound for a result of Spectral Algorithm used on the so called Non-parametric Bradley-Terry Model.¹

1 Notation

A quick summary of the notation.

1. **Random variables:** X, Y, U, V
2. **Ranges/alphabets:** $\mathcal{X}, \mathcal{Y}, \mathcal{U}, \mathcal{V}$
3. **Specific values:** x, y, u, v

For a vector $u \in \mathbb{R}^d$, we use $\|u\|_2$ to denote its ℓ_2 norm and $\|u\|_\infty$ its ℓ_∞ norm. For a matrix $A \in \mathbb{R}^{d_1 \times d_2}$, we use $\|A\|_F$ to denote its Frobenius norm and $\|A\|_{op}$ its operator/spectral norm (i.e., the largest singular value of A). For two matrices A, B of the same dimension, $\langle A, B \rangle := \text{tr}(A^\top B)$ denotes their trace inner product. The trace inner product reduces to the usual inner product between vectors for when $A, B \in \mathbb{R}^{d \times 1}$.

2 Preliminaries

Lemma 1 (Matrix Bernstein's inequality²). *Let X_1, \dots, X_N be independent, mean zero, $n \times n$ symmetric random matrices, such that $\|X_i\| \leq K$ almost surely for all i . Then, for every $t \geq 0$, we have*

$$\mathbb{P} \left\{ \left\| \sum_{i=1}^N X_i \right\|_{op} \geq t \right\} \leq 2n \exp \left(-\frac{t^2/2}{\sigma^2 + Kt/3} \right).$$

Here $\sigma^2 = \left\| \sum_{i=1}^N \mathbb{E} X_i^2 \right\|_{op}$ is the norm of the matrix variance of the sum. In particular, we can express this bound as the mixture of sub-gaussian and sub-exponential tail, just like in the scalar Bernstein's inequality:

$$\mathbb{P} \left\{ \left\| \sum_{i=1}^N X_i \right\|_{op} \geq t \right\} \leq 2n \exp \left[-c \cdot \min \left(\frac{t^2}{\sigma^2}, \frac{t}{K} \right) \right]$$

Lemma 2 (Eckart-Young-Mirsky Theorem³). *Let*

$$D = U \Sigma V^\top \in \mathbb{R}^{m \times n}, \quad m \geq n$$

be the singular value decomposition (SVD) of D and partition $U, \Sigma =: \text{diag}(\sigma_1, \dots, \sigma_m)$, and V as follows:

$$U =: [U_1 \quad U_2], \quad \Sigma =: \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix}, \quad \text{and} \quad V =: [V_1 \quad V_2]$$

¹Reference: Chatterjee, Sourav. "Matrix estimation by universal singular value thresholding." The Annals of Statistics 43.1 (2015): 177-214. Section 2.7.

²Reference: Theorem 5.4.1 of HDP-book.

³Reference: https://en.wikipedia.org/wiki/Low-rank_approximation

where U_1 is $m \times r$, Σ_1 is $r \times r$, and V_1 is $n \times r$, and $\sigma_1 \geq \dots \geq \sigma_m$. Then the rank- r matrix, obtained from the truncated singular value decomposition

$$\widehat{D}^* = U_1 \Sigma_1 V_1^\top,$$

satisfies

$$\|D - \widehat{D}^*\|_{\text{F}} = \min_{\text{rank}(\widehat{D}) \leq r} \|D - \widehat{D}\|_{\text{F}} = \sqrt{\sigma_{r+1}^2 + \dots + \sigma_m^2}$$

The minimizer \widehat{D}^* is unique if and only if $\sigma_{r+1} \neq \sigma_r$.

In words, \widehat{D}^* , given by the truncated SVD of D , is a best rank- r approximation of D .

3 Non-parametric Bradley-Terry Model

Assume an ordered set Ω with n elements ω_i and we will use ' \succ ' to denote the ordering. This ordering is unknown, but imagine the setting where one may arrange matches between pairs of items and observe the results of the matches. The results of the matches are random; if $\omega_i \succ \omega_j$, then ω_i has a higher chance of beating the opponent than ω_j does against the same opponent.

The Non-parametric Bradley-Terry model formalized the above setting. If $\omega_i \succ \omega_j$, then $\mathbb{P}(\omega_i \text{ beats } \omega_k) \geq \mathbb{P}(\omega_j \text{ beats } \omega_k)$ for any k . Denote a matrix Y^* with $Y_{ij}^* := \mathbb{P}(\omega_i \text{ beats } \omega_j)$. Further denote $Y \in \{0, 1\}^{n \times n}$ as an observed random matrix, whose entries independently follow the distribution:

$$Y_{ij} = \begin{cases} 1, & \text{with probability } pY_{ij}^*, \\ 0, & \text{with probability } p(1 - Y_{ij}^*) \\ 0, & \text{with probability } 1 - p, \end{cases}$$

where $p \in [0, 1]$. This model can be interpreted as follows: there are n teams and we want to find a ranking among them. With probability p , a match is played between a pair of teams independently. We use Y to denote the results of the matches. If ω_i beats ω_j in the match, then we write $Y_{ij} = 1$. If ω_j beats ω_i , or if a match is not played between them, then we write $Y_{ij} = 0$. Our goal is to estimate Y^* given that Y .⁴

We use a Spectral Algorithm. In particular, our estimator \widehat{Y} is given by the best rank- r approximation of $\frac{Y}{p}$, namely

$$\widehat{Y} := \arg \min_{\text{rank}(D) \leq r} \left\| \frac{Y}{p} - D \right\|_{\text{F}},$$

where we choose the rank as $r = \sqrt{np}$. It will be clear why we choose this value later in the proof of the error bound. Note that \widehat{Y} can be computed using the truncated SVD of $\frac{Y}{p}$, thanks to the Eckart-Young-Mirsky Theorem.

Our proof relies on two intermediate result. An immediate result from the Matrix Bernstein Inequality is:

Lemma 3. When $p \geq (\log n)/n$,

$$\left\| \frac{1}{p}Y - Y^* \right\|_{\text{op}} \lesssim \sqrt{\frac{n \log n}{p}}$$

holds with high probability.

The proof follows similar lines as in last lecture and is left as an exercise.

We also claim that:

⁴Given a good estimate of Y^* , one may further estimate the ranking. We will not discuss this problem in this lecture.

Claim 1. *There exists a rank- r matrix Z such that $\begin{cases} \|Z - Y^*\|_F^2 \leq \frac{n^2}{r} \\ |Z_{ij}| \leq 1, \forall i, j \end{cases}$.*

Proof For each $i = 1, \dots, n$ define the number $s_i \triangleq \sum_{j=1}^n Y_{ij}^*$. For each $l = 1, \dots, r$, define the index set $T_l \triangleq \left\{ i : s_i \in \left[\frac{n(l-1)}{r}, \frac{nl}{r} \right) \right\}$ and let $k(l) \triangleq$ first element in T_l .

For each $l = 1, \dots, r$ and all $i \in T_l$, set

$$Z_{i-} = Y_{k(l)-}^*.$$

This gives a matrix $Z \in [0, 1]^{n \times n}$ with row vectors as Z_{i-} and it has rank less than r .

For each $l = 1, \dots, r$ and each $i \in T_l$:

(1) If $\omega_i \succ \omega_{k(l)}$:

$$\begin{aligned} \sum_{j=1}^n (Y_{ij}^* - Z_{ij})^2 &= \sum_{j=1}^n (Y_{ij}^* - Y_{k(l)j}^*)^2 \\ &\leq \sum_{j=1}^n |Y_{ij}^* - Y_{k(l)j}^*| \\ &= \sum_{j=1}^n (Y_{ij}^* - Y_{k(l)j}^*) \\ &= S_i - S_{k(l)} \\ &\leq \frac{n}{r}. \end{aligned}$$

(2) if $\omega_{k(l)} \prec \omega_i$:

$$\begin{aligned} \sum_{j=1}^n (Y_{ij}^* - Z_{ij})^2 &= \sum_{j=1}^n (Y_{ij}^* - Y_{k(l)j}^*)^2 \\ &\leq \sum_{j=1}^n |Y_{ij}^* - Y_{k(l)j}^*| \\ &= - \sum_{j=1}^n (Y_{ij}^* - Y_{k(l)j}^*) \\ &= S_{k(l)} - S_i \\ &\leq \frac{n}{r}. \end{aligned}$$

Sum up over all i , we have $\|Y^* - Z\|_F^2 \leq \sum_{i=1}^n \frac{n}{r} = \frac{n^2}{r}$. ■

We are now ready to prove an error upper bound for the estimator \hat{Y} from Spectral Algorithm.

Let Z be the rank- r matrix given by Claim 1. Since \hat{Y} is best rank- r approximation of $\frac{Y}{p}$, we have

$$\left\| \frac{1}{p} Y - Z \right\|_F^2 \geq \left\| \frac{1}{p} Y - \hat{Y} \right\|_F^2 = \left\| \frac{1}{p} Y - Z \right\|_F^2 + \left\| \hat{Y} - Z \right\|_F^2 + 2 \left\langle \frac{Y}{p} - Z, Z - \hat{Y} \right\rangle.$$

Rearranging terms gives

$$\begin{aligned}
\|Z - \hat{Y}\|_F^2 &\leq 2 \left\langle \frac{Y}{p} - Z, \hat{Y} - Z \right\rangle \\
&= 2 \left\langle \frac{Y}{p} - Y^*, \hat{Y} - Z \right\rangle + 2 \left\langle Y^* - Z, \hat{Y} - Z \right\rangle \\
&\leq 2 \left\| \frac{1}{p}Y - Y^* \right\|_{op} \|\hat{Y} - Z\|_* + 2 \|Y^* - Z\|_F \|\hat{Y} - Z\|_F \\
&\leq 2 \left\| \frac{1}{p}Y - Y^* \right\|_{op} \sqrt{2r} \|\hat{Y} - Z\|_F + 2 \|Y^* - Z\|_F \|\hat{Y} - Z\|_F
\end{aligned}$$

So

$$\|Z - \hat{Y}\|_F \leq 2\sqrt{2r} \left\| \frac{1}{p}Y - Y^* \right\|_{op} + 2 \|Y^* - Z\|_F$$

Thus

$$\begin{aligned}
\|\hat{Y} - Y^*\|_F &\leq \|Z - \hat{Y}\|_F + \|Z - Y^*\|_F \\
&\leq 2\sqrt{2r} \left\| \frac{1}{p}Y - Y^* \right\|_{op} + 3 \|Z - Y^*\|_F \\
&\leq C\sqrt{r} \sqrt{\frac{n \log n}{p}} + \frac{3n}{\sqrt{r}} \quad \text{w.h.p. by Lemma 3 and Claim 1 above} \\
&\leq C\sqrt{r} \sqrt{\frac{n \log n}{p}} + \frac{3n\sqrt{\log n}}{\sqrt{r}},
\end{aligned}$$

where $C > 0$ is a constant from Lemma 3. Choosing $r = \sqrt{pn}$ (which approximately balances the two terms above), we obtain the following error bound for the spectral algorithm.

Theorem 1. *Under the above setting, we have, with high probability,*

$$\frac{1}{n^2} \|\hat{Y} - Y^*\|_F^2 \lesssim \frac{\log n}{\sqrt{np}}.$$

Note that, if $p \gtrsim \frac{\log^2 n}{n\varepsilon^2}$, then RHS $\leq \varepsilon$.