

## Lecture 6: Concentration phenomena in mean-field asymptotic

Lecturer: Song Mei

Scriber: Omer Ronen

Proof reader: Tae Joo Ahn

## 1 Concentration phenomena in mean-field asymptotic

In the **non-asymptotic** regime a typical bound is of the form

$$\mathbb{P}\left(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2/d \geq C\sqrt{\frac{d \log d/\delta}{n}}\right) \leq \delta. \quad (1)$$

In the **asymptotic** regime the picture is a bit different. Typically the risk  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2/d$  would not go to zero when we consider a regime in which  $d/n$  converges to some constant. Under such regime the risk lower bound in (1) would go to  $\infty$  and the risk would not converge to zero, but rather concentrate around its expectation, formally

1.  $\lim_{n \rightarrow \infty} \mathbb{P}(|\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2/n - \mathbb{E}\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2/n| \geq \epsilon) = 0,$
2.  $\lim_{n \rightarrow \infty} \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2/n = \text{some formula.}$

In this lecture we will focus on step 1. Step 2 will be covered in later lectures.

## 2 The Gaussian concentration inequality

**Proposition 1** (Gaussian concentration inequality). *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be an  $L$ -Lipschitz function:*

$$|f(x_1) - f(x_2)| \leq L\|x_1 - x_2\|_2, \quad \forall x_1, x_2 \in \mathbb{R}^d,$$

*and let  $\mathbf{G} = (G_1, \dots, G_d)$  be a vector of standard Gaussian random variables ( $G_i \sim_{i.i.d} N(0, 1)$ ). Then we have*

$$\mathbb{P}(|f(\mathbf{G}) - \mathbb{E}f(\mathbf{G})| \geq t) \leq 2\exp\{-t^2/(2L^2)\}.$$

This proposition can be derived from the Gaussian log Sobolev inequality coupled with the Herbst's argument (c.f. [Thm 3.25] in Ramon van Handel's notes).

## 3 $\mathbb{Z}_2$ synchronization problem

Let us define the problem, we have our true parameter  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)^\top \in \mathbb{R}^n$  with  $\theta_i \sim_{i.i.d} \text{Unif}(\mathbb{Z}_2 = \{\pm 1\})$ . We have the observation vector  $\mathbf{Y} = \frac{\lambda}{n}\boldsymbol{\theta}\boldsymbol{\theta}^\top + \mathbf{W} \in \mathbb{R}^{n \times n}$  with  $\mathbf{W} \sim \text{GOE}(n)$ . That is  $W_{ij} \sim_{i.i.d} N(0, n^{-1})$  for  $1 \leq i \leq j \leq n$  and  $W_{ii} \sim_{i.i.d}$  for  $1 \leq i \leq n$  with the restriction that  $\mathbf{W}$  is symmetric ( $W_{ij} = W_{ji}$ ).

**Remark 2.  $\boldsymbol{\theta}\boldsymbol{\theta}^\top$  scaling.** *An interesting thing to consider is the choice of scaling. We choose the scaling of  $n$  so that the operator norms of  $\frac{\lambda}{n}\boldsymbol{\theta}\boldsymbol{\theta}^\top$  and  $\mathbf{W}$  would be on the same scale (we consider to get a non-trivial behavior of the spectral estimator which involves the calculation of the leading eigenvector). Indeed  $\|\mathbf{W}\|_{op} \stackrel{2}{\approx} C$  and  $\|\boldsymbol{\theta}\boldsymbol{\theta}^\top\|_{op} \approx \|\boldsymbol{\theta}\|_2 = n$  hence the choice of scaling is  $n$ , such that the signal to noise ration is balances between the two norms.*

We observe  $\mathbf{Y}$  and try to estimate  $\boldsymbol{\theta}$ , the spectral estimator is defined to be

$$\hat{\boldsymbol{\theta}}(\mathbf{Y}) = \hat{\boldsymbol{\theta}}_{\text{spec}}(\mathbf{Y}) = \arg \max_{\boldsymbol{\sigma} \in \mathbb{S}^{n-1}(\sqrt{n})} \langle \boldsymbol{\sigma}, \mathbf{Y} \boldsymbol{\sigma} \rangle / n$$

Our first goal is to show  $\max_{\boldsymbol{\sigma} \in \mathbb{S}^{n-1}(\sqrt{n})} \langle \boldsymbol{\sigma}, \mathbf{Y} \boldsymbol{\sigma} \rangle / n \approx$  its expectation (w.h.p), which follows from the proposition below.

**Proposition 3.** *Let  $\mathcal{R}_n \subseteq \{\boldsymbol{\theta} \in \mathbb{R}^n : \|\boldsymbol{\theta}\|_2^2 n \leq 1\}$  and let  $\mathbf{Y} = \mathbf{A} + \mathbf{W} \in \mathbb{R}^{n \times n}$ , where  $\mathbf{A}$  is deterministic and  $\mathbf{W} \sim \text{GOE}(n)$ . Then  $\delta \rightarrow 0$ , we have*

$$\mathbb{P}_{\mathbf{W}} \left( \left| \sup_{\boldsymbol{\sigma} \in \mathcal{R}_n} \langle \boldsymbol{\sigma}, \mathbf{Y} \boldsymbol{\sigma} \rangle / n - \mathbb{E} \left[ \sup_{\boldsymbol{\sigma} \in \mathcal{R}_n} \langle \boldsymbol{\sigma}, \mathbf{Y} \boldsymbol{\sigma} \rangle / n \right] \right| \leq \sqrt{\frac{4 \log(2/\delta)}{n}} \right) \geq 1 - \delta$$

We will now use the Gaussian inequality to prove the above mentioned proposition

*Proof.* Let  $\mathbf{G} \in \mathbb{R}^{n \times n}$  with  $G_{ij} \sim_{i.i.d} N(0, 1)$  for  $1 \leq i, j \leq n$  (which is not symmetric).

Denote  $\tilde{\mathbf{W}} = (\mathbf{G} + \mathbf{G}^\top) / \sqrt{2n}$  then we have  $\tilde{\mathbf{W}} \stackrel{d}{=} \mathbf{W}$ .

Define  $\tilde{\mathbf{Y}} = \mathbf{A} + \tilde{\mathbf{W}}$  and

$$\begin{aligned} f(\mathbf{G}) &= \sup_{\boldsymbol{\sigma} \in \mathcal{R}_n} \langle \boldsymbol{\sigma}, \tilde{\mathbf{Y}} \boldsymbol{\sigma} \rangle / n \\ &= \sup_{\boldsymbol{\sigma} \in \mathcal{R}_n} \langle \boldsymbol{\sigma}, (\mathbf{G} + \mathbf{G}^\top) / \sqrt{2n} \boldsymbol{\sigma} \rangle / n \end{aligned}$$

We would like to show that  $f$  is a  $L$ -Lipschitz function and then apply the inequality.

Let  $\mathbf{G}_1, \mathbf{G}_2 \in \mathbb{R}^{n \times n}$  and let  $\boldsymbol{\sigma}^* = \arg \sup_{\boldsymbol{\sigma} \in \mathcal{R}_n} \langle \boldsymbol{\sigma}, \tilde{\mathbf{Y}}_1 \boldsymbol{\sigma} \rangle / n$

$$\begin{aligned} f(\mathbf{G}_1) - f(\mathbf{G}_2) &= \sup_{\boldsymbol{\sigma} \in \mathcal{R}_n} \langle \boldsymbol{\sigma}, \tilde{\mathbf{Y}}_1 \boldsymbol{\sigma} \rangle / n + \inf_{\boldsymbol{\sigma} \in \mathcal{R}_n} -\langle \boldsymbol{\sigma}, \tilde{\mathbf{Y}}_2 \boldsymbol{\sigma} \rangle / n \\ &\leq \langle \boldsymbol{\sigma}^*, \tilde{\mathbf{Y}}_1 \boldsymbol{\sigma}^* \rangle / n + -\langle \boldsymbol{\sigma}^*, \tilde{\mathbf{Y}}_2 \boldsymbol{\sigma}^* \rangle / n \\ &= \langle \mathbf{G}_1 - \mathbf{G}_2, \boldsymbol{\sigma}^* (\boldsymbol{\sigma}^*)^\top \rangle \sqrt{\frac{2}{n}} / n \\ &\leq \|\mathbf{G}_1 - \mathbf{G}_2\|_{\text{op}} \underbrace{\|\boldsymbol{\sigma}^*\|_2^2 / n}_{\leq 1} \sqrt{\frac{2}{n}} \\ &\leq \|\mathbf{G}_1 - \mathbf{G}_2\|_{\text{F}} \sqrt{\frac{2}{n}} \end{aligned}$$

Which means that  $f$  is  $\sqrt{\frac{2}{n}}$ -Lipschitz

**Remark 4. Differential  $f$  proof sketch** *To get an intuition for this proof we can consider the case where  $f$  is differentiable. We use the implicit differentiation theorem to calculate the gradient of  $f$  and evaluate it at  $\boldsymbol{\sigma}^* = \arg \sup_{\boldsymbol{\sigma} \in \mathcal{R}_n}$ , indeed*

$$\|\nabla_{\mathbf{G}} f(\mathbf{G})\|_F = \sqrt{\frac{2}{n}} \left\| \frac{\boldsymbol{\sigma}^* (\boldsymbol{\sigma}^*)^\top}{n} \right\|_F = \sqrt{\frac{2}{n}} \|\boldsymbol{\sigma}^*\|_2^2 / n = \frac{2}{n}$$

Which means  $f$  is  $\frac{2}{n}$ -Lipschitz

We now use proposition 1 with  $L = \frac{2}{n}$  and get

$$\mathbb{P}(|f(\mathbf{G}) - \mathbb{E}f(\mathbf{G})| \geq t) \leq 2 \exp\{-(nt^2/4)\}$$

we solve

$$2 \exp\{-nt^2/4\} = \delta$$

and conclude our proof with  $t = \sqrt{\frac{4 \log(2/\delta)}{n}}$  □

The next step is to show the concentration of  $\langle \hat{\theta}, \theta \rangle^2 / 2$  with  $\hat{\theta} = \arg \sup_{\sigma \in \mathbb{S}^{n-1}(\sqrt{n})} \langle \sigma, Y \sigma \rangle / n$ . We first state a result to be used:

**Remark 5. BBP phase transition** We define  $U_n(\lambda) = \sup_{\sigma \in \mathbb{S}^{n-1}(\sqrt{n})} \langle \sigma, Y \sigma \rangle / n$ . We have that

$$\lim_{n \rightarrow \infty} \mathbb{E} U_n(\lambda) = \begin{cases} 2 & , \lambda \leq 1 \\ \lambda + \frac{1}{\lambda} & , \lambda > 1 \end{cases}$$

Using implicit differentiation,

$$\begin{aligned} \frac{d}{d\lambda} U_n(\lambda) &= \frac{d}{d\lambda} \sup_{\sigma \in \mathbb{S}} \left[ \frac{\lambda}{n^2} \langle \sigma, \theta \rangle^2 + \langle \sigma, W \sigma \rangle / 2 \right] \\ &= \frac{1}{n^2} \langle \hat{\theta}, \theta \rangle^2 \end{aligned}$$

We now give a heuristic proof, for some intuition, we have shown ( $U$  is the same function as  $f$ )

$$\lim_{n \rightarrow \infty} \mathbb{P}(|U_n(\lambda) - \mathbb{E} U_n(\lambda)| \geq \epsilon) = 0$$

and we calculate ( $\approx$ ) the expectation

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \langle \theta, \hat{\theta} \rangle^2 &\approx \partial \lambda \lim_{n \rightarrow \infty} \mathbb{E}[U_n(\lambda)] \\ &= \partial \lambda \begin{cases} 2 & , \lambda \leq 1 \\ \lambda + \frac{1}{\lambda} & , \lambda > 1 \end{cases} \\ &= \begin{cases} 0 & , \lambda \leq 1 \\ 1 - \frac{1}{\lambda^2} & , \lambda > 1 \end{cases} \end{aligned}$$

We now show the deviation a bit more carefully, we define the "discrete differential of size  $\delta$ "

$$\begin{aligned} \Delta_n^+(\lambda, \delta) &= \frac{U_n(\lambda + \delta) - U_n(\lambda)}{\delta} \\ \Delta_n^-(\lambda, \delta) &= \frac{U_n(\lambda) - U_n(\lambda - \delta)}{\delta} \end{aligned}$$

Indeed

$$\begin{aligned} U_n(\lambda + \delta) &= \sup_{\sigma \in \mathbb{S}} \langle \sigma, W \sigma \rangle / n + \frac{\lambda + \delta}{n^2} \langle \sigma, \theta \rangle^2 \\ &\geq \underbrace{\frac{\langle \sigma^*, W \sigma^* \rangle}{n} + \frac{\lambda}{n^2} \langle \sigma^*, \theta \rangle^2}_{U_n(\lambda)} + \frac{\delta}{n^2} \langle \sigma^*, \theta \rangle^2 \\ &= U_n(\lambda) + \frac{\delta}{n^2} \langle \hat{\theta}, \theta \rangle^2 \end{aligned}$$

We conclude (by symmetry) the following inequality

$$\Delta_n^-(\lambda, \delta) \leq \frac{\langle \theta, \hat{\theta} \rangle^2}{n^2} \leq \Delta_n^+(\lambda, \delta) \tag{2}$$

Using BBP phase transition result we get

$$\begin{aligned}\Delta^+(\lambda, \delta) &= \frac{u(\lambda + \delta) - u(\lambda)}{\delta} \\ \Delta^-(\lambda, \delta) &= \frac{u(\lambda) - u(\lambda - \delta)}{\delta}\end{aligned}$$

for

$$u(\lambda) = \begin{cases} 2 & \lambda \leq 1 \\ \lambda + \frac{1}{\lambda} & \lambda > 1 \end{cases}$$

Next we apply the Gaussian concentration inequality

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\Delta_n^\pm - \Delta^\pm| \geq \epsilon) = 0 \quad (3)$$

Combining (3) with (2) we get

$$\lim_{n \rightarrow \infty} \mathbb{P}(\Delta^-(\lambda, \delta) - \epsilon \leq \frac{\langle \boldsymbol{\theta}, \hat{\boldsymbol{\theta}} \rangle^2}{n^2} \leq \Delta^+(\lambda, \delta) + \epsilon) = 1 \quad (4)$$

By the definition of the derivative

$$\lim_{\delta \rightarrow 0} \Delta^\pm(\delta, \lambda) = \begin{cases} 0 & \lambda \leq 1 \\ 1 - \frac{1}{\lambda^2} & \lambda > 1 \end{cases} \quad (5)$$

And finally combining (5) and (4) we get

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\langle \hat{\boldsymbol{\theta}}, \boldsymbol{\theta} \rangle^2 / n^2 - \Delta(\lambda)| \geq \epsilon) = 0$$

□

**Remark 6. General recipe to show the concentration of any  $M(\boldsymbol{\sigma})$**

1. Define the appropriate perturbed Hamiltonian  $H_h(\boldsymbol{\sigma}) = \langle \boldsymbol{\sigma}, \mathbf{Y} \boldsymbol{\sigma} \rangle / n + hM(\boldsymbol{\sigma})$
2. Show the concentration of  $\sup_{\boldsymbol{\sigma}} H_h(\boldsymbol{\sigma})$  at some  $h$  near 0 with  $H_h(\boldsymbol{\sigma})$  differentiable with respect to  $h$  at  $h = 0$
3. Obtain the result by taking derivative with respect to  $h$

## 4 Concentration of Lasso training loss

Let us consider the Lasso problem - we have a random matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$ ,  $A_{ij} \sim N(0, 1)$ . Our estimand is  $\mathbf{x}_0 \in \mathbb{R}^d$  with  $\|\mathbf{x}_0\|_2^2 / d \leq M$ . The observation vector is

$$\mathbf{y} = \mathbf{A} \mathbf{x}_0 + \boldsymbol{\epsilon}$$

with  $\boldsymbol{\epsilon} \sim N(\mathbf{0}_n, \tau^2 \mathbf{I}_n)$ .

We wish to consider two lasso estimators

- $\hat{\mathbf{x}}_1 = \arg \min_{\mathbf{x}} \frac{1}{\sqrt{n}} \|\mathbf{y} - \mathbf{A} \mathbf{x}\|_2 + \frac{\lambda}{d} \|\mathbf{x}\|_1$
- $\hat{\mathbf{x}}_2 = \arg \min_{\mathbf{x}} \frac{1}{n} \|\mathbf{y} - \mathbf{A} \mathbf{x}\|_2^2 + \frac{\lambda}{d} \|\mathbf{x}\|_1$

$\hat{\mathbf{x}}_2$  is the familiar Lasso problem while  $\hat{\mathbf{x}}_1$  is called "square root" lasso. We first show the concentration of  $\hat{\mathbf{x}}_1$

**Proposition 7.** Let  $\Omega \subseteq \mathbb{R}^d$  be a compact region and consider

$$\sup_{\mathbf{x} \in \Omega} \{\|\mathbf{x}\|_2^2/d\} \leq D$$

this is a compact set of radius  $\sqrt{D}$  times a constant. We define

$$^1f_\Omega(\mathbf{A}, \epsilon) = \min_{\mathbf{x} \in \Omega} \frac{1}{\sqrt{n}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 + \frac{\lambda}{d} \|\mathbf{x}\|_1$$

Then  $\exists K < \infty$  s.t  $\forall \delta > 0$  then

$$\mathbb{P}\left(|f_\Omega(\mathbf{A}, \epsilon) - \mathbb{E}f_\Omega(\mathbf{A}, \epsilon)| \geq \underbrace{K\left(\sqrt{\frac{d}{n}(M+D)} + \tau\right)}_{\text{constant}} \sqrt{\frac{\log(2/\delta)}{n}}\right) \leq \delta$$

*Proof.* We define the "standardized" version of  $\mathbf{A}$  and  $\epsilon$ ,  $\bar{\mathbf{A}} = \sqrt{n}\mathbf{A}$ ,  $\bar{\epsilon} = \epsilon/\tau$ . Let

$$F(\bar{\mathbf{A}}, \bar{\epsilon}) = f_\Omega(\bar{\mathbf{A}}/\sqrt{n}, \bar{\epsilon} \cdot \tau) = f_\Omega(\mathbf{A}, \epsilon)$$

we want to show that  $F$  is Lipschitz and apply the Gaussian concentration inequality. By definition

$$\begin{aligned} F(\bar{\mathbf{A}}, \bar{\epsilon}) &= \min_{\mathbf{x} \in \Omega} \frac{1}{\sqrt{n}} \|\mathbf{A}(\mathbf{x}_0 - \mathbf{x}) + \epsilon\|_2 + \frac{\lambda}{n} \|\mathbf{x}\|_1 \\ &= \min_{\mathbf{x} \in \Omega} \sup_{\boldsymbol{\nu}; \|\boldsymbol{\nu}\|_2 \leq 1} \underbrace{\frac{1}{n} \langle \boldsymbol{\nu}, \bar{\mathbf{A}}(\mathbf{x}_0 - \mathbf{x}) \rangle + \frac{\tau}{\sqrt{n}} \langle \boldsymbol{\nu}, \bar{\epsilon} \rangle + \frac{\lambda}{n} \|\mathbf{x}\|_1}_{L(\bar{\mathbf{A}}, \bar{\epsilon}, \mathbf{x}, \boldsymbol{\nu})} \end{aligned}$$

Now let us consider the difference  $F(\bar{\mathbf{A}}_1, \bar{\epsilon}_1) - F(\bar{\mathbf{A}}_2, \bar{\epsilon}_2)$  and denote  $\boldsymbol{\nu}_j^*$ ,  $\mathbf{x}_j^*$  the argmax's for the optimization problem  $F$  solves with respect to  $\bar{\mathbf{A}}_j, \bar{\epsilon}_j$ .

$$\begin{aligned} F(\bar{\mathbf{A}}_1, \bar{\epsilon}_1) - F(\bar{\mathbf{A}}_2, \bar{\epsilon}_2) &= \min_{\mathbf{x} \in \Omega} \sup_{\boldsymbol{\nu}; \|\boldsymbol{\nu}\|_2 \leq 1} L(\bar{\mathbf{A}}_1, \bar{\epsilon}_1, \mathbf{x}, \boldsymbol{\nu}) + \max_{\mathbf{x} \in \Omega} \inf_{\boldsymbol{\nu}; \|\boldsymbol{\nu}\|_2 \leq 1} -L(\bar{\mathbf{A}}_2, \bar{\epsilon}_2, \mathbf{x}, \boldsymbol{\nu}) \\ &\leq \sup_{\boldsymbol{\nu}; \|\boldsymbol{\nu}\|_2 \leq 1} L(\bar{\mathbf{A}}_1, \bar{\epsilon}_1, \mathbf{x}_2^*, \boldsymbol{\nu}) + \inf_{\boldsymbol{\nu}; \|\boldsymbol{\nu}\|_2 \leq 1} -L(\bar{\mathbf{A}}_2, \bar{\epsilon}_2, \mathbf{x}_2^*, \boldsymbol{\nu}) \\ &\leq L(\bar{\mathbf{A}}_1, \bar{\epsilon}_1, \mathbf{x}_2^*, \boldsymbol{\nu}_1^*) - L(\bar{\mathbf{A}}_2, \bar{\epsilon}_2, \mathbf{x}_2^*, \boldsymbol{\nu}_1^*) \\ &= \frac{1}{n} \langle \boldsymbol{\nu}_1^*, (\bar{\mathbf{A}}_2 - \bar{\mathbf{A}}_1)(\mathbf{x}_0 - \mathbf{x}_2^*) \rangle + \frac{\tau}{\sqrt{n}} \langle \boldsymbol{\nu}_1^*, \bar{\epsilon}_1 - \bar{\epsilon}_2 \rangle \\ &\stackrel{\text{C.S.}}{\leq} \frac{1}{n} \|\bar{\mathbf{A}}_1 - \bar{\mathbf{A}}_2\|_{\text{op}} \underbrace{\|\boldsymbol{\nu}_1^*\|_2}_{\leq 1} \underbrace{\|\mathbf{x}_0 - \mathbf{x}_2^*\|_2}_{\leq \sqrt{dM} \propto \sqrt{nM}} + \frac{\tau}{\sqrt{n}} \|\bar{\epsilon}_1 - \bar{\epsilon}_2\| \underbrace{\|\boldsymbol{\nu}_1^*\|_2}_{\leq 1} \\ &\leq C(\|\bar{\mathbf{A}}_1 - \bar{\mathbf{A}}_2\|_{\text{F}} + \|\bar{\epsilon}_1 - \bar{\epsilon}_2\|_2)/\sqrt{n} \\ &\leq C' \sqrt{(\|\bar{\mathbf{A}}_1 - \bar{\mathbf{A}}_2\|_{\text{F}}^2 + \|\bar{\epsilon}_1 - \bar{\epsilon}_2\|_2^2)/n} \end{aligned}$$

Now let  $G(\bar{\mathbf{A}}, \bar{\epsilon}) \in \mathbb{R}^{nd \times n}$  we have that  $F$  is  $\frac{C'}{\sqrt{n}}$ -Lipschitz in  $G$  ( $\|G\|_2^2 = \|\bar{\mathbf{A}}_1 - \bar{\mathbf{A}}_2\|_{\text{F}}^2 + \|\bar{\epsilon}_1 - \bar{\epsilon}_2\|_2^2$ ). We obtain the conclusion by applying the Gaussian inequality.  $\square$

---

<sup>1</sup> $f_\Omega$  is the square root lasso problem, with the minimization taken over a constrained set  $\Omega$ . In future lectures we will show that w.h.p the square root lasso solution lies within  $\Omega$ .