

A. Details of σ_l and other Φ^{adv}

Assumption of σ_l in [25]. To prove Eq. (2) according to Theorem 1.5 in [76], [25] considers $f_{\widetilde{\mathbf{W}}}$ such that $\left| \|\mathbf{W}_l\|_2 - \|\widetilde{\mathbf{W}}_l\|_2 \right| \leq \frac{1}{n} \|\widetilde{\mathbf{W}}_l\|_2$, then they assume $\sigma_l = \frac{\|\widetilde{\mathbf{W}}_l\|_2}{\beta_{\widetilde{\mathbf{W}}_l}} \sigma$, where $\beta_{\widetilde{\mathbf{W}}_l} := \left(\prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2 \right)^{\frac{1}{n}}$.

Assumption of σ_l in [54]. To prove the PAC-Bayesian bound in [54] according to Theorem 1.5 in [76], [54] assumes all variances are same across layers, that is, $\sigma_l = \sigma$.

Our assumption of σ_l . We can prove Lem. 3.2 under both of above assumptions. To make the main paper more clear, we assume that $\sigma_l = \sigma$ in the main paper. And we provide the proofs of Lem. 3.2 for $\sigma_l = \sigma$ and $\sigma_l = \frac{\|\widetilde{\mathbf{W}}_l\|_2}{\beta_{\widetilde{\mathbf{W}}_l}} \sigma$ in Appendix B (the assumption of $\sigma_l = \frac{\|\widetilde{\mathbf{W}}_l\|_2}{\beta_{\widetilde{\mathbf{W}}_l}} \sigma$ includes the assumption of $\sigma_l = \sigma$).

PGM attack for Φ^{adv} . For a PGM attack with noise power ϵ given Euclidean norm $\|\cdot\|$, r iterations for attack and step size \mathcal{Z} , let $\kappa \leq \|\nabla_{\mathbf{s}'} \mathcal{L}(f_{\mathbf{W}}(\mathbf{s}''))\|$ hold for every $\mathbf{s}'' \in \{\mathcal{D} \cup \mathcal{D}'\}$ with constant $\kappa > 0$, then we get [25]

$$\Phi^{\text{adv}} = \left\{ \prod_{l=1}^n \|\mathbf{W}_l\|_2 \left(1 + \frac{\mathcal{Z}}{\kappa} \frac{1 - (2\mathcal{Z}/\kappa)^r \overline{\text{lip}}(\nabla \mathcal{L} \circ f_{\mathbf{W}})^r}{1 - (2\mathcal{Z}/\kappa) \overline{\text{lip}}(\nabla \mathcal{L} \circ f_{\mathbf{W}})} \right)^2 \sum_{l=1}^n \frac{\|\mathbf{W}_l\|_F^2}{\|\mathbf{W}_l\|_2^2}, \right. \\ \left. \left(\prod_{l=1}^n \|\mathbf{W}_l\|_2 \right) \sum_{l=1}^n \prod_{j=1}^l \|\mathbf{W}_j\|_2 \right\}^2 \sum_{l=1}^n \frac{\|\mathbf{W}_l\|_F^2}{\|\mathbf{W}_l\|_2^2}, \quad (22)$$

where

$$\overline{\text{lip}}(\nabla \mathcal{L} \circ f_{\mathbf{W}}) := \left(\prod_{l=1}^n \|\mathbf{W}_l\|_2 \right) \sum_{l=1}^n \prod_{j=1}^l \|\mathbf{W}_j\|_2$$

gives an upper bound on the Lipschitz constant of $\nabla_{\mathbf{s}} \mathcal{L}(f_{\mathbf{W}}(\mathbf{s}))$.

B. Proof of Lem. 3.2

We provide our proofs based on the proofs in [25], to be clearer about the proofs, we suggest readers go through Appendix C.2 in [25] firstly. To prove Eq. (2), [25] considers $f_{\widetilde{\mathbf{W}}}$ such that $\left| \|\mathbf{W}_l\|_2 - \|\widetilde{\mathbf{W}}_l\|_2 \right| \leq \frac{1}{n} \|\widetilde{\mathbf{W}}_l\|_2$, since $(1 + \frac{1}{n})^n \leq e$ and $\frac{1}{e} \leq (1 - \frac{1}{n})^{n-1}$, we get

$$\left(\frac{1}{e} \right)^{\frac{n-1}{n}} \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2 \leq \prod_{l=1}^n \|\mathbf{W}_l\|_2 \leq e \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2, \quad (23)$$

and for each j , we get

$$\frac{1}{\|\mathbf{W}_j\|_2} \prod_{l=1}^n \|\mathbf{W}_l\|_2 \leq \frac{e}{\|\widetilde{\mathbf{W}}_j\|_2} \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2 \quad (24)$$

and

$$\frac{1}{\|\widetilde{\mathbf{W}}_j\|_2} \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2 \leq \left(1 - \frac{1}{n} \right)^{-(n-1)} \frac{1}{\|\mathbf{W}_j\|_2} \prod_{l=1}^n \|\mathbf{W}_l\|_2 \\ \leq \frac{e}{\|\mathbf{W}_j\|_2} \prod_{l=1}^n \|\mathbf{W}_l\|_2. \quad (25)$$

Then let $\sigma_l = \sigma \left(\frac{\|\widetilde{\mathbf{W}}_l\|_2}{\beta_{\widetilde{\mathbf{W}}_l}} = 1 \right)$ or $\sigma_l = \frac{\|\widetilde{\mathbf{W}}_l\|_2}{\beta_{\widetilde{\mathbf{W}}_l}} \sigma$ and let FGM perturbs vector be

$$\delta_{\mathbf{W}}^{\text{fgm}}(\mathbf{s}) := \arg \max_{\|\delta\| \leq \epsilon} \delta^\top \nabla_{\mathbf{s}} \mathcal{L}(f_{\mathbf{W}}(\mathbf{s})). \quad (26)$$

According to Appendix C.2 Eq. (22) in [25], we get the following inequation

$$\|f_{\mathbf{W}+\mathbf{U}}(\mathbf{s} + \delta_{\mathbf{W}+\mathbf{U}}^{\text{fgm}}(\mathbf{s})) - f_{\mathbf{W}}(\mathbf{s} + \delta_{\mathbf{W}}^{\text{fgm}}(\mathbf{s}))\| \\ \leq e(B + \epsilon) \prod_{l=1}^n \|\mathbf{W}_l\|_2 \sum_{l=1}^n \frac{\|\mathbf{U}_l\|_2}{\|\mathbf{W}_l\|_2} \\ + 2e^2 \frac{\epsilon}{\kappa} \prod_{l=1}^n \|\mathbf{W}_l\|_2^2 \sum_{l=1}^n \left[\frac{\|\mathbf{U}_l\|_2}{\|\mathbf{W}_l\|_2} \right. \\ \left. + B \left(\prod_{j=1}^l \|\mathbf{W}_j\|_2 \right) \sum_{j=1}^l \frac{\|\mathbf{U}_j\|_2}{\|\mathbf{W}_j\|_2} \right]. \quad (27)$$

According to Section 1.1 in [5], we have

$$\mathbb{E} \|\mathbf{U}_l\|_2 \lesssim (1 + \sqrt{\ln h}) \|\mathbb{E}(\mathbf{U}_l^\top \mathbf{U}_l)\|_2^{\frac{1}{2}} + \|\mathbb{E}(\mathbf{U}_l \mathbf{U}_l^\top)\|_2^{\frac{1}{2}} \\ \leq c \left((1 + \sqrt{\ln h}) \|\mathbb{E}(\mathbf{U}_l^\top \mathbf{U}_l)\|_2^{\frac{1}{2}} + \|\mathbb{E}(\mathbf{U}_l \mathbf{U}_l^\top)\|_2^{\frac{1}{2}} \right), \\ \mathbb{P} \left(\left| \|\mathbf{U}_l\|_2 - \mathbb{E} \|\mathbf{U}_l\|_2 \right| \geq t \right) \leq 2e^{-t^2/2\sigma_*(\mathbf{U}_l)^2},$$

$$\sigma_*(\mathbf{U}_l) \leq \|\mathbb{E}(\mathbf{U}_l^\top \mathbf{U}_l)\|_2^{\frac{1}{2}},$$

where $c > 0$ is a universal constant. Taking a union bond over the layers, we get that, with probability $> \frac{1}{2}$, the spectral norm of \mathbf{U}_l is bounded by $(\sqrt{2 \ln(4n)} + c + c\sqrt{\ln h}) \|\mathbb{E}(\mathbf{U}_l^\top \mathbf{U}_l)\|_2^{\frac{1}{2}} + c \|\mathbb{E}(\mathbf{U}_l \mathbf{U}_l^\top)\|_2^{\frac{1}{2}}$, let $c_1 = \sqrt{2 \ln(4n)} + c + c\sqrt{\ln h}$ and $c_2 = c$, we have

$$\|\mathbf{U}_l\|_2 \leq \left(c_1 \|\mathbf{R}'_l\|_2^{\frac{1}{2}} + c_2 \|\mathbf{R}''_l\|_2^{\frac{1}{2}} \right) \sigma_l. \quad (28)$$

Thus, $\frac{\beta_{\widetilde{\mathbf{W}}_l}}{\|\widetilde{\mathbf{W}}_l\|_2} \|\mathbf{U}_l\|_2$ is bounded by $\left(c_1 \|\mathbf{R}'_l\|_2^{\frac{1}{2}} + c_2 \|\mathbf{R}''_l\|_2^{\frac{1}{2}} \right) \sigma$. Then, according to Appendix C.2 Eq.

(22) in [25], Eqs. (24) and (27), we can get

$$\begin{aligned}
& \|f_{\mathbf{w}+\mathbf{u}}(\mathbf{s} + \delta_{\mathbf{w}+\mathbf{u}}^{\text{fgm}}(\mathbf{s})) - f_{\mathbf{w}}(\mathbf{s} + \delta_{\mathbf{w}}^{\text{fgm}}(\mathbf{s}))\| \\
& \leq e^2(B + \epsilon) \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2 \sum_{l=1}^n \frac{\|\mathbf{U}_l\|_2}{\|\widetilde{\mathbf{W}}_l\|_2} \\
& \quad + 2e^5 \frac{\epsilon}{\kappa} \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2^2 \sum_{l=1}^n \left[\frac{\|\mathbf{U}_l\|_2}{\|\widetilde{\mathbf{W}}_l\|_2} \right. \\
& \quad \left. + B \left(\prod_{j=1}^l \|\widetilde{\mathbf{W}}_j\|_2 \right) \sum_{j=1}^l \frac{\|\mathbf{U}_j\|_2}{\|\widetilde{\mathbf{W}}_j\|_2} \right] \\
& \leq 2e^5(B + \epsilon) \sigma \left(\sum_{l=1}^n (c_1 \|\mathbf{R}'_l\|_2^{\frac{1}{2}} + c_2 \|\mathbf{R}''_l\|_2^{\frac{1}{2}}) \right) \\
& \quad \left\{ \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2^{\frac{n-1}{n}} + \frac{\epsilon}{\kappa} \left(\prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2^{\frac{2n-1}{n}} \right) \left(\frac{1}{B} + \sum_{l=1}^n \prod_{j=1}^l \|\widetilde{\mathbf{W}}_j\|_2 \right) \right\} \\
& \leq \frac{\gamma}{4},
\end{aligned}$$

hence we choose

$$\begin{aligned}
\sigma & = \frac{\gamma}{8e^5(B + \epsilon) \left(\sum_{l=1}^n (c_1 \|\mathbf{R}'_l\|_2^{\frac{1}{2}} + c_2 \|\mathbf{R}''_l\|_2^{\frac{1}{2}}) \right) \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2^{\frac{n-1}{n}}} \\
& \quad \cdot \frac{1}{\left(1 + \frac{\epsilon}{\kappa} \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2 \left(\frac{1}{B} + \sum_{l=1}^n \prod_{j=1}^l \|\widetilde{\mathbf{W}}_j\|_2 \right) \right)} \quad (29)
\end{aligned}$$

Then we can get

$$\begin{aligned}
\text{KL}(Q_{\text{vec}(\mathbf{W})+\mathbf{u}}\|P) & = \sum_{l=1}^n \left(\frac{\|\mathbf{W}_l\|_F^2}{2\sigma_l^2} - \ln \det \mathbf{R}_l \right) \\
& \leq \mathcal{O} \left((B + \epsilon)^2 \left(\sum_{l=1}^n (c_1 \|\mathbf{R}'_l\|_2^{\frac{1}{2}} + c_2 \|\mathbf{R}''_l\|_2^{\frac{1}{2}}) \right)^2 \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2^2 \right. \\
& \quad \left. \frac{\left(1 + \frac{\epsilon}{\kappa} \prod_{l=1}^n \|\widetilde{\mathbf{W}}_l\|_2 \sum_{l=1}^n \prod_{j=1}^l \|\widetilde{\mathbf{W}}_j\|_2 \right)^2}{\gamma^2} \sum_{l=1}^n \frac{\|\mathbf{W}_l\|_F^2}{\|\widetilde{\mathbf{W}}_l\|_2^2} \right. \\
& \quad \left. - \sum_{l=1}^n \ln \det \mathbf{R}_l \right) \\
& \leq \mathcal{O} \left((B + \epsilon)^2 \left(\sum_{l=1}^n (c_1 \|\mathbf{R}'_l\|_2^{\frac{1}{2}} + c_2 \|\mathbf{R}''_l\|_2^{\frac{1}{2}}) \right)^2 \prod_{l=1}^n \|\mathbf{W}_l\|_2^2 \right. \\
& \quad \left. \frac{\left(1 + \frac{\epsilon}{\kappa} \prod_{l=1}^n \|\mathbf{W}_l\|_2 \sum_{l=1}^n \prod_{j=1}^l \|\mathbf{W}_j\|_2 \right)^2}{\gamma^2} \sum_{l=1}^n \frac{\|\mathbf{W}_l\|_F^2}{\|\mathbf{W}_l\|_2^2} \right. \\
& \quad \left. - \sum_{l=1}^n \ln \det \mathbf{R}_l \right). \quad (30)
\end{aligned}$$

Thus, we have

$$\begin{aligned}
\mathcal{L}_{\mathcal{D}'}(f_{\mathbf{w}}) & \leq \mathcal{L}_{\gamma, S'}(f_{\mathbf{w}}) + \mathcal{O} \left(\left(\frac{-\sum_l \ln \det \mathbf{R}_l + \ln \frac{m}{\delta}}{\gamma^2 m} \right. \right. \\
& \quad \left. \left. + \frac{\Psi^{\text{adv}} \left(\sum_l (c_1 \|\mathbf{R}'_l\|_2^{\frac{1}{2}} + c_2 \|\mathbf{R}''_l\|_2^{\frac{1}{2}}) \right)^2}{\gamma^2 m} \right)^{\frac{1}{2}} \right),
\end{aligned}$$

where $\Psi^{\text{adv}} = (B + \epsilon)^2 \Phi^{\text{adv}}$. And

$$\begin{aligned}
\Phi^{\text{adv}} & = \prod_{l=1}^n \|\mathbf{W}_l\|_2^2 \left\{ 1 + \frac{\epsilon}{\kappa} \left(\prod_{l=1}^n \|\mathbf{W}_l\|_2 \right) \right. \\
& \quad \left. \cdot \sum_{l=1}^n \prod_{j=1}^l \|\mathbf{W}_j\|_2 \right\}^2 \sum_{l=1}^n \frac{\|\mathbf{W}_l\|_F^2}{\|\mathbf{W}_l\|_2^2} \quad (31)
\end{aligned}$$

for FGM attack.

Proofs for PGM attack are similar (combine Eqs. (28) and (30) and Appendix C.3 in [25]).

C. Sampling Method

We use sharpness-like method [34] to get a set of weight samples $(\mathbf{W} + \eta)$ such that $|\mathcal{L}(f_{\mathbf{w}+\eta}) - \mathcal{L}(f_{\mathbf{w}})| \leq \epsilon'$ (e.g., $\epsilon' = 0.05$ for CIFAR-10/SVHN and $\epsilon' = 0.1$ for CIFAR-100), where $\text{vec}(\eta)$ is a $\mathbf{0}$ mean Gaussian noise. To get the samples from the posteriori distribution steadily and fastly, we train the convergent network with learning rate 0.0001, noise η and 50 epochs, then collect corresponding 50 samples. As the samples are stabilized at (clean/adversarial) training loss and validation loss but with different weights, we can treat them as the samples from same (clean/adversarial) posteriori distribution and estimate the correlation matrix through these samples.

D. Proofs of Lems. 4.1, 4.2

As we assume $r_s r_{s'} \geq 0$ (above Lem. 4.1), we give the proofs with two cases ($r_s \geq 0$ and $r_s \leq 0$).

Proof for Lem. 4.1.

Let $r_s \geq 0$ and $r_{s'} \geq 0$, we get

$$\begin{aligned}
\Lambda'_{l, \max} & = \max \left(\|\mathbf{R}'_{l, S}\|_2^{\frac{1}{2}}, \|\mathbf{R}'_{l, S'}\|_2^{\frac{1}{2}} \right) \\
& = \sqrt{h(1 + (h-1) \max(r_s, r_{s'}))} \quad (32)
\end{aligned}$$

and

$$\begin{aligned}
\Lambda''_{l, \max} & = \max \left(\|\mathbf{R}''_{l, S}\|_2^{\frac{1}{2}}, \|\mathbf{R}''_{l, S'}\|_2^{\frac{1}{2}} \right) \\
& = \sqrt{h(1 + (h-1) \max(r_s, r_{s'}))}. \quad (33)
\end{aligned}$$

Thus, decreasing $\|\mathbf{R}_{l, S}\|_F^2$ and $\|\mathbf{R}_{l, S'}\|_F^2$ leads to a decline in $\Lambda'_{l, \max}$ and $\Lambda''_{l, \max}$.

Let $r_s \leq 0$ and $r_{s'} \leq 0$, we get

$$\begin{aligned}
\Lambda'_{l, \max} & = \max \left(\|\mathbf{R}'_{l, S}\|_2^{\frac{1}{2}}, \|\mathbf{R}'_{l, S'}\|_2^{\frac{1}{2}} \right) \\
& = \sqrt{h(1 - \min(r_s, r_{s'}))} \quad (34)
\end{aligned}$$

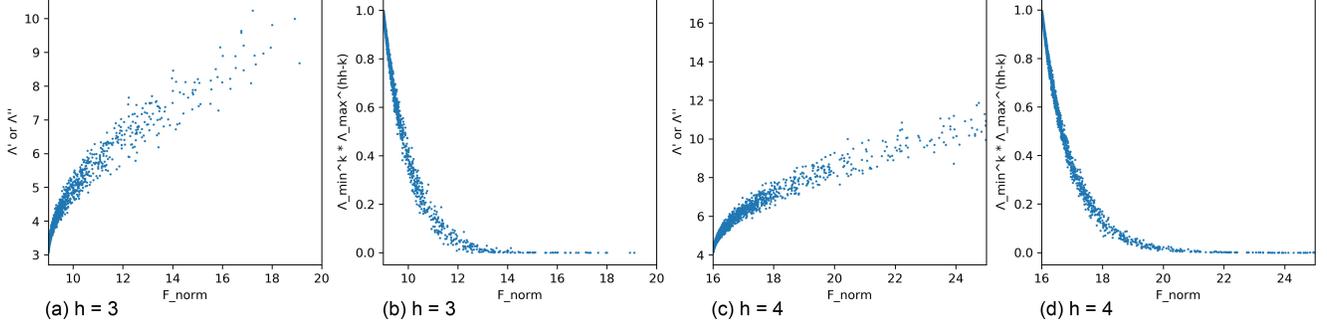


Figure 3. **(a)** We sample 10000 9-dimensional correlation matrices and demonstrate $\|\mathbf{R}_l\|_F^2$ w.r.t $\Lambda'_{l,\max}$ or $\Lambda''_{l,\max}$. **(b)** We sample 10000 9-dimensional correlation matrices and demonstrate $\|\mathbf{R}_l\|_F^2$ w.r.t $\Lambda_{l,\min}^{k_l} \Lambda_{l,\max}^{h^2-k_l}$. **(c)** We sample 10000 16-dimensional correlation matrices and demonstrate $\|\mathbf{R}_l\|_F^2$ w.r.t $\Lambda'_{l,\max}$ or $\Lambda''_{l,\max}$. **(d)** We sample 10000 16-dimensional correlation matrices and demonstrate $\|\mathbf{R}_l\|_F^2$ w.r.t $\Lambda_{l,\min}^{k_l} \Lambda_{l,\max}^{h^2-k_l}$.

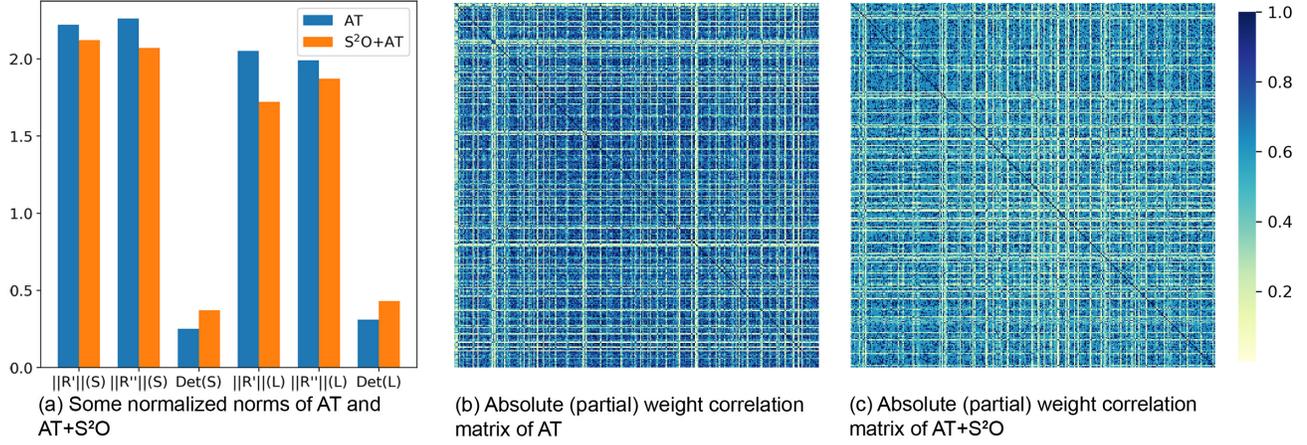


Figure 4. **(a)** shows the normalized spectral norm of \mathbf{R}'_S , \mathbf{R}''_S , and the determinant of \mathbf{R}_S , with sampling estimation (S) and Laplace approximation (L) respectively. **(b)** and **(c)** demonstrate the absolute correlation matrix of partial weights (estimate under clean data), for AT and AT+S²O respectively.

and

$$\begin{aligned} \Lambda''_{l,\max} &= \max(\|\mathbf{R}''_{l,S}\|_2^{\frac{1}{2}}, \|\mathbf{R}'_{l,S'}\|_2^{\frac{1}{2}}) \\ &= \sqrt{h(1 - \min(r_s, r_{s'}))}. \end{aligned} \quad (35)$$

Thus, decreasing $\|\mathbf{R}_{l,S}\|_F^2$ and $\|\mathbf{R}_{l,S'}\|_F^2$ leads to a decline in $\Lambda'_{l,\max}$ and $\Lambda''_{l,\max}$.

Proof for Lem. 4.2.

Let $r_s \geq r_{s'} \geq 0$, we get

$$\begin{aligned} c(r) &= \Lambda_{l,\min}^{k_l} \Lambda_{l,\max}^{h^2-k_l} \\ &= (1 - r_s)^{h^2-1} (1 + (h^2 - 1)r_s) \end{aligned} \quad (36)$$

and

$$\frac{\partial c(r)}{\partial r_s} = -h^2(h^2 - 1)r_s(1 - r_s)^{h^2-2} \leq 0, \quad (37)$$

it is easy to get $c(r)$ is negative correlated with r_s . Similarly, if $r_{s'} \geq r_s \geq 0$, we can get $c(r)$ is negative correlated with $r_{s'}$. Thus, decreasing $\|\mathbf{R}_{l,S}\|_F^2$ and $\|\mathbf{R}_{l,S'}\|_F^2$ leads to an increase in $\Lambda_{l,\min}^{k_l} \Lambda_{l,\max}^{h^2-k_l}$.

Let $r_s \leq r_{s'} \leq 0$, we get

$$\begin{aligned} c(r) &= \Lambda_{l,\min}^{k_l} \Lambda_{l,\max}^{h^2-k_l} \\ &= (1 + (h^2 - 1)r_s)(1 - r_{s'})^{h^2-1} \end{aligned} \quad (38)$$

and

$$\frac{\partial c(r)}{\partial r_s} = -h^2(h^2 - 1)r_s(1 - r_{s'})^{h^2-2} \geq 0, \quad (39)$$

it is also easy to get $c(r)$ is positive correlated with r_s . Similarly, if $r_{s'} \leq r_s \leq 0$, we can get $c(r)$ is positive correlated with $r_{s'}$. Thus, decreasing $\|\mathbf{R}_{l,S}\|_F^2$ and $\|\mathbf{R}_{l,S'}\|_F^2$ leads to an increase in $\Lambda_{l,\min}^{k_l} \Lambda_{l,\max}^{h^2-k_l}$.

E. Simulations of Lems. 4.1, 4.2 and Second-Order Statistics of Weights under Clean Data

As Fig. 3 shows, for 10000 random general 9-dimensional correlation matrices and 16-dimensional correlation matrices respectively, Lems. 4.1 and 4.2 also hold approximately.

The results in Fig. 4 also suggest that S²O can decrease the spectral norm of $\mathbf{R}'_{\mathcal{S}}$, $\mathbf{R}''_{\mathcal{S}}$ and increases the determinant of $\mathbf{R}_{\mathcal{S}}$.

F. Approximate Optimization

We use a fast approximate method to update $\mathbf{g}(\mathbf{A})$, i.e., add a penalty term to the high correlated $\mathbf{a}_{l,i}$ and $\mathbf{a}_{l,j}$ to reduce their correlation. Details are given in the code.