

SUPPLEMENTARY MATERIAL:

ASYMPTOTIC NORMALITY AND CONFIDENCE INTERVALS FOR DERIVATIVES OF 2-LAYER NEURAL NETWORK IN THE RANDOM FEATURES MODEL

We gather here the supporting propositions, the proofs and additional figures. The python code used to generate the figures is also attached. In Figure 4 on page 2 we plot $R(\psi_p, \lambda, \rho, \sigma)$ under various regimes as a supplement to Figure 3.

CONTENTS

Supplementary Material:

Asymptotic normality and confidence intervals for derivatives of 2-layer neural network in the random features model	1
Appendix A. Outline of the proof of the main result	1
Appendix B. Figure 4	2
Appendix C. Supporting propositions	3
C.1. Notation and constants	3
C.2. Asymptotic normality result	4
C.3. The limits $V(\psi_d, \psi_p, \lambda, \rho, \sigma)$ and $R(\psi_p, \lambda, \rho, \sigma)$	5
C.4. Asymptotic normality result for a general direction	7
Appendix D. Proofs	8
D.1. Proofs of Theorem 1 and 2	8
D.2. Proofs of the limits $V(\psi_d, \psi_p, \lambda, \rho, \sigma)$ and $R(\psi_p, \lambda, \rho, \sigma)$	18
D.3. Proof of Theorem C.4.1	21
References	27

APPENDIX A. OUTLINE OF THE PROOF OF THE MAIN RESULT

Since Theorem 1 can be regarded as a special case of Theorem 2, it suffices to prove Theorem 2. The formal proof of Theorem 2 is provided in Section D.1. The proof combines three supporting results using Slutsky's Theorem: the asymptotic results

$$(A.1) \quad \frac{\zeta(\mathbf{e}_j)}{(\text{Var}_j[\zeta(\mathbf{e}_j)])^{1/2}} \xrightarrow{d} N(0, 1), \quad \frac{\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2}{(\text{Var}_j[\zeta(\mathbf{e}_j)])^{1/2}} \xrightarrow{\mathbb{P}} 1 \quad \text{and} \quad \frac{\zeta(\mathbf{e}_j) - \zeta_L(\mathbf{e}_j)}{\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2} \xrightarrow{\mathbb{P}} 0.$$

The first two are proved in Proposition C.2.5 while the third is proved in Proposition C.2.6. The asymptotic results are shown to hold for most $j \in [p]$, in the sense that the current proof shows that (A.1) holds for all $j \in [p] \setminus J$ for some set J with negligible cardinality compared to p . More specifically, Proposition C.2.5 leverages the flexible central limit theorems in Bellec and Zhang [2019] to obtain asymptotic normality of $\frac{\zeta(\mathbf{e}_j)}{(\text{Var}_j[\zeta(\mathbf{e}_j)])^{1/2}}$. Proposition C.2.6 shows that the nonlinear component (induced by the nonlinear perturbation) in $\zeta(\mathbf{e}_j) - \zeta_L(\mathbf{e}_j)$ is negligible.

APPENDIX B. FIGURE 4

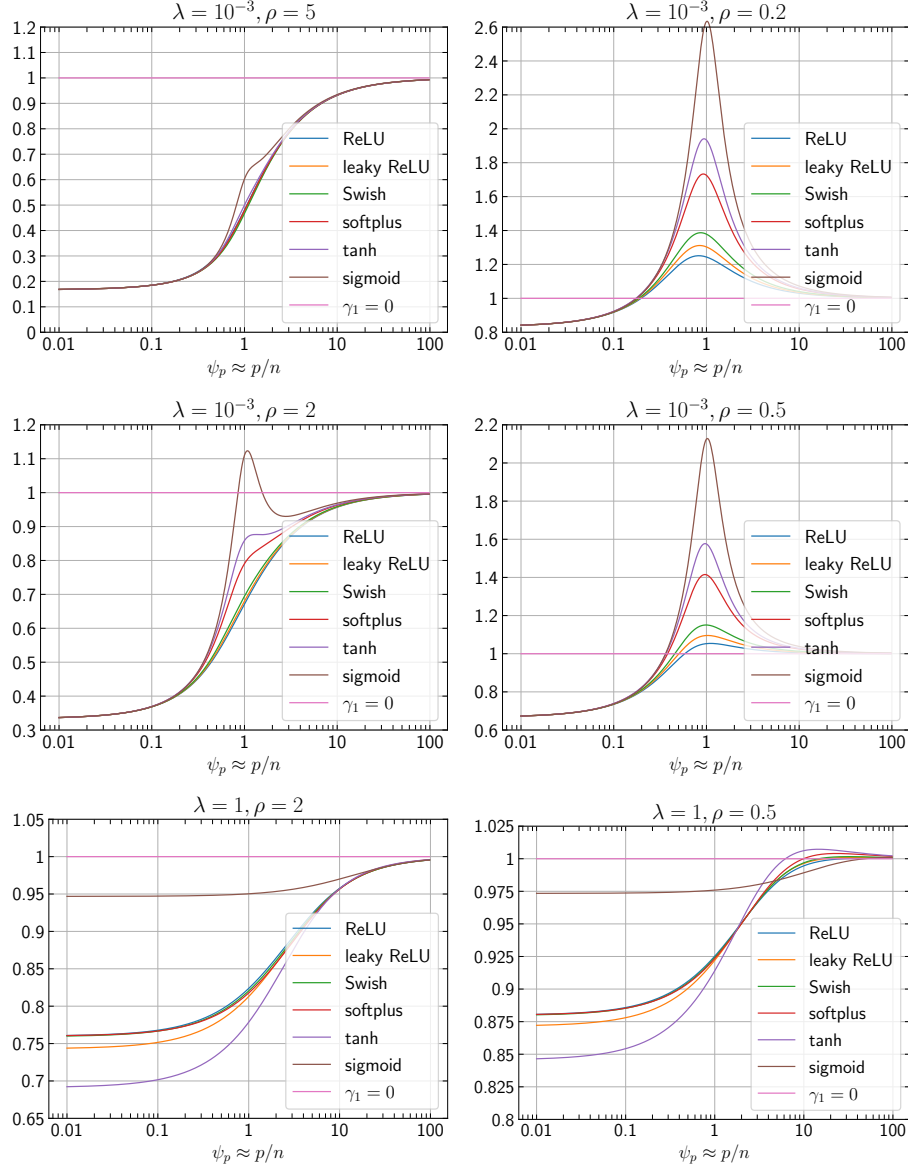


FIGURE 4. The indicator $R(\psi_p, \lambda, \rho, \sigma)$ (see Definition C.3.2) of the squared length in the infinite width. The activation functions are ReLU $\max(x, 0)$, leaky ReLU $\max(x, 0.1x)$, Swish $\frac{x}{1+e^{-x}}$, softplus $\ln(1 + \exp(x))$, tanh $\frac{e^x - e^{-x}}{e^x + e^{-x}}$ and sigmoid $\frac{1}{1+e^{-x}}$. Activation functions $\sigma(x)$ with $\gamma_1 := \mathbb{E}[\sigma'(Z)] = 0$ for $Z \sim N(0, 1)$ have limit $R = 1$ always. We consider Ridge penalty parameter $\lambda \in \{10^{-3}, 1\}$ and signal-to-noise ratio $\rho \in \{5, 2, 0.5, 0.2\}$.

APPENDIX C. SUPPORTING PROPOSITIONS

C.1. Notation and constants. The proof will require some further notation and constants.

C.1.1. Notation. Let $\|\mathbf{M}\|_F$ denote the Frobenius norm of a matrix \mathbf{M} and $\|\mathbf{M}\|_{op}$ its operator norm.

For the functions $f : \mathbb{R}^p \rightarrow \mathbb{R}$ and $G : \mathbb{R}^p \rightarrow \mathbb{R}$, denote by ∂_j the partial derivative with respect to the j -th coordinate for each $j \in [p]$, as well as

$$\begin{aligned} \mathbf{f} &:= (f(\mathbf{x}_i))^{i \in [n]} \in \mathbb{R}^n, & \mathbf{f}' &:= [\partial_j f(\mathbf{x}_i)]^{(i,j) \in [n] \times [p]} \in \mathbb{R}^{n \times p}, \\ \mathbf{G} &:= (G(\mathbf{x}_i))^{i \in [n]} \in \mathbb{R}^n, & \mathbf{G}' &:= [\partial_j G(\mathbf{x}_i)]^{(i,j) \in [n] \times [p]} \in \mathbb{R}^{n \times p} \end{aligned}$$

where $\mathbf{x}_1 \in \mathbb{R}^p, \dots, \mathbf{x}_n \in \mathbb{R}^p$ are the observed feature vectors, i.e., the rows of \mathbf{X} .

Let \mathbb{E}_j and Var_j be the conditional expectation and the conditional variance given $(\mathbf{X}_{-j}, G, \mathbf{W}, \boldsymbol{\varepsilon})$, where \mathbf{X}_{-j} is the matrix \mathbf{X} with j -th column removed. Let $\mathbb{E}_{\boldsymbol{\varepsilon}}$ be the conditional expectation given $\mathbf{X}, G, \mathbf{W}$. Let \mathbb{P} and \mathbb{E} be with respect to the total probability $\mathbb{P}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G}$. We let a_+ denote $\max(a, 0)$.

We will also consider gradients with respect to columns of \mathbf{X} or with respect to the noise $\boldsymbol{\varepsilon}$. Consider an expression $\mathbf{r} \in \mathbb{R}^q$ which is a function of $(\mathbf{X}, \boldsymbol{\varepsilon}, \mathbf{W}, G)$.

- For two given indices $(i, j) \in [n] \times [p]$, the column vector $\partial_{x_{ij}} \mathbf{r}$ has the same dimension as \mathbf{r} and denotes the partial derivative of \mathbf{r} with respect to x_{ij} while

$$((x_{i', j'})_{i' \in [n], j' \in [p]: (i', j') \neq (i, j)}, \boldsymbol{\varepsilon}, \mathbf{W}, G)$$

remain fixed. If the dimension q of \mathbf{r} is greater than 1, then $\nabla_{x_{ij}}$ acts componentwise on the components of \mathbf{r} .

- The matrix $\nabla_{\mathbf{X}_j} \mathbf{r}$ is the derivative (in the sense of the Frechet derivative) of \mathbf{r} with respect to the j -th column \mathbf{X}_j of \mathbf{X} while $(\mathbf{X}_{-j}, \boldsymbol{\varepsilon}, \mathbf{W}, G)$ (which are random variables independent of \mathbf{X}_j) remain fixed; if the dimension q of \mathbf{r} is greater than 1, then $\nabla_{\mathbf{X}_j}$ acts componentwise on the components of \mathbf{r} . For instance, the derivative of the residuals $\mathbf{r} = \mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}$ with respect to \mathbf{X}_j while $(\mathbf{X}_{-j}, \boldsymbol{\varepsilon}, \mathbf{W}, G)$ remain fixed is the $n \times n$ matrix $\nabla_{\mathbf{X}_j}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\alpha}})$ whose i -th column is $\partial_{x_{ij}}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\alpha}})$, the directional derivative with respect to the (i, j) -th entry of \mathbf{X} . We may refer to $\nabla_{\mathbf{X}_j}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\alpha}})$ as the Jacobian of the map $\mathbf{X}_j \mapsto \mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}$.
- Similarly, $\nabla_{\boldsymbol{\varepsilon}} \mathbf{r}$ is the derivative of \mathbf{r} with respect to the noise vector $\boldsymbol{\varepsilon}$ while $(\mathbf{X}, \mathbf{W}, G)$ (which are random variables independent of $\boldsymbol{\varepsilon}$) remain fixed. If the dimension q of \mathbf{r} is greater than 1, then $\nabla_{\boldsymbol{\varepsilon}}$ acts componentwise on the components of \mathbf{r} . For instance, the derivative of the residuals $\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}$ with respect to $\boldsymbol{\varepsilon}$ while $(\mathbf{X}, \mathbf{W}, G)$ remain fixed is the $n \times n$ matrix $\nabla_{\boldsymbol{\varepsilon}}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\alpha}})$ whose i -th column is the i -th entry of $\boldsymbol{\varepsilon}$. We may refer to $\nabla_{\boldsymbol{\varepsilon}}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}})$ as the Jacobian of the map $\boldsymbol{\varepsilon} \mapsto \mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}$.
- (This is only used in Section C.4) If $\mathbf{u}_0 \in \mathbb{R}^p$ has $\|\mathbf{u}_0\|_2 = 1$, we define $\mathbf{X}_0 = \mathbf{X}\mathbf{u}_0$. Then $\nabla_{\mathbf{X}_0} \mathbf{r}$ is the derivative with respect to \mathbf{X}_0 while $(\mathbf{X}(\mathbf{I}_p - \mathbf{u}_0\mathbf{u}_0^\top), \boldsymbol{\varepsilon}, \mathbf{W}, G)$, (which are random variables independent of \mathbf{X}_0) remain fixed. If $\mathbf{r} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\alpha}}$, the Jacobian $\nabla_{\mathbf{X}_0}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\alpha}})$ is the matrix whose i -th column is the derivative with respect to the i -th entry of \mathbf{X}_0 .

C.1.2. Constants. The following positive finite constants L, c_1, c_8, δ, L_2 are independent of n, p, d .

- (i) The activation function σ is Lipschitz with constant L .
- (ii) \mathbf{W} is deterministic with $\|\mathbf{W}\|_{op} < c_1$.
- (iii) For some $\delta > 0$, $\Sigma_p''(x)$ exists in $(1 - \delta, 1 + \delta)$ and $(-\delta, \delta)$ and is L_2 -Lipschitz in $(1 - \delta, 1 + \delta)$ and $(-\delta, \delta)$.
- (iv) $\max(|\Sigma_p'(1)|, |\Sigma_p''(0)|, |\Sigma_p''(1)|) < c_8$.

For the proof of the asymptotic normality result in Theorem 2, we only consider the case when \mathbf{W} is deterministic. The proof is applicable for \mathbf{W} with entries iid $N(0, 1/p)$ by Corollary 7.3.3 in Vershynin [2018] by conditioning on \mathbf{W} .

C.2. Asymptotic normality result. In this section, we present supporting propositions for the asymptotic normality result in Theorem 2. Proposition C.2.1 provides the calculation of the Jacobian matrix $\nabla_{\mathbf{X}_j}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}})$.

Proposition C.2.1 (Calculation of the Jacobian matrix). *Under model (5),*

$$\begin{aligned} \nabla_{\mathbf{X}_j}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}) = & -(\mathbf{I}_n - \mathbf{H}) \text{diag}(\sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\mathbf{W}\mathbf{e}_j)\hat{\boldsymbol{\alpha}}) \\ & - \mathbf{A}(\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} \text{diag}(\mathbf{W}\mathbf{e}_j)\sigma'(\mathbf{W}\mathbf{X}^\top) \text{diag}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}) \\ & + (\mathbf{I}_n - \mathbf{H}) [(\mathbf{e}_j^\top \boldsymbol{\beta}) \mathbf{I}_n + \text{diag}(\mathbf{G}'\mathbf{e}_j)]. \end{aligned} \quad (\text{C.1})$$

Under Definitions 1 and C.2.1, we have

$$\nabla_{\mathbf{X}_j}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}) = \mathbf{T}_0(\mathbf{e}_j) + \mathbf{T}_1(\mathbf{e}_j) + \mathbf{T}_L(\mathbf{e}_j) + \mathbf{T}_{\text{NL}}(\mathbf{e}_j). \quad (\text{C.2})$$

Definition C.2.1. Let $\zeta_L(\mathbf{e}_j)$ be in Definition 1. Let

$$\zeta(\mathbf{e}_j) = \zeta_L(\mathbf{e}_j) - \text{trace}[\mathbf{T}_{\text{NL}}(\mathbf{e}_j)], \quad (\text{C.3})$$

where

$$\begin{aligned} \mathbf{T}_{\text{NL}}(\mathbf{e}_j) &= (\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}'\mathbf{e}_j), \\ \mathbf{G}' &= [\partial_j G(\mathbf{x}_i)]^{(i,j) \in [n] \times [p]} \in \mathbb{R}^{n \times p}. \end{aligned} \quad (\text{C.4})$$

Proposition C.2.2 provides upper bounds on the components of the Jacobian matrix $\nabla_{\mathbf{X}_j}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}})$.

Proposition C.2.2. Let $\mathbf{T}_0, \mathbf{T}_1, \mathbf{T}_L, \mathbf{T}_{\text{NL}}$ be as in Definitions 1 and C.2.1. Let c_1, L, \mathbf{f}' be as in §C.1.

- (i) $\|\mathbf{A}(n\tau \mathbf{I}_d + \mathbf{A}^\top \mathbf{A})^{-1}\|_{op} \leq 1/(2\sqrt{n\tau})$.
- (ii) $\|\mathbf{I}_n - \mathbf{H}\|_{op} \leq 1$.
- (iii) $\sum_{j \in [p]} \|\mathbf{T}_0(\mathbf{e}_j) + \mathbf{T}_L(\mathbf{e}_j) + \mathbf{T}_{\text{NL}}(\mathbf{e}_j)\|_F^2 \leq 2c_1^2 L^2 n \|\hat{\boldsymbol{\alpha}}\|_2^2 + 2\|\mathbf{f}'\|_F^2$.
- (iv) $\|\mathbf{T}_1(\mathbf{e}_j)\|_F^2 \leq L^2 c_1^2 / (4n\tau) \cdot \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2$.
- (v) $\mathbb{E}\|\hat{\boldsymbol{\alpha}}\|_2^2 = O(1)$.

Proposition C.2.3 shows that the training error $\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2$ is of order at least n with large probability.

Proposition C.2.3. There exists a large event Ω such that $\mathbb{P}(\Omega^c) \leq o(\exp(-c_6 n))$ for some $c_6 > 0$ and that on Ω ,

- (i) $\|\mathbf{X}\|_F^2/n^2 \leq c_{4,n}$,
- (ii) $1/n \cdot \min_{j \in [p]} \mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq c_{2,n}$,
- (iii) $1/n \cdot \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq c_{2,n}$.

The constant c_6 is independent of n, d, p . The constants $c_{4,n}$ and $c_{2,n}$ are given in Proposition C.2.4.

Proposition C.2.4. Let

- (1) $c_{3,n} := \psi_{p,n}^{-1/2} + 1 + c^{-1/2} \rightarrow c_3 := \psi_p^{-1/2} + 1 + c^{-1/2}$ where $c > 0$ is some universal constant specified in the proof.
- (2) $c_{4,n} := \psi_{p,n} + 2\psi_{p,n}^{1/2} + 2 \rightarrow c_4 := \psi_p + 2\psi_p^{1/2} + 2$.
- (3) $F_n := 2c_1^2 L^2 \|\mathbf{X}\|_F^2/n^2 + 2\psi_{d,n}(\sigma(0))^2$.
- (4) $\bar{F}_n := 2c_1^2 L^2 (n + \|\mathbf{X}\|_F^2)/n^2 + 2\psi_{d,n}(\sigma(0))^2$.

- (5) $\bar{c}_{2,n} := (1 + 2c_1^2 L^2 n^{-1} \tau^{-1} + 2c_1^2 L^2 c_{4,n} \tau^{-1} + 2\psi_d(\sigma(0))^2 \tau^{-1})^{-1} \rightarrow \bar{c}_2 := (1 + 2c_1^2 L^2 c_4 \tau^{-1} + 2\psi_d(\sigma(0))^2 \tau^{-1})^{-1} > 0$.
 (6) $c_{2,n} := \theta_\varepsilon^2 (2\bar{c}_{2,n}/3 - n^{-1/2})_+^2 \rightarrow c_2 := \theta_\varepsilon^2 (2\bar{c}_2/3)^2 > 0$.

Then

- (i) $\mathbb{P}(\|\mathbf{X}/\sqrt{p}\|_{op} \geq c_{3,n}) \leq 2\exp(-p)$.
 (ii) $\mathbb{P}(\|\mathbf{X}\|_F^2/n^2 > c_{4,n}) \leq \exp(-n^2)$.
 (iii) $\|\sigma(\mathbf{X}\mathbf{W}^\top)\|_F^2/n^2 \leq F_n$.
 (iv) $\|\mathbf{I}_n - \mathbf{H}\|_F^2/n \geq (1 + F_n/\tau)^{-2}$.
 (v) $\min_{j \in [p]} (n^{-1/2} \mathbb{E}_j \|\mathbf{I}_n - \mathbf{H}\|_F) \geq (1 + \bar{F}_n/\tau)^{-1}$.

Proposition C.2.5 helps us prove the asymptotic normality of $\zeta(\mathbf{e}_j)$ and helps us estimate the variance term in the asymptotic normality result with the training error: Propositions C.2.5(i) and C.2.5(iii) imply the asymptotic normality of $\zeta(\mathbf{e}_j)$ for most $j \in [p]$; Propositions C.2.5(ii) and C.2.5(iii) imply that we can estimate $\text{Var}_j(\zeta(\mathbf{e}_j))$ using $\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2$ for most $j \in [p]$.

Proposition C.2.5. *Let*

- (i) $\zeta(\mathbf{e}_j)$ be as in Definition C.2.1.
 (ii) $\mathbf{r} = \mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}$.
 (iii) $\epsilon_j^2 := \mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2] / (\mathbb{E}_j [\|\mathbf{r}\|_2^2] + \mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2])$.

Then

- (i) $\mathbb{E}_j \left[\left(\frac{\zeta(\mathbf{e}_j)}{(\text{Var}_j \zeta(\mathbf{e}_j))^{1/2}} - \frac{\mathbf{X}_j^\top \mathbb{E}_j [\mathbf{r}]}{\|\mathbb{E}_j [\mathbf{r}]\|_2} \right)^2 \right] \leq 6\epsilon_j^2$.
 (ii) $\mathbb{E}_j \left[\left| \frac{\|\mathbf{r}\|_2}{(\text{Var}_j [\zeta(\mathbf{e}_j)])^{1/2}} - 1 \right| \right] \leq (1 + \sqrt{2}) \epsilon_j / (1 - 2\epsilon_j^2)_+^{1/2}$.
 (iii) $\mathbb{E} \left[\sum_{j \in [p]} \epsilon_j^2 \right] = O(1)$.

Proposition C.2.6 shows that the nonlinear component $\text{trace}(\mathbf{T}_{\text{NL}}(\mathbf{e}_j))/\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2$ of the asymptotic normality quantity $\zeta(\mathbf{e}_j)/\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2$ is negligible. Since $\zeta_L(\mathbf{e}_j) = \zeta(\mathbf{e}_j) + \text{trace}[\mathbf{T}_{\text{NL}}(\mathbf{e}_j)]$, we will then have the asymptotic normality of $\zeta_L(\mathbf{e}_j)/\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2$.

Proposition C.2.6. *Under Assumptions 1, 2, 3 and 4, we have the following convergence in probability: for all $\epsilon > 0$,*

$$(C.5) \quad \lim_{n \rightarrow +\infty} \mathbb{P} \left(\max_{j \in [p]} \left| \frac{\text{trace}((\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}' \mathbf{e}_j))}{\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2} \right| > \epsilon \right) = 0.$$

Proposition C.2.7 characterizes the expected value of the derivatives of the nonlinear perturbation function G in terms of Σ'_p and Σ''_p .

Proposition C.2.7. *For all vector $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^p$ such that the derivatives below exist,*

$$\mathbb{E}_G [\partial_{j_1} G(\mathbf{v}_1) \partial_{j_2} G(\mathbf{v}_2)] = \Sigma''_p(\mathbf{v}_1^\top \mathbf{v}_2/p) v_{1j_2} v_{2j_1}/p^2 + \Sigma'_p(\mathbf{v}_1^\top \mathbf{v}_2/p) \delta(j_1 = j_2)/p.$$

C.3. The limits $V(\psi_d, \psi_p, \lambda, \rho, \sigma)$ and $R(\psi_p, \lambda, \rho, \sigma)$. In this section we present the results related to the limits $V(\psi_d, \psi_p, \lambda, \rho, \sigma)$ and $R(\psi_p, \lambda, \rho, \sigma)$ in (18).

Assumption C.3.1 recalls the setting in Mei and Montanari [2019] that is comparable to our Gaussian setting.

Assumption C.3.1. *Let us assume that*

- (i) $(\mathbf{x}_i)_{i \in [n]} \sim^{iid} \text{Unif}(\mathbb{S}^{p-1}(\sqrt{p}))$.
 (ii) $(\mathbf{w}_k)_{k \in [d]} \sim^{iid} \text{Unif}(\mathbb{S}^{p-1}(1))$.
 (iii) $G(\mathbf{x})$ is a centered Gaussian process indexed by $\mathbf{x} \in \mathbb{S}^{p-1}(\sqrt{p})$, such that
 (i) $\mathbb{E}_G[G(\mathbf{x})] = 0$ and $\mathbb{E}_G[G(\mathbf{x}_1)G(\mathbf{x}_2)] = \Sigma_p(\mathbf{x}_1^\top \mathbf{x}_2/p)$,

- (ii) $\mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathbb{S}^{p-1}(\sqrt{p}))} [\Sigma_p(x_1/\sqrt{p})] = 0$, $\mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathbb{S}^{p-1}(\sqrt{p}))} [\Sigma_p(x_1/\sqrt{p})x_1] = 0$ and $\lim_{p \rightarrow +\infty} \Sigma_p(1) = \theta_{\text{NL}}^2$.
- (iv) $\mathbb{E}[\varepsilon_1] = 0$, $\mathbb{E}[\varepsilon_1^2] = \theta_\varepsilon^2$, $\mathbb{E}[\varepsilon_1^4] < +\infty$.
- (v) $\beta_0 \rightarrow \theta_0$, $\|\beta\|_2^2 \rightarrow \theta_\beta^2$ and $\Sigma_p(1) \rightarrow \theta_{\text{NL}}^2$. The signal-to-noise ratio $\rho = \theta_\beta^2/(\theta_\varepsilon^2 + \theta_{\text{NL}}^2)$.

Proposition C.3.1 provides the limiting squared length of our confidence intervals under Assumption C.3.1.

Proposition C.3.1. Let $L^2 := \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2/(\text{trace}(\mathbf{I}_n - \mathbf{H}))^2$. Under the model (5) in the asymptotic setting (6), Assumption C.3.1 and Definition C.3.1,

$$\text{Var}(y_1) \rightarrow \theta_\beta^2 + \theta_\varepsilon^2 + \theta_{\text{NL}}^2, \quad nL^2/\text{Var}(y_1) \xrightarrow{\mathbb{P}_{\mathbf{x}, \mathbf{w}}^{G, \varepsilon}} V.$$

Definition C.3.1 defines $V(\psi_d, \psi_p, \lambda, \rho, \sigma)$.

Definition C.3.1 (Definition 2 in Mei and Montanari [2019]). Let $Z \sim N(0, 1)$, $\rho = \theta_\beta^2/(\theta_\varepsilon^2 + \theta_{\text{NL}}^2)$,

$$\begin{aligned} \mu_1 &= \mathbb{E}[\sigma(Z)], & \mu_2 &= \mathbb{E}[(\sigma(Z))^2], & \gamma_1 &= \mathbb{E}[\sigma'(Z)], \\ \mu_*^2 &= \mu_2 - \mu_1^2 - \gamma_1^2, & \varrho^2 &= \gamma_1^2/\mu_*^2, & \bar{\lambda} &= \lambda/\mu_*^2, \\ \psi_1 &= \psi_d/\psi_p, & \psi_2 &= 1/\psi_p, & \bar{z} &= \mathbf{i}(\psi_1\psi_2\lambda)^{1/2}/\mu_*. \end{aligned}$$

Let $\Im(z)$ be the imaginary part of complex number z . Let \mathbb{C}_+ be the set of complex numbers with positive imaginary parts. Let $\nu_1, \nu_2 \in \mathbb{C}_+$ solves uniquely

$$\begin{aligned} \nu_1 &= \psi_1 \left(-\bar{z} - \nu_2 - \frac{\varrho^2 \nu_2}{1 - \varrho^2 \nu_1 \nu_2} \right)^{-1}, \\ \nu_2 &= \psi_2 \left(-\bar{z} - \nu_1 - \frac{\varrho^2 \nu_1}{1 - \varrho^2 \nu_1 \nu_2} \right)^{-1}, \end{aligned} \quad (\text{C.6})$$

under constraint $|\nu_1| \leq \psi_1/\Im(\bar{z})$ and $|\nu_2| \leq \psi_2/\Im(\bar{z})$. Let

$$\begin{aligned} \chi &= \nu_1 \nu_2, \\ \mathcal{Q} &= 1 + \psi_2^{-1} \left(\chi + \frac{\chi \varrho^2}{1 - \chi \varrho^2} \right), \\ \mathcal{L} &= \mathcal{Q} \cdot \left[\frac{\rho}{1 + \rho} \cdot \frac{1}{1 - \chi \varrho^2} + \frac{1}{1 + \rho} \right], \\ \mathcal{A}_1 &= \frac{\rho}{1 + \rho} [-\chi^2 (\chi \varrho^4 - \chi \varrho^2 + \psi_2 \varrho^2 + \varrho^2 - \chi \psi_2 \varrho^4 + 1)] \\ &\quad + \frac{1}{1 + \rho} [\chi^2 (\chi \varrho^2 - 1) (\chi^2 \varrho^4 - 2\chi \varrho^2 + \varrho^2 + 1)], \\ \mathcal{A}_0 &= -\chi^5 \varrho^6 + 3\chi^4 \varrho^4 + (\psi_1 \psi_2 - \psi_2 - \psi_1 + 1) \chi^3 \varrho^6 - 2\chi^3 \varrho^4 - 3\chi^3 \varrho^2, \\ &\quad + (\psi_1 + \psi_2 - 3\psi_1 \psi_2 + 1) \chi^2 \varrho^4 + 2\chi^2 \varrho^2 + \chi^2 + 3\psi_1 \psi_2 \chi \varrho^2 - \psi_1 \psi_2, \\ \mathcal{A} &= \mathcal{A}_1/\mathcal{A}_0, \end{aligned}$$

$$V(\psi_d, \psi_p, \lambda, \rho, \sigma) = (\mathcal{L} - \psi_1 \bar{\lambda} \mathcal{A})/\mathcal{Q}^2.$$

Random vectors uniformly distributed on the sphere $\mathbb{S}^{p-1}(\sqrt{p})$ are close to standard normal vectors (e.g., in the sense of § 3.3.3 in Vershynin [2018]). We expect the limit in Proposition C.3.1 to be valid for Gaussian $\mathbf{x}_i, \mathbf{w}_k$ as well. Proposition C.3.2 guarantees that the limit V is valid for Gaussian vectors $(\mathbf{x}_i)_{i \in [n]} \sim^{iid} N(\mathbf{0}_p, \mathbf{I}_p)$, $(\mathbf{w}_k)_{k \in [d]} \sim^{iid} N(\mathbf{0}_p, (1/p)\mathbf{I}_p)$ for the pure linear model (4).

Proposition C.3.2. Under the pure linear model (4) in the asymptotic setting (6), Assumption 1 and 2 with (ii) and Definition C.3.1,

$$\text{Var}(y_1) \rightarrow \theta_\beta^2 + \theta_\epsilon^2, \quad nL^2/\text{Var}(y_1) \xrightarrow{\mathbb{P}_{\mathbf{X}, \mathbf{W}, \epsilon}} V.$$

Proposition C.3.3 provides the $\psi_d \rightarrow +\infty$ limit of V .

Proposition C.3.3. Let V, R be in Definitions C.3.1 and C.3.2. Then

$$\lim_{\psi_d \rightarrow +\infty} V(\psi_d, \psi_p, \lambda, \rho, \sigma) = R(\psi_p, \lambda, \rho, \sigma).$$

Definition C.3.2. Let $\varrho, \rho, \psi_1, \psi_2, \bar{\lambda}$ be as in Definition C.3.1. Let

$$R(\psi_p, \lambda, \rho, \sigma) = \begin{cases} 1 & \varrho = 0 \\ (\bar{\mathcal{L}} - \bar{\lambda} \bar{\mathcal{A}}_1 / \bar{\mathcal{A}}_*) / \bar{\mathcal{Q}}^2 & \varrho \neq 0 \end{cases},$$

where

$$\begin{aligned} \bar{\lambda} &= \begin{cases} -\frac{\psi_2}{1+\psi_2\bar{\lambda}} & \varrho = 0 \\ \frac{\left(\varrho^{-2} - \frac{\psi_2-1}{1+\psi_2\bar{\lambda}}\right) - \sqrt{\left(\varrho^{-2} - \frac{\psi_2-1}{1+\psi_2\bar{\lambda}}\right)^2 + 4\frac{\psi_2}{1+\psi_2\bar{\lambda}}\varrho^{-2}}}{2} & \varrho \neq 0 \end{cases}, \\ \bar{\mathcal{Q}} &= -\bar{\lambda}\bar{\lambda}, \\ \bar{\mathcal{L}} &= (-\bar{\lambda}\bar{\lambda}) \left[\frac{\rho}{1+\rho} \frac{1}{1-\bar{\lambda}\varrho^2} + \frac{1}{1+\rho} \right], \\ \bar{\mathcal{A}}_1 &= \frac{\rho}{1+\rho} [-\bar{\lambda}^2 (\bar{\lambda}\varrho^4 - \bar{\lambda}\varrho^2 + \psi_2\varrho^2 + \varrho^2 - \bar{\lambda}\psi_2\varrho^4 + 1)] \\ &\quad + \frac{1}{1+\rho} [\bar{\lambda}^2 (\bar{\lambda}\varrho^2 - 1) (\bar{\lambda}^2\varrho^4 - 2\bar{\lambda}\varrho^2 + \varrho^2 + 1)], \\ \bar{\mathcal{A}}_* &= (\psi_2 - 1) \bar{\lambda}^3 \varrho^6 + (1 - 3\psi_2) \bar{\lambda}^2 \varrho^4 + 3\psi_2 \bar{\lambda} \varrho^2 - \psi_2. \end{aligned}$$

C.4. Asymptotic normality result for a general direction. In this section we consider a general $\mathbf{u}_0 \in \text{Unif}(\mathbb{S}^{p-1}(1))$ instead of a canonical basis \mathbf{e}_j . Defined in Definitions 1 and C.2.1, the functions $\mathbf{T}_0, \mathbf{T}_1, \mathbf{T}_L, \mathbf{T}_{NL}, \zeta_L$ and ζ are linear in \mathbf{e}_j . So we can naturally extend the functions for a general direction $\mathbf{u}_0 \in \text{Unif}(\mathbb{S}^{p-1}(1))$ by linear combinations, for e.g., $\zeta_L(\mathbf{u}_0) = \mathbf{u}_0^\top (\zeta_L(\mathbf{e}_j))^{j \in [p]}$. Extending the functions by linear combinations is equivalent to replacing \mathbf{e}_j with \mathbf{u}_0 in the definitions of the functions.

Theorem C.4.1 provides the asymptotic normality of $\zeta_L(\mathbf{u}_0)$ for a general \mathbf{u}_0 satisfying $\|\mathbf{u}_0\|_2 = 1$.

Theorem C.4.1. Let $t \in \mathbb{R}$. Under model (5), Assumption 1, 2, 3 and 4, Definition 1 and a further assumption that $\Sigma'_p(0) = O(1/p)$, we have

$$(C.7) \quad \sup_{\mathbf{u}_0 \in S_p} \left| \mathbb{P}_{\mathbf{X}, \mathbf{W}, \epsilon, G | \mathbf{u}_0} \left(\frac{\zeta_L(\mathbf{u}_0)}{\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2} \leq t \right) - \Phi(t) \right| \rightarrow 0$$

for some $S_p \subset \mathbb{S}^{p-1}(1)$ satisfying $|S_p|/|\mathbb{S}^{p-1}(1)| \geq 1 - \log(p)/p \rightarrow 1$.

The operations of taking expectations in this section are defined as follows:

- (i) Let $\mathbb{E}_0, \text{Var}_0$ denote the conditional expectation and the conditional variance given $\mathbf{X}\mathbf{Q}_0, \epsilon, \mathbf{W}, G, \mathbf{u}_0$.
- (ii) Let $\mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \epsilon, G}$ or \mathbb{E} denote the expectation with respect to the total probability.
- (iii) Let $\mathbb{E}_{\mathbf{X}, \mathbf{W}, \epsilon, G | \mathbf{u}_0}$ or $\mathbb{E}_{\mathbf{X}, \mathbf{W}, \epsilon, G}$ denote the conditional expectation given \mathbf{u}_0 .
- (iv) Let $\mathbb{E}_{\mathbf{u}_0}$ denote the conditional expectation given $\mathbf{X}, \mathbf{W}, \epsilon, G$.

Proposition C.4.1 is comparable to Proposition C.2.5 but for a general direction \mathbf{u}_0 .

Proposition C.4.1. Let

- (i) $\mathbf{u}_0 \sim \text{Unif}(\mathbb{S}^{p-1}(1))$ independent with $\mathbf{X}, \mathbf{W}, G, \epsilon$.

- (ii) $\mathbf{X}_0 := \mathbf{X} \mathbf{u}_0$.
- (iii) $\zeta(\mathbf{u}_0) := \mathbf{u}_0^\top (\zeta(\mathbf{e}_j))^{j \in [p]}$ where $\zeta(\mathbf{e}_j)$ is defined in Definition C.2.1.
- (iv) $\mathbf{r} = \mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}$.
- (v) $\epsilon_0^2 := \mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2] / (\mathbb{E}_0 [\|\mathbf{r}\|_2^2] + \mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2])$.

Then

- (i) $\mathbb{E}_0 \left[\left(\frac{\zeta(\mathbf{u}_0)}{(\text{Var}_0 \zeta(\mathbf{u}_0))^{1/2}} - \frac{\mathbf{X}_0^\top \mathbb{E}_0[\mathbf{r}]}{\|\mathbb{E}_0[\mathbf{r}]\|_2} \right)^2 \right] \leq 6\epsilon_0^2$.
- (ii) $\mathbb{E}_0 \left[\left| \frac{\|\mathbf{r}\|_2}{(\text{Var}_0[\zeta(\mathbf{u}_0)])^{1/2}} - 1 \right| \right] \leq (1 + \sqrt{2}) \epsilon_0 / (1 - 2\epsilon_0^2)_+^{1/2}$.
- (iii) $\mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G} [\epsilon_0^2] = O(1/p)$.

Proposition C.4.2 is comparable to Proposition C.2.6 but for a general direction \mathbf{u}_0 .

Proposition C.4.2. *Let Ω be as in Proposition C.2.3. Let $\delta_{i,j} = 1$ if $i = j$, 0 otherwise. Let $\bar{\Omega}_3 := \Omega \cap \{\max_{i_1, i_2 \in [n]} |\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2} / p - \delta_{i_1, i_2}| < \delta\}$ where δ is a fixed positive defined in Assumption 4. Let $\mathbf{u}_0 \sim \text{Unif}(\mathbb{S}^{p-1}(1))$ be independent of $\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G$. Under Assumptions 1, 2, 3 and 4 and a further assumption that $\Sigma'_p(0) = O(1/p)$, we have that*

$$(C.8) \quad \mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G} \left[\left(\frac{\text{trace}((\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}' \mathbf{u}_0))}{\|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2} \right)^2 I_{\Omega_3} \right] = O(1/p).$$

APPENDIX D. PROOFS

In this section we provide the proofs of our theorems and the supporting propositions.

D.1. Proofs of Theorem 1 and 2. Theorem 1 is a special case of Theorem 2 when there is no intercept nor perturbation. We prove in this section Theorem 2 based on supporting propositions.

D.1.1. Proof of Theorem 2.

Proof. Let ζ be as in Definition C.2.1. As we explain in the next paragraphs, Propositions C.2.5 and C.2.6 imply that, for a large subset $J_p \subset [p]$, for all $j \in J_p$ we have

- (i) $\frac{\zeta(\mathbf{e}_j)}{(\text{Var}_j[\zeta(\mathbf{e}_j)])^{1/2}} \xrightarrow{d} N(0, 1)$.
- (ii) $\frac{\|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2}{(\text{Var}_j[\zeta(\mathbf{e}_j)])^{1/2}} \xrightarrow{\mathbb{P}} 1$.
- (iii) $\frac{\zeta(\mathbf{e}_j) - \zeta_L(\mathbf{e}_j)}{\|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2} \xrightarrow{\mathbb{P}} 0$.

The above convergence are uniform over $j \in J_p$, where J_p is a large subset of $[p]$. So by Slutsky's Theorem we obtain the convergence in distribution of $\zeta(\mathbf{e}_j)/\|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2$ to $N(0, 1)$ and the convergence is uniform over all $j \in J_p$.

We first specify the subset $J_p \subset [p]$. Notice that (iii) in Proposition C.2.5 provides the existence of a constant $c_{11} > 0$ independent of n, p, d such that $\sum_{j \in [p]} \mathbb{E}[\epsilon_j^2] \leq c_{11}$. We can specify the large volume index set $J_p \subset [p]$ as

$$J_p := \left\{ j \in [p] : \mathbb{E}[\epsilon_j^2] \leq \frac{c_{11}}{\log(p)} \right\}.$$

Since

$$\frac{1}{p} \# \left\{ j \in [p] : \mathbb{E}[\epsilon_j^2] \geq \frac{c_{11}}{\log(p)} \right\} \leq \frac{\frac{1}{p} \sum_{j \in [p]} \mathbb{E}[\epsilon_j^2]}{c_{11}/\log(p)} \leq \frac{\log(p)}{p},$$

we have $|J_p|/p \geq 1 - \frac{\log(p)}{p}$.

We claim that (i) is provided by Proposition C.2.5 and (iii) is provided by Proposition C.2.6 directly. Now let us explain (ii) rigorously based on Proposition C.2.5 (ii). From the definition of J_p and Chebyshev's inequality, we can see that for all $\bar{\epsilon} > 0$,

$$\max_{j \in J_p} \mathbb{P}(|\epsilon_j| > \bar{\epsilon}) \leq \frac{c_{11}}{\bar{\epsilon}^2 \log(p)}.$$

Let $\mathbf{r} = \mathbf{A}\hat{\boldsymbol{\alpha}} - \mathbf{y}$, $V_j := \text{Var}_j[\zeta(\mathbf{e}_j)]$ and $U_j := \left| \|\mathbf{r}\|_2 / V_j^{1/2} - 1 \right|$. If $\epsilon_j \leq \frac{1}{2}$, by Proposition C.2.5 and simple algebra,

$$\mathbb{E}_j[U_j] \leq (1 + \sqrt{2}) \epsilon_j / (1 - 2\epsilon_j^2)^{1/2} \leq (2 + \sqrt{2}) \epsilon_j.$$

Let us consider $\bar{\epsilon} \leq \frac{1}{2}$. We let $\Omega_j(\bar{\epsilon}) := \{\mathbb{E}_j[U_j] < (2 + \sqrt{2})\bar{\epsilon}\}$. Then

$$\mathbb{P}(\Omega_j(\bar{\epsilon})) \geq \mathbb{P}(|\epsilon_j| \leq \bar{\epsilon}) \geq 1 - \frac{c_{11}}{\bar{\epsilon}^2 \log(p)}.$$

Then, letting $\mathbb{I}(\cdot) := I_{\{\cdot\}}$ be the indicator function, we have

$$\begin{aligned} \mathbb{P}(U_j > \epsilon) &:= \mathbb{E}[\mathbb{I}(U_j > \epsilon)] \\ &= \mathbb{E}[\mathbb{E}_j[\mathbb{I}(U_j > \epsilon)] I_{\Omega_j(\bar{\epsilon})}] + \mathbb{E}[\mathbb{E}_j[\mathbb{I}(U_j > \epsilon)] I_{\Omega_j^c(\bar{\epsilon})}] \\ &\leq \mathbb{E}\left[\frac{\mathbb{E}_j[U_j]}{\epsilon} I_{\Omega_j(\bar{\epsilon})}\right] + \mathbb{P}(\Omega_j^c(\bar{\epsilon})) \\ &\leq (2 + \sqrt{2}) \bar{\epsilon} / \epsilon + \frac{c_{11}}{\bar{\epsilon}^2 \log(p)}. \end{aligned}$$

Choosing $\bar{\epsilon} := \min\left(\frac{1}{\log \log(p)}, \frac{1}{2}\right)$, we have that, for all $\epsilon > 0$,

$$\lim_{p \rightarrow +\infty} \max_{j \in J_p} \mathbb{P}(U_j > \epsilon) = 0.$$

Thus we have (ii). \square

D.1.2. *Proof of Proposition C.2.1.* Propositions are restated before their proofs for convenience.

Proposition C.2.1 (Calculation of the Jacobian matrix). *Under model (5),*

$$\begin{aligned} \nabla_{\mathbf{X}_j}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}) &= -(\mathbf{I}_n - \mathbf{H}) \text{diag}(\sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\mathbf{W}\mathbf{e}_j) \hat{\boldsymbol{\alpha}}) \\ &\quad - \mathbf{A}(\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} \text{diag}(\mathbf{W}\mathbf{e}_j) \sigma'(\mathbf{W}\mathbf{X}^\top) \text{diag}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}) \\ &\quad + (\mathbf{I}_n - \mathbf{H}) [(\mathbf{e}_j^\top \boldsymbol{\beta}) \mathbf{I}_n + \text{diag}(\mathbf{G}'\mathbf{e}_j)]. \end{aligned} \quad (\text{C.1})$$

Proof. The calculation of $\nabla_{\mathbf{X}_j} \mathbf{y}$ follows directly by

$$\nabla_{\mathbf{X}_j} \mathbf{y} = \nabla_{\mathbf{X}_j} \mathbf{f} = \text{diag}(\mathbf{f}'\mathbf{e}_j) = (\mathbf{e}_j^\top \boldsymbol{\beta}) \mathbf{I}_n + \text{diag}(\mathbf{G}'\mathbf{e}_j). \quad (\text{D.1})$$

For the calculation of $\nabla_{\mathbf{X}_j}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}})$, we notice that by the KKT condition,

$$\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}.$$

Proposition C.2.1 follows by the following intermediate steps:

$$\begin{aligned} \text{(i)} \quad \nabla_{\mathbf{X}_j} [(\mathbf{I}_n - \mathbf{H})\mathbf{y}] &= (-1) \left[((\partial_{x_{i_2j}} \mathbf{H})\mathbf{y})_{i_1} \right]_{i_1, i_2 \in [n]} + (\mathbf{I}_n - \mathbf{H}) \nabla_{\mathbf{X}_j} \mathbf{y}. \\ \text{(ii)} \end{aligned}$$

$$\begin{aligned} \partial_{x_{i_2j}} \mathbf{H} &= (\mathbf{I}_n - \mathbf{H})(\partial_{x_{i_2j}} \mathbf{A})(\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} \mathbf{A}^\top \\ &\quad + \mathbf{A}(\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} (\partial_{x_{i_2j}} \mathbf{A})^\top (\mathbf{I}_n - \mathbf{H}). \end{aligned}$$

$$\text{(iii)} \quad \partial_{x_{i_2j}} \mathbf{A} = \text{diag}(\mathbf{e}_{i_2}) \sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\mathbf{W}\mathbf{e}_j).$$

The above intermediate steps can be seen from the followings.

(i) By the chain rule for multiplication.

(ii) Notice that $\mathbf{H} := \mathbf{A}(\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} \mathbf{A}$. For the inverse matrix, we use the fact that

$$\partial_{x_{i_2 j}} \left[(\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} \right] = (-1)(\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} (\partial_{x_{i_2 j}} (\mathbf{A}^\top \mathbf{A})) (\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1}.$$

(iii) We notice that $\mathbf{A} = \sigma(\mathbf{X}\mathbf{W}^\top)$. So that by the chain rule,

$$\partial_{x_{i_2 j}} \mathbf{A} = [\delta_{i=i_2} \sigma'(\mathbf{x}_i^\top \mathbf{w}_k) w_{kj}]^{i \in [n], k \in [d]} = \text{diag}(\mathbf{e}_{i_2}) \sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\mathbf{W}\mathbf{e}_j).$$

Combining the intermediate steps above, noticing that

$$\begin{aligned} & \left[\left((\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{e}_{i_2}) \sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\mathbf{W}\mathbf{e}_j) (\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} \mathbf{A}^\top \mathbf{y} \right)_{i_1} \right]^{i_1, i_2 \in [n]} \\ &= (\mathbf{I}_n - \mathbf{H}) \text{diag} \left(\sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\mathbf{W}\mathbf{e}_j) \hat{\boldsymbol{\alpha}} \right), \end{aligned}$$

and

$$\begin{aligned} & \left[\left(\mathbf{A}(\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} \text{diag}(\mathbf{W}\mathbf{e}_j) \sigma'(\mathbf{W}\mathbf{X}^\top) \text{diag}(\mathbf{e}_{i_2}) (\mathbf{I}_n - \mathbf{H}) \mathbf{y} \right)_{i_1} \right]^{i_1, i_2 \in [n]} \\ &= \mathbf{A}(\mathbf{A}^\top \mathbf{A} + n\tau \mathbf{I}_n)^{-1} \text{diag}(\mathbf{W}\mathbf{e}_j) \sigma'(\mathbf{W}\mathbf{X}^\top) \text{diag}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}), \end{aligned}$$

we have our calculation. \square

D.1.3. Proof of Proposition C.2.2.

Proposition C.2.2. Let $\mathbf{T}_0, \mathbf{T}_1, \mathbf{T}_L, \mathbf{T}_{NL}$ be as in Definitions 1 and C.2.1. Let c_1, L, \mathbf{f}' be as in §C.1.

- (i) $\|\mathbf{A}(n\tau \mathbf{I}_d + \mathbf{A}^\top \mathbf{A})^{-1}\|_{op} \leq 1/(2\sqrt{n\tau})$.
- (ii) $\|\mathbf{I}_n - \mathbf{H}\|_{op} \leq 1$.
- (iii) $\sum_{j \in [p]} \|\mathbf{T}_0(\mathbf{e}_j) + \mathbf{T}_L(\mathbf{e}_j) + \mathbf{T}_{NL}(\mathbf{e}_j)\|_F^2 \leq 2c_1^2 L^2 n \|\hat{\boldsymbol{\alpha}}\|_2^2 + 2\|\mathbf{f}'\|_F^2$.
- (iv) $\|\mathbf{T}_1(\mathbf{e}_j)\|_F^2 \leq L^2 c_1^2 / (4n\tau) \cdot \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2$.
- (v) $\mathbb{E}\|\hat{\boldsymbol{\alpha}}\|_2^2 = O(1)$.

Proof. (i) Let $\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$ be the SVD of \mathbf{A} where $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times d}$ has diagonal elements the singular values σ_l for $l \in [\min(n, d)]$. Then the singular values of $\mathbf{A}(n\tau \mathbf{I}_d + \mathbf{A}^\top \mathbf{A})^{-1}$ are either $\sigma_l/(n\tau + \sigma_l^2)$ or 0. We notice that

$$\sigma_l/(n\tau + \sigma_l^2) \leq 1/(2\sqrt{n\tau}).$$

So that the operator norm is no more than $1/(2\sqrt{n\tau})$.

- (ii) Let $\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$ be the SVD of \mathbf{A} where $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times d}$ has diagonal elements the singular values σ_l for $l \in [\min(n, d)]$. Then the singular values of $\mathbf{I}_n - \mathbf{H}$ are either $(n\tau)/(n\tau + \sigma_l^2)$ or 1, no more than 1.

(iii)

$$\begin{aligned} \sum_{j \in [p]} \|\mathbf{T}_0(\mathbf{e}_j) + \mathbf{T}_L(\mathbf{e}_j) + \mathbf{T}_{NL}(\mathbf{e}_j)\|_F^2 &\leq \sum_{j \in [p]} \|\mathbf{I}_n - \mathbf{H}\|_{op}^2 \left\| \left(\sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right) \mathbf{e}_j \right\|_2^2 \\ &\leq \sum_{j \in [p]} \left\| \left(\sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right) \mathbf{e}_j \right\|_2^2 \\ &= \left\| \sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right\|_F^2 \\ &\leq 2 \left\| \sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} \right\|_F^2 + 2 \|\mathbf{f}'\|_F^2 \\ &\leq 2 \|\mathbf{W}\|_{op}^2 \left\| \sigma'(\mathbf{X}\mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \right\|_F^2 + 2 \|\mathbf{f}'\|_F^2 \\ &\leq 2c_1^2 L^2 n \|\hat{\boldsymbol{\alpha}}\|_2^2 + 2 \|\mathbf{f}'\|_F^2. \end{aligned}$$

We used $\|\mathbf{I}_n - \mathbf{H}\|_{\text{op}} \leq 1$ in the above display.

(iv)

$$\begin{aligned} \|\mathbf{T}_1(\mathbf{e}_j)\|_F^2 &\leq \|\mathbf{A}(n\tau\mathbf{I}_d + \mathbf{A}^\top \mathbf{A})^{-1}\|_{\text{op}}^2 \cdot \|\text{diag}(\mathbf{W}\mathbf{e}_j)\sigma'(\mathbf{W}\mathbf{X}^\top) \text{diag}(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}})\|_F^2 \\ &\leq 1/(4n\tau) \cdot L^2 \cdot \|(\mathbf{W}\mathbf{e}_j)(\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}})^\top\|_F^2 \\ &= L^2/(4n\tau) \cdot \|\mathbf{W}\mathbf{e}_j\|_2^2 \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \\ &\leq L^2/(4n\tau) \cdot \|\mathbf{W}\|_{\text{op}}^2 \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \\ &\leq L^2 c_1^2/(4n\tau) \cdot \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2. \end{aligned}$$

We used $\|\mathbf{A}(n\tau\mathbf{I}_d + \mathbf{A}^\top \mathbf{A})^{-1}\|_{\text{op}}^2 \leq 1/(4n\tau)$ in the above display.

(v) From the KKT condition, we have

$$n\tau\|\hat{\boldsymbol{\alpha}}\|_2^2 = (\mathbf{A}\hat{\boldsymbol{\alpha}})^\top (\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}) = \mathbf{y}^\top \mathbf{H}^\top (\mathbf{I}_n - \mathbf{H})\mathbf{y}.$$

Let the singular values of \mathbf{A} be σ_l for $l \in [\min(n, d)]$, then the singular value of $\mathbf{H}^\top (\mathbf{I}_n - \mathbf{H})$ are $\frac{\sigma_l^2}{n\tau + \sigma_l^2} \cdot \frac{n\tau}{n\tau + \sigma_l^2} \in [0, 1/4]$. So that $\|\hat{\boldsymbol{\alpha}}\|_2^2 \leq 1/(4n\tau) \cdot \|\mathbf{y}\|_2^2$. Taking expectation we have

$$\begin{aligned} \mathbb{E}\|\hat{\boldsymbol{\alpha}}\|_2^2 &\leq 1/(4n\tau) \cdot \mathbb{E}\|\beta_0 \mathbf{1}_n + \mathbf{X}\boldsymbol{\beta} + \mathbf{G} + \boldsymbol{\varepsilon}\|_2^2 \\ &= 1/(4\tau) \cdot (\beta_0^2 + \|\boldsymbol{\beta}\|_2^2 + \mathbb{E}[\Sigma_p(\|\mathbf{x}_1\|_2^2/p)] + \theta_\varepsilon^2) = O(1). \end{aligned}$$

□

D.1.4. Proof of Proposition C.2.3.

Proposition C.2.3. *There exists a large event Ω such that $\mathbb{P}(\Omega^c) \leq o(\exp(-c_6 n))$ for some $c_6 > 0$ and that on Ω ,*

- (i) $\|\mathbf{X}\|_F^2/n^2 \leq c_{4,n}$,
- (ii) $1/n \cdot \min_{j \in [p]} \mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq c_{2,n}$,
- (iii) $1/n \cdot \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq c_{2,n}$.

The constant c_6 is independent of n, d, p . The constants $c_{4,n}$ and $c_{2,n}$ are given in Proposition C.2.4.

Proof. We first notice that by the KKT condition for the ridge regression type estimator $\hat{\boldsymbol{\alpha}}$, we have

$$\mathbb{E}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 = \mathbb{E}_\varepsilon \|(\mathbf{I}_n - \mathbf{H})\mathbf{y}\|_2^2 \geq \mathbb{E}_\varepsilon \|(\mathbf{I}_n - \mathbf{H})\boldsymbol{\varepsilon}\|_2^2 = \theta_\varepsilon^2 \|\mathbf{I}_n - \mathbf{H}\|_F^2.$$

Next we show that $\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|$ is concentrated around $\mathbb{E}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|$: By KKT condition, $\nabla_\varepsilon (\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}) = \mathbf{I}_n - \mathbf{H}$. Since $\|\mathbf{I}_n - \mathbf{H}\|_{\text{op}} \leq 1$, $\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|$ is 1-Lipschitz in ε (see also Bellec and Tsybakov [2017] for general results of this kind). By the triangle inequality and the independence between \mathbf{X}_j and ε , if $u(\varepsilon, \mathbf{X}_j)$ is a function that is 1-Lipschitz with respect to ε for every value of \mathbf{X}_j , then

$$|\mathbb{E}_j u(\varepsilon, \mathbf{X}_j) - \mathbb{E}_j u(\tilde{\varepsilon}, \mathbf{X}_j)| \leq \mathbb{E}_j |u(\varepsilon, \mathbf{X}_j) - u(\tilde{\varepsilon}, \mathbf{X}_j)| \leq \|\varepsilon - \tilde{\varepsilon}\|_2 \mathbb{E}_j [1] = \|\varepsilon - \tilde{\varepsilon}\|_2.$$

This implies that $\mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2$ is also 1-Lipschitz in ε . By the concentration inequality for Lipschitz functions of a standard normal random vector (See Theorem 5.2.2 in Vershynin [2018] or Theorem 5.5 in Boucheron et al. [2013]) applied to the mappings $\varepsilon \mapsto \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2$ and $\varepsilon \mapsto \mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2$, we have that for some universal constant $c_5 > 0$ and for all $t > 0$,

$$\begin{aligned} \mathbb{P}(\mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 \geq \mathbb{E}_\varepsilon \mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 + \sqrt{n}\theta_\varepsilon t) &\geq 1 - 2\exp(-c_5 n t^2), \\ \mathbb{P}(\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 \geq \mathbb{E}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 + \sqrt{n}\theta_\varepsilon t) &\geq 1 - 2\exp(-c_5 n t^2). \end{aligned}$$

By the union bound, the following event occurs with probability at least $1 - 2(p+1)\exp(-c_5 n t^2)$:

$$\bigcap_{j \in [p]} \left\{ \mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 \geq \mathbb{E}_\varepsilon \mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 - \sqrt{n}\theta_\varepsilon t, \quad \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 \geq \mathbb{E}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 - \sqrt{n}\theta_\varepsilon t \right\}.$$

Since $\mathbf{X}_j, \mathbf{X}_{-j}, \varepsilon$ are independent, we can exchange the order of expectation by Fubini's Theorem,

$$\begin{aligned}
\mathbb{E}_\varepsilon \mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 - \sqrt{n}\theta_\varepsilon t &= \mathbb{E}_j \mathbb{E}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 - \sqrt{n}\theta_\varepsilon t, \\
&= \mathbb{E}_j \left[(\mathbb{E}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 - \text{Var}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2)^{1/2} \right] - \sqrt{n}\theta_\varepsilon t, \\
&\geq \theta_\varepsilon \mathbb{E}_j \left[(\|\mathbf{I}_n - \mathbf{H}\|_F^2 - 1)_+^{1/2} \right] - \sqrt{n}\theta_\varepsilon t \\
&\geq \theta_\varepsilon \mathbb{E}_j \|\mathbf{I}_n - \mathbf{H}\|_F - (\theta_\varepsilon + \sqrt{n}\theta_\varepsilon t) \\
&= \theta_\varepsilon [\mathbb{E}_j \|\mathbf{I}_n - \mathbf{H}\|_F - \sqrt{nt} - 1].
\end{aligned}$$

The first inequality above is due to the fact that $\mathbb{E}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq \theta_\varepsilon^2 \|\mathbf{I}_n - \mathbf{H}\|_F^2$ and the fact that $\text{Var}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 \leq \theta_\varepsilon^2$ by the Gaussian Poincaré Inequality [Boucheron et al., 2013, Theorem 3.20] with respect to the 1-Lipschitz mapping $\varepsilon \mapsto \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2$. The second last inequality above is due to the fact that for any two positive real number a, b , $(a^2 - b^2)_+^{1/2} \geq a - b$. By a similar argument,

$$\mathbb{E}_\varepsilon \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2 - \sqrt{n}\theta_\varepsilon t \geq \theta_\varepsilon [\|\mathbf{I}_n - \mathbf{H}\|_F - \sqrt{nt} - 1].$$

From Proposition C.2.4 we can have for all $t > 0$ and some universal constant $c_5 > 0$,

$$\begin{aligned}
\mathbb{P} \left(1/n \cdot \min_{j \in [p]} \mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq \theta_\varepsilon^2 ((1 + \bar{F}_n/\tau)^{-1} - t - 1/\sqrt{n})_+^2 \right) &\geq 1 - 2p \exp(-c_5 n t^2). \\
\mathbb{P} \left(1/n \cdot \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq \theta_\varepsilon^2 ((1 + \bar{F}_n/\tau)^{-1} - t - 1/\sqrt{n})_+^2 \right) &\geq 1 - \exp(-c_5 n t^2).
\end{aligned}$$

Consider the intersection of the above events with the event $\{\|\mathbf{X}\|_F^2/n^2 \leq c_{4,n}\}$. We notice that on that intersection,

$$\begin{aligned}
\bar{F}_n &\leq 2c_1^2 L^2/n + 2c_1^2 L^2 c_{4,n} + 2\psi_{d,n}(\sigma(0))^2 \\
&\leq 2c_1^2 L^2 c_4 + 2\psi_d(\sigma(0))^2 + o(1), \\
(1 + \bar{F}_n/\tau)^{-1} &\geq \bar{c}_{2,n} = \bar{c}_2 + o(1).
\end{aligned}$$

Taking $t = \bar{c}_{2,n}/3$, we obtain

$$\begin{aligned}
\mathbb{P} \left(1/n \cdot \min_{j \in [p]} \mathbb{E}_j \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq \theta_\varepsilon^2 (2\bar{c}_{2,n}/3 - 1/\sqrt{n})_+^2 \right) &\geq 1 - 2p \exp(-9^{-1} c_5 \bar{c}_{2,n}^2 n) \\
&\quad - \exp(-n^2) \\
\text{as well as } \mathbb{P} \left(1/n \cdot \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq \theta_\varepsilon^2 (2\bar{c}_{2,n}/3 - 1/\sqrt{n})_+^2 \right) &\geq 1 - \exp(-9^{-1} c_5 \bar{c}_{2,n}^2 n) \\
&\quad - \exp(-n^2).
\end{aligned}$$

□

D.1.5. Proof of Proposition C.2.4.

Proposition C.2.4. *Let*

- (1) $c_{3,n} := \psi_{p,n}^{-1/2} + 1 + c^{-1/2} \rightarrow c_3 := \psi_p^{-1/2} + 1 + c^{-1/2}$ where $c > 0$ is some universal constant specified in the proof.
- (2) $c_{4,n} := \psi_{p,n} + 2\psi_{p,n}^{1/2} + 2 \rightarrow c_4 := \psi_p + 2\psi_p^{1/2} + 2$.
- (3) $F_n := 2c_1^2 L^2 \|\mathbf{X}\|_F^2/n^2 + 2\psi_{d,n}(\sigma(0))^2$.
- (4) $\bar{F}_n := 2c_1^2 L^2 (n + \|\mathbf{X}\|_F^2)/n^2 + 2\psi_{d,n}(\sigma(0))^2$.
- (5) $\bar{c}_{2,n} := (1 + 2c_1^2 L^2 n^{-1} \tau^{-1} + 2c_1^2 L^2 c_{4,n} \tau^{-1} + 2\psi_{d,n}(\sigma(0))^2 \tau^{-1})^{-1} \rightarrow \bar{c}_2 := (1 + 2c_1^2 L^2 c_4 \tau^{-1} + 2\psi_d(\sigma(0))^2 \tau^{-1})^{-1} > 0$.
- (6) $c_{2,n} := \theta_\varepsilon^2 (2\bar{c}_{2,n}/3 - n^{-1/2})_+^2 \rightarrow c_2 := \theta_\varepsilon^2 (2\bar{c}_2/3)^2 > 0$.

Then

- (i) $\mathbb{P}\left(\|\mathbf{X}/\sqrt{p}\|_{\text{op}} \geq c_{3,n}\right) \leq 2\exp(-p).$
- (ii) $\mathbb{P}(\|\mathbf{X}\|_F^2/n^2 > c_{4,n}) \leq \exp(-n^2).$
- (iii) $\|\sigma(\mathbf{X}\mathbf{W}^\top)\|_F^2/n^2 \leq F_n.$
- (iv) $\|\mathbf{I}_n - \mathbf{H}\|_F^2/n \geq (1 + F_n/\tau)^{-2}.$
- (v) $\min_{j \in [p]} (n^{-1/2} \mathbb{E}_j \|\mathbf{I}_n - \mathbf{H}\|_F) \geq (1 + \bar{F}_n/\tau)^{-1}.$

Proof. (i) Corollary 7.3.3 in [Vershynin \[2018\]](#) provides the high probability upper bound for the operator norm of random matrix,

$$\mathbb{P}(\|\mathbf{X}\|_{\text{op}} \geq \sqrt{n} + \sqrt{p} + t) \leq 2\exp(-ct^2)$$

for some universal constant $c > 0$. Taking $ct^2 = p$, we have

$$\mathbb{P}\left(\|\mathbf{X}/\sqrt{p}\|_{\text{op}} \geq \psi_{p,n}^{-1/2} + 1 + c^{-1/2}\right) \leq 2\exp(-p).$$

- (ii) By Lemma 1 in [Laurent and Massart \[2000\]](#), we have the concentration for $\|\mathbf{X}\|_F^2$, the chi-square random-variable with degree of freedom np , as follows: For any $x > 0$,

$$\mathbb{P}(\|\mathbf{X}\|_F^2 - np \geq 2\sqrt{xn p} + 2x) \leq \exp(-x).$$

Here we take $x = n^2$ for simplicity of proof.

- (iii) If σ is L -Lipschitz, then

$$\begin{aligned} |\sigma(x)| &\leq L|x| + |\sigma(0)|, \\ \implies (\sigma(x))^2 &\leq 2L^2x^2 + 2(\sigma(0))^2. \end{aligned}$$

Taking $x = \mathbf{x}_i^\top \mathbf{w}_k$ and summing over $(i, k) \in [n] \times [d]$ we have

$$\begin{aligned} \|\sigma(\mathbf{X}\mathbf{W}^\top)\|_F^2/n^2 &\leq 2L^2\|\mathbf{X}\mathbf{W}^\top\|_F^2/n^2 + 2\psi_{d,n}(\sigma(0))^2 \\ &\leq 2c_1^2L^2\|\mathbf{X}\|_F^2/n^2 + 2\psi_{d,n}(\sigma(0))^2. \end{aligned}$$

- (iv) Let σ_l , $l \in [\min(n, d)]$ be the singular values of \mathbf{A} . If $n > d$, we define $\sigma_i = 0$ for $i \in (d, n]$. We notice that by the SVD of \mathbf{H} ,

$$\begin{aligned} \|\mathbf{I}_n - \mathbf{H}\|_F^2 &= \sum_{i \in [n]} \left(\frac{1}{1 + \sigma_i^2/(n\tau)} \right)^2 \\ &\geq n \left(\frac{1}{1 + \sum_{l \in \min(n, d)} \sigma_l^2/(n^2\tau)} \right)^2 \\ &= n \left(1 + \|\mathbf{A}\|_F^2/(n^2\tau) \right)^{-2} \end{aligned}$$

where the inequality is due to Jensen's Inequality.

- (v) This is by the convexity of $x \mapsto \frac{1}{1+x}$ on \mathbb{R}^+ and Jensen's Inequality, $\mathbb{E}_j(\frac{1}{1+X}) \geq \frac{1}{1+\mathbb{E}_j X}$ for any random variable X supported on \mathbb{R}^+ . We then notice that $\mathbb{E}_j[\|\mathbf{X}\|_F^2] \leq n + \|\mathbf{X}\|_F^2$. \square

D.1.6. Proof of Proposition [C.2.5](#).

Proposition C.2.5. *Let*

- (i) $\zeta(\mathbf{e}_j)$ be as in Definition [C.2.1](#).
- (ii) $\mathbf{r} = \mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}$.
- (iii) $\epsilon_j^2 := \mathbb{E}_j[\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2] / (\mathbb{E}_j[\|\mathbf{r}\|_2^2] + \mathbb{E}_j[\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2])$.

Then

$$(i) \quad \mathbb{E}_j \left[\left(\frac{\zeta(\mathbf{e}_j)}{(\text{Var}_j \zeta(\mathbf{e}_j))^{1/2}} - \frac{\mathbf{X}_j^\top \mathbb{E}_j[\mathbf{r}]}{\|\mathbb{E}_j[\mathbf{r}]\|_2} \right)^2 \right] \leq 6\epsilon_j^2.$$

$$\begin{aligned}
(ii) \quad & \mathbb{E}_j \left[\left| \frac{\|\mathbf{r}\|_2}{(\text{Var}_j[\zeta(\mathbf{e}_j)])^{1/2}} - 1 \right| \right] \leq (1 + \sqrt{2}) \epsilon_j / (1 - 2\epsilon_j^2)_+^{1/2}. \\
(iii) \quad & \mathbb{E} \left[\sum_{j \in [p]} \epsilon_j^2 \right] = O(1).
\end{aligned}$$

Proof. (i) The first inequality follows directly from Theorem 2.1 in [Bellec and Zhang \[2019\]](#).
(ii) By Second Order Stein's lemma applied to the mapping $\mathbf{X}_j \mapsto \mathbf{r}$ (cf. Theorem 2.1 in [Bellec and Zhang \[2018\]](#)), we have

$$V := \text{Var}_j [\zeta(\mathbf{e}_j)] = \mathbb{E}_j [\|\mathbf{r}\|_2^2] + \mathbb{E}_j [\text{trace}[(\nabla_{\mathbf{X}_j} \mathbf{r})^2]].$$

By Second Order Stein's lemma applied to the mapping $\mathbf{X}_j \mapsto \mathbf{r} - \mathbb{E}_j [\mathbf{r}]$, we have

$$\begin{aligned}
\bar{V} &:= \text{Var}_j \left(\mathbf{X}_j^\top (\mathbf{r} - \mathbb{E}_j [\mathbf{r}]) - \text{trace} [\nabla_{\mathbf{X}_j} \mathbf{r}] \right) \\
&= \mathbb{E}_j [\|\mathbf{r} - \mathbb{E}_j [\mathbf{r}]\|_2^2 + \text{trace}[(\nabla_{\mathbf{X}_j} \mathbf{r})^2]] \\
&= \text{Var}_j [\zeta(\mathbf{e}_j)] - \|\mathbb{E}_j [\mathbf{r}]\|_2^2.
\end{aligned}$$

By the fact that $2ab \in [-a^2 - b^2, a^2 + b^2]$ for two real a and b ,

$$|\mathbb{E}_j [\text{trace}[(\nabla_{\mathbf{X}_j} \mathbf{r})^2]]| \leq \mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2].$$

By Gaussian Poincaré Inequality [[Boucheron et al., 2013](#), Theorem 3.20] applied to the mapping $\mathbf{X}_j \mapsto r_i(\mathbf{X}_j)$ for each $i \in [n]$,

$$\mathbb{E}_j [\|\mathbf{r} - \mathbb{E}_j [\mathbf{r}]\|_2^2] \leq \mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2].$$

So that $\bar{V} \leq 2\mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2]$. We also notice that

$$\begin{aligned}
\mathbb{E}_j \left[\left| \frac{\|\mathbf{r}\|_2}{V^{1/2}} - 1 \right| \right] &= V^{-1/2} \mathbb{E}_j \left[\left| \|\mathbf{r}\|_2 - V^{1/2} \right| \right] \\
&\leq V^{-1/2} \mathbb{E}_j [|\|\mathbf{r}\|_2 - \|\mathbb{E}_j [\mathbf{r}]\|_2|] + V^{-1/2} \left| \|\mathbb{E}_j [\mathbf{r}]\|_2 - V^{1/2} \right|.
\end{aligned}$$

By Jensen's Inequality,

$$\begin{aligned}
\mathbb{E}_j [|\|\mathbf{r}\|_2 - \|\mathbb{E}_j [\mathbf{r}]\|_2|] &\leq \left(\mathbb{E}_j \left[(\|\mathbb{E}_j [\mathbf{r}]\|_2 - \|\mathbf{r}\|_2)^2 \right] \right)^{1/2} \\
&\leq (\mathbb{E}_j [\|\mathbf{r} - \mathbb{E}_j [\mathbf{r}]\|_2^2])^{1/2} \\
&\leq \sqrt{\mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2]}.
\end{aligned}$$

The last inequality above follows by the Gaussian Poincaré Inequality [[Boucheron et al., 2013](#), Theorem 3.20] applied to the mapping $\mathbf{X}_j \mapsto r_i(\mathbf{X}_j)$ for each $i \in [n]$. By $|a - b| \leq \sqrt{a^2 + b^2}$ for two real $a, b > 0$,

$$\left| \|\mathbb{E}_j [\mathbf{r}]\|_2 - V^{1/2} \right| \leq \bar{V}^{1/2} \leq \sqrt{2} \sqrt{\mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2]}.$$

So that

$$\begin{aligned}
\mathbb{E}_j \left[\left| \frac{\|\mathbf{r}\|_2}{V^{1/2}} - 1 \right| \right] &\leq (1 + \sqrt{2}) \left(\frac{\mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2]}{V} \right)^{1/2} \\
&\leq (1 + \sqrt{2}) \left(\frac{\mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2]}{(\mathbb{E}_j [\|\mathbf{r}\|_2^2] - \mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2])_+} \right)^{1/2} \\
&\leq (1 + \sqrt{2}) \epsilon_j / (1 - 2\epsilon_j^2)_+^{1/2},
\end{aligned}$$

where $\epsilon_j^2 = \mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2] / (\mathbb{E}_j [\|\mathbf{r}\|_2^2] + \mathbb{E}_j [\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2])$. In the above we let $a/0 = +\infty$ by convention.

- (iii) We first recall a large probability event Ω , defined in Proposition C.2.3. Since $1 = I_\Omega + I_{\Omega^c}$ and $\sum_{j \in [p]} \epsilon_j^2 \leq p$,

$$\mathbb{E} \left[\sum_{j \in [p]} \epsilon_j^2 \right] \leq \mathbb{E} \left[\sum_{j \in [p]} \epsilon_j^2 I_\Omega \right] + o(p \exp(-c_6 n)).$$

Let $\mathbf{T}_0, \mathbf{T}_1, \mathbf{T}_L, \mathbf{T}_{NL}$ be as in Proposition C.2.1. Then $\nabla_{\mathbf{X}_j} \mathbf{r} = \mathbf{T}_0 + \mathbf{T}_1 + \mathbf{T}_L + \mathbf{T}_{NL}$,

$$\|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2 \leq 2\|\mathbf{T}_0 + \mathbf{T}_L + \mathbf{T}_{NL}\|_F^2 + 2\|\mathbf{T}_1\|_F^2.$$

We have

$$\mathbb{E} \left[\sum_{j \in [p]} \epsilon_j^2 I_\Omega \right] \leq \mathbb{E} \left[\sum_{j \in [p]} \frac{\mathbb{E}_j \|\mathbf{T}_0 + \mathbf{T}_L + \mathbf{T}_{NL}\|_F^2 + \mathbb{E}_j \|\mathbf{T}_1\|_F^2}{\mathbb{E}_j \|\mathbf{r}\|_2^2 + \mathbb{E}_j \|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2} I_\Omega \right]$$

By Proposition C.2.2 and Proposition C.2.3,

$$\begin{aligned} \mathbb{E} \left[\sum_{j \in [p]} \frac{\mathbb{E}_j \|\mathbf{T}_1(\mathbf{e}_j)\|_F^2}{\mathbb{E}_j \|\mathbf{r}\|_2^2 + \mathbb{E}_j \|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2} I_\Omega \right] &\leq \mathbb{E} \left[\sum_{j \in [p]} \frac{L^2 c_1^2 / (4n\tau) \cdot \mathbb{E}_j \|\mathbf{r}\|_2^2}{\mathbb{E}_j \|\mathbf{r}\|_2^2 + \mathbb{E}_j \|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2} \right] \\ &\leq (1/4) L^2 c_1^2 \psi_{p,n} \tau^{-1}, \\ &= (1/4) L^2 c_1^2 \psi_p \tau^{-1} + o(1). \end{aligned}$$

$$\begin{aligned} \mathbb{E} \left[\sum_{j \in [p]} \frac{\mathbb{E}_j \|\mathbf{T}_0(\mathbf{e}_j) + \mathbf{T}_L(\mathbf{e}_j) + \mathbf{T}_{NL}(\mathbf{e}_j)\|_F^2}{\mathbb{E}_j \|\mathbf{r}\|_2^2 + \mathbb{E}_j \|\nabla_{\mathbf{X}_j} \mathbf{r}\|_F^2} I_\Omega \right] &\leq 1/(c_{2,n} n) \cdot \mathbb{E} \left[\sum_{j \in [p]} \|\mathbf{T}_0 + \mathbf{T}_L + \mathbf{T}_{NL}\|_F^2 \right] \\ &\leq 1/(c_{2,n} n) \cdot \left(2c_1^2 L^2 n \mathbb{E} [\|\hat{\boldsymbol{\alpha}}\|_2^2] + 2\mathbb{E} [\|\mathbf{f}'\|_F^2] \right) \\ &\leq 2c_1^2 c_{2,n}^{-1} L^2 \mathbb{E} [\|\hat{\boldsymbol{\alpha}}\|_2^2] + 2c_{2,n}^{-1} \mathbb{E} [\|\boldsymbol{\beta} + \nabla G(\mathbf{x}_1)\|_2^2]. \end{aligned}$$

Proposition C.2.2 and Assumptions 1 and 3 provide us $\mathbb{E} [\|\hat{\boldsymbol{\alpha}}\|_2^2] = O(1)$ and $\mathbb{E} [\|\boldsymbol{\beta} + \nabla G(\mathbf{x}_1)\|_2^2] = O(1)$. Combining the above we have $\mathbb{E} [\sum_{j \in [p]} \epsilon_j^2] \leq O(1)$.

□

D.1.7. Proof of Proposition C.2.6.

Proposition C.2.6. Under Assumptions 1, 2, 3 and 4, we have the following convergence in probability: for all $\epsilon > 0$,

$$(C.5) \quad \lim_{n \rightarrow +\infty} \mathbb{P} \left(\max_{j \in [p]} \left| \frac{\text{trace}((\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}' \mathbf{e}_j))}{\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2} \right| > \epsilon \right) = 0.$$

Proof. Let us first define

- (i) $\mathbf{v}_n = \mathbf{1}_n - (H_{11}, H_{22}, \dots, H_{nn})^\top$.
- (ii) $q_j := \text{trace}(\mathbf{T}_{NL}(\mathbf{e}_j)) = \text{trace}[(\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}' \mathbf{e}_j)] = \mathbf{v}_n^\top [\mathbf{G}']_j$.
- (iii) $\mathbf{r} := \mathbf{A}\hat{\boldsymbol{\alpha}} - \mathbf{y}$.
- (iv) Ω be in Proposition C.2.3 such that
 - (i) $\Omega := \{1/n \cdot \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 \geq c_{2,n}\}$.
 - (ii) $\mathbb{P}(\Omega^c) \leq o(\exp(-c_6 n))$ for some constant $c_6 > 0$.
- (v) $\bar{\Omega}_3$ be such that (cf. Corollary 2.8.3 and Proposition 2.5.2 (i) in Vershynin [2018])

(i)

$$\begin{aligned} \bar{\Omega}_3 = & \left\{ \max_{i_1, i_2 \in [n]} |\mathbf{x}_{i_1}^T \mathbf{x}_{i_2}/p - \delta_{i_1, i_2}| < \delta \right\} \cap \left\{ \max_{j_1, j_2 \in [p]} |\mathbf{X}_{j_1}^\top \mathbf{X}_{j_2}/n - \delta_{j_1, j_2}| < \delta \right\} \\ & \cap \left\{ \max_{i \in [n], j \in [p]} |x_{ij}| \leq p^{1/4} \right\} \end{aligned}$$

(ii) $\mathbb{P}(\bar{\Omega}_3^c) \leq o(\exp(-c_{10}n^{1/2}))$ for some universal constant $c_{10} > 0$.(vi) $\Omega_3 := \Omega \cap \bar{\Omega}_3$.

By Chebyshev's inequality,

$$\begin{aligned} \mathbb{P} \left(\max_{j \in [p]} \frac{|q_j|}{\|\mathbf{r}\|_2} > \epsilon \right) & \leq \mathbb{P} \left(\left\{ \max_{j \in [p]} \frac{|q_j|}{\|\mathbf{r}\|_2} > \epsilon \right\} \cap \Omega_3 \right) + \mathbb{P}(\Omega_3^c) \\ & \leq \mathbb{P} \left(\left\{ \max_{j \in [p]} \frac{|q_j|}{n^{-1/2}c_{2,n}^{1/2}} > \epsilon \right\} \cap \Omega_3 \right) + \mathbb{P}(\Omega_3^c) \\ & \leq \mathbb{P} \left(\left\{ \max_{j \in [p]} \frac{|q_j|}{n^{-1/2}c_{2,n}^{1/2}} > \epsilon \right\} \cap \bar{\Omega}_3 \right) + \mathbb{P}(\Omega_3^c) \\ & \leq \mathbb{P} \left(\left\{ \max_{j \in [p]} \frac{|q_j|}{n^{-1/2}c_{2,n}^{1/2}} I_{\bar{\Omega}_3} > \epsilon \right\} \right) + \mathbb{P}(\Omega_3^c) \\ & \leq \frac{\mathbb{E} [\max_{j \in [p]} q_j^2 I_{\bar{\Omega}_3}]}{n\epsilon^2 c_{2,n}} + o(1). \end{aligned}$$

So to show our proposition, it suffices to show that

$$\mathbb{E} \left[\max_{j \in [p]} q_j^2 I_{\bar{\Omega}_3} \right] = o(n).$$

Letting $\mathbf{v}_n = \mathbf{1}_n - (H_{11}, H_{22}, \dots, H_{nn})^\top$, we notice that by some algebra,

$$\begin{aligned} \mathbb{E} \left[\max_{j \in [p]} q_j^2 I_{\bar{\Omega}_3} \right] & = \mathbb{E} \left[\max_{j \in [p]} (\text{trace}[(\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}' \mathbf{e}_j)])^2 I_{\bar{\Omega}_3} \right] \\ & = \mathbb{E} \left[\max_{j \in [p]} (\mathbf{v}_n^\top [\mathbf{G}']_j)^2 I_{\bar{\Omega}_3} \right] \\ & = \mathbb{E} \left[\max_{j \in [p]} \mathbf{v}_n^\top \mathbb{E}_G [[\mathbf{G}']_j ([\mathbf{G}']_j)^\top] \mathbf{v}_n I_{\bar{\Omega}_3} \right] \end{aligned}$$

Let

$$\mathbf{M}(\mathbf{e}_j) := \mathbb{E}_G [[\mathbf{G}']_j ([\mathbf{G}']_j)^\top].$$

From Proposition C.2.7, on $\bar{\Omega}_3$, the (i_1, i_2) -th element of the above matrix is

$$m_{i_1, i_2}(\mathbf{e}_j) = \Sigma_p''(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) x_{i_1 j} x_{i_2 j}/p^2 + \Sigma_p'(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p)/p.$$

Let $\delta_{i_1, i_2} = 1$ if $i_1 = i_2$, 0 otherwise. We look at Taylor expansions around δ_{i_1, i_2} , and do some arrangement as following: for $i_1, i_2 \in [n]$ and $|\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2}| \leq \delta$,

$$\begin{aligned}\Sigma'_p(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) &= \Sigma'_p(\delta_{i_1, i_2}) + \Sigma''_p(\kappa_{i_1, i_2})(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2}), \\ &= \Sigma'_p(\delta_{i_1, i_2}) + \Sigma''_p(0)(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2}) \\ &\quad + (\Sigma''_p(\kappa_{i_1, i_2}) - \Sigma''_p(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2})(1 - \delta_{i_1, i_2}) \\ &\quad + (\Sigma''_p(\kappa_{i_1, i_2}) - \Sigma''_p(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2})\delta_{i_1, i_2}, \\ \Sigma''_p(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p)x_{i_1j}x_{i_2j} &= \Sigma''_p(0)x_{i_1j}x_{i_2j} + (\Sigma''_p(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) - \Sigma''_p(0))x_{i_1j}x_{i_2j}(1 - \delta_{i_1, i_2}) \\ &\quad + (\Sigma''_p(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) - \Sigma''_p(0))x_{i_1j}x_{i_2j}\delta_{i_1, i_2},\end{aligned}$$

where κ_{i_1, i_2} satisfies $|\kappa_{i_1, i_2} - \delta_{i_1, i_2}| \leq |\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2}|$. From this we have decomposition of $\mathbf{M}(\mathbf{e}_j) := \mathbb{E}_G [\mathbf{G}'_j [\mathbf{G}'_j]^\top]$ into several matrices with small operator norm easy to calculate. With a slight abuse of notations $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{E}, \mathbf{F}, \mathbf{G}, \mathbf{H}$, we have

$$\mathbf{M} = \mathbf{A} + \mathbf{B} + \mathbf{C} + \mathbf{D} + \mathbf{E} + \mathbf{F} + \mathbf{G} + \mathbf{H},$$

where

$$\begin{aligned}\mathbf{A} &= (\Sigma'_p(1) - \Sigma'_p(0)) \mathbf{I}_n/p, \\ \mathbf{B} &= \Sigma'_p(0) \mathbf{1}_n \mathbf{1}_n^\top/p, \\ \mathbf{C} &= \Sigma''_p(0)(\mathbf{X} \mathbf{X}^\top/p)/p, \\ \mathbf{D}_{i_1, i_2} &= (\Sigma''_p(\kappa_{i_1, i_2}) - \Sigma''_p(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) \delta_{i_1 \neq i_2}/p, \\ \mathbf{E}_{i_1, i_2} &= [\Sigma''_p(\kappa_{i_1, i_2}) \mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \Sigma''_p(0) \mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \Sigma''_p(\kappa_{i_1, i_2})] \delta_{i_1 = i_2}/p, \\ \mathbf{F} &= \Sigma''_p(0) \mathbf{X}_j \mathbf{X}_j^\top/p^2, \\ \mathbf{G}_{i_1, i_2}(\mathbf{e}_j) &= (\Sigma''_p(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) - \Sigma''_p(0)) x_{i_1j} x_{i_2j} \delta_{i_1 \neq i_2}/p^2, \\ \mathbf{H}_{i_1, i_2}(\mathbf{e}_j) &= (\Sigma''_p(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) - \Sigma''_p(0)) x_{i_1j} x_{i_2j} \delta_{i_1 = i_2}/p^2.\end{aligned}$$

It suffices to show that $q(\mathbf{N}) := \mathbb{E} [\max_{j \in [p]} \mathbf{v}_n^\top \mathbf{N}(\mathbf{e}_j) \mathbf{v}_n I_{\overline{\Omega}_3}] = o(p)$ for \mathbf{N} being from \mathbf{A} to \mathbf{H} . We notice that $\|\mathbf{I}_n - \mathbf{H}\|_{\text{op}} \leq 1$ implies $|\mathbf{v}_n, i| \leq 1$ for all $i \in [n]$. Then we have

- (i) $q(\mathbf{A}) = o(p)$ provided that $\Sigma'_p(1), \Sigma'_p(0) = o(p)$.
- (ii) $q(\mathbf{B}) = o(p)$ provided that $\Sigma'_p(0) = o(1)$.
- (iii) $q(\mathbf{C}) = o(p)$ provided that $\mathbb{E} [\|\mathbf{X}/\sqrt{p}\|_{\text{op}}] = O(1)$ and $\Sigma''_p(0) = o(p)$.
- (iv) $q(\mathbf{E}) = o(p)$ provided that $\sup_{x \in [1-\delta, 1+\delta]} \Sigma''_p(x), \Sigma''_p(0) = o(p)$.
- (v) $q(\mathbf{F}) = o(p)$ provided that $\Sigma''_p(0) = o(p)$.
- (vi) $q(\mathbf{H}) = o(p)$ provided that $\sup_{x \in [1-\delta, 1+\delta]} \Sigma''_p(x), \Sigma''_p(0) = o(p)$.

We notice that the above are true by assumptions on Σ_p and \mathbf{X} . For \mathbf{D} and \mathbf{G} , we notice the following:

$$|q(\mathbf{D})| := \left| \mathbb{E} \left[\max_{j \in [p]} \mathbf{v}_n^\top \mathbf{D}(\mathbf{e}_j) \mathbf{v}_n I_{\overline{\Omega}_3} \right] \right| \leq \mathbb{E} [|\mathbf{v}_n|^\top |\mathbf{D}| |\mathbf{v}_n| I_{\overline{\Omega}_3}] \leq \mathbf{1}_n^\top \mathbb{E} [|\mathbf{D}| I_{\overline{\Omega}_3}] \mathbf{1}_n,$$

where the absolute value operation is taken element-wise for the vector \mathbf{v}_n and matrix \mathbf{D} . By the Lipschitz assumption of Σ''_p around 0, for $i_1 \neq i_2$,

$$\mathbb{E} [|\mathbf{D}_{i_1, i_2}| I_{\overline{\Omega}_3}] \leq \mathbb{E} [|\Sigma''_p(\kappa_{i_1, i_2}) - \Sigma''_p(0)| |\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p| I_{\overline{\Omega}_3}/p] \leq L_2/p \cdot \mathbb{E} [(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p)^2] = L_2/p^2.$$

This implies $|q(\mathbf{D})| = O(1)$. For $|q(\mathbf{G})|$, we notice that

$$|q(\mathbf{G})| := \left| \mathbb{E} \left[\max_{j \in [p]} \mathbf{v}_n^\top \mathbf{G}(\mathbf{e}_j) \mathbf{v}_n I_{\overline{\Omega}_3} \right] \right| \leq \mathbb{E} \left[\max_{j \in [p]} |\mathbf{v}_n|^\top |\mathbf{G}(\mathbf{e}_j)| |\mathbf{v}_n| I_{\overline{\Omega}_3} \right] \leq \mathbb{E} \left[\max_{j \in [p]} \mathbf{1}_n^\top |\mathbf{G}(\mathbf{e}_j)| \mathbf{1}_n I_{\overline{\Omega}_3} \right]$$

where the absolute value operation is taken element-wise for the vector \mathbf{v}_n and matrix $\mathbf{G}(\mathbf{e}_j)$. By the Lipschitz assumption of Σ_p'' around 0, for $i_1 \neq i_2$,

$$\begin{aligned} \mathbb{E} \left[\max_{j \in [p]} |g_{i_1, i_2}(\mathbf{e}_j)| I_{\bar{\Omega}_3} \right] &\leq L_2 \mathbb{E} \left[\max_{j \in [p]} |(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) x_{i_1 j} x_{i_2 j}| I_{\bar{\Omega}_3} \right] / p^2 \\ &\leq L_2 \delta \cdot \mathbb{E} \left[\max_{j \in [p]} |x_{i_1 j} x_{i_2 j}| I_{\bar{\Omega}_3} \right] / p^2 \\ &\leq L_2 \delta p^{-3/2}, \end{aligned}$$

where we used that $\max_j |x_{i_1 j} x_{i_2 j}| I_{\bar{\Omega}_3} \leq p^{1/2}$. So that $|q(\mathbf{G})| = O(p^{1/2})$. Combining the above we have our proposition. \square

D.1.8. *Proof of Proposition C.2.7.*

Proposition C.2.7. *For all vector $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^p$ such that the derivatives below exist,*

$$\mathbb{E}_G [\partial_{j_1} G(\mathbf{v}_1) \partial_{j_2} G(\mathbf{v}_2)] = \Sigma_p''(\mathbf{v}_1^\top \mathbf{v}_2/p) v_{1j_2} v_{2j_1} / p^2 + \Sigma_p'(\mathbf{v}_1^\top \mathbf{v}_2/p) \delta(j_1 = j_2) / p.$$

Proof. Let $\mathbf{t}_1, \mathbf{t}_2 \in \mathbb{R}^p$. Let

$$\mathcal{G}(\mathbf{t}_1, \mathbf{t}_2) = \mathbb{E}_G [G(\mathbf{v}_1 + \mathbf{t}_1) G(\mathbf{v}_2 + \mathbf{t}_2)] = \Sigma_p((\mathbf{v}_1 + \mathbf{t}_1)^\top (\mathbf{v}_2 + \mathbf{t}_2) / p).$$

Let $j_1, j_2 \in [p]$. Let us assume that the derivatives below exists,

$$\begin{aligned} \partial_{t_{2j_2}} \partial_{t_{1j_1}} \mathcal{G}(\mathbf{t}_1, \mathbf{t}_2) &= \mathbb{E}_G [\partial_{j_1} G(\mathbf{v}_1 + \mathbf{t}_1) \partial_{j_2} G(\mathbf{v}_2 + \mathbf{t}_2)], \\ \partial_{t_{2j_2}} \partial_{t_{1j_1}} \mathcal{G}(\mathbf{t}_1, \mathbf{t}_2) &= \Sigma_p''((\mathbf{v}_1 + \mathbf{t}_1)^\top (\mathbf{v}_2 + \mathbf{t}_2) / p) (v_{1j_2} + t_{1j_2}) (v_{2j_1} + t_{2j_1}) / p^2 \\ &\quad + \Sigma_p'((\mathbf{v}_1 + \mathbf{t}_1)^\top (\mathbf{v}_2 + \mathbf{t}_2) / p) \delta(j_1 = j_2) / p. \end{aligned}$$

Then

$$\begin{aligned} \mathbb{E}_G [\partial_{j_1} G(\mathbf{v}_1) \partial_{j_2} G(\mathbf{v}_2)] &= \partial_{t_{2j_2}} \partial_{t_{1j_1}} \mathcal{G}(\mathbf{0}_p, \mathbf{0}_p) \\ &= \Sigma_p''(\mathbf{v}_1^\top \mathbf{v}_2/p) v_{1j_2} v_{2j_1} / p^2 + \Sigma_p'(\mathbf{v}_1^\top \mathbf{v}_2/p) \delta(j_1 = j_2) / p. \end{aligned}$$

\square

D.2. **Proofs of the limits $V(\psi_d, \psi_p, \lambda, \rho, \sigma)$ and $R(\psi_p, \lambda, \rho, \sigma)$.**

D.2.1. *Proof of Proposition C.3.1.*

Proposition C.3.1. *Let $L^2 := \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 / (\text{trace}(\mathbf{I}_n - \mathbf{H}))^2$. Under the model (5) in the asymptotic setting (6), Assumption C.3.1 and Definition C.3.1,*

$$\text{Var}(y_1) \rightarrow \theta_{\boldsymbol{\beta}}^2 + \theta_{\boldsymbol{\epsilon}}^2 + \theta_{\text{NL}}^2, \quad nL^2 / \text{Var}(y_1) \xrightarrow{\mathbb{P}_{\mathbf{X}, \mathbf{W}}^{G, \epsilon}} V.$$

Proof. In this proof we consider the model in Mei and Montanari [2019]. Under this model, the limit of $(1/n) \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2$ is given in Theorem 6 in Mei and Montanari [2019] as

$$\lim_{n \rightarrow +\infty} \mathbb{E} \left| (1/n) \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 - (\theta_{\boldsymbol{\beta}}^2 + \theta_{\boldsymbol{\epsilon}}^2 + \theta_{\text{NL}}^2) \left(\mathcal{L} - \frac{\psi_d \lambda}{\psi_p \mu_*^2} \mathcal{A} \right) \right| = 0.$$

The limit of $\text{Var}(y_1)$ follows by

$$\begin{aligned} \text{Var}(y_1) &= \mathbb{E} [y_1^2] - \beta_0^2 \\ &= \mathbb{E} [(\beta_0 + \mathbf{x}_1^\top \boldsymbol{\beta} + G(\mathbf{x}_1) + \varepsilon_1)^2] - \beta_0^2 \\ &= \mathbb{E} [(\mathbf{x}_1^\top \boldsymbol{\beta})^2] + \mathbb{E} [(G(\mathbf{x}_1))^2] + \mathbb{E} [\varepsilon_1^2] \\ &\rightarrow \theta_{\boldsymbol{\beta}}^2 + \theta_{\text{NL}}^2 + \theta_{\boldsymbol{\epsilon}}^2. \end{aligned}$$

The limit of $(1/n) \text{trace}(\mathbf{I}_n - \mathbf{H})$ is implied in [Mei and Montanari \[2019\]](#) by the following facts. Let us first denote $\mathbf{Z} = (1/\sqrt{p})\sigma(\mathbf{X}\mathbf{W}^\top)$. Let

$$\overline{\mathbf{Z}}(t) = (1+t) \begin{bmatrix} 0 & \mathbf{Z}^\top \\ \mathbf{Z} & 0 \end{bmatrix} \in \mathbb{R}^{(n+d) \times (n+d)}.$$

Let $u = (\psi_1\psi_2\lambda)^{1/2} \in \mathbb{R}^+$ and let \log denote the complex logarithm with branch cut on the negative real axis. Let $\lambda_i(\overline{\mathbf{Z}}(t))$ be the eigenvalues of $\overline{\mathbf{Z}}(t)$ in non-increasing order. Let

$$\mathcal{J}(u, t) = (1/p) \sum_{i \in [n+d]} \log(\lambda_i(\overline{\mathbf{Z}}(t)) - iu).$$

From Proposition 7.3 in [Mei and Montanari \[2019\]](#),

$$\partial_t \mathcal{J}(u, 0) = \frac{2}{p} \text{trace}((u^2 \mathbf{I}_d + \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Z}) = \frac{2}{p} \text{trace}(\mathbf{H}).$$

From the fact that $\overline{\mathbf{Z}}(t) = (1+t)\overline{\mathbf{Z}}(0)$, we have from the chain rule and the definition of \mathcal{J} ,

$$\partial_t \mathcal{J}(u, 0) = (1/p) \sum_{i \in [n+d]} (\lambda_i(\overline{\mathbf{Z}}(0)) - iu)^{-1} \cdot (iu) + \psi_1 + \psi_2.$$

From Proposition 7.2 and Step 2 in Lemma C.1. in [Mei and Montanari \[2019\]](#),

$$\lim_{n \rightarrow +\infty} \mathbb{E} \left| (1/p) \sum_{i \in [n+d]} (\lambda_i(\overline{\mathbf{Z}}(0)) - iu)^{-1} - (\nu_1 + \nu_2)/\mu_* \right| = 0.$$

So that we have

$$\lim_{n \rightarrow +\infty} \mathbb{E} |(1/n) \text{trace}(\mathbf{I}_n - \mathbf{H}) - (1/2)(1 - \psi_p((\nu_1 + \nu_2)/\mu_*)u\mathbf{i} - \psi_d)| = 0.$$

We notice by the definition of u and by (C.6),

$$\begin{aligned} 1 - \psi_p((\nu_1 + \nu_2)/\mu_*)u\mathbf{i} - \psi_d &= 1 + \psi_p(\nu_1 + \nu_2)(-\bar{z}) - \psi_d \\ &= 1 + \psi_p \left(\psi_1 + \psi_2 + 2\chi + \frac{2\varrho^2\chi}{1 - \varrho^2\chi} \right) - \psi_d \\ &= 2 + \psi_p(2\chi + \frac{2\varrho^2\chi}{1 - \varrho^2\chi}). \end{aligned}$$

This implies

$$\lim_{n \rightarrow +\infty} \mathbb{E} |(1/n) \text{trace}(\mathbf{I}_n - \mathbf{H}) - \mathcal{Q}| = 0.$$

Combining the above, we have the limit of $nL^2/\text{Var}(y_1)$. □

D.2.2. Sketch of Proof of Proposition C.3.2.

Proposition C.3.2. *Under the pure linear model (4) in the asymptotic setting (6), Assumption 1 and 2 with (ii) and Definition C.3.1,*

$$\text{Var}(y_1) \rightarrow \theta_\beta^2 + \theta_\varepsilon^2, \quad nL^2/\text{Var}(y_1) \xrightarrow{\mathbb{P}^{\mathbf{x}, \mathbf{w}, \varepsilon}} V.$$

Proof. In this sketch of proof, we consider the setting for the pure linear model,

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon,$$

with Gaussian \mathbf{x}_i and \mathbf{w}_k under Assumptions 1 and 2 with (ii). The limit of $(1/n) \text{trace}(\mathbf{I}_n - \mathbf{H})$ satisfies

$$\lim_{n \rightarrow +\infty} \mathbb{E} |(1/n) \text{trace}(\mathbf{I}_n - \mathbf{H}) - \mathcal{Q}| = 0,$$

under the same reasoning in Section D.2.1, which also holds for $\mathbf{x}_i, \mathbf{w}_k$ being Gaussian.

So it suffices to show that the limit of $(1/n) \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2$ is the same as that in Section D.2.1, that is,

$$\lim_{n \rightarrow +\infty} \mathbb{E} \left| (1/n) \|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2 - (\theta_{\boldsymbol{\beta}}^2 + \theta_{\boldsymbol{\varepsilon}}^2) \left(\mathcal{L} - \frac{\psi_d \lambda}{\psi_p \mu_*^2} \mathcal{A} \right) \right| = 0.$$

First, by Remark 8 in Mei and Montanari [2019], we can consider $\boldsymbol{\beta} \in \mathbb{S}^{p-1}(\|\boldsymbol{\beta}\|_2^2)$ independent of $\mathbf{X}, \boldsymbol{\varepsilon}, \mathbf{W}$, instead of $\boldsymbol{\beta}$ being deterministic, since the training error $\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2$ as a function of $\boldsymbol{\beta}$ is invariant in distribution after orthogonal rotation of $\boldsymbol{\beta}$. So letting \mathbb{E} denote $\mathbb{E}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, \boldsymbol{\beta} \sim \text{Unif}(\mathbb{S}^{p-1}(\|\boldsymbol{\beta}\|_2^2))}$, by the fact that $\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}} = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$, we have

$$\begin{aligned} (1/n) \mathbb{E} [\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2] &= (1/n) \mathbb{E} [\|(\mathbf{I}_n - \mathbf{H})\mathbf{y}\|_2^2] \\ &= (1/n) \mathbb{E} [\|(\mathbf{I}_n - \mathbf{H})(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})\|_2^2] \\ &= (1/n) \left[\mathbb{E} [\text{trace}(\widetilde{\mathbf{H}}\mathbf{X}\boldsymbol{\beta}\boldsymbol{\beta}^\top \mathbf{X}^\top)] + \mathbb{E} [\text{trace}(\widetilde{\mathbf{H}}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top)] \right] \\ &= (1/n) \left[\theta_{\boldsymbol{\beta}}^2 \mathbb{E} [\text{trace}(\widetilde{\mathbf{H}}\mathbf{X}\mathbf{X}^\top/p)] + \theta_{\boldsymbol{\varepsilon}}^2 \mathbb{E} [\text{trace}(\widetilde{\mathbf{H}})] \right] \end{aligned}$$

where

$$\widetilde{\mathbf{H}} := (\mathbf{I}_n - \mathbf{H})^2 = (n\tau)^2 (\mathbf{A}\mathbf{A}^\top + n\tau\mathbf{I}_n)^{-2}.$$

We notice that the expected value of traces above can be calculated by the quantities in (228) in Mei and Montanari [2019]. We also claim that the calculations of the traces in the (228) are the same for the random vectors on the spheres and the Gaussian random vectors, cf. Section C in the Appendix of Mei and Montanari [2019], especially Lemma C.1, Proposition 7.2 and Lemma C.7 there. So that $\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2^2$ has the same limit as that in Theorem 6 in Mei and Montanari [2019]. \square

D.2.3. Proof of Proposition C.3.3.

Proposition C.3.3. *Let V, R be in Definitions C.3.1 and C.3.2. Then*

$$\lim_{\psi_d \rightarrow +\infty} V(\psi_d, \psi_p, \lambda, \rho, \sigma) = R(\psi_p, \lambda, \rho, \sigma).$$

Proof. In this proof we let \rightarrow denote the convergence when $\psi_d \rightarrow +\infty$. We recall all the notations in Definition C.3.1 and C.3.2. We notice that $|\chi| \leq |\nu_1| \cdot |\nu_2| \leq 1/\bar{\lambda}^2$ is always bounded when $\psi_d \rightarrow +\infty$. From (C.6),

$$\begin{aligned} \nu_1(-\bar{z}) &= \psi_1 + \chi + \frac{\varrho^2 \chi}{1 - \varrho^2 \chi}, \\ \nu_2(-\bar{z}) &= \psi_2 + \chi + \frac{\varrho^2 \chi}{1 - \varrho^2 \chi}. \end{aligned} \tag{D.2}$$

The quantities ν_1, ν_2 are purely imaginary with positive imaginary part, so that χ is a negative real. In fact,

$$\nu_1 - \nu_2 = (\psi_1 - \psi_2)/(-\bar{z}) = \frac{\psi_1 - \psi_2}{(\psi_1 \psi_2 \lambda)^{1/2}} \mu_* \mathbf{i}$$

is purely imaginary. We can specify the real and the imaginary parts of ν_1, ν_2 ,

$$\nu_1 = a + (b + c)\mathbf{i}, \quad \nu_2 = a + b\mathbf{i}, \quad c = \frac{\psi_1 - \psi_2}{(\psi_1 \psi_2 \lambda)^{1/2}} \mu_*, \quad \chi = (a^2 - b^2 - bc) + (2ab + ac)\mathbf{i}.$$

From the fact that $\mathcal{L}, \mathcal{Q}, \mathcal{A}$ are real numbers, we deduce that χ is real number. So that $a(2b + c) = 0$. From the fact that $\nu_1, \nu_2 \in \mathbb{C}_+$, we have $2b + c > 0$ so that $a = 0$. This shows that ν_1, ν_2 are purely

imaginary numbers and $\chi < 0$. Next, from (D.2), we have

$$\begin{aligned}\chi(-\psi_1\psi_2\bar{\lambda}) &= \left(\psi_1 + \chi + \frac{\varrho^2\chi}{1-\varrho^2\chi}\right) \left(\psi_2 + \chi + \frac{\varrho^2\chi}{1-\varrho^2\chi}\right), \\ \implies \psi_1 &= \frac{-\left(\chi + \frac{\varrho^2\chi}{1-\varrho^2\chi}\right) \left(\psi_2 + \chi + \frac{\varrho^2\chi}{1-\varrho^2\chi}\right)}{\psi_2 + \chi + \frac{\varrho^2\chi}{1-\varrho^2\chi} + \chi\psi_2\bar{\lambda}}.\end{aligned}$$

When $\psi_1 \rightarrow +\infty$ and $\varrho \neq 0$ and $\psi_2 \neq 0$ fixed and χ is bounded, negative, we have

$$\begin{aligned}\psi_2 + \chi + \frac{\varrho^2\chi}{1-\varrho^2\chi} + \chi\psi_2\bar{\lambda} &\rightarrow 0, \\ \implies \chi^2 + \left(\frac{\psi_2-1}{1+\psi_2\bar{\lambda}} - \varrho^{-2}\right) \chi - \frac{\psi_2}{1+\psi_2\bar{\lambda}} \varrho^{-2} &\rightarrow 0, \\ \implies \chi \rightarrow \bar{\chi} &:= \frac{\left(\varrho^{-2} - \frac{\psi_2-1}{1+\psi_2\bar{\lambda}}\right) - \sqrt{\left(\varrho^{-2} - \frac{\psi_2-1}{1+\psi_2\bar{\lambda}}\right)^2 + 4\frac{\psi_2}{1+\psi_2\bar{\lambda}}\varrho^{-2}}}{2}.\end{aligned}$$

When $\varrho = 0$, we define $\bar{\chi} := -\frac{\psi_2}{1+\psi_2\bar{\lambda}}$ and we have $\chi \rightarrow \bar{\chi}$ still holds.

The quantities have limits

$$\begin{aligned}\mathcal{Q} \rightarrow \bar{\mathcal{Q}} &:= \psi_2^{-1}(-\bar{\chi}\psi_2\bar{\lambda}) = -\bar{\chi}\bar{\lambda}, \\ \mathcal{L} \rightarrow \bar{\mathcal{L}} &:= (-\bar{\chi}\bar{\lambda}) \left[\frac{\rho}{1+\rho} \frac{1}{1-\bar{\chi}\varrho^2} + \frac{1}{1+\rho} \right], \\ \mathcal{A}_1 \rightarrow \bar{\mathcal{A}}_1 &:= \frac{\rho}{1+\rho} [-\bar{\chi}^2 (\bar{\chi}\varrho^4 - \bar{\chi}\varrho^2 + \psi_2\varrho^2 + \varrho^2 - \bar{\chi}\psi_2\varrho^4 + 1)] \\ &\quad + \frac{1}{1+\rho} [\bar{\chi}^2 (\bar{\chi}\varrho^2 - 1) (\bar{\chi}^2\varrho^4 - 2\bar{\chi}\varrho^2 + \varrho^2 + 1)], \\ \mathcal{A}_0/\psi_1 \rightarrow \bar{\mathcal{A}}_* &:= (\psi_2 - 1) \bar{\chi}^3 \varrho^6 + (1 - 3\psi_2) \bar{\chi}^2 \varrho^4 + 3\psi_2 \bar{\chi} \varrho^2 - \psi_2.\end{aligned}$$

So that

$$V = (\mathcal{L} - \psi_1\bar{\lambda}\mathcal{A})/\mathcal{Q}^2 \rightarrow R = (\bar{\mathcal{L}} - \bar{\lambda}\bar{\mathcal{A}}_1/\bar{\mathcal{A}}_*)/\bar{\mathcal{Q}}^2.$$

Furthermore, by simple algebra, we have that when $\varrho = 0$, $(\bar{\mathcal{L}} - \bar{\lambda}\bar{\mathcal{A}}_1/\bar{\mathcal{A}}_*)/\bar{\mathcal{Q}}^2 = 1$. \square

D.3. Proof of Theorem C.4.1.

Theorem C.4.1. *Let $t \in \mathbb{R}$. Under model (5), Assumption 1, 2, 3 and 4, Definition 1 and a further assumption that $\Sigma'_p(0) = O(1/p)$, we have*

$$(C.7) \quad \sup_{\mathbf{u}_0 \in S_p} \left| \mathbb{P}_{\mathbf{X}, \mathbf{W}, \varepsilon, G | \mathbf{u}_0} \left(\frac{\zeta_L(\mathbf{u}_0)}{\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2} \leq t \right) - \Phi(t) \right| \rightarrow 0$$

for some $S_p \subset \mathbb{S}^{p-1}(1)$ satisfying $|S_p|/|\mathbb{S}^{p-1}(1)| \geq 1 - \log(p)/p \rightarrow 1$.

D.3.1. Proof of Theorem C.4.1.

Proof. Let ζ be as in Definition C.2.1. Provided with Propositions C.4.1 and C.4.2, we will show that, for a large proportion of $\mathbf{u}_0 \in \mathbb{S}^{p-1}$, given \mathbf{u}_0 ,

- (i) $\frac{\zeta(\mathbf{u}_0)}{(\text{Var}_0[\zeta(\mathbf{u}_0)])^{1/2}} \xrightarrow{d} N(0, 1)$.
- (ii) $\frac{\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2}{(\text{Var}_0[\zeta(\mathbf{u}_0)])^{1/2}} \xrightarrow{\mathbb{P}} 1$.
- (iii) $\frac{\zeta(\mathbf{u}_0) - \zeta_L(\mathbf{u}_0)}{\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2} \xrightarrow{\mathbb{P}} 0$.

The above convergence are uniform over $\mathbf{u}_0 \in S_p$, where S_p is a large subset of $\mathbb{S}^{p-1}(1)$. So by Slutsky's Theorem we have our proposition.

We specify S_p as follows. Let $\mathbf{v}_n = \mathbf{1}_n - (H_{11}, H_{22}, \dots, H_{nn})^\top$, $\mathbf{r} := \mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}$ and $q_0 := \text{trace}(\mathbf{T}_{\text{NL}}(\mathbf{u}_0)) = \text{trace}[(\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}'\mathbf{u}_0)] = \mathbf{v}_n^\top \mathbf{G}'\mathbf{u}_0$. Notice that (iii) in Proposition C.4.1 and Proposition C.4.2 provide the existence of a constant $c_9 > 0$ independent of n, p, d such that

$$\max \left(\mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G} [\epsilon_0^2], \mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G} \left[q_0^2 / \|\mathbf{r}\|_2^2 I_{\Omega_3} \right] \right) \leq c_9/p$$

We specify the large volume index set $S_p \subset \mathbb{S}^{p-1}(1)$ as

$$S_p := \left\{ \mathbf{u}_0 \in \mathbb{S}^{p-1} : \max \left(\mathbb{E}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} [\epsilon_0^2], \mathbb{E}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} \left[q_0^2 / \|\mathbf{r}\|_2^2 I_{\Omega_3} \right] \right) \leq \frac{2c_9}{\log(p)} \right\}.$$

Since

$$\begin{aligned} \mathbb{P}_{\mathbf{u}_0}(S_p^c) &= \mathbb{P}_{\mathbf{u}_0} \left(\mathbb{E}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} [\epsilon_0^2] \geq \frac{2c_9}{\log(p)} \text{ or } \mathbb{E}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} \left[q_0^2 / \|\mathbf{r}\|_2^2 I_{\Omega_3} \right] \geq \frac{2c_9}{\log(p)} \right) \\ &\leq \frac{\mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G} [\epsilon_0^2] + \mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G} \left[q_0^2 / \|\mathbf{r}\|_2^2 I_{\Omega_3} \right]}{2c_9} \log(p) \leq \log(p)/p, \end{aligned}$$

the relative volume $|S_p|/|\mathbb{S}^{p-1}(1)| \geq 1 - \log(p)/p \rightarrow 1$ as $p \rightarrow +\infty$.

We notice that (i) can be directly obtained from Proposition C.4.1 (i), by noticing that $(\mathbf{X}\mathbf{u}_0)^\top \mathbb{E}_0[\mathbf{r}] / \|\mathbb{E}_0[\mathbf{r}]\|_2 \sim N(0, 1)$.

For (iii), we notice that for any $\epsilon > 0$,

$$\begin{aligned} \sup_{\mathbf{u}_0 \in S_p} \mathbb{P}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} \left(\left| \frac{\zeta(\mathbf{u}_0) - \zeta_L(\mathbf{u}_0)}{\|\mathbf{y} - \mathbf{A}\hat{\boldsymbol{\alpha}}\|_2} \right| > \epsilon \right) &= \sup_{\mathbf{u}_0 \in S_p} \mathbb{P}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} \left(\left| \frac{q_0}{\|\mathbf{r}\|_2} \right| > \epsilon \right) \\ &\leq \sup_{\mathbf{u}_0 \in S_p} \mathbb{P}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} \left(\left| \frac{q_0}{\|\mathbf{r}\|_2} \right| I_{\Omega_3} > \epsilon \right) \\ &\quad + \mathbb{P}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G}(I_{\Omega_3}^c) \\ &\leq \frac{2c_9}{\log(p)\epsilon^2} + o(-\exp(-\min(c_6, c_{11})n)) \\ &= o(1). \end{aligned}$$

For (ii), we first recall Proposition C.4.1 (ii). Let $\epsilon > 0$ be a fixed number. From the definition of S_p combined with Chebyshev's inequality, we can see that for all $\bar{\epsilon} > 0$,

$$\sup_{\mathbf{u}_0 \in S_p} \mathbb{P}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} (|\epsilon_0| > \bar{\epsilon}) \leq \frac{2c_9}{\bar{\epsilon}^2 \log(p)}.$$

Let $\mathbf{u}_0 \in S_p$. Letting $\mathbf{r} = \mathbf{A}\hat{\boldsymbol{\alpha}} - \mathbf{y}$, $V_0 := \text{Var}_0[\zeta(\mathbf{u}_0)]$ and $U_0 := \left| \|\mathbf{r}\|_2 / V_0^{1/2} - 1 \right|$, we have that if $\epsilon_0 \leq \frac{1}{2}$,

$$\mathbb{E}_0[U_0] \leq (1 + \sqrt{2}) \epsilon_0 / (1 - 2\epsilon_0^2)^{1/2} \leq (2 + \sqrt{2})\epsilon_0.$$

Now let us focus on $\bar{\epsilon} < 1/2$. We let $\Omega_0(\bar{\epsilon}) := \{\mathbb{E}_0[U_0] < (2 + \sqrt{2})\bar{\epsilon}\}$. Then

$$\mathbb{P}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0}(\Omega_0(\bar{\epsilon})) = \mathbb{P}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} \left(\mathbb{E}_0[U_0] < (2 + \sqrt{2})\bar{\epsilon} \right) \geq \mathbb{P}_{\mathbf{X}, \mathbf{W}, \boldsymbol{\varepsilon}, G | \mathbf{u}_0} (|\epsilon_0| \leq \bar{\epsilon}) \geq 1 - \frac{2c_9}{\bar{\epsilon}^2 \log(p)}.$$

Then, letting $\mathbb{I}(\cdot) := I_{\{\cdot\}}$ be the indicator function, we have

$$\begin{aligned}
\mathbb{P}_{\mathbf{X}, \mathbf{W}, \varepsilon, G | \mathbf{u}_0}(U_0 > \epsilon) &:= \mathbb{E}_{\mathbf{X}, \mathbf{W}, \varepsilon, G | \mathbf{u}_0} [\mathbb{I}(U_0 > \epsilon)] \\
&= \mathbb{E}_{\mathbf{X}, \mathbf{W}, \varepsilon, G | \mathbf{u}_0} [\mathbb{E}_0 [\mathbb{I}(U_0 > \epsilon)] I_{\Omega_0(\bar{\epsilon})}] + \mathbb{E}_{\mathbf{X}, \mathbf{W}, \varepsilon, G | \mathbf{u}_0} [\mathbb{E}_0 [\mathbb{I}(U_0 > \epsilon)] I_{\Omega_0^c(\bar{\epsilon})}] \\
&\leq \mathbb{E}_{\mathbf{X}, \mathbf{W}, \varepsilon, G | \mathbf{u}_0} \left[\frac{\mathbb{E}_0 [U_0]}{\epsilon} I_{\Omega_0(\bar{\epsilon})} \right] + \mathbb{P}_{\mathbf{X}, \mathbf{W}, \varepsilon, G | \mathbf{u}_0}(\Omega_0^c(\bar{\epsilon})) \\
&\leq \left(2 + \sqrt{2}\right) \bar{\epsilon} / \epsilon + \frac{2c_9}{\bar{\epsilon}^2 \log(p)}.
\end{aligned}$$

Choosing $\bar{\epsilon} := \min \left(\frac{1}{\log \log(p)}, \frac{1}{2} \right)$, we have that, for all $\epsilon > 0$,

$$\lim_{p \rightarrow +\infty} \sup_{\mathbf{u}_0 \in S_p} \mathbb{P}_{\mathbf{X}, \mathbf{W}, \varepsilon, G | \mathbf{u}_0}(U_0 > \epsilon) = 0.$$

Thus we have (ii). □

D.3.2. Proof of Proposition C.4.1.

Proposition C.4.1. *Let*

- (i) $\mathbf{u}_0 \sim \text{Unif}(\mathbb{S}^{p-1}(1))$ independent with $\mathbf{X}, \mathbf{W}, G, \varepsilon$.
- (ii) $\mathbf{X}_0 := \mathbf{X} \mathbf{u}_0$.
- (iii) $\zeta(\mathbf{u}_0) := \mathbf{u}_0^\top (\zeta(\mathbf{e}_j))^{j \in [p]}$ where $\zeta(\mathbf{e}_j)$ is defined in Definition C.2.1.
- (iv) $\mathbf{r} = \mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}$.
- (v) $\epsilon_0^2 := \mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2] / (\mathbb{E}_0 [\|\mathbf{r}\|_2^2] + \mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2])$.

Then

- (i) $\mathbb{E}_0 \left[\left(\frac{\zeta(\mathbf{u}_0)}{(\text{Var}_0 \zeta(\mathbf{u}_0))^{1/2}} - \frac{\mathbf{X}_0^\top \mathbb{E}_0[\mathbf{r}]}{\|\mathbb{E}_0[\mathbf{r}]\|_2} \right)^2 \right] \leq 6\epsilon_0^2$.
- (ii) $\mathbb{E}_0 \left[\left| \frac{\|\mathbf{r}\|_2}{(\text{Var}_0[\zeta(\mathbf{u}_0)])^{1/2}} - 1 \right| \right] \leq (1 + \sqrt{2}) \epsilon_0 / (1 - 2\epsilon_0^2)_+^{1/2}$.
- (iii) $\mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \varepsilon, G} [\epsilon_0^2] = O(1/p)$.

Proof. (i) The proof is the same as that of Proposition C.2.5 (i).

(ii) The proof is the same as that of Proposition C.2.5 (ii).

(iii) For a general vector \mathbf{u}_0 satisfying $\|\mathbf{u}_0\|_2 = 1$, the propositions discussed before for canonical basis \mathbf{e}_j can be updated as follows.

Step 1. We show that Proposition C.2.2 (iii) and (iv) can be replaced with

$$p\mathbb{E}_{\mathbf{u}_0} \|\mathbf{T}_0(\mathbf{u}_0) + \mathbf{T}_L(\mathbf{u}_0) + \mathbf{T}_{NL}(\mathbf{u}_0)\|_F^2 \leq 2c_1^2 L^2 n \|\hat{\boldsymbol{\alpha}}\|_2^2 + 2\|\mathbf{f}'\|_F^2.$$

and

$$\|\mathbf{T}_1(\mathbf{u}_0)\|_F^2 \leq L^2 c_1^2 / (4n\tau) \cdot \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2^2.$$

We will use the fact that $\mathbb{E}_{\mathbf{u}_0} \mathbf{u}_0 \mathbf{u}_0^\top = (1/p) \mathbf{I}_p$ and $\|\mathbf{u}_0\|_2 = 1$.

$$\begin{aligned}
p\mathbb{E}_{\mathbf{u}_0} \|\mathbf{T}_0(\mathbf{u}_0) + \mathbf{T}_L(\mathbf{u}_0) + \mathbf{T}_{NL}(\mathbf{u}_0)\|_F^2 &\leq p\mathbb{E}_{\mathbf{u}_0} \|\mathbf{I}_n - \mathbf{H}\|_{\text{op}}^2 \left\| \left(\sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right) \mathbf{u}_0 \right\|_2^2 \\
&\leq p\mathbb{E}_{\mathbf{u}_0} \left\| \left(\sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right) \mathbf{u}_0 \right\|_2^2 \\
&= \left\| \sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right\|_F^2 \\
&\leq 2 \left\| \sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} \right\|_F^2 + 2 \|\mathbf{f}'\|_F^2 \\
&\leq 2 \|\mathbf{W}\|_{\text{op}}^2 \left\| \sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \right\|_F^2 + 2 \|\mathbf{f}'\|_F^2 \\
&\leq 2c_1^2 L^2 n \|\hat{\boldsymbol{\alpha}}\|_2^2 + 2 \|\mathbf{f}'\|_F^2.
\end{aligned}$$

We used $\|\mathbf{I}_n - \mathbf{H}\|_{\text{op}} \leq 1$ in the second inequality in the above display. In the equality above, We used the fact that $\mathbb{E}_{\mathbf{u}_0} \mathbf{u}_0 \mathbf{u}_0^\top = (1/p) \mathbf{I}_p$ and that

$$\begin{aligned} & \mathbb{E}_{\mathbf{u}_0} \left\| \left(\sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right) \mathbf{u}_0 \right\|_2^2 \\ &= \mathbb{E}_{\mathbf{u}_0} \text{trace} \left(\left(\sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right) \mathbf{u}_0 \mathbf{u}_0^\top \left(\sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right)^\top \right) \\ &= \text{trace} \left(\left(\sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right) \mathbb{E}_{\mathbf{u}_0} [\mathbf{u}_0 \mathbf{u}_0^\top] \left(\sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right)^\top \right) \\ &= (1/p) \left\| \sigma'(\mathbf{X} \mathbf{W}^\top) \text{diag}(\hat{\boldsymbol{\alpha}}) \mathbf{W} - \mathbf{f}' \right\|_F^2. \end{aligned}$$

Noticing that $\|\mathbf{u}_0\|_2 = 1$, we have

$$\begin{aligned} \|\mathbf{T}_1(\mathbf{u}_0)\|_F^2 &\leq \|\mathbf{A}(n\tau \mathbf{I}_d + \mathbf{A}^\top \mathbf{A})^{-1}\|_{\text{op}}^2 \cdot \|\text{diag}(\mathbf{W} \mathbf{u}_0) \sigma'(\mathbf{W} \mathbf{X}^\top) \text{diag}(\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}})\|_F^2 \\ &\leq 1/(4n\tau) \cdot L^2 \cdot \|(\mathbf{W} \mathbf{u}_0)(\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}})^\top\|_F^2 \\ &= L^2/(4n\tau) \cdot \|\mathbf{W} \mathbf{u}_0\|_2^2 \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2^2 \\ &\leq L^2/(4n\tau) \cdot \|\mathbf{W}\|_{\text{op}}^2 \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2^2 \\ &\leq L^2 c_1^2/(4n\tau) \cdot \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2^2. \end{aligned}$$

Step 2. By the proof of Proposition C.2.3, there exists a class of large events $\{\Omega(\mathbf{u}_0)\}_{\mathbf{u}_0 \in \mathbb{S}^{p-1}(1)}$ such that $\mathbb{E}_{\mathbf{X}, \boldsymbol{\varepsilon} | \mathbf{W}, G, \mathbf{u}_0} (I_{\Omega(\mathbf{u}_0)}) \geq 1 - o(\exp(-c_{12}n))$ for some constant $c_{12} > 0$ independent of \mathbf{u}_0 and that on $\Omega(\mathbf{u}_0)$,

$$\mathbb{E}_0 \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2^2 \geq n \cdot c_{2,n}.$$

The starting point of the construction is as follows: Since the mapping $\boldsymbol{\varepsilon} \mapsto \mathbb{E}_0 \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2$ is 1-Lipschitz, by Theorem 5.2.2 in Vershynin [2018], for some universal constant $c_5 > 0$ and for all $t > 0$, there exists an event $\Omega_2(\mathbf{u}_0, t)$ such that:

- (i) On $\Omega_2(\mathbf{u}_0, t)$, $\mathbb{E}_0 \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2 \geq \mathbb{E}_\varepsilon \mathbb{E}_0 \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2 - \sqrt{n} \theta_\varepsilon t$
- (ii) $\mathbb{P}(\Omega_2(\mathbf{u}_0, t)) \geq 1 - 2 \exp(-c_5 n t^2)$.

By the reasoning in the proof of Proposition C.2.3,

$$\mathbb{E}_\varepsilon \mathbb{E}_0 \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2 - \sqrt{n} \theta_\varepsilon t \geq \theta_\varepsilon [\mathbb{E}_0 \|\mathbf{I}_n - \mathbf{H}\|_F - \sqrt{nt} - 1].$$

From Proposition C.2.4 and its proof,

$$\mathbb{E}_0 \|\mathbf{I}_n - \mathbf{H}\|_F / \sqrt{n} \geq \mathbb{E}_0 (1 + F_n/\tau)^{-1} \geq (1 + [2c_1^2 L^2 \mathbb{E}_0 \|\mathbf{X}\|_F^2 / n^2 + 2\psi_{d,n}(\sigma(0))^2] / \tau)^{-1}.$$

Since $\mathbf{I}_n = \mathbf{u}_0 \mathbf{u}_0^\top + \mathbf{Q}_0$, $\mathbf{u}_0^\top = \mathbf{u}_0^\top + \mathbf{u}_0^\top \mathbf{Q}_0$, $\mathbf{u}_0^\top \mathbf{Q}_0 = \mathbf{0}_n^\top$.

$$\begin{aligned} \mathbb{E}_0 \|\mathbf{X}\|_F^2 &= \mathbb{E}_0 \|\mathbf{X}_0 \mathbf{u}_0^\top + \mathbf{X} \mathbf{Q}_0\|_F^2 \\ &= \mathbb{E}_0 [\|\mathbf{X}_0 \mathbf{u}_0^\top\|_F^2] + \|\mathbf{X} \mathbf{Q}_0\|_F^2 + 2\mathbb{E}_0 \text{trace}(\mathbf{X}_0 \mathbf{u}_0^\top \mathbf{Q}_0^\top \mathbf{X}^\top) \\ &= n + \|\mathbf{X} \mathbf{Q}_0\|_F^2 + 0 \leq n + \|\mathbf{X}\|_F^2. \end{aligned}$$

So that we have

$$\mathbb{E}_0 \|\mathbf{I}_n - \mathbf{H}\|_F / \sqrt{n} \geq (1 + \bar{F}_n/\tau)^{-1}$$

where \bar{F}_n is as given in Proposition C.2.4. Following the rest of the proof of Proposition C.2.3, we will have a desired large event $\Omega(\mathbf{u}_0)$ for $\mathbf{u}_0 \in \mathbb{S}^{p-1}(1)$.

We notice that, letting $\mathbb{E} := \mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \varepsilon, G}$, we have

$$\begin{aligned}
\mathbb{E} [\epsilon_0^2] &= \mathbb{E} \left[\frac{\mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2]}{\mathbb{E}_0 [\|\mathbf{r}\|_2^2] + \mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2]} \right] \\
&\leq \mathbb{E} \left[\frac{\mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2]}{\mathbb{E}_0 [\|\mathbf{r}\|_2^2] + \mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2]} I_{\Omega(\mathbf{u}_0)} \right] + \mathbb{P}(\Omega(\mathbf{u}_0)^c) \\
&\leq \mathbb{E} \left[\frac{\mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2]}{\mathbb{E}_0 [\|\mathbf{r}\|_2^2] + \mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2]} I_{\Omega(\mathbf{u}_0)} \right] + o(\exp(-c_{12}n)) \\
&\leq \mathbb{E} \left[\frac{\mathbb{E}_0 [2 \|\mathbf{T}_0(\mathbf{u}_0) + \mathbf{T}_L(\mathbf{u}_0) + \mathbf{T}_{NL}(\mathbf{u}_0)\|_F^2]}{\mathbb{E}_0 [\|\mathbf{r}\|_2^2] + \mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2]} I_{\Omega(\mathbf{u}_0)} \right] \\
&\quad + \mathbb{E} \left[\frac{\mathbb{E}_0 [2 \|\mathbf{T}_1(\mathbf{u}_0)\|_F^2]}{\mathbb{E}_0 [\|\mathbf{r}\|_2^2] + \mathbb{E}_0 [\|\nabla_{\mathbf{X}_0} \mathbf{r}\|_F^2]} I_{\Omega(\mathbf{u}_0)} \right] + o(\exp(-c_{12}n)) \\
&\leq (2/n) c_{2,n}^{-1} \mathbb{E} [\|\mathbf{T}_0(\mathbf{u}_0) + \mathbf{T}_L(\mathbf{u}_0) + \mathbf{T}_{NL}(\mathbf{u}_0)\|_F^2] \\
&\quad + L^2 c_1^2 / (2n\tau) + o(\exp(-c_{12}n)) \\
&\leq (2/n) c_{2,n}^{-1} (1/p) \left(2c_1^2 L^2 n \mathbb{E} [\|\hat{\boldsymbol{\alpha}}\|_2^2] + 2\mathbb{E} [\|\mathbf{f}'\|_F^2] \right) \\
&\quad + L^2 c_1^2 / (2n\tau) + o(\exp(-c_{12}n)) \\
&= O(1/p).
\end{aligned}$$

The last two inequalities above are due to Step 1. The fact that $\mathbb{E} [\|\hat{\boldsymbol{\alpha}}\|_2^2] = O(1)$ is provided in Proposition C.2.2. Assumption 1 and 3 provide us $\mathbb{E} [\|\mathbf{f}'\|_F^2] / n = \mathbb{E} [\|\boldsymbol{\beta} + \nabla G(\mathbf{x}_1)\|_2^2] = O(1)$. \square

D.3.3. Proof of Proposition C.4.2.

Proposition C.4.2. *Let Ω be as in Proposition C.2.3. Let $\delta_{i,j} = 1$ if $i = j$, 0 otherwise. Let $\bar{\Omega}_3 := \Omega \cap \{\max_{i_1, i_2 \in [n]} |\mathbf{x}_{i_1}^T \mathbf{x}_{i_2} / p - \delta_{i_1, i_2}| < \delta\}$ where δ is a fixed positive defined in Assumption 4. Let $\mathbf{u}_0 \sim \text{Unif}(\mathbb{S}^{p-1}(1))$ be independent of $\mathbf{X}, \mathbf{W}, \varepsilon, G$. Under Assumptions 1, 2, 3 and 4 and a further assumption that $\Sigma'_p(0) = O(1/p)$, we have that*

$$(C.8) \quad \mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, \varepsilon, G} \left[\left(\frac{\text{trace}((\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}' \mathbf{u}_0))}{\|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2} \right)^2 I_{\Omega_3} \right] = O(1/p).$$

Proof. Let us first define

- (i) $\mathbf{v}_n = \mathbf{1}_n - (H_{11}, H_{22}, \dots, H_{nn})^\top$.
- (ii) $q_0 := \text{trace}(\mathbf{T}_{NL}(\mathbf{u}_0)) = \text{trace}[(\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}' \mathbf{u}_0)] = \mathbf{v}_n^\top \mathbf{G}' \mathbf{u}_0$.
- (iii) $\mathbf{r} := \mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}$.
- (iv) Ω be in Proposition C.2.3 such that
 - (i) $\Omega := \{1/n \cdot \|\mathbf{y} - \mathbf{A} \hat{\boldsymbol{\alpha}}\|_2^2 \geq c_{2,n}\}$.
 - (ii) $\mathbb{P}(\Omega^c) \leq o(\exp(-c_6 n))$ for some constant $c_6 > 0$.
- (v) $\bar{\Omega}_3$ be such that (cf. Corollary 2.8.3 in Vershynin [2018])
 - (i)

$$\bar{\Omega}_3 = \left\{ \max_{i_1, i_2 \in [n]} |\mathbf{x}_{i_1}^T \mathbf{x}_{i_2} / p - \delta_{i_1, i_2}| < \delta \right\}$$

- (ii) $\mathbb{P}(\bar{\Omega}_3^c) \leq o(\exp(-c_{10}n))$ for some universal constant $c_{10} > 0$.
 (vi) $\Omega_3 := \Omega \cap \bar{\Omega}_3$.

We abbreviate $\mathbb{E} := \mathbb{E}_{\mathbf{u}_0, \mathbf{X}, \mathbf{W}, G, \varepsilon}$. The order of $\|\mathbf{r}\|_2$ is specified on event Ω so that

$$\mathbb{E} \left[q_0^2 / \|\mathbf{r}\|_2^2 \cdot I_{\Omega_3} \right] \leq (1/n) c_{2,n}^{-1} \mathbb{E} \left[q_0^2 I_{\bar{\Omega}_3} \right].$$

By some algebra,

$$\begin{aligned} \mathbb{E} \left[q_0^2 I_{\bar{\Omega}_3} \right] &= \mathbb{E} \left[\left(\text{trace} \left[(\mathbf{I}_n - \mathbf{H}) \text{diag}(\mathbf{G}' \mathbf{u}_0) \right] \right)^2 I_{\bar{\Omega}_3} \right] \\ &= \mathbb{E} \left[(\mathbf{v}_n^\top \mathbf{G}' \mathbf{u}_0)^2 I_{\bar{\Omega}_3} \right] \\ &= \mathbb{E} \left[\mathbf{v}_n^\top \mathbb{E}_G \left[\mathbf{G}' \mathbb{E}_{\mathbf{u}_0} [\mathbf{u}_0 \mathbf{u}_0^\top] \mathbf{G}'^\top \right] \mathbf{v}_n I_{\bar{\Omega}_3} \right] \\ &= (1/p) \mathbb{E} \left[\mathbf{v}_n^\top \mathbb{E}_G \left[\mathbf{G}' \mathbf{G}'^\top \right] \mathbf{v}_n I_{\bar{\Omega}_3} \right], \end{aligned}$$

where we used $\mathbb{E}_{\mathbf{u}_0} [\mathbf{u}_0 \mathbf{u}_0^\top] = (1/p) \mathbf{I}_p$. So it suffices to show that

$$\mathbb{E} \left[\mathbf{v}_n^\top \mathbb{E}_G \left[\mathbf{G}' \mathbf{G}'^\top \right] \mathbf{v}_n I_{\bar{\Omega}_3} \right] = O(p).$$

Let $\mathbf{M} := \mathbb{E}_G \left[\mathbf{G}' \mathbf{G}'^\top \right]$. From Proposition C.2.7, on $\bar{\Omega}_3$, the (i_1, i_2) -th element of the above matrix is

$$m_{i_1, i_2} = \Sigma_p''(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) (\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p)/p + \Sigma_p'(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p).$$

We look at Taylor expansions around δ_{i_1, i_2} , and do some arrangement as following: for $i_1, i_2 \in [n]$ and $|\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2}| \leq \delta$,

$$\begin{aligned} \Sigma_p'(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) &= \Sigma_p'(\delta_{i_1, i_2}) + \Sigma_p''(\kappa_{i_1, i_2})(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2}), \\ &= \Sigma_p'(\delta_{i_1, i_2}) + \Sigma_p''(0)(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2}) \\ &\quad + (\Sigma_p''(\kappa_{i_1, i_2}) - \Sigma_p''(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2})(1 - \delta_{i_1, i_2}) \\ &\quad + (\Sigma_p''(\kappa_{i_1, i_2}) - \Sigma_p''(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2})\delta_{i_1, i_2}, \\ \Sigma_p''(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p)(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p)/p &= \Sigma_p''(0)(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p)/p + (\Sigma_p''(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) - \Sigma_p''(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p)(1 - \delta_{i_1, i_2})/p \\ &\quad + (\Sigma_p''(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) - \Sigma_p''(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p)\delta_{i_1, i_2}/p, \end{aligned}$$

where κ_{i_1, i_2} satisfies $|\kappa_{i_1, i_2} - \delta_{i_1, i_2}| \leq |\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \delta_{i_1, i_2}|$. From this we have decomposition of $\mathbf{M} := \mathbb{E}_G \left[\mathbf{G}' \mathbf{G}'^\top \right]$ into several matrices with small operator norm easy to calculate. With a slight abuse of notations $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{E}, \mathbf{F}, \mathbf{G}, \mathbf{H}$, we have

$$\mathbf{M} = \mathbf{A} + \mathbf{B} + \mathbf{C} + \mathbf{D} + \mathbf{E} + \mathbf{F} + \mathbf{G} + \mathbf{H},$$

where

$$\begin{aligned} \mathbf{A} &= (\Sigma_p'(1) - \Sigma_p'(0)) \mathbf{I}_n, \\ \mathbf{B} &= \Sigma_p'(0) \mathbf{1}_n \mathbf{1}_n^\top, \\ \mathbf{C} &= \Sigma_p''(0)(\mathbf{X} \mathbf{X}^\top / p), \\ \mathbf{D}_{i_1, i_2} &= (\Sigma_p''(\kappa_{i_1, i_2}) - \Sigma_p''(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) \delta_{i_1 \neq i_2}, \\ \mathbf{E}_{i_1, i_2} &= [\Sigma_p''(\kappa_{i_1, i_2}) \mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \Sigma_p''(0) \mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p - \Sigma_p''(\kappa_{i_1, i_2})] \delta_{i_1 = i_2}, \\ \mathbf{F} &= \Sigma_p''(0)(\mathbf{X} \mathbf{X}^\top / p)/p, \\ \mathbf{G}_{i_1, i_2} &= (\Sigma_p''(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) - \Sigma_p''(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) \delta_{i_1 \neq i_2}/p, \\ \mathbf{H}_{i_1, i_2} &= (\Sigma_p''(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) - \Sigma_p''(0))(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}/p) \delta_{i_1 = i_2}/p. \end{aligned}$$

It suffices to show that $q(\mathbf{N}) := \mathbb{E} [\mathbf{v}_n^\top \mathbf{N} \mathbf{v}_n I_{\bar{\Omega}_3}] = O(p)$ for \mathbf{N} being from \mathbf{A} to \mathbf{H} . We notice that $\|\mathbf{I}_n - \mathbf{H}\|_{\text{op}} \leq 1$ implies $|v_{n,i}| \leq 1$ for all $i \in [n]$. Then we have

- (i) $q(\mathbf{A}) = O(p)$ provided that $\Sigma'_p(1), \Sigma'_p(0) = O(1)$.
- (ii) $q(\mathbf{B}) = O(p)$ provided that $\Sigma'_p(0) = O(1/p)$.
- (iii) $q(\mathbf{C}) = O(p)$ provided that $\mathbb{E} [\|\mathbf{X}/\sqrt{p}\|_{\text{op}}] = O(1)$ and $\Sigma''_p(0) = O(1)$.
- (iv) $q(\mathbf{E}) = O(p)$ provided that $\sup_{x \in [1-\delta, 1+\delta]} \Sigma''_p(x), \Sigma''_p(0) = O(1)$.
- (v) $q(\mathbf{F}) = O(p)$ provided that $\Sigma''_p(0) = O(1)$.
- (vi) $q(\mathbf{H}) = O(p)$ provided that $\sup_{x \in [1-\delta, 1+\delta]} \Sigma''_p(x), \Sigma''_p(0) = O(p)$.

We notice that the above are true by assumptions on Σ_p and \mathbf{X} . For \mathbf{D} and \mathbf{G} , we notice the following:

$$|q(\mathbf{D})| := |\mathbb{E} [\mathbf{v}_n^\top \mathbf{D} \mathbf{v}_n I_{\bar{\Omega}_3}]| \leq \mathbb{E} [|\mathbf{v}_n|^\top |\mathbf{D}| |\mathbf{v}_n| I_{\bar{\Omega}_3}] \leq \mathbf{1}_n^\top \mathbb{E} [|\mathbf{D}| I_{\bar{\Omega}_3}] \mathbf{1}_n,$$

where the absolute value operation is taken element-wise for the vector \mathbf{v}_n and matrix \mathbf{D} . By the Lipschitz assumption of Σ''_p around 0, for $i_1 \neq i_2$,

$$\mathbb{E} [d_{i_1, i_2} | I_{\bar{\Omega}_3}] \leq \mathbb{E} [|(\Sigma''_p(\kappa_{i_1, i_2}) - \Sigma''_p(0))| |\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2} / p| I_{\bar{\Omega}_3}] \leq L_2 \cdot \mathbb{E} [(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2} / p)^2] = L_2 / p.$$

This implies $|q(\mathbf{D})| = O(p)$. For $|q(\mathbf{G})|$, we notice that

$$|q(\mathbf{G})| := |\mathbb{E} [\mathbf{v}_n^\top \mathbf{G} \mathbf{v}_n I_{\bar{\Omega}_3}]| \leq \mathbb{E} [|\mathbf{v}_n|^\top |\mathbf{G}| |\mathbf{v}_n| I_{\bar{\Omega}_3}] \leq \mathbf{1}_n^\top \mathbb{E} [|\mathbf{G}| I_{\bar{\Omega}_3}] \mathbf{1}_n,$$

where the absolute value operation is taken element-wise for the vector \mathbf{v}_n and matrix \mathbf{G} . By the Lipschitz assumption of Σ''_p around 0, for $i_1 \neq i_2$,

$$\mathbb{E} [g_{i_1, i_2} | I_{\bar{\Omega}_3}] \leq L_2 \mathbb{E} [(\mathbf{x}_{i_1}^\top \mathbf{x}_{i_2} / p)^2] / p \leq L_2 / p^2.$$

So that $|q(\mathbf{G})| = O(1)$. Combining the above we have our proposition. \square

REFERENCES

- Pierre C Bellec and Alexandre B Tsybakov. Bounds on the prediction error of penalized least squares estimators with convex penalty. In *Modern Problems of Stochastic Analysis and Statistics, Selected Contributions In Honor of Valentin Konakov*. Springer, 2017. URL <https://arxiv.org/pdf/1609.06675.pdf>.
- Pierre C Bellec and Cun-Hui Zhang. Second order stein: Sure for sure and other applications in high-dimensional inference, 2018.
- Pierre C Bellec and Cun-Hui Zhang. Second order poincaré inequalities and de-biasing arbitrary convex regularizers when $p/n \rightarrow \gamma$, 2019.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 2013.
- Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.
- Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve, 2019.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.