

# Asymptotic normality of robust $M$ -estimators with convex penalty\*

Pierre C. Bellec, Yiwei Shen and Cun-Hui Zhang

*Department of Statistics, Rutgers University*  
*e-mail: [pcb71@stat.rutgers.edu](mailto:pcb71@stat.rutgers.edu)*

**Abstract:** This paper develops asymptotic normality results for individual coordinates of robust  $M$ -estimators with convex penalty in high-dimensions, where the dimension  $p$  is at most of the same order as the sample size  $n$ , i.e.,  $p/n \leq \gamma$  for some fixed constant  $\gamma > 0$ . The asymptotic normality requires a bias correction and holds for most coordinates of the  $M$ -estimator for a large class of loss functions including the Huber loss and its smoothed versions regularized with a strongly convex penalty.

The asymptotic variance that characterizes the width of the resulting confidence intervals is estimated with data-driven quantities. This estimate of the variance adapts automatically to low ( $p/n \rightarrow 0$ ) or high ( $p/n \leq \gamma$ ) dimensions and does not involve the proximal operators seen in previous works on asymptotic normality of  $M$ -estimators. For the Huber loss, the estimated variance has a simple expression involving an effective degrees-of-freedom as well as an effective sample size. The case of the Huber loss with Elastic-Net penalty is studied in details and a simulation study confirms the theoretical findings. The asymptotic normality results follow from Stein formulae for high-dimensional random vectors on the sphere developed in the paper which are of independent interest.

**Keywords and phrases:** Robust estimation,  $M$ -estimator, asymptotic normality, confidence Intervals, high-dimensional statistics, bias-correction, Stein's formula.

Received June 2021.

## 1. Introduction

### 1.1. Robust inference

In his seminal paper on robustness, Huber [20] introduced  $M$ -estimators for an unknown location parameter  $\mu \in \mathbb{R}$  from observations  $Y_i = \mu + \varepsilon_i$ ,  $i = 1, \dots, n$ , where  $\varepsilon_i$  are iid noise random variables distributed as a mixture  $F = (1 - \epsilon)N(0, 1) + \epsilon H$  with  $H$  being the distribution of the contaminated samples, possibly chosen by an adversary. Given a differentiable loss function  $\rho : \mathbb{R} \rightarrow \mathbb{R}$  and its derivative  $\psi = \rho'$ , Huber defined  $M$ -estimators as minimizers  $\hat{\mu} = \operatorname{argmin}_{b \in \mathbb{R}} \sum_{i=1}^n \rho(Y_i - b)$ , or equivalently as solutions to  $\sum_{i=1}^n \psi(Y_i - b) = 0$ . Huber [20] went on to prove consistency and asymptotic normality of such  $M$ -estimators, obtaining among other results that if  $\rho$  is convex and  $\psi$  is absolutely

---

\*P.C. Bellec was partially supported by the NSF Grants DMS-1811976 and DMS-1945428. C.-H. Zhang was partially supported by the NSF Grants DMS-1721495, IIS-1741390, CCF-1934924, DMS-2052949 and DMS-2210850.

continuous, then under mild assumptions on  $F$ , the convergence  $\hat{\mu} \rightarrow \mu$  in probability holds as well as the asymptotic normality

$$n^{1/2}(\hat{\mu} - \mu) \rightarrow^d N(0, \mathbb{E}[\psi^2(\varepsilon_1)] / \mathbb{E}[\psi'(\varepsilon_1)]^2).$$

Huber's  $M$ -estimators were extended to regression models, where a design matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$  is observed together with responses  $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$  where  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are the rows of  $\mathbf{X}$  and  $\varepsilon_1, \dots, \varepsilon_n$  are possibly contaminated noise random variables as in the previous paragraph. For fixed or slowly growing dimension  $p$  as  $n \rightarrow +\infty$ , consistency and asymptotic normality of  $M$ -estimators of the form  $\hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n \rho(y_i - \mathbf{x}_i^\top \mathbf{b})$  were obtained, see [21, Section 7] or [27] among others. Explicitly, if  $\mathbf{e}_j \in \mathbb{R}^p$  is a canonical basis vector and one is interested in the asymptotic normality of  $\hat{\beta}_j - \beta_j$  for the purpose of confidence intervals, [27] shows that

$$(\mathbf{e}_j^\top (\mathbf{X}^\top \mathbf{X} / n)^{-1} \mathbf{e}_j)^{-1/2} \sqrt{n}(\hat{\beta}_j - \beta_j) \rightarrow^d N\left(0, \frac{\mathbb{E}[\psi^2(\varepsilon_1)]}{\mathbb{E}[\psi'(\varepsilon_1)]^2}\right) \quad (1.1)$$

if  $(p \log n)^{3/2} / n \rightarrow 0$  and under mild assumptions on  $\mathbf{X}$ . As in the location parameter problem of the previous paragraph, the asymptotic variance is characterized by the ratio  $\mathbb{E}[\psi^2(\varepsilon_1)] / \mathbb{E}[\psi'(\varepsilon_1)]^2$ .

The last decade has seen striking developments of similar asymptotic normality results in high-dimensions, where  $p/n \rightarrow \gamma$  for some constant  $\gamma < 1$ , cf. [17, 3, 25, 15, 16]. In terms of asymptotic normality, these works show that if  $\mathbf{X}$  has iid  $N(0, \boldsymbol{\Sigma})$  rows, the unregularized  $M$ -estimator  $\hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n \rho(y_i - \mathbf{x}_i^\top \mathbf{b})$  satisfies asymptotic normality of the form

$$(\mathbf{e}_j^\top \boldsymbol{\Sigma}^{-1} \mathbf{e}_j)^{-1/2} \sqrt{p}(\hat{\beta}_j - \beta_j) \rightarrow^d N(0, r^2) \quad (1.2)$$

where  $r > 0$  is a deterministic constant that captures the high-dimensionality of the problem [17, Lemma 1]. The constant  $r > 0$  is determined by solving a system of nonlinear equations with two unknowns: In the unregularized setting, [17, S2] describes this system of nonlinear equations with unknowns  $(r, c)$  as

$$\begin{cases} \mathbb{E}[1 - [\operatorname{prox}_c(\rho)]'(\varepsilon_1 + rZ)] = \gamma, \\ \mathbb{E}[(\varepsilon_1 + rZ - [\operatorname{prox}_c(\rho)](\varepsilon_1 + rZ))^2] = \gamma r^2 \end{cases} \quad \text{as } p/n \rightarrow \gamma, \quad (1.3)$$

where  $Z \sim N(0, 1)$  is independent of  $\varepsilon_1$ , and  $\operatorname{prox}_c(\rho)(t) = \operatorname{argmin}_{u \in \mathbb{R}} \rho(u) + (t - u)^2 / (2c)$  denotes the proximal operator of the convex function  $t \rightarrow c\rho(t)$  with derivative  $[\operatorname{prox}_c(\rho)]'(t)$ . The optimality conditions

$$c^{-1}(t - [\operatorname{prox}_c(\rho)](t)) = \psi([\operatorname{prox}_c(\rho)](t))$$

of the proximal minimization problem leads to the expressions

$$c^{-2}\gamma r^2 = \mathbb{E}[\psi([\operatorname{prox}_c(\rho)](\varepsilon_1 + rZ))^2] \quad \text{and} \quad c^{-1}\gamma = \mathbb{E}\left[\frac{d}{dt}\psi([\operatorname{prox}_c(\rho)](t))\right]_{t=\varepsilon_1 + rZ}$$

for the solutions  $(r, c)$  to the above system. Hence (1.2) can be rewritten as

$$(\mathbf{e}_j^\top \boldsymbol{\Sigma}^{-1} \mathbf{e}_j)^{-1/2} \sqrt{n}(\hat{\beta}_j - \beta_j) \rightarrow^d N\left(0, \frac{\mathbb{E}[\psi([\text{prox}_c(\rho)](\varepsilon_1 + rZ))^2]}{\mathbb{E}[\frac{d}{dt}\psi([\text{prox}_c(\rho)](t))|_{t=\varepsilon_1+rZ}]^2}\right), \quad (1.4)$$

see, e.g., [15, Theorem 4.1 and Corollary 4.6]. These results embody that when  $p$  and  $n$  are of the same order, the asymptotic variance in (1.1) must be modified to account for the high-dimensionality of the problem by (a) replacing  $\psi$  in the numerator and  $\psi'$  in the denominator by their compositions with the proximal operator  $\text{prox}_c(\rho)$ , and (b) adding the extra Gaussian term  $rZ$  to the initial noise  $\varepsilon_1$ . The distribution of  $\varepsilon_1 + rZ$  is sometimes referred to as the effective noise. The Gaussian assumption can be relaxed and some of the above results are still valid if  $\mathbf{X}$  has iid centered entries with variance one [25, 16]. Despite the subtle introduction of the proximal operator and the constants  $(r, c)$ , it is remarkable that the informal ratio  $\frac{\text{average}[\psi^2]}{\text{average}[(d/dt)\psi]^2}$  unifies the results (1.1) and (1.4) in both low and high-dimensions.

All results of the previous section are applicable when  $p/n \rightarrow \gamma$  with  $\gamma < 1$ . For  $\gamma > 1$  regularization is required to ensure the uniqueness of  $\hat{\beta}$ , for instance through an additive penalty which leads to regularized  $M$ -estimators of the form

$$\hat{\beta} = \underset{\mathbf{b} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \rho(y_i - \mathbf{x}_i^\top \mathbf{b}) + g(\mathbf{b}) \quad (1.5)$$

for some convex penalty function  $g: \mathbb{R}^p \rightarrow \mathbb{R}$ . The case of Ridge regularization with  $g(\mathbf{b}) = \tau \|\mathbf{b}\|_2^2/2$  for some constant  $\tau > 0$  is treated in [25, 16]. In this case, the two solutions  $(r, c)$  of a system of two nonlinear equations similar to (1.3) characterize the error  $\|\hat{\beta} - \beta\|_2$ , asymptotic normality and asymptotic variance of  $\sqrt{n}((1+a)\hat{\beta}_j - \beta_j)$  where  $a$  is a constant capturing the bias induced by regularization [16, Proposition 3.30] and  $a$  is a function of  $(\gamma, r, c)$ . Thrampoulidis et al. [29] characterize the error  $\|\hat{\beta} - \beta\|_2$  for a large class of  $(\rho, g)$  pairs using a technique known as the Convex Gordon Min-Max theorem pioneered by [28], and the recent paper [19] on Approximate Message Passing focused on  $g(\mathbf{b}) = \lambda \|\mathbf{b}\|_1$  and  $\rho$  either the Huber loss or the absolute value. Little is known, however, on asymptotic normality of the regularized estimators (1.5) for penalty functions different from the square norm  $\mathbf{b} \mapsto \tau \|\mathbf{b}\|_2^2$ . The theories developed in [29, 19] do not readily provide asymptotic normality results and regularized  $M$ -estimators of the form (1.5) lack confidence interval capabilities. One goal of the present paper is to fill this gap.

A separate line of research develops asymptotic normality results and confidence intervals for regularized least-squares estimators of the form

$$\hat{\beta} = \underset{\mathbf{b} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{2n} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{b})^2 + g(\mathbf{b}) = \underset{\mathbf{b} \in \mathbb{R}^p}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2/(2n) + g(\mathbf{b}) \quad (1.6)$$

where  $\mathbf{X}$  has rows  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . Early results studied the Lasso with  $g(\mathbf{b}) = \lambda \|\mathbf{b}\|_1$  [33, 22, 30] under a sparsity condition  $s \log(p) = o(\sqrt{n})$  where  $s = \|\beta\|_0$ , or

Ridge regression [11]. For the Lasso the sparsity condition was later improved to  $s \log^2(p)/n \rightarrow 0$  [24], to  $s \log(p/s)/n \rightarrow 0$  [8] and  $p/n \rightarrow \gamma \in (0, \infty)$  with  $s \lesssim n/\log(p/s)$  ([23, 26] for isotropic Gaussian design and [13] [7, Theorem 3.2] for non-isotropic Gaussian design). For penalty functions beyond the Lasso and Ridge regularization, [12, Proposition 4.3(iii)] provides asymptotic normality on average over the coordinates for permutation invariant penalty function  $g$  in (1.6), and [7, Theorem 3.1] proves asymptotic normality for individual coordinates of (1.6) under a strong convexity assumption. A high-level message of these works is that one must de-bias the regularized estimator (1.6) in order to obtain asymptotic normality at the  $\sqrt{n}$ -adjusted rate and construct confidence intervals. In the  $p/n \rightarrow \gamma$  regime that is the focus of the present paper, this bias correction takes the following form. Under a strong convexity assumption and for  $\mathbf{X}$  with iid  $N(\mathbf{0}_p, \Sigma)$  rows, [7] proves that for most coordinates  $j = 1, \dots, p$ ,

$$\frac{(n - \hat{\text{df}})(\hat{\beta}_j - \beta_j) + \mathbf{e}_j^\top \Sigma^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \hat{\beta})}{\|\mathbf{y} - \mathbf{X} \hat{\beta}\|_2} \Omega_{jj}^{-1/2} \rightarrow^d N(0, 1) \quad (1.7)$$

where  $\Omega_{jj} = \mathbf{e}_j^\top \Sigma^{-1} \mathbf{e}_j$  and  $\hat{\text{df}}$  is the effective degrees of freedom of  $\hat{\beta}$  defined as the Jacobian of the map  $\mathbf{y} \mapsto \mathbf{X} \hat{\beta}$  for fixed  $\mathbf{X}$ . For  $\Sigma = \mathbf{I}_p$  and consequently  $\Omega_{jj} = 1$ , a similar bias correction proportional to  $\mathbf{e}_j^\top \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \hat{\beta})$  is visible in the asymptotic normality result [12, Proposition 4.3(iii)], although there the coefficients  $(n - \hat{\text{df}})$  and  $\|\mathbf{y} - \mathbf{X} \hat{\beta}\|_2$  in (1.7) are replaced with deterministic scalar counterparts obtained by solving a system of nonlinear equations of a similar nature as (1.3). Another goal of the present paper is to equip robust  $M$ -estimators (1.5), with  $\rho$  different than the squared loss, with de-biasing and asymptotic normality results similar to the previous display, allowing for general robust loss functions  $\rho$  coupled with general convex penalty functions  $g$ .

## 1.2. Contributions

Our contributions are the following.

1. We provide de-biasing and asymptotic normality results for robust  $M$ -estimators with convex penalty functions when  $p$  and  $n$  are of the same order. This leads to confidence intervals for the  $j$ -th coordinate  $\beta_j$  of the unknown coefficient vector  $\beta$ . Asymptotic normality holds for a large class of robust loss functions, including the Huber loss and its smoothed versions.
2. Although this bias correction required for asymptotic normality resembles the one-step estimators recommended in the theory of classical  $M$ -estimator to improve efficiency (e.g., [31, Eq. (1.11)]), a notable difference from the low-dimensional case is the requirement of a degrees-of-freedom adjustment to amplify the one-step correction. For the squared loss, this degrees-of-freedom adjustment takes the form of multiplication by  $(n - \hat{\text{df}})$  in (1.7); one contribution of this paper is to identify the degrees-of-freedom

adjustment that leads to asymptotic normality for robust and regularized M-estimators, beyond the squared loss.

3. The asymptotic variance is estimated by random, data-driven quantities, as opposed to the deterministic scalars  $(r, c)$  that determine the asymptotic variance for unregularized estimators in (1.4). The fact that the asymptotic variance is estimated by observable quantities makes this results more suitable for confidence intervals (case in point: computing the solution  $(r, c)$  of (1.3) and the asymptotic variance in (1.4) requires the knowledge of the noise distribution subject to contamination). The asymptotic normality result takes the form

$$\hat{V}^{-1/2} \Omega_{jj}^{-1/2} \sqrt{n}(\hat{\beta}_j - \beta_j) + [\text{observable bias correction}] \approx N(0, 1)$$

where the data-driven variance estimate  $\hat{V}$  again is a ratio of the form  $\frac{\text{average}[\psi^2]}{\text{average}[(d/dt)\psi]^2}$  for a particular sense of average to be defined in (2.14) below. Interestingly, the expression for this average and  $\hat{V}$  does not involve the proximal mapping in (1.4). This informal statement will be made precise in Section 2.4 below.

4. In order to derive these new asymptotic normality results, we develop in Appendix B new identities for random vectors uniformly distributed on the Euclidean sphere. Although the argument of the present paper for asymptotic normality is closely related to that used in [7] for the squared loss, this previous theory for the squared loss for functions of standard normal vector does not extend to robust loss functions due to the lack of strong convexity of  $\rho$  for robust losses, and consequently the lack of explicit lower bounds on  $\frac{1}{n} \sum_{i=1}^n \psi(y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}})^2$ . Developing these new identities and the corresponding asymptotic normality results for random vectors uniformly distributed on the sphere is a crucial step to overcome the lack of global strong convexity of  $\rho$  for robust losses and to obtain the asymptotic normality results. These new identities provide novel Stein formulae for random vectors on the sphere and may be used more broadly for elliptical distributions.

## 2. Model and main results

### 2.1. Model and assumptions

We consider a linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (2.1)$$

where  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ , with a regularized M-estimator

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \rho(y_i - \mathbf{x}_i^\top \mathbf{b}) + g(\mathbf{b}), \quad (2.2)$$

where  $\rho : \mathbb{R} \rightarrow \mathbb{R}$  is the loss and  $g : \mathbb{R}^p \rightarrow \mathbb{R}$  is the penalty.

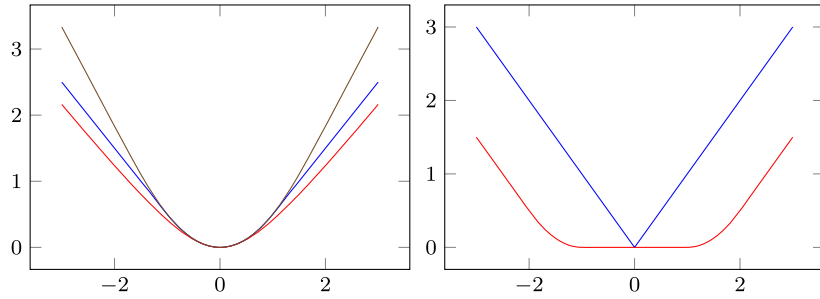


TABLE 1

Left: robust loss functions satisfying Assumption A: the Huber loss  $x \mapsto H(x) = \int_0^{|x|} \min(t, 1) dt$ , its smoothed versions  $x \mapsto \sqrt{1+x^2}$  and  $x \mapsto \frac{x^2}{2} I\{|x| \leq 1\} + (\frac{1}{6} - \frac{|x|}{2} + x^2 - \frac{|x|^3}{6}) I\{|x| \in (1, 2)\} + (\frac{7}{6} + \frac{3|x|}{2}) I\{|x| \geq 2\}$ . Right: two loss functions that do not satisfy Assumption A: the absolute deviation loss  $x \mapsto |x|$  and the 1-insensitive loss  $x \mapsto H(|x| - 1)_+$  where  $H(\cdot)$  is the Huber loss.

**Assumption A** (Assumptions on the loss). Let  $\rho : \mathbb{R} \rightarrow \mathbb{R}$  be convex and continuously differentiable on  $\mathbb{R}$ , with derivative  $\psi = \rho'$  being  $L$ -Lipschitz with

$$K^2 \leq \psi'(x) + \psi(x)^2 \quad \text{for almost every } x \in \mathbb{R} \quad (2.3)$$

for some positive constant  $K > 0$  independent of  $n, p$ .

Two families of robust losses that do not satisfy Assumption A are non-differentiable losses such as the least absolute deviations  $\rho(x) = |x|$ , and  $\delta$ -insensitive losses such as  $\rho(x) = (|x| - \delta)_+^2$  as  $\psi(x)^2 + \psi'(x) = 0$  in a neighborhood of 0. Assumption A is verified by the Huber loss  $\rho(x) = \int_0^{|x|} \min(1, t) dt$  with  $K = 1$ , as well as by any smooth version of the Huber loss, for instance  $\rho(x) = \sqrt{1+x^2}$  with  $K^2 = \frac{23}{27} \approx 0.852$ . The one-sided logistic loss  $\rho(x) = \log(1 + e^x)$  also satisfies Assumption A.

**Assumption B** (Strong convexity of  $g$ ). For some constant  $\tau > 0$  independent of  $n, p$ , the penalty  $g : \mathbb{R}^p \rightarrow \mathbb{R}$  is  $\tau$ -strongly convex in the sense that  $\mathbf{x} \mapsto g(\mathbf{x}) - \tau \|\mathbf{x}\|_2^2/2$  is convex.

Some useful characterizations of strong convexity (Assumption B) are the following. Throughout,  $\partial g(\mathbf{b})$  denotes the subdifferential of  $g$  at  $\mathbf{b} \in \mathbb{R}^p$ . Then  $g$  is  $\tau$ -strongly convex if and only if

$$g(\mathbf{a}) - g(\mathbf{b}) \geq \mathbf{u}^\top (\mathbf{a} - \mathbf{b}) + (\tau/2) \|\mathbf{a} - \mathbf{b}\|_2^2 \quad \text{for all } \mathbf{u} \in \partial g(\mathbf{a}) \text{ and } \mathbf{a}, \mathbf{b} \in \mathbb{R}^p. \quad (2.4)$$

Similarly  $g$  is  $\tau$ -strongly convex if and only for any  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$

$$(\mathbf{u} - \mathbf{v})^\top (\mathbf{a} - \mathbf{b}) \geq \tau \|\mathbf{a} - \mathbf{b}\|_2^2 \quad \text{for all } \mathbf{u} \in \partial g(\mathbf{a}), \mathbf{v} \in \partial g(\mathbf{b}). \quad (2.5)$$

As  $\psi$  in Assumption A is increasing, Assumption A implies that  $\rho$  is  $K^2/2$ -strongly convex in the interval  $\{x \in \mathbb{R} : \psi(x)^2 \leq K^2/2\}$ , and conversely if  $\rho$  is  $\mu$ -strongly convex in the interval  $\{x \in \mathbb{R} : \psi(x)^2 \leq C\}$  then Assumption A is satisfied with  $K^2 = \min(\mu, C)$ .

**Assumption C.** The rows of the design matrix  $\mathbf{X}$  are iid  $N(\mathbf{0}, \Sigma)$  random vectors and all the eigenvalues of  $\Sigma \in \mathbb{R}^{p \times p}$  are in  $[\kappa, 1/\kappa]$  for some constant  $\kappa \in (0, 1)$  independent of  $n, p$ . The noise  $\varepsilon$  is independent of  $\mathbf{X}$  and admits a density with respect to the Lebesgue measure.

**Assumption D.**  $p/n \leq \gamma$  for some constant  $\gamma > 0$  independent of  $n, p$ .

## 2.2. Notation

Throughout the paper  $\Phi(t) = (2\pi)^{-1/2} \int_{-\infty}^t \exp(-u^2/2) du$  is the standard normal cumulative distribution function. Furthermore,  $[n] = \{1, \dots, n\}$  and we use the notation

$$\boldsymbol{\psi} = (\psi(y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}))_{i \in [n]}, \quad \boldsymbol{\psi}' = (\psi'(y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}))_{i \in [n]}, \quad \mathbf{h} = \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}. \quad (2.6)$$

For each  $j \in [p]$ , let  $\mathbf{e}_j$  denote the  $j$ -th vector in the standard basis of  $\mathbb{R}^p$ , and let

$$\Omega_{jj} = \mathbf{e}_j^\top \Sigma^{-1} \mathbf{e}_j, \quad \mathbf{z}_j = \mathbf{X} \Sigma^{-1} \Omega_{jj}^{-1} \mathbf{e}_j, \quad \mathbf{Q}_j = \mathbf{I}_p - \Sigma^{-1} \Omega_{jj}^{-1} \mathbf{e}_j \mathbf{e}_j^\top. \quad (2.7)$$

We remark that the above definition implies the following properties:

- $\mathbf{X} = \mathbf{X} \mathbf{Q}_j + \mathbf{z}_j \mathbf{e}_j^\top$  and  $\mathbf{X} \boldsymbol{\beta} = \mathbf{X} \mathbf{Q}_j \boldsymbol{\beta} + \mathbf{z}_j \beta_j$ .
- $\mathbf{z}_j \sim N(\mathbf{0}, \Omega_{jj}^{-1} \mathbf{I}_n)$  is independent of  $\mathbf{X} \mathbf{Q}_j$  (cf. Proposition D.1).
- Under Assumption C,  $\Omega_{jj} \in [\kappa, 1/\kappa]$ .

By construction of  $\mathbf{z}_j$  and  $\mathbf{Q}_j$ , the response  $\mathbf{y}$  can be decomposed as  $\mathbf{y} = \beta_j \mathbf{z}_j + \mathbf{X} \mathbf{Q}_j \boldsymbol{\beta} + \varepsilon$  where  $\beta_j \in \mathbb{R}$  is the scalar parameter of interest for a fixed covariate  $j \in [p]$ , the vector  $\mathbf{Q}_j \boldsymbol{\beta}$  is a high-dimensional nuisance parameter and  $\varepsilon$  is independent noise. Under the additional assumption of  $\varepsilon_i \sim N(0, 1)$ ,  $\Omega_{jj}^{-1}$  is the Fisher information for the estimation of  $\beta_j$ .

In the proof, it will be useful to treat  $\boldsymbol{\psi} = \boldsymbol{\psi}(\varepsilon, \mathbf{X})$  as a map from  $\mathbb{R}^{n \times (p+1)} \rightarrow \mathbb{R}^n$ , formally defined as

$$\begin{aligned} \hat{\boldsymbol{\beta}}(\varepsilon, \mathbf{X}) &= \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} \sum_{i \in [n]} \frac{\rho(\varepsilon_i - \mathbf{x}_i^\top (\mathbf{b} - \boldsymbol{\beta}))}{n} + g(\mathbf{b}), \\ \boldsymbol{\psi}(\varepsilon, \mathbf{X}) &= \boldsymbol{\psi}(\varepsilon + \mathbf{X} \boldsymbol{\beta} - \mathbf{X} \hat{\boldsymbol{\beta}}(\varepsilon, \mathbf{X})). \end{aligned} \quad (2.8)$$

Since  $(\boldsymbol{\beta}, \varepsilon)$  are unknown, we cannot compute the derivatives of  $\boldsymbol{\psi}$  a priori. However, for a fixed  $\mathbf{X}$ , the quantity  $\boldsymbol{\psi}(\mathbf{y} - \mathbf{X} \boldsymbol{\beta}, \mathbf{X})$  is observable since all terms in  $(\boldsymbol{\beta}, \varepsilon)$  cancel out (indeed  $\boldsymbol{\psi}(\mathbf{y} - \mathbf{X} \boldsymbol{\beta}, \mathbf{X})$  is simply  $\boldsymbol{\psi}(\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}})$  with  $\hat{\boldsymbol{\beta}}$  in (2.2)). We can thus define the observable matrix of size  $n \times n$

$$[\nabla_{\mathbf{y}} \boldsymbol{\psi}]^\top \stackrel{\text{def}}{=} (\partial / \partial \mathbf{y}) \boldsymbol{\psi}(\mathbf{y} - \mathbf{X} \boldsymbol{\beta}, \mathbf{X}) \quad (2.9)$$

holding  $\mathbf{X}$  fixed, at every point  $\mathbf{y} \in \mathbb{R}^n$  where  $\mathbf{y} \mapsto \boldsymbol{\psi}(\mathbf{y} - \mathbf{X} \boldsymbol{\beta}, \mathbf{X})$  is differentiable. By Proposition C.4 below, the map  $\mathbf{y} \mapsto \boldsymbol{\psi}(\mathbf{y} - \mathbf{X} \boldsymbol{\beta}, \mathbf{X})$  is  $L$ -Lipschitz

and  $\nabla_{\mathbf{y}}\psi$  exists at Lebesgue almost every  $\mathbf{y}$ , and with probability one since  $\mathbf{y}$  has continuous distribution under Assumption C. Furthermore, the gradient (2.9) at the currently observed data  $(\mathbf{y}, \mathbf{X})$  does not depend on any unknown quantity. It can be computed from  $(\mathbf{y}, \mathbf{X})$  either by finding a closed form expression for (2.9) for a given penalty function, or by approximation using finite difference or other numerical methods (e.g., [4, Section 2.7] in the same regularized M-estimation setting as the present paper). Note that in the concurrent paper [5] that focuses on asymptotic normality of residuals and estimation of the generalization error, the  $n \times n$  matrix (2.9) is denoted  $\mathbf{V}$  [5, (2.3)], and  $\mathbf{V} = \text{diag}(\psi')(\mathbf{I}_n - \mathbf{X}\hat{\mathbf{A}}\mathbf{X}^T \text{diag}(\psi'))$  for the matrix  $\hat{\mathbf{A}} \in \mathbb{R}^{p \times p}$  defined in [5, Theorem 1].

### 2.3. Main result

In the following result, we consider a sequence of integer pairs  $(n, p)$ , regression problems (2.1) and  $M$ -estimator (2.2), without explicit reference to their dependence on  $n$  or  $p$ . For instance, one can think of  $p = p_n$  as a nondecreasing function of  $n$  and  $(g, \beta, \hat{\beta}, \rho)$  are also implicitly indexed by  $n$  with values possibly changing with  $n$ .

**Theorem 2.1** (Asymptotic Normality result for  $M$ -estimator). *Consider the linear model (2.1) and the  $M$ -estimator  $\hat{\beta}$  in (2.2). Assume  $\mathbb{E}[\|\Sigma^{1/2}\mathbf{h}\|_2^2] \leq \mathcal{R} < +\infty$ . Let Assumptions A, B, C and D be fulfilled for constants  $\mathcal{R}, \tau, \kappa, K, L, \gamma > 0$  independent of  $n, p$ . Define the map  $\mathbf{y} \mapsto \psi(\mathbf{y} - \mathbf{X}\beta, \mathbf{X})$  and its Jacobian  $[\nabla_{\mathbf{y}}\psi]^\top$  in (2.9) holding  $\mathbf{X}$  fixed. For each  $j \in [p]$  let*

$$\xi_j = \|\psi\|^{-1} [\psi^\top \mathbf{z}_j - n^{-1} \|\mathbf{z}_j\|_2^2 (\beta_j - \hat{\beta}_j) \text{tr}(\nabla_{\mathbf{y}}\psi)], \quad (2.10)$$

$$\xi'_j = \|\psi\|^{-1} [\psi^\top \mathbf{z}_j - \Omega_{jj}^{-1} (\beta_j - \hat{\beta}_j) \text{tr}(\nabla_{\mathbf{y}}\psi)]. \quad (2.11)$$

Then for any positive sequence  $(a_p)$  with  $\lim_{p \rightarrow +\infty} a_p = +\infty$ ,

$$\max_{j \in J_{n,p}} \left| \mathbb{P}(\Omega_{jj}^{1/2} \xi_j \leq t) - \Phi(t) \right| + \left| \mathbb{P}(\Omega_{jj}^{1/2} \xi'_j \leq t) - \Phi(t) \right| \rightarrow 0, \quad (2.12)$$

for some  $J_{n,p} \subseteq [p]$  satisfying  $|J_{n,p}|/p \geq 1 - a_p/p$ .

The proof is given in Appendix A.2. To interpret the above result, we remark that for any slowly increasing sequence  $a_p$  such as  $a_p = \log p$  or  $a_p = \log \log p$ , the asymptotic normality (2.12) holds uniformly over all coordinates  $j = 1, \dots, p$  except  $a_p$  of them, so that both  $\xi_j$  and  $\xi'_j$  are asymptotically pivotal for an overwhelming majority of  $\beta_j$ . Another interpretation is given in the following Corollary: For any given precision threshold  $v > 0$ , there exist at most  $a_*$  coordinates  $j = 1, \dots, p$  such that  $|\mathbb{P}(\Omega_{jj}^{1/2} \xi_j \leq t) - \Phi(t)| \geq v$  where  $a_*$  is a constant independent of  $n, p$ .

**Corollary 2.2.** *Let the setting and assumptions of Theorem 2.1 be fulfilled. For any arbitrarily small constant  $v > 0$  independent of  $n, p$ , define*

$$J_{n,p}^v = \left\{ j \in [p] : \left| \mathbb{P}(\Omega_{jj}^{1/2} \xi_j \leq t) - \Phi(t) \right| > v \right\}.$$



Then,  $\sup_{n,p} |J_{n,p}^v| \leq a_*$  for a certain constant  $a_*$  depending on  $\{v, \tau, \mathcal{R}, \gamma, L, K, t\}$  only. In other words, for any  $(n, p)$  with  $p/n \leq \gamma$  there are at most a constant number of coordinates  $j = 1, \dots, p$  such that  $|\mathbb{P}(\Omega_{jj}^{1/2} \xi_j \leq t) - \Phi(t)| > v$ . The same conclusion holds with  $\xi_j$  replaced by  $\xi'_j$ .

*Proof of Corollary 2.2.* We proceed by contradiction. If the claim does not hold, there exists a constant  $v_* > 0$  such that  $|J_{n,p}^{v_*}| \geq 2a_p$  for an unbounded sequence  $a_p$ . By extracting a subsequence if necessary, we may assume without loss of generality that  $a_p$  is monotonically increasing with  $a_p \rightarrow +\infty$ . By Theorem 2.1 there exists  $J_{n,p} \subset [p]$  such that (2.12) holds. By definition of  $J_{n,p}^{v_*}$  and  $J_{n,p}$ , we have  $J_{n,p} \cap J_{n,p}^{v_*} = \emptyset$  for  $p$  large enough. This implies that  $p \geq |J_{n,p}| + |J_{n,p}^{v_*}| \geq (p - a_p) + 2a_p$  for  $p$  large enough, a contradiction.  $\square$

*Remark 2.1.* Theorem 2.1 requires the assumption  $\mathbb{E}[\|\Sigma^{1/2}(\hat{\beta} - \beta)\|_2^2] \leq \mathcal{R} < +\infty$  for some constant  $\mathcal{R}$ . This assumes that the expected risk of  $\hat{\beta}$  is bounded, which holds true under the following additional assumptions:

- The penalty is minimized at 0:  $\mathbf{0} \in \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} g(\mathbf{b})$ ;
- The loss is Lipschitz:  $\sup_{t \in \mathbb{R}} |\psi(t)| \leq \mathcal{L}$ ;
- The noiseless signal has bounded variance:  $\mathbb{E}[(\mathbf{x}_i^\top \beta)^2] = \|\Sigma^{1/2} \beta\|^2 \leq \mathcal{V}$ .

Above,  $\mathcal{L}, \mathcal{V} > 0$  are constants. Then by the KKT conditions,  $\mathbf{X}^\top \psi \in n \partial g(\hat{\beta})$  and thanks to (2.5) from Assumption B and the first bullet above,

$$(\hat{\beta} - \mathbf{0})^\top \mathbf{X}^\top \psi \in n(\hat{\beta} - \mathbf{0})^\top (\partial g(\hat{\beta}) - \partial g(\mathbf{0})) \geq n\tau \|\hat{\beta}\|^2.$$

It follows by the Cauchy-Schwarz inequality that  $n\tau \|\hat{\beta}\| \leq \|\mathbf{X}\|_{op} \sqrt{n} \mathcal{L}$ . Thanks to Assumption C we thus obtain

$$\begin{aligned} \mathbb{E}[\|\Sigma^{1/2}(\hat{\beta} - \beta)\|_2^2] &\leq 2\mathbb{E}[\|\Sigma^{1/2}\beta\|_2^2] + 2\mathbb{E}[\|\Sigma^{1/2}\hat{\beta}\|_2^2] \\ &\leq 2\mathcal{V} + 2(\mathcal{L}^2 \kappa^{-2} \tau^{-2}) \mathbb{E}[\|\mathbf{X}\Sigma^{-1/2}\|_{op}^2/n]. \end{aligned}$$

The upper bound (C.1) completes the proof that we can find a suitable constant  $\mathcal{R}$  depending on  $\gamma, \tau, \kappa, \mathcal{L}, \mathcal{V}$  only.

*Remark 2.2.* The set  $J_{n,p}$  in (2.12) excludes a few coordinates  $j \in \{1, \dots, p\} \setminus J_{n,p}$ . A natural question is whether the exclusion of some coordinates is necessary or an artefact of the current proofs. By Section 3.7 in [7], there are known examples where a “variance spike” occurs in at most a finite number of coordinates  $j \in \{1, \dots, p\}$ : for these few coordinates, the asymptotic variance of  $\xi_j$  or  $\xi'_j$  is strictly larger than 1 and the convergence in distribution  $\xi_j \rightarrow^d N(0, 1)$  cannot hold.

## 2.4. Data-driven variance estimate

Except for at most a constant number of coordinates  $j \in [p]$ , the approximation

$$\hat{V}^{-1/2} \Omega_{jj}^{-1/2} \sqrt{n}(\hat{\beta}_j - \beta_j) + [\text{bias correction}] \approx N(0, 1) \quad (2.13)$$

holds, where the bias correction is observable and determined by the first term in (2.11) and the data-driven variance estimate is

$$\hat{V} = \frac{\|\boldsymbol{\psi}\|_2^2/n}{(\text{tr}(\nabla_{\mathbf{y}}\boldsymbol{\psi})/n)^2} = \frac{n^{-1} \sum_{i=1}^n \psi_i^2}{(n^{-1} \sum_{i=1}^n (\partial/\partial y_i)\psi_i)^2}. \quad (2.14)$$

which characterizes the length of confidence intervals for  $\beta_j$ . This ratio of an average of  $\psi^2$  by a squared average of a derivative of  $\psi$  mirrors the robust asymptotic results in (1.1) and (1.4) as discussed in the introduction. Confidence intervals can be readily constructed from Theorem 2.1 or the informal approximation (2.13): a 95%-confidence interval for  $\beta_j$  is given by  $\hat{\beta}_j + (\Omega_{jj}\hat{V}/n)^{1/2}([\text{bias correction}] \pm 1.96)$ , that is,

$$\left[ \hat{\beta}_j + \frac{\Omega_{jj}\boldsymbol{\psi}^\top \mathbf{z}_j}{\text{tr}[\nabla_{\mathbf{y}}\boldsymbol{\psi}]} - \left( \frac{\Omega_{jj}\hat{V}}{n} \right)^{1/2} 1.96, \hat{\beta}_j + \frac{\Omega_{jj}\boldsymbol{\psi}^\top \mathbf{z}_j}{\text{tr}[\nabla_{\mathbf{y}}\boldsymbol{\psi}]} + \left( \frac{\Omega_{jj}\hat{V}}{n} \right)^{1/2} 1.96 \right].$$

In contrast with the asymptotic variance in (1.4), the above  $\hat{V}$  involves observable quantities. In particular, while the asymptotic variance in (1.4) depends on the distribution of the noise through the solutions  $(r, c)$  of the system (1.3), the knowledge of the noise distribution is not required to compute  $\hat{V}$  and construct confidence intervals for  $\beta_j$ .

Theorem 2.1 is valid for  $p/n \leq \gamma$ , without requiring a specific limit for the ratio  $p/n$ . Theorem 2.1 is also valid for  $p = o(n)$ , so that Theorem 2.1 and the estimated asymptotic variance (2.14) unifies both low- and high-dimensional asymptotic normality results.

For the Huber loss

$$H(u) = \int_0^{|u|} \min(1, t) dt = \begin{cases} u^2/2 & \text{if } |u| \leq 1, \\ |u| - 1/2 & \text{if } |u| > 1, \end{cases} \quad (2.15)$$

the quantity  $\text{tr}[\nabla_{\mathbf{y}}\boldsymbol{\psi}]$  has a simpler form: By the chain rule

$$\text{tr}[\nabla_{\mathbf{y}}\boldsymbol{\psi}] = \text{tr}[\text{diag}(\boldsymbol{\psi}')] - \text{tr}[\text{diag}(\boldsymbol{\psi}')(\partial/\partial \mathbf{y})\mathbf{X}\hat{\boldsymbol{\beta}}]$$

where the differentiation is understood holding  $\mathbf{X}$  fixed as in (2.9). With  $\hat{n} = \text{tr}[\text{diag}(\boldsymbol{\psi}')] - \text{df}$  the number of observations such that  $y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$  falls in the range where the Huber loss is quadratic and  $\hat{\text{df}} = \text{tr}[\text{diag}(\boldsymbol{\psi}')(\partial/\partial \mathbf{y})\mathbf{X}\hat{\boldsymbol{\beta}}]$  representing the degrees-of-freedom of the M-estimator  $\hat{\boldsymbol{\beta}}$ , the quantity  $\text{tr}[\nabla_{\mathbf{y}}\boldsymbol{\psi}]$  appearing in the denominator of  $\hat{V}$  simplifies to  $\text{tr}[\nabla_{\mathbf{y}}\boldsymbol{\psi}] = \hat{n} - \hat{\text{df}}$ . In this case, the asymptotic normality for  $\xi_j'$  in Theorem 2.1 takes the form

$$\Omega_{jj}^{1/2} \xi_j' = \frac{(\hat{n} - \hat{\text{df}})(\hat{\beta}_j - \beta_j) + \Omega_{jj}\mathbf{z}_j^\top \boldsymbol{\psi}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{\|\boldsymbol{\psi}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})\|_2} \Omega_{jj}^{-1/2} \rightarrow^d N(0, 1) \quad (2.16)$$

uniformly over all  $j \in J_{n,p}$ . This extends the asymptotic normality result (1.7) to the Huber loss with the following important modifications: the sample size

$n$  is replaced by  $\hat{n}$  and the residuals  $\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  are replaced by  $\boldsymbol{\psi} = \psi(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ . The variance  $\hat{V}$  that determines the length of the confidence interval resulting from (2.13) presents a trade-off among  $\|\boldsymbol{\psi}\|^2/n$ , an effective sample size  $\hat{n}$  and the degrees-of-freedom  $\hat{\text{df}}$ : For confidence intervals with small length, the  $M$ -estimator  $\hat{\boldsymbol{\beta}}$  should have small residuals measured as  $\|\psi(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})\|_2$ , small degrees-of-freedom  $\hat{\text{df}}$ , and large effective sample size  $\hat{n}$ .

### 2.5. Example

This section specializes the above results to scaled versions of the Huber loss (2.15) and the Elastic-Net penalty  $g(\mathbf{b}) = \lambda\|\mathbf{b}\|_1 + \tau\|\mathbf{b}\|_2^2/2$  for tuning parameters  $\lambda, \tau > 0$ . We consider the  $M$ -estimator

$$\hat{\boldsymbol{\beta}} = \underset{\mathbf{b} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \sigma^2 H(\sigma^{-1}(y_i - \mathbf{x}_i^\top \mathbf{b})) + \lambda\|\mathbf{b}\|_1 + \tau\|\mathbf{b}\|_2^2/2, \quad (2.17)$$

which corresponds to the scaled Huber loss  $\rho(u) = \sigma^2 H(\sigma^{-1}u)$  where  $\sigma > 0$  is a scaling parameter. Since the derivative  $H'$  is 1-Lipschitz, so is  $\psi = \rho'$ . Furthermore,  $\psi'(u) = H''(\sigma^{-1}u) = 1$  for  $|u| \leq \sigma$  and  $|\psi(u)| = \sigma|H'(\sigma^{-1}u)| = \sigma$  for  $|u| > \sigma$ , so that  $\min_{x \in \mathbb{R}} [\psi^2(x) + \psi'(x)] \geq \min(1, \sigma^2)$  and Assumption A is satisfied with  $L = 1$  and  $K^2 = \min(1, \sigma^2)$ . Assumption B is also satisfied as the penalty is the sum of the  $\ell_1$  norm plus  $\tau\|\mathbf{b}\|^2/2$ . The quantity  $\operatorname{tr}[\nabla_{\mathbf{y}} \boldsymbol{\psi}]$  appearing in Theorem 2.1 in the denominator of estimated variance  $\hat{V}$  is computed in [4, Proposition 2.3]: For almost every  $(\mathbf{X}, \mathbf{y})$ ,

$$\begin{aligned} \operatorname{tr}[\nabla_{\mathbf{y}} \boldsymbol{\psi}] &= \hat{n} - \hat{\text{df}}, \\ \hat{\text{df}} &= \operatorname{tr}[\operatorname{diag}(\boldsymbol{\psi}') \mathbf{X}_{\hat{S}} (\mathbf{X}_{\hat{S}}^\top \operatorname{diag}(\boldsymbol{\psi}') \mathbf{X}_{\hat{S}} + n\tau \mathbf{I}_{|\hat{S}|})^{-1} \mathbf{X}_{\hat{S}}^\top \operatorname{diag}(\boldsymbol{\psi}')], \end{aligned}$$

where  $\hat{n} = \operatorname{tr}[\operatorname{diag}(\boldsymbol{\psi}')] = |\{i \in [n] : \psi'(y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}) = 1\}|$  is the number of observations  $i = 1, \dots, n$  such that  $y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$  falls in the range where  $\rho$  is quadratic,  $\hat{S} = \{j \in [p] : \hat{\beta}_j \neq 0\}$  and  $\mathbf{X}_{\hat{S}} \in \mathbb{R}^{n \times |\hat{S}|}$  is the submatrix of  $\mathbf{X}$  containing columns indexed in  $\hat{S}$ .

### 2.6. Simulation study

We provide here simulations to showcase the asymptotic normality in the Huber loss and Elastic-Net penalty example of the previous section.

We set  $n = 200, p = 300$ , and generate  $\varepsilon_i \sim N(0, 1)$  and the coordinates of  $\boldsymbol{\beta}$  from iid Bernoulli variables with parameter 0.1. We compute 1000 simulations of the Z-score  $\Omega_{jj}^{1/2} \xi'_j$  from Theorem 2.1 for  $j = 1$  for the  $M$ -estimator with the Huber loss and the Elastic-Net penalty (2.17) with  $\sigma = 1$  and the four combinations  $(\lambda, \tau) \in \{n^{-1/2}, 2n^{-1/2}\} \times \{0, 0.1\}$ . The covariance matrix  $\boldsymbol{\Sigma}$  is set as  $\mathbf{R}^\top \mathbf{R}/(2p)$  where  $\mathbf{R} \in \mathbb{R}^{2p \times p}$  has iid Rademacher entries;  $\boldsymbol{\Sigma}$  is generated once and is the same across the 1000 simulations. The average value and standard

$(\lambda, \tau)$	$(n^{-1/2}, 0.1)$	$(n^{-1/2}, 0)$	$(2n^{-1/2}, 0.1)$	$(2n^{-1/2}, 0)$
$\hat{n}$	81.2 $\pm$ 5.8	111.9 $\pm$ 6.4	38.3 $\pm$ 4.6	49.3 $\pm$ 5.5
$\hat{\text{df}}$	53.6 $\pm$ 3.9	81.4 $\pm$ 5.1	14.2 $\pm$ 2.3	24.8 $\pm$ 4.1
$ \hat{S} $	96.4 $\pm$ 6.8	81.4 $\pm$ 5.1	26.5 $\pm$ 4.6	24.8 $\pm$ 4.1
$\sqrt{\hat{V}/n}$	0.44 $\pm$ 0.06	0.37 $\pm$ 0.07	0.56 $\pm$ 0.09	0.55 $\pm$ 0.11
$\Omega_{jj}^{1/2} \xi'_j$				

TABLE 2

Averages and standard errors over 1000 simulations of  $(\hat{n}, \hat{\text{df}}, |\hat{S}|, \hat{V}^{1/2}n^{-1/2})$ , as well as histograms and QQ-plots for  $\Omega_{jj}^{1/2} \xi'_j$  given in (2.11). The coordinate  $j$  is always  $j = 1$ . The noise  $\varepsilon_i$  is set as iid standard Cauchy. The red line on the QQ-plots is the diagonal line with equation  $x = y$ .

$(\lambda, \tau)$	$(n^{-1/2}, 0.1)$	$(n^{-1/2}, 0)$	$(2n^{-1/2}, 0.1)$	$(2n^{-1/2}, 0)$
$\hat{n}$	90.0 $\pm$ 5.8	125.7 $\pm$ 6.4	42.5 $\pm$ 4.8	55.4 $\pm$ 5.7
$\hat{\text{df}}$	57.8 $\pm$ 3.8	82.1 $\pm$ 4.9	16.5 $\pm$ 2.4	28.1 $\pm$ 4.3
$ \hat{S} $	97.8 $\pm$ 6.5	82.1 $\pm$ 4.9	29.7 $\pm$ 4.5	28.1 $\pm$ 4.3
$\sqrt{\hat{V}/n}$	0.37 $\pm$ 0.04	0.25 $\pm$ 0.04	0.51 $\pm$ 0.08	0.48 $\pm$ 0.09
$\Omega_{jj}^{1/2} \xi'_j$				

TABLE 3

Averages and standard errors over 1000 simulations of  $(\hat{n}, \hat{\text{df}}, |\hat{S}|, \hat{V}^{1/2}n^{-1/2})$ , as well as histograms and QQ-plots for  $\Omega_{jj}^{1/2} \xi'_j$  given in (2.11). The coordinate  $j$  is always  $j = 1$ . The noise  $\varepsilon_i$  is set as iid t-distribution with degree of freedom 2. The red line on the QQ-plots is the diagonal line with equation  $x = y$ .

error over the 1000 simulations of  $\hat{n}$ ,  $\hat{\text{df}}$ ,  $|\hat{S}|$  and  $\hat{V}^{1/2}n^{-1/2}$  are presented in Table 2 for  $\varepsilon$  with iid standard Cauchy components and Table 3 for  $\varepsilon$  with iid components from the t-distribution with 2 degrees of freedom, together with

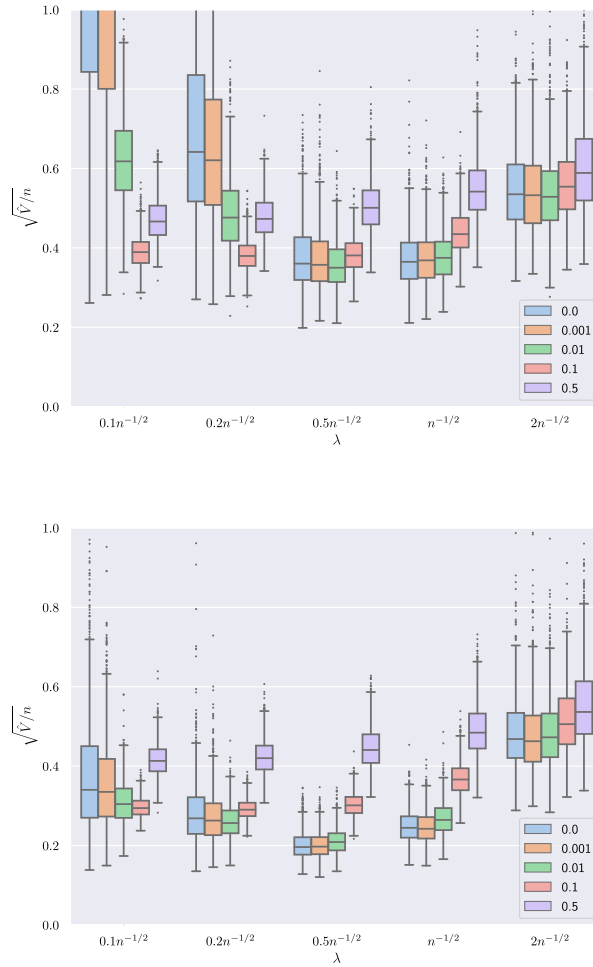


FIG 1. Boxplots of simulated  $\hat{V}^{1/2}n^{-1/2}$ . The simulation setup is described in Section 2.6. The top plot corresponds to iid standard Cauchy noise  $\varepsilon_i$ . The noise in the second plot is iid from the  $t$ -distribution with 2 degrees of freedom. Different colors correspond to different values of  $\tau \in \{0, 10^{-3}, 10^{-2}, 0.1, 0.5\}$ . In the  $x$ -axis,  $\lambda$  takes values in  $\{0.1n^{-1/2}, 0.2n^{-1/2}, 0.5n^{-1/2}, n^{-1/2}, 2n^{-1/2}\}$ .

histograms and QQ-plots against standard normal quantiles of  $\Omega_{jj}^{1/2}\xi'_j$  in (2.16).

The quantity  $\hat{V}^{1/2}n^{-1/2} = \|\psi\|_2/(\hat{n} - \hat{\text{df}})$  featured in the boxplots of Figure 1 characterizes the length of our confidence intervals in (2.13). Computing the values  $\hat{V}^{1/2}$  for different tuning parameters lets the practitioner pick the tuning parameters leading to the smallest confidence interval width, although this process amounts to the construction of multiple confidence intervals and warrants a Bonferroni multiple testing correction.

The histograms and QQ-plots in Table 2 and Table 3 confirm the normality

of  $\xi'_j$  for these two heavy-tailed continuous noise distributions.

Since  $\hat{n} - \hat{\mathbf{d}}\mathbf{f}$  appears in the denominator of  $\hat{V}$ , the length of confidence intervals can be large if  $\hat{n} - \hat{\mathbf{d}}\mathbf{f}$  is nearly zero and the length is infinite if  $\hat{n} - \hat{\mathbf{d}}\mathbf{f} = 0$ . This explains the large values and large variances observed in the boxplots of Figure 1 for small tuning parameters.

## Appendix A: Proof of the main result

Throughout,  $\phi_{\min}(\mathbf{M})$  and  $\phi_{\max}(\mathbf{M})$  denote the smallest and largest eigenvalues of a positive definite matrix  $\mathbf{M}$ . We also recall the important notation defined in (2.6)-(2.7):  $\mathbf{h} = \hat{\beta} - \beta$  is the error vector,  $\psi$  and  $\psi'$  are the random vectors in  $\mathbb{R}^n$  with coordinates  $(\psi(y_i - \mathbf{x}_i^\top \hat{\beta}))_{i \in [n]}$  and  $(\psi'(y_i - \mathbf{x}_i^\top \hat{\beta}))_{i \in [n]}$  respectively. Finally,

$$\Omega_{jj} = \mathbf{e}_j^\top \Sigma^{-1} \mathbf{e}_j, \quad \mathbf{z}_j = \mathbf{X} \Sigma^{-1} \Omega_{jj}^{-1} \mathbf{e}_j, \quad \mathbf{Q}_j = \mathbf{I}_p - \Sigma^{-1} \Omega_{jj}^{-1} \mathbf{e}_j \mathbf{e}_j^\top$$

and  $\mathbf{z}_j \sim N(\mathbf{0}, \Omega_{jj}^{-1} \mathbf{I}_n)$  is independent of  $\mathbf{X} \mathbf{Q}_j$  (cf. Proposition D.1).

### A.1. Supporting propositions

The proof of Theorem 2.1 relies on the two intermediary results given below. Proposition A.1 will be proved in Appendix B and Lemma A.2 in Appendix C.

**Proposition A.1.** *Let  $n \geq 3$ ,  $R > 0$  and  $\mathbf{z} \sim \text{Unif}(\mathbb{S}^{n-1}(R))$ , where  $\mathbb{S}^{n-1}(R)$  is sphere of radius  $R$  in  $\mathbb{R}^n$ . Assume either:*

- $\mathbf{f}$  is locally Lipschitz on  $\mathbb{S}^{n-1}(R)$  with  $\mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})\|_F^2] < \infty$ .
- $\mathbf{f}(\mathbf{z})$  is of the form  $\tilde{\mathbf{f}}(\mathbf{z})/\|\tilde{\mathbf{f}}(\mathbf{z})\|_2$  where  $\tilde{\mathbf{f}}: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is locally Lipschitz and satisfies  $\mathbb{P}(\|\tilde{\mathbf{f}}(\mathbf{z})\|_2 \neq 0) = 1$  and  $\mathbb{E}[\|\tilde{\mathbf{f}}(\mathbf{z})\|_2^{-2} \|\nabla \tilde{\mathbf{f}}(\mathbf{z})\|_F^2] < +\infty$ .

Define  $\mathbf{P}_z^\perp = \mathbf{I}_n - \mathbf{z} \mathbf{z}^\top / \|\mathbf{z}\|_2^2$  and  $\xi_{\mathbf{f}}(\mathbf{z}) = \mathbf{f}(\mathbf{z})^\top \mathbf{z} - R^2 n^{-1} \text{tr}(\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp)$ . Then

$$\mathbb{E}[(\xi_{\mathbf{f}}(\mathbf{z}) - \mathbb{E}[\mathbf{f}(\mathbf{z})]^\top \mathbf{z})^2] \leq 2R^4(n^2 - 2n)^{-1} \mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp\|_F^2], \quad (\text{A.1})$$

$$\mathbb{E}[\|\mathbf{f}(\mathbf{z})\|_2 - \|\mathbb{E}[\mathbf{f}(\mathbf{z})]\|_2]^2 \leq R^2(n - 2)^{-1} \mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp\|_F^2]. \quad (\text{A.2})$$

**Lemma A.2.** *Let Assumptions A, B and C be fulfilled. Let  $\eta_n = \sqrt{2 \log(n)/n} + n^{-1/2}$ , define the events*

$$\mathcal{E}_j = \{\|\mathbf{X} \mathbf{Q}_j \Sigma^{-1/2} n^{-1/2}\|_{op} \leq 1 + \sqrt{p/n} + \eta_n\} \cap \{|n^{-1/2} \Omega_{jj}^{1/2} \|\mathbf{z}_j\|_2 - 1| \leq \eta_n$$

for each  $j \in [p]$ , define  $u \geq 0$  by

$$u_* = [K^2(1 - 1/n)(L\tau^{-1}\|\Sigma\|_{op}(2 + \sqrt{p/n} + 2\eta_n)^2 + 1)^{-1} - 3L/n]_+,$$

and let  $\mathbb{E}_j$  be the conditional expectation given  $(\|\mathbf{z}_j\|_2, \mathbf{X}\mathbf{Q}_j, \varepsilon)$ . Then, when  $u_* > 0$

$$\sup_{\delta > 0} \mathbb{E} \left[ \frac{n}{\|\boldsymbol{\psi}\|_2^2 + \delta} \sum_{j \in [p]} I_{\mathcal{E}_j} h_j^2 \right]^{1/2} \leq \frac{[(1 + \sqrt{\frac{p}{n}})^2 + \frac{1}{n}]^{1/2}}{\phi_{\min}(\boldsymbol{\Sigma})^{1/2} (1 - \eta_n)_+^2 u_*} + \frac{[\frac{2p}{n\tau} + \mathbb{E}[\|\mathbf{h}\|_2^2]]^{1/2}}{u_*^{1/2}}. \quad (\text{A.3})$$

## A.2. Proof of the main result

*Proof of Theorem 2.1.* Since  $\boldsymbol{\psi} \neq \mathbf{0}_n$  for almost every  $\mathbf{X}$  by Proposition D.2,  $\xi_j$  is well-defined with  $\mathbb{P}$ -probability 1. By Lemma A.2 and the monotone convergence theorem as  $\delta \rightarrow 0$  for the left-hand side (A.3), when  $u_* > 0$

$$\mathbb{E} \left[ \frac{n}{\|\boldsymbol{\psi}\|_2^2} \sum_{j \in [p]} I_{\mathcal{E}_j} h_j^2 \right] \leq \frac{[(1 + \sqrt{\frac{p}{n}})^2 + \frac{1}{n}]^{1/2}}{\phi_{\min}(\boldsymbol{\Sigma})^{1/2} (1 - \eta_n)_+^2 u_*} + \frac{[\frac{2p}{n\tau} + \mathbb{E}[\|\mathbf{h}\|_2^2]]^{1/2}}{u_*^{1/2}}.$$

Under our assumptions,  $\|\boldsymbol{\Sigma}\|_{\text{op}} \leq 1/\kappa < +\infty$  and  $\mathbb{E}[\|\boldsymbol{\Sigma}^{1/2} \mathbf{h}\|_2^2] \leq \mathcal{R} < +\infty$ , so that there exists some finite constant  $\mathcal{N} > 0$  and  $\mathcal{A} < +\infty$  independent of  $n, p$ , such that for  $n \geq \mathcal{N}$ ,

$$\sum_{j \in [p]} \mathbb{E} \left[ n \|\boldsymbol{\psi}\|_2^{-2} I_{\mathcal{E}_j} h_j^2 \right] \leq \mathcal{A} < +\infty. \quad (\text{A.4})$$

By Markov's inequality with respect to the uniform distribution on  $[p] = \{1, \dots, p\}$ , the set

$$J_{n,p} := \{j \in [p] : \mathbb{E} [n I_{\mathcal{E}_j} \|\boldsymbol{\psi}\|_2^{-2} h_j^2] \leq \mathcal{A}/a_p\} \quad \text{satisfies} \quad |J_{n,p}|/p \geq 1 - a_p/p. \quad (\text{A.5})$$

While the function  $\phi(\varepsilon, \mathbf{X})$  in (2.8) is formally  $\mathbb{R}^n \times \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ , in this paragraph it is useful to only consider variations in  $\mathbf{z}_j$ , holding a fixed value of  $(\varepsilon, \mathbf{X}\mathbf{Q}_j)$ . Since  $(\varepsilon, \mathbf{X}\mathbf{Q}_j)$  is independent of  $\mathbf{z}_j$ , the conditional probability distribution of  $\mathbf{z}_j$  given  $(\varepsilon, \mathbf{X}\mathbf{Q}_j)$  is still  $N(\mathbf{0}, \Omega_{jj}^{-1} \mathbf{I}_n)$ . To this end and with a slight abuse of notation, for a given, fixed value of  $(\varepsilon, \mathbf{X}\mathbf{Q}_j)$  we view  $\boldsymbol{\psi}$  as a function of  $\mathbf{z}_j$  only,

$$\boldsymbol{\psi} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \boldsymbol{\psi} : \mathbf{z}_j \mapsto \boldsymbol{\psi}(\mathbf{z}_j) = \boldsymbol{\psi}(\varepsilon, \mathbf{X}\mathbf{Q}_j + \mathbf{z}_j \mathbf{e}_j^\top) \quad (\text{A.6})$$

and we denote its Jacobian by  $\nabla \boldsymbol{\psi}(\mathbf{z}_j)^\top$  at any point  $\mathbf{z}_j$  where (A.6) is Fréchet differentiable. Next, we argue conditionally on  $(\mathbf{X}\mathbf{Q}_j, \|\mathbf{z}_j\|_2)$ : Since  $\mathbf{z}_j/\|\mathbf{z}_j\|_2$  is independent of  $(\|\mathbf{z}_j\|_2, \mathbf{X}\mathbf{Q}_j)$  and by rotational invariance of the Gaussian distribution, conditionally on  $(\|\mathbf{z}_j\|_2, \mathbf{X}\mathbf{Q}_j)$  the vector  $\mathbf{z}_j$  is uniformly distributed on the sphere  $\mathbb{S}^{n-1}(\|\mathbf{z}_j\|_2)$ . Let  $\mathbb{E}_j$  denote the conditional expectation of  $\mathbf{z}_j$  given  $(\mathbf{X}\mathbf{Q}_j, \|\mathbf{z}_j\|_2)$ . By Proposition C.4, the above function is locally Lipschitz. From Proposition C.4 (iii) we have that

$$I_{\mathcal{E}_j} \mathbb{E}_j [\|\boldsymbol{\psi}\|_2^{-2} \|\nabla \boldsymbol{\psi}(\mathbf{z}_j)\|_F^2] \leq L(n\tau)^{-1} + I_{\mathcal{E}_j} \mathbb{E}_j [nL^2 \|\boldsymbol{\psi}\|_2^{-2} h_j^2] \quad (\text{A.7})$$

and the right-hand side is finite with probability one with respect to  $(\mathbf{X}\mathbf{Q}_j, \|\mathbf{z}_j\|_2)$  thanks to (A.4) and Tonelli's theorem for non-negative measurable functions. By Proposition D.2, conditional on almost every  $(\mathbf{X}\mathbf{Q}_j, \|\mathbf{z}_j\|_2)$ ,  $\boldsymbol{\psi}(\mathbf{z}_j) \neq \mathbf{0}_n$  for almost every  $\mathbf{z}_j \in \mathbb{R}$ . We are now in position to apply Proposition A.1 with  $\mathbf{f} = \boldsymbol{\psi}/\|\boldsymbol{\psi}\|_2$ ,  $\mathbf{z} = \mathbf{z}_j$  and

$$\xi_{\mathbf{f}}(\mathbf{z}_j) := \|\boldsymbol{\psi}\|_2^{-1} (\boldsymbol{\psi}^\top \mathbf{z}_j - (\|\mathbf{z}_j\|_2^2/n) \operatorname{tr}[\mathbf{P}_{\boldsymbol{\psi}}^\perp (\nabla \boldsymbol{\psi}(\mathbf{z}_j))^\top \mathbf{P}_{\mathbf{z}_j}^\perp]),$$

where  $\mathbf{P}_{\mathbf{v}}^\perp := \mathbf{I}_n - \mathbf{P}_{\mathbf{v}}$  with  $\mathbf{P}_{\mathbf{v}} := \mathbf{v}\mathbf{v}^\top/\|\mathbf{v}\|_2^2$  for any  $\mathbf{v} \in \mathbb{R}^n$ . By Proposition A.1 and (A.7),

$$\begin{aligned} & I_{\mathcal{E}_j} \mathbb{E}_j[(\xi_{\mathbf{f}}(\mathbf{z}_j) - \mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]^\top \mathbf{z}_j)^2] \\ & \leq 2I_{\mathcal{E}_j} \|\mathbf{z}_j\|_2^4 (n^2 - 2n)^{-1} \mathbb{E}_j[\|(\nabla \mathbf{f}(\mathbf{z}_j))^\top \mathbf{P}_{\mathbf{z}_j}^\perp\|_F^2] \\ & \leq 2(1 + \eta_n)^4 (1 - 2/n)^{-1} I_{\mathcal{E}_j} \Omega_{jj}^{-2} \mathbb{E}_j[\|\boldsymbol{\psi}\|_2^{-2} \|\nabla \boldsymbol{\psi}(\mathbf{z}_j)\|_F^2] \\ & \leq 2(1 + \eta_n)^4 (1 - 2/n)^{-1} \kappa^{-2} (L(n\tau)^{-1} + I_{\mathcal{E}_j} \mathbb{E}_j[nL^2 \|\boldsymbol{\psi}\|_2^{-2} h_j^2]) \end{aligned}$$

and

$$\begin{aligned} & I_{\mathcal{E}_j} (\|\mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]\|_2 - 1)^2 \\ & \leq I_{\mathcal{E}_j} \|\mathbf{z}_j\|_2^2 (n - 2)^{-1} \mathbb{E}_j[\|(\nabla \mathbf{f}(\mathbf{z}_j))^\top \mathbf{P}_{\mathbf{z}_j}^\perp\|_F^2] \\ & \leq (1 + \eta_n)^2 (1 - 2/n)^{-1} I_{\mathcal{E}_j} \Omega_{jj}^{-1} \mathbb{E}_j[\|\boldsymbol{\psi}\|_2^{-2} \|\nabla \boldsymbol{\psi}(\mathbf{z}_j)\|_F^2] \\ & \leq (1 + \eta_n)^2 (1 - 2/n)^{-1} \kappa^{-1} (L(n\tau)^{-1} + I_{\mathcal{E}_j} \mathbb{E}_j[nL^2 \|\boldsymbol{\psi}\|_2^{-2} h_j^2]), \end{aligned}$$

where the upper bounds follow from  $I_{\mathcal{E}_j} \|\mathbf{z}_j\|_2^4 \leq (1 + \eta_n)^4 n^2 \Omega_{jj}^{-2}$ . Taking  $\mathbb{E}$  on both sides, we obtain  $\max_{j \in J_{n,p}} \mathbb{E}[I_{\mathcal{E}_j} (\xi_{\mathbf{f}}(\mathbf{z}_j) - \mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]^\top \mathbf{z}_j)^2 + I_{\mathcal{E}_j} \|\mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]\|_2 - 1] \rightarrow 0$ . Thanks to  $\min_{j \in [p]} \mathbb{P}(\mathcal{E}_j) \rightarrow 1$  by Lemma C.1 this implies that both  $\|\mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]\|_2 - 1$  and  $|\xi_{\mathbf{f}}(\mathbf{z}_j) - \mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]^\top \mathbf{z}_j|$  converge in probability to 0 uniformly over  $j \in J_{n,p}$ .

We now study the asymptotic distribution of  $\xi_{\mathbf{f}}(\mathbf{z}_j)$ . Since  $\mathbf{z}_j \sim N(\mathbf{0}_n, \Omega_{jj}^{-1} \mathbf{I}_n)$  is independent with  $\mathbf{X}\mathbf{Q}_j$  by Proposition D.1, without loss of generality, we can assume that  $\mathbf{z}_j = \|\mathbf{z}_j\|_2 \boldsymbol{\zeta}_j / \|\boldsymbol{\zeta}_j\|_2$  for some  $\boldsymbol{\zeta}_j \sim N(\mathbf{0}_n, \mathbf{I}_n)$  independent of  $(\mathbf{X}\mathbf{Q}_j, \|\mathbf{z}_j\|_2)$ . Then  $\mathbb{E}_j$  coincides with the conditional expectation of  $\boldsymbol{\zeta}_j$  given  $(\mathbf{X}\mathbf{Q}_j, \|\mathbf{z}_j\|_2)$ . After some rearrangement,

$$\xi_{\mathbf{f}}(\mathbf{z}_j) = \xi_{\mathbf{f}}(\mathbf{z}_j) - \mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]^\top \mathbf{z}_j + \|\mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]\|_2 \frac{\|\mathbf{z}_j n^{-1/2}\|_2}{\|\boldsymbol{\zeta}_j n^{-1/2}\|_2} \left( \frac{\mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]}{\|\mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]\|_2} \right)^\top \boldsymbol{\zeta}_j.$$

Uniformly over  $j \in J_{n,p}$ , we have (i) the limit  $|\xi_{\mathbf{f}}(\mathbf{z}_j) - \mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]^\top \mathbf{z}_j| \rightarrow 0$  and the limit  $\|\mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]\|_2 \rightarrow 1$  both in probability, (ii)  $\frac{\mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]^\top}{\|\mathbb{E}_j[\mathbf{f}(\mathbf{z}_j)]\|_2} \boldsymbol{\zeta}_j \sim N(0, 1)$  and (iii) the limit  $\|\Omega_{jj}^{1/2} \mathbf{z}_j\|_2 / \|\boldsymbol{\zeta}_j\|_2 \rightarrow 1$  in probability. By Slutsky's Theorem,  $\Omega_{jj}^{1/2} \xi_{\mathbf{f}}(\mathbf{z}_j)$  converges in distribution to  $N(0, 1)$  uniformly over  $j \in J_{n,p}$ . That is, for any  $t \in \mathbb{R}$ ,  $\max_{j \in J_{n,p}} |\mathbb{P}(\Omega_{jj}^{1/2} \xi_{\mathbf{f}}(\mathbf{z}_j) \leq t) - \Phi(t)| \rightarrow 0$ .

It remains to relate  $\xi_{\mathbf{f}}(\mathbf{z}_j)$  to  $\xi_j$  defined in (2.10). As the term  $\|\boldsymbol{\psi}\|_2^{-1} \mathbf{z}_j^\top \boldsymbol{\psi}$  present in both  $\xi_j$  and  $\xi_{\mathbf{f}}(\mathbf{z}_j)$  cancel out, we have the decomposition

$$\|\boldsymbol{\psi}\|_2 \|\mathbf{z}_j\|_2^{-2} n (\xi_j - \xi_{\mathbf{f}}(\mathbf{z}_j))$$



$$= (\hat{\beta}_j - \beta_j) \text{tr}[\nabla_{\mathbf{y}} \boldsymbol{\psi}^\top] + \text{tr}(\mathbf{P}_{\boldsymbol{\psi}}^\perp (\nabla \boldsymbol{\psi}(\mathbf{z}_j))^\top \mathbf{P}_{\mathbf{z}_j}^\perp) \quad (\text{A.8})$$

$$= \text{tr}[(\hat{\beta}_j - \beta_j) \nabla_{\mathbf{y}} \boldsymbol{\psi}^\top + \nabla \boldsymbol{\psi}(\mathbf{z}_j)^\top \mathbf{P}_{\boldsymbol{\psi}}^\perp] \quad (\text{A.9})$$

$$+ (\hat{\beta}_j - \beta_j) \text{tr}[\mathbf{P}_{\boldsymbol{\psi}} \nabla_{\mathbf{y}} \boldsymbol{\psi}^\top] + \text{tr}[\mathbf{P}_{\boldsymbol{\psi}}^\perp \nabla \boldsymbol{\psi}(\mathbf{z}_j)^\top \mathbf{P}_{\mathbf{z}_j}]. \quad (\text{A.10})$$

The matrix inside the trace in (A.8) is zero thanks to (C.24). It follows that only the two terms in (A.10) remain, hence

$$\mathbb{E}[I_{\mathcal{E}_j}(\xi_j - \xi_{\mathbf{f}}(\mathbf{z}_j))^2] \leq \frac{2}{n} \mathbb{E}[I_{\mathcal{E}_j} \|\mathbf{z}_j\|_2^2 \|\boldsymbol{\psi}\|_2^{-2} ((\hat{\beta}_j - \beta_j)^2 \|\nabla_{\mathbf{y}} \boldsymbol{\psi}\|_{\text{op}}^2 + \|\nabla \boldsymbol{\psi}(\mathbf{z}_j)\|_{\text{op}}^2)].$$

Since  $I_{\mathcal{E}_j} \|\mathbf{z}_j\|_2^2 / n \leq \Omega_{jj}^{-1} (1 + \eta_n)^2$  and  $\|\nabla_{\mathbf{y}} \boldsymbol{\psi}\|_{\text{op}} \leq L$  by Proposition C.4 (i) with  $\mathbf{X} = \widetilde{\mathbf{X}}$ , the first term is bounded from above by  $L \Omega_{jj}^{-1} (1 + \eta_n)^2 \mathbb{E}[I_{\mathcal{E}_j} \|\boldsymbol{\psi}\|_2^{-2} h_j^2]$  which converges to 0 uniformly over  $j \in J_{n,p}$  by definition of  $J_{n,p}$  in (A.5). The second term also converges to 0 uniformly over  $j \in J_{n,p}$  thanks to (A.7). Thus  $\max_{j \in J_{n,p}} \mathbb{E}[I_{\mathcal{E}_j}(\xi_j - \xi_{\mathbf{f}}(\mathbf{z}_j))^2] \rightarrow 0$  which implies that  $|\xi_j - \xi_{\mathbf{f}}(\mathbf{z}_j)| \rightarrow 0$  in probability uniformly over  $j \in J_{n,p}$  and Slutsky's theorem completes the proof of (2.12) for  $\xi_j$ .

To prove a similar result for  $\xi'_j$ , it is enough to prove  $\Omega_j^{1/2} |\xi'_j - \xi_j| \xrightarrow{\mathbb{P}} 0$  uniformly over  $j \in J_{n,p}$  by Slutsky's theorem. As

$$|\xi_j - \xi'_j| = |\text{tr}[\nabla_{\mathbf{y}} \boldsymbol{\psi}]| \|\boldsymbol{\psi}\|_2^{-1} |h_j| \left| \|\mathbf{z}_j\|_2^2 / n - \Omega_{jj}^{-1} \right|$$

and  $|\text{tr}[\nabla_{\mathbf{y}} \boldsymbol{\psi}]| \leq nL$ , by the Cauchy-Schwarz inequality we find

$$\mathbb{E}[I_{\mathcal{E}_j} |\xi_j - \xi'_j|] \leq \Omega_{jj}^{-1} nL \mathbb{E}[I_{\mathcal{E}_j} h_j^2 \|\boldsymbol{\psi}\|_2^{-2}]^{1/2} \mathbb{E}[(\Omega_{jj} \|\mathbf{z}_j\|_2^2 / n - 1)^2]^{1/2}.$$

Since  $\mathbb{E}[(\Omega_{jj} \|\mathbf{z}_j\|_2^2 / n - 1)^2]^{1/2} = \sqrt{2/n}$  and  $\Omega_{jj} \in [\kappa, 1/\kappa]$ , the previous display converges to 0 uniformly over  $j \in J_{n,p}$  by definition of  $J_{n,p}$  in (A.5).  $\square$

## Appendix B: Stein formulae on the sphere

The goal of this section is to prove Proposition A.1 and to develop Stein formulae for random vectors  $\mathbf{z}$  uniformly distributed on the sphere.

Let  $\mathbb{S}^{n-1}(R)$  be the sphere in  $\mathbb{R}^n$  with center  $\mathbf{0}$  and radius  $R > 0$ . We say that  $\mathbf{z}$  is uniformly distributed in  $\mathbb{S}^{n-1}(R)$  and write  $\mathbf{z} \sim \text{Unif}(\mathbb{S}^{n-1}(R))$  if  $\mathbf{z}$  is equal in distribution to  $R\boldsymbol{\zeta}/\|\boldsymbol{\zeta}\|_2$  where  $\boldsymbol{\zeta} \sim N(\mathbf{0}, \mathbf{I}_n)$ . We first develop Stein's formulae with respect to  $\mathbf{z} \sim \text{Unif}(\mathbb{S}^{n-1}(R))$  for functions  $\mathbf{f} : \mathbf{z} \in \mathbb{S}^{n-1}(R) \mapsto \mathbf{f}(\mathbf{z}) \in \mathbb{R}^n$  in Sobolev spaces over  $\mathbb{S}^{n-1}(R)$ .

We derive next Stein formulae for functions in Sobolev spaces over  $\mathbb{S}^{n-1}(R)$ . One possible construction of such Sobolev spaces is obtained by completion of the space of infinitely differentiable functions  $\mathbb{S}^{n-1}(R) \rightarrow \mathbb{R}$  with respect to the desired Sobolev norm as follows. Here,  $\mathbb{S}^{n-1}(R)$  is viewed as a compact Riemannian manifold equipped with the canonical metric (the metric induced as a submanifold of  $\mathbb{R}^n$  equipped with the Euclidean metric). As it will be convenient for compatibility with the rest of the paper to conserve the partial derivatives with respect to the canonical basis in  $\mathbb{R}^n$ , we adopt the

following notation. For a smooth function  $h : \mathbb{S}^{n-1}(R) \rightarrow \mathbb{R}$  and  $\Omega \subset \mathbb{R}^n$  an open neighborhood of  $\mathbb{S}^{n-1}(R)$ , define the smooth function  $\check{h} : \Omega \rightarrow \mathbb{R}$  by  $\check{h}(\mathbf{x}) = h(R\mathbf{x}/\|\mathbf{x}\|_2)$  and define the gradient of  $h$  as that of  $\check{h}$ , i.e., for  $\mathbf{x} \in \mathbb{S}^{n-1}(R)$ , set  $\nabla h(\mathbf{x}) = ((\partial/\partial x_1)\check{h}(\mathbf{x}), \dots, (\partial/\partial x_n)\check{h}(\mathbf{x}))$ . For every  $\mathbf{x} \in \mathbb{S}^{n-1}$ , the gradient  $\nabla h(\mathbf{x})$  belongs to the hyperplane orthogonal to  $\mathbf{x}$  which is the tangent space of  $\mathbb{S}^{n-1}$  at  $\mathbf{x}$ . Furthermore if  $\gamma : \mathbb{R} \rightarrow \mathbb{S}^{n-1}(R)$  is a smooth curve with  $\gamma(0) = \mathbf{x}$ ,  $(d/dt)\gamma(t)|_{t=0} = \mathbf{v}$  then  $\mathbf{v}^\top \mathbf{x} = 0$  and  $\nabla h(\mathbf{x})^\top \mathbf{v} = (d/dt)h(\gamma(t))|_{t=0}$ . For such smooth function  $h$  and equipped with its gradient, for  $\alpha \in \{1, 2\}$  we define the Sobolev norm

$$\|h\|_{1,\alpha} = \mathbb{E}[|h(\mathbf{z})|^\alpha]^{1/\alpha} + \mathbb{E}[\|\nabla h(\mathbf{z})\|_2^\alpha]^{1/\alpha}, \quad \mathbf{z} \sim \text{Unif}(\mathbb{S}^{n-1}(R))$$

and the Sobolev space  $W^{1,\alpha}(\mathbb{S}^{n-1}(R))$  as the completion of the space of smooth functions over  $\mathbb{S}^{n-1}(R)$  with respect to the above norm. This definition is equivalent to the definition given in [18, Definition 2.1]. By [18, Proposition 2.3], since the manifold  $\mathbb{S}^{n-1}(R)$  is compact the resulting Sobolev spaces do not depend on the chosen metric. Equivalently, the Sobolev space  $W^{1,\alpha}(\mathbb{S}^{n-1}(R))$  is also the completion with respect to the above norm of the space of once continuously differentiable functions on  $\mathbb{S}^{n-1}(R)$ .

If  $h$  is locally Lipschitz on  $\mathbb{S}^{n-1}(R)$  (i.e., every point has a neighborhood in  $\mathbb{S}^{n-1}(R)$  on which  $h$  is Lipschitz), then again by considering an open neighborhood  $\Omega \subset \mathbb{R}^n$  of  $\mathbb{S}^{n-1}(R)$ , the function  $\check{h}(\mathbf{x}) = h(R\mathbf{x}/\|\mathbf{x}\|_2)$  is locally Lipschitz on  $\Omega$ . Thus, in this case and by Rademacher's theorem,  $\nabla h(\mathbf{z})$  is well defined almost everywhere in  $\mathbb{S}^{n-1}(R)$  as the gradient of  $\check{h}(\mathbf{x})$ , and  $h \in W^{1,\alpha}(\mathbb{S}^{n-1}(R))$  if and only if  $\mathbb{E}[\|\nabla h(\mathbf{z})\|_2^\alpha] < \infty$ . For example,  $h \in W^{1,\alpha}(\mathbb{S}^{n-1}(R))$  when  $h$  is  $L$ -Lipschitz on  $\mathbb{S}^{n-1}(R)$ .

Finally, for  $\alpha \in \{1, 2\}$  define  $W^{1,\alpha}(\mathbb{S}^{n-1}(R))^n$  as the space of  $\mathbb{R}^n$  valued functions  $\mathbf{f} = (f_1, \dots, f_n)$  with all coordinates  $f_i$  belonging to  $W^{1,\alpha}(\mathbb{S}^{n-1}(R))$ , equipped with the norm

$$\|\mathbf{f}\|_{1,\alpha} = \mathbb{E}[\|\mathbf{f}(\mathbf{z})\|_2^\alpha]^{1/\alpha} + \mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})\|_F^\alpha]^{1/\alpha}$$

where the gradient  $\nabla \mathbf{f}$  is the matrix in  $\mathbb{R}^{n \times n}$  with columns  $\nabla f_1, \dots, \nabla f_n$ ,

**Lemma B.1** (Stein's formula on the sphere). *Let  $n \geq 3$ ,  $R > 0$  and  $\mathbf{z} \sim \text{Unif}(\mathbb{S}^{n-1}(R))$ . Let  $\mathbf{P}_\mathbf{z}^\perp = \mathbf{I}_n - \mathbf{z}\mathbf{z}^\top/\|\mathbf{z}\|_2^2$  for  $\mathbf{z} \neq \mathbf{0}$ . Then, for all  $\mathbf{f} = (f_1, \dots, f_n) \in W^{1,1}(\mathbb{S}^{n-1}(R))^n$ ,*

$$\mathbb{E}[\mathbf{f}(\mathbf{z})^\top \mathbf{z}] = (n-1)^{-1} R^2 \mathbb{E}[\text{tr}((\nabla \mathbf{f}(\mathbf{z}))^\top \mathbf{P}_\mathbf{z}^\perp)], \quad (\text{B.1})$$

where  $\nabla \mathbf{f} = (\nabla f_1, \dots, \nabla f_n)$ . For all  $\mathbf{f} = (f_1, \dots, f_n) \in W^{1,2}(\mathbb{S}^{n-1}(R))^n$  we also have

$$\mathbb{E}[\|\mathbf{f}(\mathbf{z}) - \mathbb{E}[\mathbf{f}(\mathbf{z})]\|_2^2] \leq (n-2)^{-1} R^2 \mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp\|_F^2], \quad (\text{B.2})$$

$$\begin{aligned} & \mathbb{E}[(nR^{-2} \mathbf{f}(\mathbf{z})^\top \mathbf{z} - \text{tr}[\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp])^2] \\ &= nR^{-2} \mathbb{E}[\|\mathbf{f}\|_2^2] + (1 - 2/n)^{-1} \mathbb{E}[\text{tr}[(\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp)^2]] - \frac{2}{n-2} \mathbb{E}[\text{tr}(\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp)^2] \end{aligned} \quad (\text{B.3})$$

$$\leq nR^{-2} \mathbb{E}[\|\mathbf{f}\|_2^2] + (1 - 2/n)^{-1} \mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp\|_F^2]. \quad (\text{B.4})$$

Note that (B.2) is the classical Poincaré inequality on the sphere. A proof is provided for completeness.

*Proof of Lemma B.1.* As the operators  $T$  and  $T_{ij}$  defined by  $T\mathbf{f}(\mathbf{z}) = \mathbf{f}(\mathbf{z})^\top \mathbf{z}$  and  $T_{ij}\mathbf{f}(\mathbf{z}) = [\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp]_{ij}$  for every  $i, j = 1, \dots, n$  are all continuous linear mappings from  $W^{1,\alpha}(\mathbb{S}^{n-1}(R))^n$  to the  $L_\alpha$  space with the norm  $\|\mathbf{f}\|_{L_\alpha} = (\mathbb{E}[\|\mathbf{f}(\mathbf{z})\|^\alpha])^{1/\alpha}$ , we assume without loss of generality that all coordinates of  $\mathbf{f}$  are continuously differentiable. Indeed, if  $(\mathbf{f}^{(q)})_{q \geq 1}$  is a sequence of smooth functions over the sphere converging to  $\mathbf{f}$  in  $W^{1,\alpha}(\mathbb{S}^{n-1}(R))^n$  for  $\alpha = 1$  and (B.1) holds for  $\mathbf{f}^{(q)}$  then (B.1) also holds for the limit by continuity; the same argument applies with  $\alpha = 2$  for (B.2)-(B.3)-(B.4).

Let  $\boldsymbol{\zeta} \sim N(\mathbf{0}, \mathbf{I}_n)$ . Assume without loss of generality  $\mathbf{z} = R\boldsymbol{\zeta}/\|\boldsymbol{\zeta}\|_2$  as  $R\boldsymbol{\zeta}/\|\boldsymbol{\zeta}\|_2 \sim \text{Unif}(\mathbb{S}^{n-1}(R))$ . Let  $\phi(t)$  be a continuously differentiable function in  $\mathbb{R}$  with  $\phi(t) = 0$  for  $t \leq 0$  and  $\phi(t) = 1$  for  $t \geq 1$ . For  $\delta > 0$  define  $\boldsymbol{\varphi}_\delta(\mathbf{x}) = \phi(\|\mathbf{x}\|_2/\delta)\mathbf{f}(R\mathbf{x}/\|\mathbf{x}\|_2)$ . As  $\|\boldsymbol{\zeta}\|_2$  is independent of  $\mathbf{z}$  and  $\boldsymbol{\varphi}_\delta(\mathbf{x})$  has uniformly bounded gradient, the first order Stein formula for  $\boldsymbol{\varphi}_\delta$  yields

$$\mathbb{E}[\phi(\|\boldsymbol{\zeta}\|_2/\delta)\|\boldsymbol{\zeta}\|_2]\mathbb{E}[\mathbf{f}(\mathbf{z})^\top \mathbf{z}] = \mathbb{E}[R\boldsymbol{\varphi}_\delta(\boldsymbol{\zeta})^\top \boldsymbol{\zeta}] = R\mathbb{E}[\text{tr}(\nabla \boldsymbol{\varphi}_\delta(\boldsymbol{\zeta}))].$$

By the product and chain rules,  $\nabla \boldsymbol{\varphi}_\delta(\mathbf{x})$  is given by

$$\phi'(\|\mathbf{x}\|_2/\delta)(\mathbf{x}/(\|\mathbf{x}\|_2\delta))\mathbf{f}(R\mathbf{x}/\|\mathbf{x}\|_2)^\top + \phi(\|\mathbf{x}\|_2/\delta)(R/\|\mathbf{x}\|_2)\mathbf{P}_\mathbf{x}^\perp(\nabla \mathbf{f}(R\mathbf{x}/\|\mathbf{x}\|_2))$$

with  $\phi'(t) = (d/dt)\phi(t)$  and  $\phi'(0) = \phi'(1) = 0$ . As  $\sup_{\delta > 0, \mathbf{x} \in \mathbb{R}^n} \|\mathbf{x}\|_2 \|\nabla \boldsymbol{\varphi}_\delta(\mathbf{x})\|_F \leq C < \infty$  and  $\mathbb{E}[\|\boldsymbol{\zeta}\|_2^{-2}] = 1/(n-2) < \infty$ , the dominated convergence theorem gives

$$\mathbb{E}[\mathbf{f}(\mathbf{z})^\top \mathbf{z}] = \frac{R \lim_{\delta \rightarrow 0} \mathbb{E}[\text{tr}(\nabla \boldsymbol{\varphi}_\delta(\boldsymbol{\zeta}))]}{\lim_{\delta \rightarrow 0} \mathbb{E}[\phi(\|\boldsymbol{\zeta}\|_2/\delta)\|\boldsymbol{\zeta}\|_2]} = \frac{R^2 \mathbb{E}[1/\|\boldsymbol{\zeta}\|_2] \mathbb{E}[\text{tr}(\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp)]}{\mathbb{E}[\|\boldsymbol{\zeta}\|_2]},$$

which yields (B.1) due to  $\mathbb{E}[\|\boldsymbol{\zeta}\|_2]/\mathbb{E}[\|\boldsymbol{\zeta}\|_2^{-1}] = (n-1)$ .

Next, as the exchange of limit and expectation is allowed when  $\boldsymbol{\varphi}_\delta \rightarrow \boldsymbol{\varphi}_{0+} = \mathbf{f}$ , the Gaussian Poincaré inequality [9, Theorem 1.6.4] yields

$$\mathbb{E}[\|\mathbf{f}(\mathbf{z}) - \mathbb{E}[\mathbf{f}(\mathbf{z})]\|_2^2] \leq \lim_{\delta \rightarrow 0+} \mathbb{E}[\|\nabla \boldsymbol{\varphi}_\delta(\boldsymbol{\zeta})\|_F^2] = \mathbb{E}[R^2 \|\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp\|_F^2 \|\boldsymbol{\zeta}\|_2^{-2}].$$

Since  $\mathbb{E}[\|\boldsymbol{\zeta}\|_2^{-2}] = 1/(n-2)$  and  $(\|\boldsymbol{\zeta}\|_2, \mathbf{z})$  are independent, we obtain (B.2). Finally by applying the Second Order Stein formula of [6] to  $\boldsymbol{\varphi}_\delta(\boldsymbol{\zeta})$  we find by dominated convergence

$$\begin{aligned} & \mathbb{E}[(R^{-1}\|\boldsymbol{\zeta}\|_2 \mathbf{z}^\top \mathbf{f}(\mathbf{z}) - R\|\boldsymbol{\zeta}\|_2^{-1} \text{tr}[\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp])^2] \\ &= \lim_{\delta \rightarrow 0+} \mathbb{E}[(\boldsymbol{\zeta}^\top \boldsymbol{\varphi}_\delta(\boldsymbol{\zeta}) - \text{tr}[\nabla \boldsymbol{\varphi}_\delta(\boldsymbol{\zeta}))]^2] \\ &= \lim_{\delta \rightarrow 0+} \mathbb{E}[\|\boldsymbol{\varphi}_\delta(\boldsymbol{\zeta})\|_2^2] + \mathbb{E} \text{tr}[(\nabla \boldsymbol{\varphi}_\delta(\boldsymbol{\zeta}))^\top]^2 \\ &= \mathbb{E}[\|\mathbf{f}(\mathbf{z})\|_2^2] + R^2 \mathbb{E} \text{tr}[(\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp)^2] \mathbb{E}[\|\boldsymbol{\zeta}\|_2^{-2}] \\ &= \mathbb{E}[\|\mathbf{f}(\mathbf{z})\|_2^2] + R^2 \mathbb{E} \text{tr}[(\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_\mathbf{z}^\perp)^2]/(n-2), \end{aligned}$$

where the last equality follows from the independence of  $\mathbf{z}$  and  $\|\boldsymbol{\zeta}\|_2$ . We now simplify the left-hand side in order to get rid of  $\|\boldsymbol{\zeta}\|_2$ . By expanding the square, using the independence of  $(\mathbf{z}, \|\boldsymbol{\zeta}\|_2)$  and  $\mathbb{E}[\|\boldsymbol{\zeta}\|_2^{-2}] = \frac{1}{n-2}$ , the above display is equal to

$$\begin{aligned} & \begin{cases} R^{-2} \mathbb{E}[\|\boldsymbol{\zeta}\|_2^2] \mathbb{E}[(\mathbf{z}^\top \mathbf{f}(\mathbf{z}))^2] - 2 \mathbb{E}[\mathbf{z}^\top \mathbf{f}(\mathbf{z}) \operatorname{tr}[\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp]] \\ + R^2 \mathbb{E}[\|\boldsymbol{\zeta}\|_2^{-2}] \mathbb{E}[(\operatorname{tr}[\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp])^2] \end{cases} \\ &= \begin{cases} n R^{-2} \mathbb{E}[(\mathbf{z}^\top \mathbf{f}(\mathbf{z}))^2] - 2 \mathbb{E}[\mathbf{z}^\top \mathbf{f}(\mathbf{z}) \operatorname{tr}[\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp]] \\ + R^2 (n-2)^{-1} \mathbb{E}[(\operatorname{tr}[\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp])^2] \end{cases} \\ &= \begin{cases} \mathbb{E}[(\frac{\sqrt{n}}{R} \mathbf{z}^\top \mathbf{f}(\mathbf{z}) - \frac{R}{\sqrt{n}} \operatorname{tr}[\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp])^2] \\ + R^2 ((n-2)^{-1} - n^{-1}) \mathbb{E}[(\operatorname{tr}[\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp])^2]. \end{cases} \end{aligned}$$

Since  $1/(n-2) - 1/n = 2n^{-1}(n-2)^{-1}$ , we obtain (B.3) after multiplying by  $n/R^2$ . The proof is complete since (B.4) follows directly from (B.3) by the Cauchy-Schwarz inequality.  $\square$

*Proof of Proposition A.1.* If  $\mathbf{f}$  is locally Lipschitz and  $\mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})^\top\|_F^2] < \infty$ , then  $\mathbf{f} \in W^{1,2}(\mathbb{S}^{n-1}(R))$ . We consider the mapping  $\mathbf{f}(\mathbf{z}) - \mathbb{E}[\mathbf{f}(\mathbf{z})]$  rather than  $\mathbf{f}(\mathbf{z})$  in (B.3) and (B.4). Multiplying  $R^4 n^{-2}$  on both sides of the inequality (B.4), it provides

$$\begin{aligned} & \mathbb{E}[(\xi_{\mathbf{f}}(\mathbf{z}) - \mathbb{E}[\mathbf{f}(\mathbf{z})]^\top \mathbf{z})^2] \\ & \leq R^2 n^{-1} \mathbb{E}[\|\mathbf{f} - \mathbb{E}[\mathbf{f}(\mathbf{z})]\|_2^2] + R^4 (n^2 - 2n)^{-1} \mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp\|_F^2] \\ & \leq 2R^4 (n^2 - 2n)^{-1} \mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp\|_F^2], \end{aligned}$$

where the second inequality follows from (B.2). By the triangle inequality and (B.2),

$$\begin{aligned} \mathbb{E}[|\|\mathbf{f}(\mathbf{z})\|_2 - \|\mathbb{E}[\mathbf{f}(\mathbf{z})]\|_2|^2] & \leq \mathbb{E}[\|\mathbf{f}(\mathbf{z}) - \mathbb{E}[\mathbf{f}(\mathbf{z})]\|_2^2] \\ & \leq R^2 (n-2)^{-1} \mathbb{E}[\|\nabla \mathbf{f}(\mathbf{z})^\top \mathbf{P}_z^\perp\|_F^2]. \end{aligned}$$

If  $\mathbf{f}(\mathbf{z}) = \tilde{\mathbf{f}}(\mathbf{z})/\|\tilde{\mathbf{f}}(\mathbf{z})\|_2$  with locally Lipschitz  $\tilde{\mathbf{f}}$  on  $\mathbb{S}^{n-1}(R)$  then  $\tilde{\mathbf{f}}(\mathbf{z})/(\delta \vee \|\tilde{\mathbf{f}}(\mathbf{z})\|_2)$  is locally Lipschitz for  $\delta > 0$  and converges to  $\mathbf{f}$  in  $W^{1,2}(\mathbb{S}^{n-1}(R))^n$  as  $\delta \rightarrow 0+$  when  $\mathbb{P}(\|\mathbf{f}(\mathbf{z})\|_2 \neq 0) = 1$  and  $\mathbb{E}[\|\tilde{\mathbf{f}}(\mathbf{z})\|_2^{-2} \|\nabla \tilde{\mathbf{f}}(\mathbf{z})\|_F^2] < +\infty$ . Thus, the proof still applies.  $\square$

## Appendix C: Bounds on $(\hat{\beta}_j - \beta_j)^2 \|\psi\|_2^{-2}$

The goal of this section is to prove Lemma A.2.

### C.1. High probability events $\mathcal{E}_j$

**Lemma C.1** (high probability of  $\mathcal{E}_j$ ). *Assume that  $\mathbf{X}$  has iid  $N(\mathbf{0}, \boldsymbol{\Sigma})$  rows. Then*

$$\mathbb{E}[\|\mathbf{X} \boldsymbol{\Sigma}^{-1/2}\|_{op}^2] \leq (\sqrt{n} + \sqrt{p})^2 + 1. \quad (\text{C.1})$$

Furthermore, with  $\eta_n = \sqrt{2\log(n)/n} + n^{-1/2}$  and for the events

$$\mathcal{E}_j = \{\|\mathbf{X}\mathbf{Q}_j\boldsymbol{\Sigma}^{-1/2}n^{-1/2}\|_{\text{op}} \leq 1 + \sqrt{p/n} + \eta_n\} \cap \{|n^{-1/2}\Omega_{jj}^{1/2}\|\mathbf{z}_j\|_2 - 1| \leq \eta_n\},$$

we have  $\mathbb{P}(\cap_{j \in [p]} \mathcal{E}_j) \geq 1 - (p + 1/2)n^{-1}(\pi \log(n))^{-1/2}$ .

*Proof of Lemma C.1.* Let us first notice that  $\mathbf{X}\boldsymbol{\Sigma}^{-1/2}$  is a random Gaussian matrix with iid standard normal entries. Theorem 7.3.1 in [32] provides that  $\mathbb{E}[\|\mathbf{X}\boldsymbol{\Sigma}^{-1/2}\|_{\text{op}}] \leq \sqrt{n} + \sqrt{p}$ . Since the operator norm of a matrix is a 1-Lipschitz function of the entries of the matrix, by the Gaussian Poincaré inequality [10, Theorem 3.20],  $\text{Var}(\|\mathbf{X}\boldsymbol{\Sigma}^{-1/2}\|_{\text{op}}) \leq 1$ . This proves (C.1).

By Theorem II.13 in [14], we have for  $t > 0$ ,

$$\mathbb{P}(\|\mathbf{X}\boldsymbol{\Sigma}^{-1/2}n^{-1/2}\|_{\text{op}} \geq 1 + \sqrt{p/n} + t) \leq \Phi(-t\sqrt{n}),$$

$$\mathbb{P}(|\Omega_{jj}^{1/2}\|\mathbf{z}_jn^{-1/2}\|_2 - 1| \geq n^{-1/2} + t) \leq 2\Phi(-t\sqrt{n}).$$

Since  $\mathbf{X}\mathbf{Q}_j\boldsymbol{\Sigma}^{-1/2} = \mathbf{X}\boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\Sigma}^{1/2}\mathbf{Q}_j\boldsymbol{\Sigma}^{-1/2})$  and  $\|\boldsymbol{\Sigma}^{1/2}\mathbf{Q}_j\boldsymbol{\Sigma}^{-1/2}\|_{\text{op}} \leq 1$ ,

$$\|\mathbf{X}\mathbf{Q}_j\boldsymbol{\Sigma}^{-1/2}\|_{\text{op}} \leq \|\mathbf{X}\boldsymbol{\Sigma}^{-1/2}\|_{\text{op}}$$

for all  $j \in [p]$ . Next, using a union bound over  $j \in [p]$  when  $t = \sqrt{2\log(n)/n}$  we have  $\mathbb{P}(\cap_{j \in [p]} \mathcal{E}_j) \geq 1 - (2p + 1)\Phi(-\sqrt{2\log(n)})$ . We conclude the proof using  $\Phi(-t) \leq (2\pi)^{-1/2} \exp(-t^2/2)/t$  for  $t > 0$ , which provides

$$\mathbb{P}(\cap_{j \in [p]} \mathcal{E}_j) \geq 1 - (p + 1/2)n^{-1}(\pi \log(n))^{-1/2}. \quad \square$$

*Remark C.1.* Our specific construction of  $\mathcal{E}_j$  satisfies the following properties:

1.  $I_{\mathcal{E}_j}$  is a function of  $\|\mathbf{X}\mathbf{Q}_j\|_{\text{op}}$  and  $\|\mathbf{z}_j\|_2$  only, consequently the event  $\mathcal{E}_j$  is independent of  $\mathbf{z}_j/\|\mathbf{z}_j\|_2$ .
2.  $\|\mathbf{X}\boldsymbol{\Sigma}^{-1/2}n^{-1/2}\|_{\text{op}} I_{\mathcal{E}_j} \leq 1 + \sqrt{p/n} + \eta_n + (1 + \eta_n) = 2 + \sqrt{p/n} + 2\eta_n$ .

### C.2. Twice continuously differentiable penalty

**Lemma C.2.** Let  $L, \tau$  be such that Assumptions A and B are fulfilled. Further assume that the Hessian matrix  $\mathbf{G} = (\nabla^2 g(\mathbf{b}))|_{\mathbf{b}=\hat{\mathbf{b}}}$  of  $g$  at  $\hat{\mathbf{b}}$  exists and is symmetric, and define

$$\mathbf{M} = (\mathbf{X}^\top \text{diag}(\boldsymbol{\psi}')\mathbf{X} + n\mathbf{G})^{-1}, \quad (\text{C.2})$$

$$\mathbf{V} = \text{diag}(\boldsymbol{\psi}') - \text{diag}(\boldsymbol{\psi}')\mathbf{X}\mathbf{M}\mathbf{X}^\top \text{diag}(\boldsymbol{\psi}'). \quad (\text{C.3})$$

Then with the partial order  $\mathbf{A} \preceq \mathbf{B}$  if and only if the matrix  $\mathbf{B} - \mathbf{A}$  is positive semi-definite, we have

$$\|\text{diag}(\boldsymbol{\psi}')\mathbf{X}\mathbf{M}\|_{\text{op}} \leq (1/2)L^{1/2}(n\tau)^{-1/2}, \quad (\text{C.4})$$

$$\mathbf{M} \preceq (n\tau)^{-1}\mathbf{I}_p, \quad (\text{C.5})$$

$$(L\tau^{-1}\|\mathbf{X}n^{-1/2}\|_{\text{op}}^2 + 1)^{-1} \text{diag}(\boldsymbol{\psi}') \preceq \mathbf{V} \preceq \text{diag}(\boldsymbol{\psi}') \preceq L\mathbf{I}_n. \quad (\text{C.6})$$

*Proof of Lemma C.2.* Throughout the proof we use the notation

$$\mathbf{B} = \text{diag}(\boldsymbol{\psi}')^{1/2} \mathbf{X} n^{-1/2} \mathbf{G}^{-1/2}.$$

By some algebra, we have  $\mathbf{B}(\mathbf{B}^\top \mathbf{B} + \mathbf{I}_p)^{-1} \mathbf{B}^\top = \text{diag}(\boldsymbol{\psi}')^{1/2} \mathbf{X} \mathbf{M} \mathbf{X}^\top \text{diag}(\boldsymbol{\psi}')^{1/2}$ . For an upper bound of  $\text{diag}(\boldsymbol{\psi}') \mathbf{X} \mathbf{M}$ , we notice

$$\begin{aligned} \text{diag}(\boldsymbol{\psi}') \mathbf{X} \mathbf{M} &= \text{diag}(\boldsymbol{\psi}') \mathbf{X} (\mathbf{X}^\top \text{diag}(\boldsymbol{\psi}') \mathbf{X} + n \mathbf{G})^{-1} \\ &= \text{diag}(\boldsymbol{\psi}')^{1/2} \mathbf{B} (\mathbf{B}^\top \mathbf{B} + \mathbf{I}_p)^{-1} n^{-1/2} \mathbf{G}^{-1/2}. \end{aligned}$$

For any matrix  $\mathbf{B}$ ,  $\|\mathbf{B}(\mathbf{B}^\top \mathbf{B} + \mathbf{I}_p)^{-1}\|_{\text{op}} \leq \max_{t \geq 0} t/(t^2 + 1) = 1/2$ . By strongly convexity of  $g$ ,  $\|\mathbf{G}^{-1/2}\|_{\text{op}} \leq \tau^{-1/2}$ . Since  $\psi$  is  $L$ -Lipschitz,  $\|\text{diag}(\boldsymbol{\psi}')^{1/2}\|_{\text{op}} \leq L^{1/2}$ . Combining those upper bounds, we obtain (C.4). For the upper bound of  $\mathbf{M}$ , we notice  $\mathbf{M} = n^{-1} \mathbf{G}^{-1/2} (\mathbf{B}^\top \mathbf{B} + \mathbf{I}_p)^{-1} \mathbf{G}^{-1/2} \preceq n^{-1} \mathbf{G}^{-1} \preceq (n\tau)^{-1} \mathbf{I}_p$ . This gives (C.5). For lower and upper bounds on  $\mathbf{V}$ , we first notice that by definition of  $\mathbf{B}$ ,  $\|\mathbf{B}\|_{\text{op}} \leq L^{1/2} \tau^{-1/2} \|\mathbf{X} n^{-1/2}\|_{\text{op}}$ , thus

$$(L\tau^{-1} \|\mathbf{X} n^{-1/2}\|_{\text{op}}^2 + 1)^{-1} \mathbf{I}_n \preceq \mathbf{I}_n - \mathbf{B}(\mathbf{B}^\top \mathbf{B} + \mathbf{I}_p)^{-1} \mathbf{B}^\top \preceq \mathbf{I}_n.$$

Since  $\mathbf{V} = \text{diag}(\boldsymbol{\psi}')^{1/2} (\mathbf{I}_n - \mathbf{B}(\mathbf{B}^\top \mathbf{B} + \mathbf{I}_p)^{-1} \mathbf{B}^\top) \text{diag}(\boldsymbol{\psi}')^{1/2}$ , we have

$$(L\tau^{-1} \|\mathbf{X} n^{-1/2}\|_{\text{op}}^2 + 1)^{-1} \text{diag}(\boldsymbol{\psi}') \preceq \mathbf{V} \preceq \text{diag}(\boldsymbol{\psi}').$$

By the  $L$ -Lipschitz property of  $\psi$ , we have  $\text{diag}(\boldsymbol{\psi}') \preceq L \mathbf{I}_n$ . Thus (C.6) holds.  $\square$

**Proposition C.3.** Assume that  $g$  is strongly convex with parameter  $\tau > 0$  and  $\psi = \rho'$  is  $L$ -Lipschitz. Let  $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}(\boldsymbol{\varepsilon}, \mathbf{X})$ ,  $\boldsymbol{\psi} = \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X})$  be as in (2.8) and set  $\mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X}) = \hat{\boldsymbol{\beta}}(\boldsymbol{\varepsilon}, \mathbf{X}) - \boldsymbol{\beta}$ . Define  $\nabla_{\mathbf{z}} \mathbf{h} = (\partial/\partial \mathbf{z}) \mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X} + \mathbf{z} \mathbf{a}^\top)|_{\mathbf{z}=\mathbf{0}}$  and  $\nabla_{\mathbf{z}} \boldsymbol{\psi} = (\partial/\partial \mathbf{z}) \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X} + \mathbf{z} \mathbf{a}^\top)|_{\mathbf{z}=\mathbf{0}}$  for fixed  $\mathbf{a}$ ,  $\boldsymbol{\varepsilon}$  and  $\mathbf{X}$ . Let  $\mathbf{P}_{\mathbf{x}}^\perp = \mathbf{I}_n - \mathbf{x} \mathbf{x}^\top / \|\mathbf{x}\|_2^2$  for  $\mathbf{x} \neq \mathbf{0}$ . Then

(i) For fixed  $\boldsymbol{\varepsilon}$ ,  $\mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X})$  and  $\boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X})$  are Lipschitz in  $\mathbf{X}$  on every compact subset of  $\mathbb{R}^{n \times p}$ .

For (ii) and (iii), additionally assume that  $g$  is twice continuously differentiable.

(ii) For almost every  $\mathbf{X}$  and every  $\mathbf{a} \in \mathbb{R}^p$ ,

$$(\nabla_{\mathbf{z}} \mathbf{h})^\top = -(\mathbf{h}^\top \mathbf{a}) \mathbf{M} \mathbf{X}^\top \text{diag}(\boldsymbol{\psi}') + \mathbf{M} \mathbf{a} \boldsymbol{\psi}^\top, \quad (\text{C.7})$$

$$(\nabla_{\mathbf{z}} \boldsymbol{\psi})^\top = -(\mathbf{h}^\top \mathbf{a}) \mathbf{V} - \text{diag}(\boldsymbol{\psi}') \mathbf{X} \mathbf{M} \mathbf{a} \boldsymbol{\psi}^\top. \quad (\text{C.8})$$

(iii) For almost every  $\mathbf{X}$  and every  $\mathbf{a} \in \mathbb{R}^p$ , if  $\boldsymbol{\psi} \neq \mathbf{0}$ , then

$$(\nabla_{\mathbf{z}} (\boldsymbol{\psi} / \|\boldsymbol{\psi}\|_2))^\top = \|\boldsymbol{\psi}\|_2^{-1} \mathbf{P}_{\boldsymbol{\psi}}^\perp [ -(\mathbf{h}^\top \mathbf{a}) \mathbf{V} - \text{diag}(\boldsymbol{\psi}') \mathbf{X} \mathbf{M} \mathbf{a} \boldsymbol{\psi}^\top ] \quad (\text{C.9})$$

We remark that in view of (A.6),  $\nabla_{\mathbf{z}} \boldsymbol{\psi} = \nabla \boldsymbol{\psi}(\mathbf{z}_j)$  when  $\mathbf{a} = \mathbf{e}_j$ .

*Proof.* (i) We refer our readers to the proof of Proposition C.4.

(ii) As the functions  $\mathbf{h}$  and  $\boldsymbol{\psi}$  are Lipschitz on every compact, their Fréchet derivatives exist almost everywhere by Rademacher's theorem. Moreover, the chain rule holds almost everywhere by [34, Theorem 2.1.11]. Let  $\mathbf{G}, \mathbf{V}$  and  $\mathbf{M}$  be as in Lemma C.2. By differentiating these KKT conditions  $\mathbf{X}^\top \boldsymbol{\psi}(\mathbf{X}) = n(\nabla g(\mathbf{b}))|_{\mathbf{b}=\boldsymbol{\beta}+\mathbf{h}}$ , and by the chain rule, we obtain that for almost every  $\mathbf{X}$ ,

$$\begin{aligned} n\mathbf{G}(\nabla_{\mathbf{z}}\mathbf{h})^\top &= \mathbf{a}\boldsymbol{\psi}^\top + \mathbf{X}^\top (\nabla_{\mathbf{z}}\boldsymbol{\psi})^\top, \\ (\nabla_{\mathbf{z}}\boldsymbol{\psi})^\top &= \text{diag}(\boldsymbol{\psi}')\{-\mathbf{a}^\top \mathbf{h}\mathbf{I}_n - \mathbf{X}(\nabla_{\mathbf{z}}\mathbf{h})^\top\}. \end{aligned}$$

Solving the above system of equations gives (C.7) and (C.8).

(iii) Since the map  $\mathbf{v} \mapsto \mathbf{v}/\|\mathbf{v}\|_2$  with  $\mathbf{v} \in \mathbb{R}^n$  has Fréchet derivative  $\|\mathbf{v}\|_2^{-1}\mathbf{P}_{\mathbf{v}}^\perp$  at every point  $\mathbf{v} \neq \mathbf{0} \in \mathbb{R}^n$ , by chain rule, (C.9) follows for almost every  $\mathbf{X}$  if  $\boldsymbol{\psi}(\mathbf{X}) \neq \mathbf{0}_n$ .  $\square$

*Proof of Lemma A.2 when  $g$  is twice continuously differentiable.* Here, we further assume that  $g$  is twice differentiable so that  $\mathbf{V}, \mathbf{M}$  in Lemma C.2 are well-defined.

By Proposition C.4, the map  $\mathbf{z}_j \mapsto (\mathbf{h}, \boldsymbol{\psi})$  is locally Lipschitz, thus the map of the product  $\mathbf{z}_j \mapsto h_j \boldsymbol{\psi}(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1}$  is also locally Lipschitz. By Proposition D.2, without loss of generality, we can assume that  $\boldsymbol{\psi} \neq \mathbf{0}_n$  at almost every point  $\mathbf{z}_j \in \mathbb{R}^n$ . By the first order Stein's formula on the sphere (B.1) for  $(n-1)\|\mathbf{z}_j\|_2^{-2}\mathbb{E}_j[\mathbf{z}_j^\top \mathbf{f}(\mathbf{z}_j)]$  with  $\mathbf{f}(\mathbf{z}_j) = h_j \boldsymbol{\psi}(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1}$ , we have

$$\mathbb{E}_j[h_j^2(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \text{tr}(\mathbf{V}\mathbf{P}_{\mathbf{z}_j}^\perp)] \quad (\text{C.10})$$

$$= -\mathbb{E}_j[h_j(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \boldsymbol{\psi}^\top \mathbf{z}_j](n-1)\|\mathbf{z}_j\|_2^{-2} \quad (\text{C.11})$$

$$-2\mathbb{E}_j[h_j(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \mathbf{e}_j^\top (\text{diag}(\boldsymbol{\psi}')\mathbf{X}\mathbf{M})^\top \mathbf{P}_{\mathbf{z}_j}^\perp \boldsymbol{\psi}] \quad (\text{C.12})$$

$$+ \mathbb{E}_j[(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \mathbf{e}_j^\top \mathbf{M}\mathbf{e}_j \boldsymbol{\psi}^\top \mathbf{P}_{\mathbf{z}_j}^\perp \boldsymbol{\psi}] \quad (\text{C.13})$$

$$+ 2\mathbb{E}_j[h_j^2(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-2} \boldsymbol{\psi}^\top \mathbf{V}\mathbf{P}_{\mathbf{z}_j}^\perp \boldsymbol{\psi}] \quad (\text{C.14})$$

$$+ 2\mathbb{E}_j[h_j(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-2} \boldsymbol{\psi}^\top \text{diag}(\boldsymbol{\psi}')\mathbf{X}\mathbf{M}\mathbf{e}_j \boldsymbol{\psi}^\top \mathbf{P}_{\mathbf{z}_j}^\perp \boldsymbol{\psi}]. \quad (\text{C.15})$$

For the terms (C.12)-(C.14), by Lemma C.2 and  $\|\mathbf{P}_{\mathbf{z}_j}^\perp\|_{\text{op}} \leq 1$  we find

$$(\text{C.12}) \vee (\text{C.15}) \leq \sqrt{L/(\tau n)}\mathbb{E}_j[|h_j|(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1/2}] \quad \text{by (C.4)}$$

$$\leq \frac{1}{2}(\tau n)^{-1} + \frac{L}{2}\mathbb{E}_j[h_j^2(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1}],$$

$$(\text{C.13}) \leq (\tau n)^{-1} \quad \text{by (C.5),}$$

$$(\text{C.14}) \leq 2L\mathbb{E}_j[h_j^2(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1}] \quad \text{by (C.6).}$$

By leaving term (C.11) unchanged and using the above inequalities,

$$\begin{aligned} (\text{C.10}) &\leq -(n-1)\|\mathbf{z}_j\|_2^{-2}\mathbb{E}_j[h_j(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \boldsymbol{\psi}^\top \mathbf{z}_j] \\ &\quad + 2(\tau n)^{-1} \\ &\quad + 3L\mathbb{E}_j[h_j^2(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1}]. \end{aligned} \quad (\text{C.16})$$

We now bound (C.10) from below. Let

$$U_j = (L\tau^{-1}(\|\mathbf{X}\mathbf{Q}_j n^{-1/2}\|_{\text{op}} + \|\mathbf{z}_j n^{-1/2}\|_2)^2 + 1)^{-1}.$$

By  $(L\tau^{-1}\|\mathbf{X}n^{-1/2}\|_{\text{op}}^2 + 1)^{-1} \text{diag}(\boldsymbol{\psi}') \preceq \mathbf{V} \preceq \text{diag}(\boldsymbol{\psi}') \preceq L\mathbf{I}_n$  in (C.6) and by  $\|\mathbf{X}\|_{\text{op}} \leq \|\mathbf{X}\mathbf{Q}_j\|_{\text{op}} + \|\mathbf{z}_j\|_2$  for all  $j \in [p]$ , we have that  $U_j \text{diag}(\boldsymbol{\psi}') \preceq \mathbf{V} \preceq L\mathbf{I}_n$ .

Since  $U_j \text{diag}(\boldsymbol{\psi}') \preceq \mathbf{V}$  and  $K^2\mathbf{I}_n \preceq \text{diag}(\boldsymbol{\psi}') + \text{diag}\{\psi_i^2, i = 1, \dots, n\}$  both holds, multiplying the second inequality by  $U_j$  and summing yields  $K^2U_j\mathbf{I}_n \preceq U_j \text{diag} \psi_i^2 + \mathbf{V}$ . Multiplying both sides  $\mathbf{P}_{\mathbf{z}_j}^\perp$  to the left and to the right and using that  $\mathbf{P}_{\mathbf{z}_j}^\perp \preceq \mathbf{I}_n$  and  $\text{tr} \mathbf{P}_{\mathbf{z}_j}^\perp = n - 1$  we find

$$K^2U_j\mathbf{P}_{\mathbf{z}_j}^\perp \preceq U_j \text{diag} \psi_i^2 + \mathbf{P}_{\mathbf{z}_j}^\perp \mathbf{V} \mathbf{P}_{\mathbf{z}_j}^\perp, \quad K^2U_j(n-1) \leq U_j\|\boldsymbol{\psi}\|^2 + \text{tr}[\mathbf{P}_{\mathbf{z}_j}^\perp \mathbf{V}].$$

Multiplying by  $(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1}h_j^2$  and taking the conditional expectation  $\mathbb{E}_j$ ,

$$\begin{aligned} K^2U_j(n-1)\mathbb{E}_j[(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1}h_j^2] &\leq \mathbb{E}_j[h_j^2] + \mathbb{E}_j[\text{tr}[\mathbf{P}_{\mathbf{z}_j}^\perp \mathbf{V}]h_j^2(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1}] \\ &= \mathbb{E}_j[h_j^2] + (\text{C.10}). \end{aligned} \quad (\text{C.17})$$

Note that by definition of  $u_*$  and the events in Lemma C.1, we have when  $u_* > 0$ ,

$$\mathcal{E}_j \subset \{nu_* \leq K^2U_j(n-1) - 3L\} \cap \{\Omega_{jj}\|\mathbf{z}_j\|_2^2 \geq n(1-\eta_n)_+^2\}. \quad (\text{C.18})$$

Then combining (C.16) with the previous display, multiplying both sides by  $I_{\mathcal{E}_j}$  and summing over  $j \in [p]$  we find

$$nu_* \sum_{j=1}^p \mathbb{E}_j \left[ \frac{I_{\mathcal{E}_j} h_j^2}{\|\boldsymbol{\psi}\|_2^2 + \delta} \right] \leq \frac{2p}{\tau n} + \sum_{j=1}^p \mathbb{E}_j [I_{\mathcal{E}_j} h_j^2] - (n-1) \mathbb{E}_j \left[ \frac{\boldsymbol{\psi}^\top \mathbf{z}_j I_{\mathcal{E}_j} \mathbf{e}_j^\top \mathbf{h}}{(\|\boldsymbol{\psi}\|_2^2 + \delta) \|\mathbf{z}_j\|_2^2} \right].$$

Taking expectations  $\mathbb{E}$ , letting  $\text{diag}\{I_{\mathcal{E}_j}\}$  denote the diagonal matrix with the  $j$ -th diagonal element  $I_{\mathcal{E}_j}$ , we find

$$\begin{aligned} nu_* \mathbb{E}[(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \mathbf{h}^\top \text{diag}\{I_{\mathcal{E}_j}\} \mathbf{h}] &- 2p/(\tau n) - \mathbb{E}[\|\mathbf{h}\|_2^2] \\ &\leq -(n-1) \mathbb{E}[(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \boldsymbol{\psi}^\top \mathbf{X} \boldsymbol{\Sigma}^{-1} \text{diag}\{\Omega_{jj}^{-1} I_{\mathcal{E}_j} \|\mathbf{z}_j\|_2^{-2}\} \mathbf{h}] \\ &\leq (n-1) \mathbb{E}[\|\mathbf{X} \boldsymbol{\Sigma}^{-1} \text{diag}\{\Omega_{jj}^{-1} I_{\mathcal{E}_j} \|\mathbf{z}_j\|_2^{-2}\}\|_{\text{op}}^2]^{\frac{1}{2}} \mathbb{E}[(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \mathbf{h}^\top \text{diag}\{I_{\mathcal{E}_j}\} \mathbf{h}]^{\frac{1}{2}} \\ &\leq (1-\eta_n)_+^{-2} \mathbb{E}[\|\mathbf{X} \boldsymbol{\Sigma}^{-1}\|_{\text{op}}^2]^{\frac{1}{2}} \mathbb{E}[(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \mathbf{h}^\top \text{diag}\{I_{\mathcal{E}_j}\} \mathbf{h}]^{\frac{1}{2}} \end{aligned}$$

where the second inequality follows from the Cauchy-Schwarz inequality and  $(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1/2} \|\boldsymbol{\psi}\|_2 \leq 1$ , and the third inequality follows from  $\Omega_{jj}^{-1} I_{\mathcal{E}_j} (n-1) \|\mathbf{z}_j\|_2^{-2} \leq (1-\eta_n)_+^{-2}$  thanks to (C.18). This implies that  $x = (n \mathbb{E}[(\|\boldsymbol{\psi}\|_2^2 + \delta)^{-1} \mathbf{h}^\top \text{diag}\{I_{\mathcal{E}_j}\} \mathbf{h}])^{1/2}$  satisfies  $Ax^2 + Bx + C \leq 0$  where the polynomial coefficients are

$$A = u_*, \quad C = -2p(\tau n)^{-1} - \mathbb{E}[\|\mathbf{h}\|_2^2] \quad \text{and} \quad B = -(1-\eta_n)_+^{-2} \mathbb{E}[\|\mathbf{X} \boldsymbol{\Sigma}^{-1} n^{-1}\|_{\text{op}}^2]^{1/2}.$$



As  $A > 0$ , inequality  $Ax^2 + Bx + C \leq 0$  implies that  $x$  lies between the two real roots of the polynomial  $AX^2 + BX + C$ . In particular,  $x$  is smaller than the largest root, i.e.,  $x \leq (-B + \sqrt{B^2 - 4AC})/(2A) \leq |B|/A + (|C|/A)^{1/2}$ . Here,  $|C|/A = (2p(n\tau)^{-1} + \mathbb{E}[\|\mathbf{h}\|_2^2])/u_*$  and

$$|B|/A \leq u_*^{-1}(1 - \eta_n)_+^{-2} \|\Sigma^{-1/2}\|_{\text{op}} \mathbb{E}[\|\mathbf{X}\Sigma^{-1/2}n^{-1/2}\|_{\text{op}}^2]^{1/2}.$$

The upper bound (C.1) on  $\mathbb{E}[\|\mathbf{X}\Sigma^{-1/2}\|_{\text{op}}^2]$  completes the proof.  $\square$

### C.3. Non-smooth strongly convex penalty $g$

#### C.3.1. Almost everywhere differentiability

In this section, we provide the almost everywhere existence of the Jacobian matrices. We also notice that if our penalty  $g$  is not twice differentiable, the matrices  $\mathbf{M}$  and  $\mathbf{V}$  in Lemma C.2 are not well defined. In this case we do not have explicit formula for the Jacobian matrices  $(\partial/\partial \mathbf{z}_j)\psi$  and  $(\partial/\partial \mathbf{z}_j)\mathbf{h}$  such as those in terms of  $\mathbf{V}$ ,  $\text{diag}(\psi')$ ,  $\mathbf{M}$  in Proposition C.3 (ii) and (iii). In this section we provide upper bounds of the norms of these Jacobian matrices without using Proposition C.3.

**Proposition C.4.** *Let  $\rho : \mathbb{R} \rightarrow \mathbb{R}$  be convex and continuously differentiable with derivative  $\psi = \rho'$  being  $L$ -Lipschitz. Let  $g : \mathbb{R}^p \rightarrow \mathbb{R}$  be strongly convex with parameter  $\tau > 0$ . Let  $\mathbf{X}, \tilde{\mathbf{X}} \in \mathbb{R}^{n \times p}$ ,  $\boldsymbol{\varepsilon}, \tilde{\boldsymbol{\varepsilon}} \in \mathbb{R}^n$  and correspondingly,*

$$\begin{aligned} \hat{\boldsymbol{\beta}}(\boldsymbol{\varepsilon}, \mathbf{X}) &= \underset{\mathbf{b} \in \mathbb{R}^p}{\text{argmin}} \sum_{i \in [n]} \frac{\rho(\varepsilon_i - \mathbf{x}_i^\top (\mathbf{b} - \boldsymbol{\beta}))}{n} + g(\mathbf{b}), & \tilde{\boldsymbol{\beta}} &= \hat{\boldsymbol{\beta}}(\tilde{\boldsymbol{\varepsilon}}, \tilde{\mathbf{X}}), \\ \mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X}) &= \hat{\boldsymbol{\beta}}(\boldsymbol{\varepsilon}, \mathbf{X}) - \boldsymbol{\beta}, & \tilde{\mathbf{h}} &= \mathbf{h}(\tilde{\boldsymbol{\varepsilon}}, \tilde{\mathbf{X}}), \\ \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X}) &= \psi(\boldsymbol{\varepsilon} - \mathbf{X}\mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X})), & \tilde{\boldsymbol{\psi}} &= \boldsymbol{\psi}(\tilde{\boldsymbol{\varepsilon}}, \tilde{\mathbf{X}}). \end{aligned} \quad (\text{C.19})$$

Then (i)

$$n\tau \|\mathbf{h} - \tilde{\mathbf{h}}\|_2^2 + L^{-1} \|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2^2 \leq (\mathbf{h} - \tilde{\mathbf{h}})^\top (\mathbf{X} - \tilde{\mathbf{X}})^\top \boldsymbol{\psi} + (\boldsymbol{\varepsilon} - \tilde{\boldsymbol{\varepsilon}} + \tilde{\mathbf{X}}\mathbf{h} - \mathbf{X}\mathbf{h})^\top (\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}). \quad (\text{C.20})$$

(ii) The map  $(\boldsymbol{\varepsilon}, \mathbf{X}) \mapsto (\mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X}), \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X}))$  is Lipschitz on every compact of  $\mathbb{R}^n \times \mathbb{R}^{n \times p}$ .

(iii) For any  $\boldsymbol{\varepsilon} \in \mathbb{R}^n$  fixed, and for any  $\boldsymbol{\eta} \in \mathbb{R}^n$  and  $\mathbf{a} \in \mathbb{R}^p$

$$\begin{aligned} n\tau \|\hat{\boldsymbol{\beta}}(\boldsymbol{\varepsilon}, \mathbf{X} + \boldsymbol{\eta}\mathbf{a}^\top) - \hat{\boldsymbol{\beta}}(\boldsymbol{\varepsilon}, \mathbf{X})\|_2^2 + L^{-1} \|\boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X} + \boldsymbol{\eta}\mathbf{a}^\top) - \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X})\|_2^2 \\ \leq (n\tau)^{-1} \|\mathbf{a}\|_2^2 (\boldsymbol{\eta}^\top \boldsymbol{\psi})^2 + L(\mathbf{h}^\top \mathbf{a})^2 \|\boldsymbol{\eta}\|_2^2. \end{aligned}$$

Furthermore,

$$\sum_{i \in [n]} \|\boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X} + \mathbf{e}_i \mathbf{a}^\top) - \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X})\|_2^2 \leq L(n\tau)^{-1} \|\mathbf{a}\|_2^2 \|\boldsymbol{\psi}\|_2^2 + nL^2(\mathbf{h}^\top \mathbf{a})^2. \quad (\text{C.21})$$

(iv) If  $\boldsymbol{\eta} \in \mathbb{R}^n$  is such that  $\boldsymbol{\eta}^\top \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X}) = 0$  then  $\boldsymbol{\psi}(\boldsymbol{\varepsilon} + \mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X}))^\top \mathbf{e}_j \boldsymbol{\eta}, \mathbf{X} + \boldsymbol{\eta} \mathbf{e}_j^\top) = \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X})$ , so that if  $\boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X})$  is Fréchet differentiable at  $(\boldsymbol{\varepsilon}, \mathbf{X})$  then

$$\sum_{i=1}^n \eta_i \left[ \frac{\partial \boldsymbol{\psi}}{\partial x_{ij}}(\boldsymbol{\varepsilon}, \mathbf{X}) + (\mathbf{e}_j^\top \mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X})) \frac{\partial \boldsymbol{\psi}}{\partial \varepsilon_i}(\boldsymbol{\varepsilon}, \mathbf{X}) \right] = 0 \text{ provided that } \boldsymbol{\eta}^\top \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X}) = 0. \quad (\text{C.22})$$

The content of the above proposition appeared in [4, Proposition 5.2] with variables  $(\mathbf{y}, \mathbf{X})$  instead of  $(\boldsymbol{\varepsilon}, \mathbf{X})$ . It follows from strong convexity and the KKT conditions of  $\hat{\boldsymbol{\beta}}$  and  $\tilde{\boldsymbol{\beta}}$ . Its proof is provided below for completeness. An application of the above Proposition C.4 to normalized  $\boldsymbol{\psi}$  yields the following.

**Corollary C.5.** *Under the conditions of Proposition C.4 and with the notation of Proposition C.3, at a point where  $\|\boldsymbol{\psi}\|_2 > 0$  and  $\boldsymbol{\psi}$  is Fréchet differentiable,*

$$\left\| \nabla_{\mathbf{z}} \left( \frac{\boldsymbol{\psi}}{\|\boldsymbol{\psi}\|_2} \right) \right\|_F^2 \leq \frac{L \|\mathbf{a}\|_2^2}{n\tau} + \frac{nL^2(\mathbf{h}^\top \mathbf{a})^2}{\|\boldsymbol{\psi}\|_2^2}, \quad (\text{C.23})$$

and with the  $\nabla_{\mathbf{y}} \boldsymbol{\psi}$  in (2.9)

$$\left( \nabla_{\mathbf{z}} \left( \frac{\boldsymbol{\psi}}{\|\boldsymbol{\psi}\|_2} \right) + (\mathbf{a}^\top \mathbf{h}) \frac{\nabla_{\mathbf{y}} \boldsymbol{\psi}}{\|\boldsymbol{\psi}\|_2} \right)^\top \mathbf{P}_\psi^\perp = \mathbf{0}. \quad (\text{C.24})$$

*Proof.* For (C.23), by the chain rule

$$\|\nabla_{\mathbf{z}}(\boldsymbol{\psi}/\|\boldsymbol{\psi}\|_2)^\top\|_F^2 = \|\boldsymbol{\psi}\|_2^{-2} \|\mathbf{P}_\psi^\perp (\nabla_{\mathbf{z}} \boldsymbol{\psi})^\top\|_F^2 \leq \|\boldsymbol{\psi}\|_2^{-2} \|\nabla_{\mathbf{z}} \boldsymbol{\psi}\|_F^2.$$

By definition of the Frobenius norm  $\|\nabla_{\mathbf{z}} \boldsymbol{\psi}\|_F^2 = \sum_{i \in [n]} \|\partial \boldsymbol{\psi} / \partial z_i\|_2^2$ . By (C.21) with  $\mathbf{a}$  replaced by  $t\mathbf{a}$  and taking the limit as  $t \rightarrow 0$  we obtain  $\|\nabla_{\mathbf{z}}(\boldsymbol{\psi}/\|\boldsymbol{\psi}\|_2)\|_F^2 \leq L(n\tau)^{-1} \|\mathbf{a}\|_2^2 + nL^2(\mathbf{h}^\top \mathbf{a})^2 \|\boldsymbol{\psi}\|_2^{-2}$ .

For (C.24), we have

$$\left( \left( \frac{\partial}{\partial \boldsymbol{\eta}} + (\mathbf{a}^\top \mathbf{h}) \frac{\partial}{\partial \mathbf{y}} \right) \frac{\boldsymbol{\psi}}{\|\boldsymbol{\psi}\|_2} \right) \mathbf{P}_\psi^\perp = \left( \frac{\mathbf{P}_\psi^\perp}{\|\boldsymbol{\psi}\|_2} \left( \frac{\partial \boldsymbol{\psi}}{\partial \boldsymbol{\eta}} + (\mathbf{a}^\top \mathbf{h}) \frac{\partial \boldsymbol{\psi}}{\partial \mathbf{y}} \right) \right) \mathbf{P}_\psi^\perp = \mathbf{0}$$

due to  $\boldsymbol{\psi}(\boldsymbol{\varepsilon} + \mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X}))^\top \mathbf{a}(\mathbf{P}_\psi^\perp \boldsymbol{\eta}), \mathbf{X} + (\mathbf{P}_\psi^\perp \boldsymbol{\eta}) \mathbf{a}^\top) = \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X})$  by (C.20) as in the proof of Proposition C.4 (iv). Note that if  $F(\mathbf{y}, \mathbf{X})$  and  $G(\boldsymbol{\varepsilon}, \mathbf{X})$  are functions such that  $G(\boldsymbol{\varepsilon}, \mathbf{X}) = F(\boldsymbol{\varepsilon} + \mathbf{X}\boldsymbol{\beta}, \mathbf{X})$  and  $F(\mathbf{y}, \mathbf{X}) = G(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}, \mathbf{X})$  then  $(\partial/\partial y_i)F(\mathbf{y}, \mathbf{X}) = (\partial/\partial \varepsilon_i)G(\boldsymbol{\varepsilon}, \mathbf{X})$  whenever  $F$  is Fréchet differentiable at  $(\mathbf{y}, \mathbf{X})$  and  $G$  is Fréchet differentiable at  $(\boldsymbol{\varepsilon}, \mathbf{X})$  (i.e., translation by a constant in the variables does not change the derivatives).  $\square$

*Proof of Proposition C.4 (i).* Let  $\partial g(\cdot)$  denote the subdifferential of  $g$ . The KKT conditions read  $\mathbf{X}^\top \boldsymbol{\psi} \in n\partial g(\hat{\boldsymbol{\beta}})$  and  $\widetilde{\mathbf{X}}^\top \widetilde{\boldsymbol{\psi}} \in n\partial g(\tilde{\boldsymbol{\beta}})$ . Taking the difference and by  $\tau$ -strong convexity of  $g$ , we have

$$\begin{aligned} n\tau \|\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}\|_2^2 &\leq (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})^\top (\mathbf{X}^\top \boldsymbol{\psi} - \widetilde{\mathbf{X}}^\top \widetilde{\boldsymbol{\psi}}) \\ &= (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})^\top [(\mathbf{X} - \widetilde{\mathbf{X}})^\top \boldsymbol{\psi} + \widetilde{\mathbf{X}}^\top (\boldsymbol{\psi} - \widetilde{\boldsymbol{\psi}})] \end{aligned}$$

For the second term,

$$\begin{aligned} & (\hat{\beta} - \tilde{\beta})^\top \widetilde{\mathbf{X}}^\top (\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}) \\ &= (\widetilde{\mathbf{X}}\mathbf{h} - \widetilde{\mathbf{X}}\tilde{\mathbf{h}})^\top (\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}) \\ &= -(\boldsymbol{\varepsilon} - \mathbf{X}\mathbf{h} - (\tilde{\boldsymbol{\varepsilon}} - \widetilde{\mathbf{X}}\tilde{\mathbf{h}}))^\top (\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}) + (\boldsymbol{\varepsilon} - \tilde{\boldsymbol{\varepsilon}} + \widetilde{\mathbf{X}}\mathbf{h} - \mathbf{X}\mathbf{h})^\top (\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}). \end{aligned}$$

Since  $\boldsymbol{\psi}$  is non-decreasing and  $L$ -Lipschitz,  $L^{-1}\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2^2 \leq (\boldsymbol{\varepsilon} - \mathbf{X}\mathbf{h} - (\tilde{\boldsymbol{\varepsilon}} - \widetilde{\mathbf{X}}\tilde{\mathbf{h}}))^\top (\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}})$  holds. Combining the above displays we obtain (C.20).  $\square$

*Proof of Proposition C.4 (ii).* For fixed values of  $(\boldsymbol{\varepsilon}, \mathbf{X}, \mathbf{h}, \boldsymbol{\psi})$ , inequality (C.20) divided by  $1 + \|\mathbf{h} - \tilde{\mathbf{h}}\|_2 + \|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2$  implies that  $(\tilde{\mathbf{h}}, \tilde{\boldsymbol{\psi}}) \rightarrow (\mathbf{h}, \boldsymbol{\psi})$  as  $(\tilde{\boldsymbol{\varepsilon}}, \widetilde{\mathbf{X}}) \rightarrow (\boldsymbol{\varepsilon}, \mathbf{X})$ , hence the function  $(\boldsymbol{\varepsilon}, \mathbf{X}) \mapsto (\mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X}), \boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X}))$  is everywhere continuous. This implies that  $S(K) = \sup_{(\boldsymbol{\varepsilon}, \mathbf{X}) \in K} ((n\tau)^{-1}\|\boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{X})\|_2^2 + L\|\mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X})\|_2^2)$  is finite for any compact  $K \subset \mathbb{R}^n \times \mathbb{R}^{n \times p}$ . The Cauchy-Schwarz inequality on the right hand side of (C.20) gives for any  $(\boldsymbol{\varepsilon}, \mathbf{X}), (\tilde{\boldsymbol{\varepsilon}}, \widetilde{\mathbf{X}}) \in K$

$$\begin{aligned} & n\tau\|\mathbf{h} - \tilde{\mathbf{h}}\|_2^2 + L^{-1}\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2^2 \\ & \leq \|\mathbf{X} - \widetilde{\mathbf{X}}\|_{\text{op}}(\|\mathbf{h} - \tilde{\mathbf{h}}\|_2\|\boldsymbol{\psi}\|_2 + \|\mathbf{h}\|_2\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2) + \|\boldsymbol{\varepsilon} - \tilde{\boldsymbol{\varepsilon}}\|_2\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2 \\ & \leq \|\mathbf{X} - \widetilde{\mathbf{X}}\|_{\text{op}}(n\tau\|\mathbf{h} - \tilde{\mathbf{h}}\|_2^2 + L^{-1}\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2^2)^{\frac{1}{2}}S(K)^{\frac{1}{2}} + \|\boldsymbol{\varepsilon} - \tilde{\boldsymbol{\varepsilon}}\|_2\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2, \end{aligned}$$

This implies that  $(\boldsymbol{\varepsilon}, \mathbf{X}) \mapsto (\mathbf{h}, \boldsymbol{\psi})$  is Lipschitz on  $K$ .  $\square$

*Proof of Proposition C.4 (iii).* Combined with (C.20) with  $\boldsymbol{\varepsilon} = \tilde{\boldsymbol{\varepsilon}}$ , we have

$$\begin{aligned} & n\tau\|\mathbf{h} - \tilde{\mathbf{h}}\|_2^2 + L^{-1}\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2^2 \\ & \leq -(\mathbf{h} - \tilde{\mathbf{h}})^\top \mathbf{a}(\boldsymbol{\eta}^\top \boldsymbol{\psi}) + (\mathbf{h}^\top \mathbf{a})\boldsymbol{\eta}^\top (\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}) \\ & \leq \|\mathbf{h} - \tilde{\mathbf{h}}\|_2\|\mathbf{a}\|_2\|\boldsymbol{\eta}^\top \boldsymbol{\psi}\| + \|\mathbf{h}^\top \mathbf{a}\|\|\boldsymbol{\eta}\|_2\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2 \\ & \leq (n\tau\|\mathbf{h} - \tilde{\mathbf{h}}\|_2^2 + L^{-1}\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2^2)^{1/2}((n\tau)^{-1}\|\mathbf{a}\|_2^2(\boldsymbol{\eta}^\top \boldsymbol{\psi})^2 + L(\mathbf{h}^\top \mathbf{a})^2\|\boldsymbol{\eta}\|_2^2)^{1/2} \end{aligned}$$

so that the first inequality holds. Taking summation over  $\boldsymbol{\eta} = \mathbf{e}_i$  for  $i \in [n]$  gives (C.21).  $\square$

*Proof of Proposition C.4 (iv).* For  $\tilde{\boldsymbol{\varepsilon}} = \boldsymbol{\varepsilon} + \mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X})^\top \mathbf{e}_j \boldsymbol{\eta}$  and  $\widetilde{\mathbf{X}} = \mathbf{X} + \boldsymbol{\eta} \mathbf{e}_j^\top$  we have  $(\mathbf{X} - \widetilde{\mathbf{X}})^\top \boldsymbol{\psi} = 0$  thanks to  $\boldsymbol{\eta}^\top \boldsymbol{\psi} = 0$  as well as  $\boldsymbol{\varepsilon} - \tilde{\boldsymbol{\varepsilon}} + (\widetilde{\mathbf{X}} - \mathbf{X})\mathbf{h} = 0$ . Hence the two terms in the right-hand side of (C.20) are 0 and  $\|\boldsymbol{\psi} - \tilde{\boldsymbol{\psi}}\|_2^2 = 0$ . Identity (C.22) then follows by definition of the Fréchet differentiability.  $\square$

### C.3.2. Approximation using smooth penalty $\tilde{g}^\nu$

**Lemma C.6** (Approximation of strongly convex functions). *Let  $g : \mathbb{R}^p \rightarrow \mathbb{R}$  be strongly convex with constant  $\tau \geq 0$ . Then for every  $\nu > 0$  there exists a real-analytic strongly convex function  $g_\nu : \mathbb{R}^p \rightarrow \mathbb{R}$  with constant  $\tau$  such that  $g_\nu - \nu \leq g \leq g_\nu$ .*

*Proof.* Since  $g$  is proper, i.e.,  $-\infty \notin g(\mathbb{R}^p)$  and  $\{\mathbf{b} \in \mathbb{R}^p \mid g(\mathbf{b}) < +\infty\} \neq \emptyset$ , by Proposition 10.8 in [2],  $g$  is strongly convex with constant  $\tau \geq 0$  if and only if  $f := g - (\tau/2)\|\cdot\|_2^2$  is convex. By Theorem 1 in [1], there exists a function  $f_\nu : \mathbb{R}^p \rightarrow \mathbb{R}$  real-analytic and convex that satisfies  $f_\nu - \nu \leq f \leq f_\nu$ . The conclusion follows by letting  $g_\nu := f_\nu + (\tau/2)\|\cdot\|_2^2$ .  $\square$

**Lemma C.7.** *Let  $\rho : \mathbb{R} \rightarrow \mathbb{R}$  be convex and continuously differentiable with derivative  $\psi = \rho'$  being  $L$ -Lipschitz. Let  $g, \tilde{g} : \mathbb{R}^p \rightarrow \mathbb{R}$  be strongly convex with parameter  $\tau, \tilde{\tau} \geq 0$ . Let  $\|g - \tilde{g}\|_\infty = \max_{\mathbf{x} \in \mathbb{R}^p} |g(\mathbf{x}) - \tilde{g}(\mathbf{x})|$ . For  $\mathbf{b} \in \mathbb{R}^p$ , let  $L(\mathbf{b}; g) = \frac{1}{n} \sum_{i \in [n]} \rho(y_i - \mathbf{x}_i^\top \mathbf{b}) + g(\mathbf{b})$  and define*

$$\hat{\beta} = \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} L(\mathbf{b}; g), \quad \tilde{\beta} = \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} L(\mathbf{b}; \tilde{g}), \quad \psi = \psi(\mathbf{y} - \mathbf{X}\hat{\beta}), \quad \tilde{\psi} = \psi(\mathbf{y} - \mathbf{X}\tilde{\beta}).$$

*Then inequality  $((\tau + \tilde{\tau})/2)\|\tilde{\beta} - \hat{\beta}\|_2^2 + (nL)^{-1}\|\psi - \tilde{\psi}\|_2^2 \leq 2\|g - \tilde{g}\|_\infty$  holds.*

*Proof of Lemma C.7.* Denote by  $\partial g(\mathbf{b})$  subdifferential of  $g$  at  $\mathbf{b} \in \mathbb{R}^p$ . The KKT conditions read

$$(1/n)\mathbf{X}^\top \psi \in \partial g(\hat{\beta}), \quad (1/n)\mathbf{X}^\top \tilde{\psi} \in \partial \tilde{g}(\tilde{\beta}).$$

By the definition of the strongly convexity, the above display implies that

$$(\tau/2)\|\tilde{\beta} - \hat{\beta}\|_2^2 + (\tilde{\beta} - \hat{\beta})^\top (1/n)\mathbf{X}^\top \psi \leq g(\tilde{\beta}) - g(\hat{\beta}),$$

$$(\tilde{\tau}/2)\|\hat{\beta} - \tilde{\beta}\|_2^2 + (\hat{\beta} - \tilde{\beta})^\top (1/n)\mathbf{X}^\top \tilde{\psi} \leq \tilde{g}(\hat{\beta}) - \tilde{g}(\tilde{\beta}).$$

Summing over the above displays, we have

$$\begin{aligned} \frac{\tau + \tilde{\tau}}{2}\|\hat{\beta} - \tilde{\beta}\|_2^2 + (\hat{\beta} - \tilde{\beta})^\top (1/n)\mathbf{X}^\top (\tilde{\psi} - \psi) &\leq g(\tilde{\beta}) - \tilde{g}(\tilde{\beta}) + \tilde{g}(\hat{\beta}) - g(\hat{\beta}) \\ &\leq 2\|g - \tilde{g}\|_\infty. \end{aligned}$$

We notice that the second term in the left hand side is can be rewritten

$$\frac{1}{n} \langle \mathbf{y} - \mathbf{X}\hat{\beta} - (\mathbf{y} - \mathbf{X}\tilde{\beta}), \psi(\mathbf{y} - \mathbf{X}\hat{\beta}) - \psi(\mathbf{y} - \mathbf{X}\tilde{\beta}) \rangle = \sum_{i=1}^n \frac{(a_i - b_i)(\psi(a_i) - \psi(b_i))}{n}$$

where  $a_i = y_i - \mathbf{x}_i^\top \hat{\beta}$  and  $b_i = y_i - \mathbf{x}_i^\top \tilde{\beta}$ . Since  $\psi$  non-decreasing and  $L$ -Lipschitz, inequality  $(a_i - b_i)(\psi(a_i) - \psi(b_i)) = |a_i - b_i| |\psi(a_i) - \psi(b_i)| \geq \frac{1}{L}(\psi(a_i) - \psi(b_i))^2$  completes the proof.  $\square$

*Proof of Lemma A.2 for  $g$   $\tau$ -strongly convex but not continuously differentiable.* In this proof, we approximate the non-smooth  $g$  with smooth  $g$  using Lemma C.6. Let  $g$  be strongly convex with constant  $\tau > 0$ , not necessarily twice differentiable. By Lemma C.6, for all  $\nu > 0$ , there exists  $\tilde{g}^\nu$  strongly convex with constant  $\tau > 0$  such that  $\|\tilde{g}^\nu - g\|_\infty \leq \nu$ . Let  $\tilde{\beta}^\nu = \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} L(\mathbf{b}; \tilde{g}^\nu)$  and  $\tilde{\psi}^\nu = \psi(\mathbf{y} - \mathbf{X}\tilde{\beta}^\nu)$  be as in Lemma C.7. For any  $\delta > 0$ ,

$$\mathbb{E} \left[ \frac{n}{\|\tilde{\psi}^\nu\|_2^2 + \delta} \sum_{j \in [p]} I_{\mathcal{E}_j} \tilde{h}_j^2 \right]^{1/2} \leq \frac{[(1 + \sqrt{\frac{p}{n}})^2 + \frac{1}{n}]^{1/2}}{\phi_{\min}(\Sigma)^{1/2}(1 - \eta_n)_+^2 u_*} + \frac{[\frac{2p}{n\tau} + \mathbb{E}[\|\tilde{\mathbf{h}}^\nu\|_2^2]]^{1/2}}{u_*^{1/2}} \quad (\text{C.25})$$

by (A.3) since  $\tilde{g}^\nu$  is twice continuously differentiable. By Lemma C.7, we have

$$\tau \|\tilde{\beta}^\nu - \hat{\beta}\|_2^2 + (nL)^{-1} \|\tilde{\psi}^\nu - \psi\|_2^2 \leq 2\nu.$$

This implies that, as  $\nu \rightarrow 0+$ , the pointwise convergence  $(\tilde{\mathbf{h}}_j^\nu, \tilde{\psi}^\nu) \rightarrow (\mathbf{h}, \psi)$  holds. By the dominated convergence theorem, a sufficient condition that (C.25) holds with  $(\tilde{\mathbf{h}}^\nu, \tilde{\psi}^\nu)$  replaced by its pointwise limit  $(\mathbf{h}, \psi)$  inside the two expectations in (C.25) is that  $\mathbb{E} \sup_{\nu \in (0,1)} \|\tilde{\mathbf{h}}^\nu\|_2^2 < +\infty$ . By Lemma C.7,

$$\|\tilde{\mathbf{h}}^\nu\|_2^2 \leq 2\|\tilde{\mathbf{h}}^\nu - \mathbf{h}\|_2^2 + 2\|\mathbf{h}\|_2^2 < (2\nu/\tau) + 2\|\mathbf{h}\|_2^2 \quad (\text{C.26})$$

which provides integrability of  $\sup_{\nu \in (0,1)} \|\tilde{\mathbf{h}}^\nu\|_2^2$  as  $\mathbb{E}[\|\Sigma^{1/2}\mathbf{h}\|_2^2] < +\infty$  when the right-hand side of (A.3) is finite.  $\square$

## Appendix D: Auxiliary propositions

### D.1. Decomposition of the design matrix into independent components

**Proposition D.1** (Independence between  $\mathbf{X}(\mathbf{I}_p - \mathbf{b}\mathbf{a}^\top)$  and  $\mathbf{X}\mathbf{b}$ ). *Let each row  $\mathbf{x}_i$  of  $\mathbf{X} \in \mathbb{R}^{n \times p}$  satisfy that  $\mathbf{x}_i \sim^{iid} N(\mathbf{0}, \Sigma)$ . Then for any deterministic vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ ,  $\Sigma\mathbf{b} = (\mathbf{b}^\top \Sigma \mathbf{b})\mathbf{a}$  holds if and only if  $\mathbf{X}(\mathbf{I}_p - \mathbf{b}\mathbf{a}^\top)$  is independent with  $\mathbf{X}\mathbf{b}$ . Furthermore, if the above holds and the inverse matrix  $\Sigma^{-1}$  exists, then  $(\mathbf{b}^\top \Sigma \mathbf{b})(\mathbf{a}^\top \Sigma^{-1} \mathbf{a}) = 1$ .*

*Proof.* From the fact that  $(\mathbf{X}(\mathbf{I}_p - \mathbf{b}\mathbf{a}^\top), \mathbf{X}\mathbf{b})$  can be represented in a linear transformation of  $n \times p$  iid  $N(0, 1)$  random variable, the pair is distributed in a multivariate normal distribution. Since the rows of  $\mathbf{X}$  are independent, the independence between  $\mathbf{X}(\mathbf{I}_p - \mathbf{b}\mathbf{a}^\top)$  and  $\mathbf{X}\mathbf{b}$  reduces to the independence between  $\mathbf{x}_i^\top(\mathbf{I}_p - \mathbf{b}\mathbf{a}^\top)$  and  $\mathbf{x}_i^\top \mathbf{b}$  for each  $i \in [n]$ , which holds if and only if the two random quantities are uncorrelated in the sense that

$$\mathbb{E}[(\mathbf{x}_i^\top(\mathbf{I}_p - \mathbf{b}\mathbf{a}^\top))(\mathbf{x}_i^\top \mathbf{b})] = \mathbb{E}[\mathbf{b}^\top \mathbf{x}_i \mathbf{x}_i^\top (\mathbf{I}_p - \mathbf{b}\mathbf{a}^\top)] = \mathbf{b}^\top \Sigma (\mathbf{I}_p - \mathbf{b}\mathbf{a}^\top) = \mathbf{0}.$$

If the inverse  $\Sigma^{-1}$  exists, the above display implies  $(\mathbf{b}^\top \Sigma \mathbf{b})(\mathbf{a}^\top \Sigma^{-1} \mathbf{a}) = 1$ .  $\square$

### D.2. $\psi$ at the residuals is almost surely nonzero

**Proposition D.2.** *If Assumptions A, B and C hold then  $\mathbb{P}(\psi(\mathbf{y} - \mathbf{X}\hat{\beta}) \neq 0) = 1$ .*

*Proof of Proposition D.2.* If  $\psi(\mathbf{y} - \mathbf{X}\hat{\beta}) = 0$  then  $\hat{\beta}$  must be a minimizer of the penalty  $g$ . Let  $\mathbf{b}_0$  be a minimizer of  $g$ , which is unique by strong convexity.

Our assumption on the convexity of  $\rho$  implies that  $\psi(x)$  is non-decreasing in  $x \in \mathbb{R}$ . Combined with our assumption  $\psi'(x) + \psi^2(x) \geq K^2 > 0$  for every point  $x \in \mathbb{R}$ , this implies that  $\psi(x) = 0$  at no more than one point in  $\mathbb{R}$ . (Otherwise,

there exists an open interval on which  $\psi(x) = 0$  and  $\psi'(x) + \psi^2(x) = 0$ . A contradiction then follows.)

Thus we have  $\mathbb{P}(\psi(\mathbf{y} - \mathbf{X}\hat{\beta}) = \mathbf{0}) \leq \mathbb{P}(\psi(\varepsilon - \mathbf{X}(\mathbf{b}_0 - \beta)) = \mathbf{0}) = 0$  as  $\varepsilon - \mathbf{X}(\mathbf{b}_0 - \beta)$  has continuous distribution by Assumption C and  $\{x \in \mathbb{R} : \psi(x) = 0\}$  has zero Lebesgue measure.  $\square$

## References

- [1] Daniel Azagra. Global and fine approximation of convex functions. *Proceedings of the London Mathematical Society*, 107(4):799–824, 2013. [MR3108831](#)
- [2] Heinz H Bauschke and Patrick L Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2017. [MR3616647](#)
- [3] Derek Bean, Peter J Bickel, Noureddine El Karoui, and Bin Yu. Optimal m-estimation in high-dimensional regression. *Proceedings of the National Academy of Sciences*, 110(36):14563–14568, 2013.
- [4] Pierre C Bellec. Out-of-sample error estimate for robust m-estimators with convex penalty. *arXiv preprint arXiv:2008.11840*, 2020.
- [5] Pierre C Bellec and Yiwei Shen. Derivatives and residual distribution of regularized m-estimators with application to adaptive tuning. In *Conference on Learning Theory*, pages 1912–1947. PMLR, 2022.
- [6] Pierre C Bellec and Cun-Hui Zhang. Second order stein: Sure for sure and other applications in high-dimensional inference. *arXiv preprint arXiv:1811.04121*, 2018. [MR4319234](#)
- [7] Pierre C Bellec and Cun-Hui Zhang. De-biasing convex regularized estimators and interval estimation in linear models. *arXiv preprint arXiv:1912.11943*, 2019.
- [8] Pierre C Bellec and Cun-Hui Zhang. De-biasing the lasso with degrees-of-freedom adjustment. *arXiv:1902.08885*, 2019. URL <https://arxiv.org/pdf/1902.08885.pdf>. [MR4389062](#)
- [9] Vladimir Igorevich Bogachev. *Gaussian measures*. Number 62. American Mathematical Soc., 1998. [MR1642391](#)
- [10] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013. [MR3185193](#)
- [11] Peter Bühlmann et al. Statistical significance in high-dimensional linear models. *Bernoulli*, 19(4):1212–1242, 2013. [MR3102549](#)
- [12] Michael Celentano and Andrea Montanari. Fundamental barriers to high-dimensional regression with convex penalties. *arXiv preprint arXiv:1903.10603*, 2019. [MR4382013](#)
- [13] Michael Celentano, Andrea Montanari, and Yuting Wei. The lasso with general gaussian designs with applications to hypothesis testing. *arXiv preprint arXiv:2007.13716*, 2020.
- [14] Kenneth R Davidson and Stanislaw J Szarek. Local operator theory, ran-

- dom matrices and banach spaces. *Handbook of the geometry of Banach spaces*, 1(317-366):131, 2001. [MR1863696](#)
- [15] David Donoho and Andrea Montanari. High dimensional robust m-estimation: Asymptotic variance via approximate message passing. *Probability Theory and Related Fields*, 166(3-4):935–969, 2016. [MR3568043](#)
- [16] Nouredine El Karoui. On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probability Theory and Related Fields*, 170(1-2):95–175, 2018. [MR3748322](#)
- [17] Nouredine El Karoui, Derek Bean, Peter J Bickel, Chingway Lim, and Bin Yu. On robust regression with high-dimensional predictors. *Proceedings of the National Academy of Sciences*, 110(36):14557–14562, 2013.
- [18] Emmanuel Hebey. *Sobolev spaces on Riemannian manifolds*, volume 1635. Springer Science & Business Media, 1996. [MR1481970](#)
- [19] Hanwen Huang. Asymptotic risk and phase transition of  $l_1$ -penalized robust estimator. *Ann. Statist.*, 48(5):3090–3111, 10 2020. URL <https://doi.org/10.1214/19-AOS1923>. [MR4152636](#)
- [20] Peter J. Huber. Robust estimation of a location parameter. *Ann. Math. Statist.*, 35(1):73–101, 03 1964. URL <https://doi.org/10.1214/aoms/1177703732>. [MR0161415](#)
- [21] Peter J Huber. *Robust statistics*, volume 523. John Wiley & Sons, 2004. [MR0606374](#)
- [22] Adel Javanmard and Andrea Montanari. Confidence intervals and hypothesis testing for high-dimensional regression. *The Journal of Machine Learning Research*, 15(1):2869–2909, 2014. [MR3277152](#)
- [23] Adel Javanmard and Andrea Montanari. Hypothesis testing in high-dimensional regression under the gaussian random design model: Asymptotic theory. *IEEE Transactions on Information Theory*, 60(10):6522–6554, 2014. [MR3265038](#)
- [24] Adel Javanmard and Andrea Montanari. Debiasing the lasso: Optimal sample size for gaussian designs. *The Annals of Statistics*, 46(6A):2593–2622, 2018. [MR3851749](#)
- [25] Nouredine El Karoui. Asymptotic behavior of unregularized and ridge-regularized high-dimensional robust regression estimators: rigorous results. *arXiv preprint arXiv:1311.2445*, 2013.
- [26] Léo Miolane and Andrea Montanari. The distribution of the lasso: Uniform control over sparse balls and adaptive parameter tuning. *arXiv preprint arXiv:1811.01212*, 2018. [MR4319252](#)
- [27] Stephen Portnoy. Asymptotic behavior of  $m$  estimators of  $p$  regression parameters when  $p^2/n$  is large; ii. normal approximation. *The Annals of Statistics*, pages 1403–1417, 1985. [MR0811499](#)
- [28] Mihailo Stojnic. A framework to characterize performance of lasso algorithms. *arXiv preprint arXiv:1303.7291*, 2013.
- [29] Christos Thrampoulidis, Ehsan Abbasi, and Babak Hassibi. Precise error analysis of regularized  $m$ -estimators in high dimensions. *IEEE Transactions on Information Theory*, 64(8):5592–5628, 2018. [MR3832326](#)

- [30] Sara Van de Geer, Peter Bühlmann, Ya'acov Ritov, and Ruben Dezeure. On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202, 2014. [MR3224285](#)
- [31] AW Van der Vaart. Estimating a real parameter in a class of semiparametric models. *The Annals of Statistics*, pages 1450–1474, 1988. [MR0964933](#)
- [32] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018. [MR3837109](#)
- [33] Cun-Hui Zhang and Stephanie S Zhang. Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):217–242, 2014. [MR3153940](#)
- [34] William P. Ziemer. *Weakly Differentiable Functions: Sobolev Spaces and Functions of Bounded Variation*. Graduate Texts in Mathematics. Springer, 1989. [MR1014685](#)