

Derivatives and residual distribution of regularized M-estimators with application to adaptive tuning

Pierre C. Bellec PIERRE.BELLECC@RUTGERS.EDU and **Yiwei Shen** YS573@STAT.RUTGERS.EDU
Department of Statistics, Rutgers University, Piscataway, NJ 08854, USA.

Editors: Po-Ling Loh and Maxim Raginsky

Abstract

This paper studies M-estimators with gradient-Lipschitz loss function regularized with convex penalty in linear models with Gaussian design matrix and arbitrary noise distribution. A practical example is the robust M-estimator constructed with the Huber loss and the Elastic-Net penalty and the noise distribution has heavy-tails. Our main contributions are three-fold. (i) We provide general formulae for the derivatives of regularized M-estimators $\hat{\beta}(\mathbf{y}, \mathbf{X})$ where differentiation is taken with respect to both \mathbf{y} and \mathbf{X} ; this reveals a simple differentiability structure shared by all convex regularized M-estimators. (ii) Using these derivatives, we characterize the distribution of the residual $r_i = y_i - \mathbf{x}_i^\top \hat{\beta}$ in the intermediate high-dimensional regime where dimension and sample size are of the same order. (iii) Motivated by the distribution of the residuals, we propose a novel adaptive criterion to select tuning parameters of regularized M-estimators. The criterion approximates the out-of-sample error up to an additive constant independent of the estimator, so that minimizing the criterion provides a proxy for minimizing the out-of-sample error. The proposed adaptive criterion does not require the knowledge of the noise distribution or of the covariance of the design. Simulated data confirms the theoretical findings, regarding both the distribution of the residuals and the success of the criterion as a proxy of the out-of-sample error. Finally our results reveal new relationships between the derivatives of $\hat{\beta}(\mathbf{y}, \mathbf{X})$ and the effective degrees of freedom of the M-estimator, which are of independent interest.

Keywords: Robust estimation, M-estimator, Adaptive tuning, High-dimensional statistics, Residual distribution.

1. Introduction

This paper studies properties of robust estimators in linear models $\mathbf{y} = \mathbf{X}\beta^* + \varepsilon$ with response $\mathbf{y} \in \mathbb{R}^n$, unknown regression vector β^* and \mathbf{X} is a design matrix with n rows $\mathbf{x}_1, \dots, \mathbf{x}_n$. Each row \mathbf{x}_i being a high-dimensional feature vector in \mathbb{R}^p , centered and normally distributed with covariance Σ , and each ε_i is independent of \mathbf{X} with continuous distribution. Throughout, let $\hat{\beta} = \hat{\beta}(\mathbf{y}, \mathbf{X})$ be a regularized M -estimator given as a solution of the convex minimization problem

$$\hat{\beta}(\mathbf{y}, \mathbf{X}) = \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \rho(y_i - \mathbf{x}_i^\top \mathbf{b}) + g(\mathbf{b}) \tag{1.1}$$

where $\rho : \mathbb{R} \rightarrow \mathbb{R}$ is a convex data-fitting loss function and $g : \mathbb{R}^p \rightarrow \mathbb{R}$ a convex penalty. We may write $\hat{\beta}_{\rho, g}(\mathbf{y}, \mathbf{X})$ for (1.1) to emphasize the dependence on the loss-penalty pair (ρ, g) ; if the argument (\mathbf{y}, \mathbf{X}) is dropped then $\hat{\beta}$ is implicitly understood at the observed that (\mathbf{y}, \mathbf{X}) . Typical examples of losses include the square loss $\rho(u) = u^2/2$, the Huber loss $H(u) = \int_0^{|u|} \min(1, t) dt$ or its scaled version $\rho = \Lambda^2 H(u/\Lambda)$ for some tuning parameter $\Lambda > 0$, while typical examples of penalty functions include the Elastic-Net $g(\mathbf{b}) = \lambda \|\mathbf{b}\|_1 + \mu \|\mathbf{b}\|^2/2$ for tuning parameters $\lambda, \mu \geq 0$.

The paper introduces the following criterion to select a loss-penalty pair (ρ, g) with small out-of-sample error $\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|^2$: for a given set of candidate loss-penalty pairs $\{(\rho, g)\}$ and the corresponding M -estimator $\hat{\beta}_{\rho, g}$ in (1.1), select the pair (ρ, g) that minimizes the criterion

$$\text{Crit}(\rho, g) = \left\| \mathbf{r} + \frac{\hat{\text{d}}\mathbf{f}}{\text{tr}[\mathbf{V}]} \psi(\mathbf{r}) \right\|^2 \text{ with } \begin{cases} \mathbf{r} = \mathbf{y} - \mathbf{X}\hat{\beta}_{\rho, g} & \in \mathbb{R}^n, \\ \hat{\text{d}}\mathbf{f} = \text{tr}[\mathbf{X}(\partial/\partial\mathbf{y})\hat{\beta}_{\rho, g}] & \in \mathbb{R}, \\ \mathbf{V} = \text{diag}\{\psi'(\mathbf{r})\}(\mathbf{I}_n - \mathbf{X}(\partial/\partial\mathbf{y})\hat{\beta}_{\rho, g}) & \in \mathbb{R}^{n \times n} \end{cases} \quad (1.2)$$

where $\text{tr}[\cdot]$ is the trace, $\psi : \mathbb{R} \rightarrow \mathbb{R}$ is the derivative of ρ , ψ' the derivative of ψ and we extend ψ and ψ' to functions $\mathbb{R}^n \rightarrow \mathbb{R}^n$ by componentwise application of the univariate function of the same symbol. Above, $(\partial/\partial\mathbf{y})\hat{\beta}_{\rho, g} \in \mathbb{R}^{p \times n}$ denotes the Jacobian of (1.1) with respect to \mathbf{y} for \mathbf{X} fixed, at the observed data (\mathbf{y}, \mathbf{X}) . As we will see while studying particular examples, for pairs (ρ, g) commonly used in robust high-dimensional statistics such as the square loss, Huber loss with the ℓ_1 -penalty or Elastic-Net penalty, the ratio $\hat{\text{d}}\mathbf{f}/\text{tr}[\mathbf{V}]$ in (1.2) admits simple, closed-form expressions. The criterion (1.2) has an appealing adaptivity property: it does not require any knowledge of the noise ε or its distribution, nor any knowledge of the covariance Σ of the design.

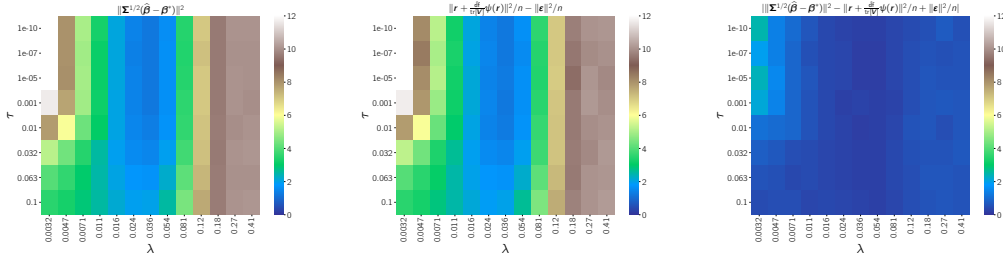


Figure 1: Heatmaps for $\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|^2$, its approximation $\|\mathbf{r} + (\hat{\text{d}}\mathbf{f}/\text{tr}[\mathbf{V}])\psi(\mathbf{r})\|^2/n - \|\varepsilon\|^2/n$ and the approximation error $\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|^2 - \|\mathbf{r} + (\hat{\text{d}}\mathbf{f}/\text{tr}[\mathbf{V}])\psi(\mathbf{r})\|^2/n - \|\varepsilon\|^2/n$ for the Huber loss and Elastic-Net penalty on a grid of tuning parameters (λ, τ) where $\lambda \in [0.0032, 0.41]$ and $\tau \in [10^{-10}, 0.1]$. Each cell is the average over 100 repetitions. See Section 6 for more details.

1.1. Contributions

1. The end goal of this paper is to provide theoretical justification and theoretical guarantees for the criterion (1.2) in the high-dimensional regime where the ratio p/n has a finite limit and \mathbf{X} has anisotropic Gaussian distribution. The theoretical results will justify the approximation

$$\left\| \mathbf{r} + \left(\hat{\text{d}}\mathbf{f} / \text{tr}[\mathbf{V}] \right) \psi(\mathbf{r}) \right\|^2 / n \approx \|\varepsilon\|^2 / n + \|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|^2. \quad (1.3)$$

Figure 1 illustrates the accuracy of (1.3) on simulated data. To study the criterion (1.2) and derive the approximation (1.3), we develop novel results of independent interest regarding M -estimators in (1.1):

2. The paper derives general formula for the derivatives $(\partial/\partial y_i)\hat{\beta}$ and $(\partial/\partial x_{ij})\hat{\beta}$. This sheds light on the differentiability structure of M -estimators for general loss-penalty pairs: for any ρ, g

with g strongly convex, there exists $\widehat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ depending on (\mathbf{y}, \mathbf{X}) such that for almost every (\mathbf{y}, \mathbf{X}) ,

$$(\partial/\partial y_i)\widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) = \widehat{\mathbf{A}}\mathbf{X}^\top \mathbf{e}_i \psi'(r_i), \quad (\partial/\partial x_{ij})\widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) = \widehat{\mathbf{A}}\mathbf{e}_j \psi(r_i) - \widehat{\mathbf{A}}\mathbf{X}^\top \mathbf{e}_i \psi'(r_i)\widehat{\beta}_j,$$

for $r_i = y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}$, $\forall i \in [n], j \in [p]$ where $\mathbf{e}_j \in \mathbb{R}^p$ and $\mathbf{e}_i \in \mathbb{R}^n$ are canonical basis vectors.

3. The paper obtains a stochastic representation for the residual $y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}$ for some fixed $i = 1, \dots, n$, extending some results of [El Karoui et al. \(2013\)](#) on unregularized M -estimators to penalized ones as in (1.1). In short, for each $i = 1, \dots, n$ the i -th residual satisfies $r_i = y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}$

$$r_i + \frac{\text{d}f}{\text{tr}[\mathbf{V}]} \psi(r_i) \approx \varepsilon_i + Z_i \|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\| \quad (1.4)$$

where $Z_i \sim N(0, 1)$ is independent of ε_i . This stochastic representation is the motivation for the criterion (1.2) as the amplitude of the normal part in the right-hand side is proportional to the out-of-sample error $\|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|$ that we wish to minimize, while the variance of the noise ε_i does not depend on the choice of (ρ, g) .

Simulated data in Figure 2 confirms that the stochastic representation for the i -th residual $r_i = y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}$ is accurate. Our working assumption throughout the paper is the following.

Assumption A For constants $\gamma, \mu > 0$ independent of n, p we have $p/n \leq \gamma$, the loss $\rho : \mathbb{R} \rightarrow \mathbb{R}$ is convex with a unique minimizer at 0, continuously differentiable and its derivative $\psi = \rho'$ is 1-Lipschitz. The design matrix \mathbf{X} has iid $N(\mathbf{0}, \boldsymbol{\Sigma})$ rows for some invertible covariance $\boldsymbol{\Sigma}$ and the noise ε is independent of \mathbf{X} with continuous distribution. The penalty $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is μ -strongly convex w.r.t. $\boldsymbol{\Sigma}$ in the sense that $\mathbf{b} \mapsto g(\mathbf{b}) - (\mu/2)\mathbf{b}^\top \boldsymbol{\Sigma} \mathbf{b}$ is convex in $\mathbf{b} \in \mathbb{R}^p$.

Throughout the paper, we consider a sequence (say, indexed by n) of regression problems with $p, \boldsymbol{\beta}^*, \boldsymbol{\Sigma}$ and the loss-penalty pair (ρ, g) depending implicitly on n . For some deterministic sequence (a_n) , the stochastically bounded notation $O_P(a_n)$ in this context may hide constants depending on γ, μ only, that is, $O_P(a_n)$ denotes a sequence of random variables W_n such that for any $\epsilon > 0$ there exists K depending on (ϵ, γ, μ) satisfying $\mathbb{P}(|W_n| \geq K a_n) \leq \epsilon$.

Since Assumption A requires $p/n \leq \gamma$, the Bolzano-Weierstrass theorem lets us extract a subsequence of regression problems such that $p/n \rightarrow \gamma'$ along this subsequence, for some constant γ' . This is the asymptotic regime we have in mind throughout the paper, although our results do not require a specific limit for the ratio p/n . For some results, we will require the following additional assumption which is satisfied by robust loss functions and penalty that shrink towards 0.

Assumption B The penalty is minimized at $\mathbf{0}$, that is, $g(\mathbf{0}) = \min_{\mathbf{b} \in \mathbb{R}^p} g(\mathbf{b})$; the loss is Lipschitz as in $|\psi| \leq M$ for some constant M independent of n, p ; the signal is bounded as in $\|\boldsymbol{\Sigma}^{1/2} \boldsymbol{\beta}^*\|^2 \leq M$.

1.2. Related works

The context of the present work is the study of M -estimators in the regime $\frac{p}{n}$ has a finite limit. This literature pioneered in [Bayati and Montanari \(2012\)](#); [El Karoui et al. \(2013\)](#); [Donoho and Montanari \(2016\)](#); [Stojnic \(2013\)](#) typically describes the subtle behavior of $\widehat{\boldsymbol{\beta}}$ in this regime by solving a system of nonlinear equations. This system depends on a prior distribution for the components of $\boldsymbol{\beta}^*$, and

either depends on the covariance Σ (Celentano et al., 2020; Dobriban and Wager, 2018) or assume $\Sigma = \mathbf{I}_p$ (Bayati and Montanari, 2012; Thrampoulidis et al., 2018; Celentano and Montanari, 2019, among many others). Solutions to the nonlinear system are a powerful tool to understand $\hat{\beta}$ in theory, e.g., to characterize the deterministic limit of $\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|$, see e.g., the general results in Celentano and Montanari (2019) for the square loss and Thrampoulidis et al. (2018) for general loss-penalty pairs. However, since the system and its solution depend on unobservable quantities (Σ and prior on β^*), the system solution is not directly usable for practical purposes such as parameter tuning.

The present work distinguishes itself from most of this literature as the goal is to describe the behavior of $\hat{\beta}$ using observable quantities that only depend on the data (\mathbf{y}, \mathbf{X}) (and not unobservable ones such as Σ or a prior distribution on β^* that appear in the aforementioned nonlinear system of equations). As we will see this view lets us perform adaptive tuning of parameters in a fully adaptive manner using the criterion (1.2). The criterion (1.2) appeared in previous works for the square loss only: Bayati et al. (2013); Miolane and Montanari (2018) studied (1.2) for the Lasso with $\Sigma = \mathbf{I}_p$ and (Bellec, 2020, Section 3) for the square loss and general penalty (note that for the square loss $\rho(u) = u^2/2$, (1.2) reduces to $n^2\|\mathbf{r}\|^2/(n - \hat{\text{df}})^2$ due to $\psi(u) = u$ and $\text{tr}[\mathbf{V}] = n - \hat{\text{df}}$. The property $\psi(u) = u$ of the square loss hides the subtle interplay between \mathbf{r} , $\psi(\mathbf{r})$, $\hat{\text{df}}$ and $\text{tr}[\mathbf{V}]$ in (1.2) for ρ different than the square loss).

A criterion different from (1.2) is studied in Bayati et al. (2013); Miolane and Montanari (2018) (for the Lasso and $\Sigma = \mathbf{I}_p$) and Bellec (2020) (for general loss-penalty pairs), with the purpose of estimating the out-of-sample error $\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|^2$. In the case of an M-estimator and with the notation in (1.2), this criterion is the right-hand side of the approximation

$$\|\Sigma^{1/2}(\hat{\beta} - \beta^*)\|^2 \approx \text{tr}[\mathbf{V}]^{-2}(\|\psi(\mathbf{r})\|^2(2\hat{\text{df}} - p) + \|\Sigma^{-1/2}\mathbf{X}^\top\psi(\mathbf{r})\|^2).$$

That criterion from Bayati et al. (2013); Miolane and Montanari (2018); Bellec (2020) has the drawback to require the knowledge of the covariance Σ , and is thus not readily usable unless Σ is known or can be consistently estimated. On the other hand, the criterion (1.2) is fully adaptive: it does not depend on Σ , and can thus be used even if Σ cannot itself be consistently estimated. Another line of work (Rad et al., 2020; Xu et al., 2021; Rad and Maleki, 2020) proposed the ALO criterion

$$\sum_{i=1}^n \left(r_i + \frac{H_{ii}}{V_{ii}} \psi(r_i) \right)^2 \tag{1.5}$$

(when specialized to linear models), where \mathbf{V} is the matrix defined in (1.2) and $H_{ii} = \mathbf{x}_i^\top \frac{\partial}{\partial y_i} \hat{\beta}(\mathbf{y}, \mathbf{X})$ in the notation of the present paper. This criterion differs from (1.2) proposed in the present paper, in that (1.2) replaces H_{ii} and V_{ii} by their respective averages, $\hat{\text{df}}/n = \frac{1}{n} \sum_{i=1}^n H_{ii}$ and $\text{tr}[\mathbf{V}]/n = \frac{1}{n} \sum_{i=1}^n V_{ii}$. This extra averaging step lets us prove, for non-smooth penalty functions, theoretical guarantees for selecting a loss-penalty pair by minimizing the criterion (1.2) (cf. Sections 4 and 5 below). On the other hand, we are not aware of similar theoretical guarantees for selecting a loss-penalty pair based on (1.5), since the theoretical analysis of (1.5) is so far restricted to twice continuously differentiable loss and penalty functions, with a uniform upper bound on the Lipschitz constant of the second derivatives (Rad and Maleki, 2020, Assumption 6 required in Theorem 3 and Corollary 1). This rules out the Elastic-Net and other non-smooth penalty functions typically used for high-dimensional data, as well as the Huber loss which is not twice continuously

differentiable. The criterion (1.2) of the present paper thus improves upon (1.5) since it enjoys theoretical guarantees for non-smooth penalty functions and the Huber loss.

This work leverages probabilistic results on functions of standard normal random variables [Bellec and Zhang \(2019\)](#)([Bellec, 2020](#), §6, §7) which are consequences of Stein's formula [Stein \(1981\)](#). Consequently, the main limitation of our work is that it currently requires Gaussian design for the probabilistic results, although simulations in Appendix G suggest that the results hold for more general distributions, including design with Rademacher entries. On the other hand, the differentiability result (2.1) is deterministic and does not rely on any probabilistic assumption.

2. Differentiability of regularized M-estimators

The first step towards the study of the criterion (1.2) is to justify the almost sure existence of the derivatives of $\hat{\beta}$ that appear in (1.2) through the scalar $\hat{d}f$ and the matrix \mathbf{V} in (1.2). Although the criterion (1.2) only involves the derivatives of $\hat{\beta}(\mathbf{y}, \mathbf{X})$ with respect to \mathbf{y} for a fixed \mathbf{X} , the proof of our results rely on the interplay between the derivatives with respect to \mathbf{y} and with respect to \mathbf{X} : this *differentiability structure* of M-estimators is the content of the following result.

Theorem 1 *Let Assumption A be fulfilled. For almost every (\mathbf{y}, \mathbf{X}) the map $(\mathbf{y}, \mathbf{X}) \mapsto \hat{\beta}(\mathbf{y}, \mathbf{X})$ is differentiable at (\mathbf{y}, \mathbf{X}) and there exists a matrix $\hat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ depending on (\mathbf{X}, \mathbf{y}) with $\|\Sigma^{1/2} \hat{\mathbf{A}} \Sigma^{1/2}\|_{op} \leq (n\mu)^{-1}$ s.t.*

$$\begin{aligned} (\partial/\partial y_i) \hat{\beta}(\mathbf{y}, \mathbf{X}) &= \hat{\mathbf{A}} \mathbf{X}^\top \mathbf{e}_i \psi'(r_i), \\ (\partial/\partial x_{ij}) \hat{\beta}(\mathbf{y}, \mathbf{X}) &= \hat{\mathbf{A}} \mathbf{e}_j \psi(r_i) - \hat{\mathbf{A}} \mathbf{X}^\top \mathbf{e}_i \psi'(r_i) \hat{\beta}_j, \end{aligned} \quad \text{where } r_i = y_i - \mathbf{x}_i^\top \hat{\beta}, \quad (2.1)$$

$\mathbf{e}_i \in \mathbb{R}^n$, $\mathbf{e}_j \in \mathbb{R}^p$ are canonical basis vectors, $\psi := \rho'$ and ψ' denote the derivatives. Furthermore,

$$\hat{d}f = \text{tr}[\mathbf{X}(\partial/\partial \mathbf{y}) \hat{\beta}] = \text{tr}[\mathbf{X} \hat{\mathbf{A}} \mathbf{X}^\top \text{diag}\{\psi'(\mathbf{r})\}], \quad (2.2)$$

$$\mathbf{V} = \text{diag}\{\psi'(\mathbf{r})\}(\mathbf{I}_n - \mathbf{X}(\partial/\partial \mathbf{y}) \hat{\beta}) = \text{diag}\{\psi'(\mathbf{r})\} - \text{diag}\{\psi'(\mathbf{r})\} \mathbf{X} \hat{\mathbf{A}} \mathbf{X}^\top \text{diag}\{\psi'(\mathbf{r})\} \quad (2.3)$$

satisfy $0 \leq \hat{d}f \leq n$ and $0 \leq \text{tr}[\mathbf{V}] \leq n$.

Since the same matrix $\hat{\mathbf{A}}$ appears in both the derivatives with respect to y_i and to x_{ij} , (2.1) provides relationship between $(\partial/\partial y_i) \hat{\beta}$ and $(\partial/\partial x_{ij}) \hat{\beta}$, for instance $(\partial/\partial x_{ij}) \hat{\beta} = \hat{\mathbf{A}} \mathbf{e}_j \psi(r_i) - \hat{\beta}_j (\partial/\partial y_i) \hat{\beta}$. Although the matrix $\hat{\mathbf{A}}$ is not explicit for arbitrary loss-penalty pair, closed-form expressions are available for particular examples such as the Elastic-Net penalty as discussed in Section 6.

Remark 2 *For the square loss $\rho(u) = u^2/2$, the differentiability formulae (2.1) reduce to*

$$\begin{aligned} (\partial/\partial y_l) \hat{\beta}(\mathbf{y}, \mathbf{X}) &= \hat{\mathbf{A}} \mathbf{X}^\top \mathbf{e}_l, \\ (\partial/\partial x_{ij}) \hat{\beta}(\mathbf{y}, \mathbf{X}) &= \hat{\mathbf{A}} \mathbf{e}_j (y_i - \mathbf{x}_i^\top \hat{\beta}) - \hat{\mathbf{A}} \mathbf{X}^\top \mathbf{e}_i \hat{\beta}_j \end{aligned} \quad (2.4)$$

for almost every (\mathbf{y}, \mathbf{X}) and some matrix $\hat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ depending on (\mathbf{y}, \mathbf{X}) , since in this case $\psi' = 1$.

In the simple case where g is twice continuously differentiable, (2.1) follows (Bellec and Zhang, 2019) with

$$\widehat{\mathbf{A}} = (\mathbf{X}^\top \text{diag}\{\psi'(\mathbf{r})\} \mathbf{X} + n \nabla^2 g(\widehat{\boldsymbol{\beta}}))^{-1} \quad (2.5)$$

by differentiating the KKT conditions $\mathbf{X}^\top \psi(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}) = n \nabla g(\widehat{\boldsymbol{\beta}})$. To illustrate why this is true, provided that $\widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X})$ is differentiable, if $(\mathbf{y}(t), \mathbf{X}(t))$ are smooth perturbations of (\mathbf{y}, \mathbf{X}) with $(\mathbf{y}(0), \mathbf{X}(0)) = (\mathbf{y}, \mathbf{X})$ and $\frac{d}{dt}(\mathbf{y}(t), \mathbf{X}(t))|_{t=0} = (\dot{\mathbf{y}}, \dot{\mathbf{X}})$, differentiation of $\mathbf{X}(t)^\top \psi(\mathbf{y}(t) - \mathbf{X}(t)\widehat{\boldsymbol{\beta}}(\mathbf{y}(t), \mathbf{X}(t))) = n \nabla g(\widehat{\boldsymbol{\beta}}(\mathbf{y}(t), \mathbf{X}(t)))$ at $t = 0$ and the chain rule yields

$$\dot{\mathbf{X}}^\top \psi(\mathbf{r}) - \mathbf{X}^\top \text{diag}\{\psi'(\mathbf{r})\}(\dot{\mathbf{y}} - \dot{\mathbf{X}}\widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X})) = \widehat{\mathbf{A}}^{-1} \frac{d}{dt} \widehat{\boldsymbol{\beta}}(\mathbf{y}(t), \mathbf{X}(t))|_{t=0}$$

with $\widehat{\mathbf{A}}$ in (2.5). This gives (2.1) if the penalty g is twice-differentiable. Theorem 1 reveals that for *arbitrary* convex penalty functions including non-differentiable ones, the differentiability structure (2.1) always holds, as in the case of twice differentiable penalty g , even for penalty functions such as $g(\mathbf{b}) = \mu \|\mathbf{b}\|^2/2 + \lambda \|\text{mat}(\mathbf{b})\|_{\text{nuc}}$ where $\text{mat} : \mathbb{R}^p \rightarrow \mathbb{R}^{d_1 \times d_2}$ is a linear isomorphism to the space of $d_1 \times d_2$ matrices and $\|\cdot\|_{\text{nuc}}$ is the nuclear norm: in this case by Theorem 1 there exists a matrix $\widehat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ such that (2.1) holds although no closed-form expression for $\widehat{\mathbf{A}}$ is known.

The representation (2.1) is a powerful tool as it provides explicit derivatives of quantities of interest such as $\mathbf{r} = \mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}$, $\|\psi(\mathbf{r})\|^2$ or $\|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2$. These explicit derivatives can then be used in probabilistic identities and inequalities that involve derivatives, for instance Stein's formulae (Stein, 1981), the Gaussian Poincaré inequality (Boucheron et al., 2013, Theorem 3.20), or normal approximations (Chatterjee, 2009; Bellec and Zhang, 2019).

Remark 3 *Similar derivative formulae hold if an intercept is included in the minimization, as in*

$$(\widehat{\beta}_0(\mathbf{y}, \mathbf{X}), \widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X})) = \underset{b_0 \in \mathbb{R}, \mathbf{b} \in \mathbb{R}^p}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \rho(y_i - b_0 - \mathbf{x}_i^\top \mathbf{b}) + g(\mathbf{b}) \quad (2.6)$$

Let Assumption A be fulfilled, and assume further $\|\psi'(\mathbf{r})\|_2 > 0$ with $\mathbf{r} := \mathbf{y} - \mathbf{1}_n \widehat{\beta}_0 - \mathbf{X}^\top \widehat{\boldsymbol{\beta}}$ where $\mathbf{1}_n = (1, \dots, 1)^\top \in \mathbb{R}^n$. For almost every (\mathbf{y}, \mathbf{X}) the map $(\mathbf{y}, \mathbf{X}) \mapsto \widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X})$ is differentiable at (\mathbf{y}, \mathbf{X}) , and there exists $\widehat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ depending on (\mathbf{y}, \mathbf{X}) with $\|\boldsymbol{\Sigma}^{1/2} \widehat{\mathbf{A}} \boldsymbol{\Sigma}^{1/2}\|_{\text{op}} \leq (n\mu)^{-1}$ such that

$$(\partial/\partial y_i) \widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) = \widehat{\mathbf{A}} \mathbf{X}^\top \boldsymbol{\Psi}' \mathbf{e}_i, \quad (\partial/\partial x_{ij}) \widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) = \widehat{\mathbf{A}} \mathbf{e}_j \psi(r_i) - \widehat{\mathbf{A}} \mathbf{X}^\top \boldsymbol{\Psi}' \mathbf{e}_i \widehat{\beta}_j, \quad (2.7)$$

where $\mathbf{e}_i \in \mathbb{R}^n, \mathbf{e}_j \in \mathbb{R}^p$ are canonical basis vectors, $\psi = \rho'$ and $\boldsymbol{\Psi}' := \text{diag}\{\psi'(\mathbf{r})\} - \psi'(\mathbf{r})\psi'(\mathbf{r})^\top / \sum_{i \in [n]} \psi'(r_i)$.

3. Distribution of individual residuals

We now turn to the distribution of a single residual $r_i = y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}$ for some fixed observation $i \in \{1, \dots, n\}$ (for instance, fix $i = 1$). By leveraging the differentiability structure (2.1) and the normal approximation from Bellec and Zhang (2019), the following result provides a clear picture of the distribution of r_i .

Theorem 4 *Let Assumption A be fulfilled and let $\widehat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ be given by Theorem 1. Then for every $i = 1, \dots, n$ there exists $Z_i \sim N(0, 1)$ independent of ε_i such that*

$$\left| \left(r_i + \text{tr}[\boldsymbol{\Sigma} \widehat{\mathbf{A}}] \psi(r_i) \right) - \left(\varepsilon_i + \|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\| Z_i \right) \right| \leq O_P(n^{-1/4})(|\psi(\varepsilon_i)| + \|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|) \quad (3.1)$$

Furthermore, if ε_i has a fixed distribution F , there exists a bivariate variable $(\tilde{\varepsilon}_i^n, \tilde{Z}_i^n)$ converging in distribution to the product measure $F \otimes N(0, 1)$ such that

$$r_i + \text{tr}[\boldsymbol{\Sigma} \widehat{\mathbf{A}}] \psi(r_i) = \tilde{\varepsilon}_i^n + \|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\| \tilde{Z}_i^n. \quad (3.2)$$

If ε_i has a fixed distribution F and Assumption B holds then $|\psi(\varepsilon_i)| + \|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^)\| = O_P(1)$.*

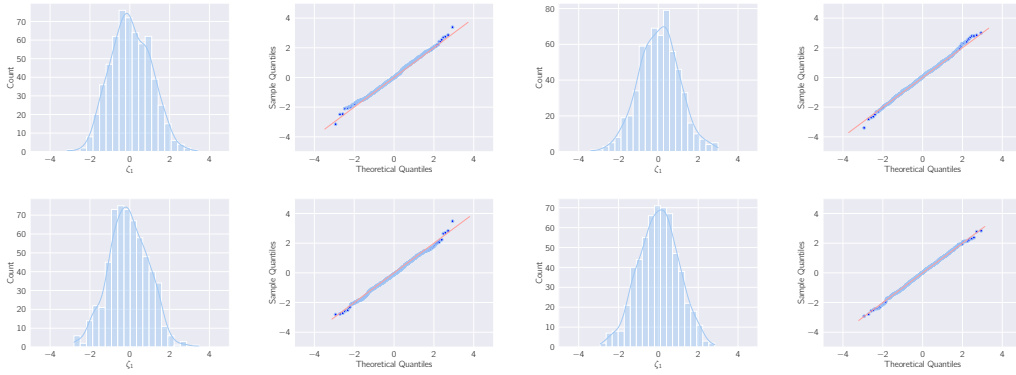


Figure 2: Histogram and QQ-plot for ζ_1 in (3.3) under Huber Elastic-Net regression for different choices of tuning parameters (λ, τ) . Left Top: $(0.036, 10^{-10})$, Right Top: $(0.054, 0.01)$, Left Bottom: $(0.036, 0.01)$, Right Bottom: $(0.024, 0.1)$. Each figure contains 600 data points generated with anisotropic design matrix and iid ε_i from the t -distribution with 2 degrees of freedom. A detailed setup is provided in Section 6.

Theorem 4 is a formal statement regarding the informal normal approximation

$$\zeta_i := \frac{r_i + \text{tr}[\boldsymbol{\Sigma} \widehat{\mathbf{A}}] \psi(r_i) - \varepsilon_i}{\|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|} \approx N(0, 1). \quad (3.3)$$

Simulations in Figure 2 confirm the normality of ζ_i for the Huber loss with Elastic-Net penalty and four combinations of tuning parameters. For the square loss $\rho(u) = u^2/2$, because $\psi(u) = u$, asymptotic normality of the residuals hold in the following form.

Theorem 5 *Let Assumption A hold with $\rho(u) = u^2/2$ and $\varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. Then for $i = 1$,*

$$\frac{(1 + \text{tr}[\boldsymbol{\Sigma} \widehat{\mathbf{A}}])(y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}})}{(\sigma^2 + \|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2)^{1/2}} \rightarrow^d N(0, 1) \quad \text{as } n \rightarrow +\infty. \quad (3.4)$$

It is informative to provide a sketch of the proof of Theorem 4 to explain the appearance of $\psi(r_i)$ and $\text{tr}[\Sigma\hat{\mathbf{A}}]$ in the normal approximation results (3.1) and (3.3). A variant of the normal approximation of Bellec and Zhang (2019) proved in the supplement states that for a differentiable function $\mathbf{f} : \mathbb{R}^q \rightarrow \mathbb{R}^q \setminus \{\mathbf{0}\}$ and $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}_q)$, there exists $Z \sim N(0, 1)$ such that

$$\mathbb{E} \left[\left| \frac{\mathbf{f}(\mathbf{z})^\top \mathbf{z} - \sum_{k=1}^q (\partial/\partial z_k) f_k(\mathbf{z})}{\|\mathbf{f}(\mathbf{z})\|} - Z \right|^2 \right] \leq C_1 \mathbb{E} \left[\frac{\sum_{k=1}^q \|(\partial/\partial z_k) \mathbf{f}(\mathbf{z})\|^2}{\|\mathbf{f}(\mathbf{z})\|^2} \right]. \quad (3.5)$$

Some technical hurdles aside, the proof sketch is the following: Apply the previous display to $q = p$, $\mathbf{z} = \Sigma^{-1/2} \mathbf{x}_i$ conditionally on $(\varepsilon, (\mathbf{x}_l)_{l \in [n] \setminus \{i\}})$ and to $\mathbf{f}(\mathbf{z}) = \Sigma^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)$ in the simple case where $\boldsymbol{\beta}^* = \mathbf{0}$ (this amounts to performing a change of variable by translation of $\hat{\boldsymbol{\beta}}$ to $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*$). Then the right-hand side of the previous display is negligible in probability compared to Z , and in the left-hand side $\mathbf{f}(\mathbf{z})^\top \mathbf{z} = \mathbf{x}_i^\top (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)$ and $\sum_{k=1}^q (\partial/\partial z_k) f_k(\mathbf{z}) \approx \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(r_i)$ as the second term in (2.1) is negligible. This completes the sketch of the proof of (3.3).

Proximal operator representation. From the above asymptotic normality results, a stochastic representation for the i -th residual $r_i = y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$ can be obtained as follows: With $\text{prox}[t\rho](u)$ the proximal operator of $x \mapsto t\rho(x)$ defined as the unique solution $z \in \mathbb{R}$ of equation $z + t\psi(z) = u$,

$$r_i = y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}} = \text{prox}[\hat{t}\rho](\tilde{\varepsilon}_i^n + \|\Sigma^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\| \tilde{Z}_i^n) \quad \text{with } \hat{t} = \text{tr}[\Sigma\hat{\mathbf{A}}]$$

where $(\tilde{\varepsilon}_i^n, \tilde{Z}_i^n)$ converges in distribution to the product measure $F \otimes N(0, 1)$ where F is the law of ε_i .

4. A proxy of the out-of-sample error if Σ is known

The approximations of the previous sections for $r_i + \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(r_i)$ and the fact that ε_i is independent of $Z_i \sim N(0, 1)$ in (3.1) suggest that $(r_i + \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(r_i))^2 \approx \varepsilon_i^2 + \|\Sigma^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2 Z_i^2$; and averaging over $\{1, \dots, n\}$ one can hope for the approximation $\|\mathbf{r} + \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(\mathbf{r})\|^2/n \approx \|\varepsilon\|^2/n + \|\Sigma^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2$. The following result makes this heuristic precise.

Theorem 6 *Let Assumption A be fulfilled and $\hat{\mathbf{A}}$ be given by Theorem 1. Then*

$$\|\Sigma^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2 + \|\varepsilon\|^2/n = \|\mathbf{r} + \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(\mathbf{r})\|^2/n + O_P(n^{-1/2}) \text{Rem},$$

where $\text{Rem} := \|\Sigma^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2 + \frac{1}{n} \|\psi(\mathbf{r})\|^2 + (\|\Sigma^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2 + \frac{1}{n} \|\psi(\mathbf{r})\|^2)^{1/2} \|\frac{1}{\sqrt{n}}\varepsilon\|$. Thus

$$\|\Sigma^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2 + \|\varepsilon\|^2/n = (1 + O_P(n^{-1/2})) \|\mathbf{r} + \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(\mathbf{r})\|^2/n.$$

Theorem 6 provides a first candidate, $\|\mathbf{r} + \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(\mathbf{r})\|^2/n$ to estimate

$$\|\Sigma^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2 + \|\varepsilon\|^2/n. \quad (4.1)$$

Estimation of (4.1) is useful as $\|\varepsilon\|^2/n$ is independent of the choice of the estimator $\hat{\boldsymbol{\beta}}$ and in particular independent of the chosen loss-penalty pair in (1.1). Given two or more estimators (1.1), choosing the one with smallest $\|\mathbf{r} + \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(\mathbf{r})\|^2$ is thus a good proxy for minimizing the out-of-sample error.

Corollary 7 Let $\widehat{\beta}, \widetilde{\beta}$ be two M -estimators (1.1) Assumption A with loss-penalty pair (ρ, g) and $(\widetilde{\rho}, \widetilde{g})$ respectively. Assume that both satisfy Assumption A and let $\psi = \rho'$ and $\widetilde{\psi} = \widetilde{\rho}'$. Let $\mathbf{r} = \mathbf{y} - \mathbf{X}\widehat{\beta}, \widetilde{\mathbf{r}} = \mathbf{y} - \mathbf{X}\widetilde{\beta}$ be the residuals, $\widehat{\mathbf{A}}, \widetilde{\mathbf{A}}$ be the corresponding matrices of size $p \times p$ given by Theorem 1. Further assume that both estimators satisfy Assumption B and that ε has iid coordinates independent with $\mathbb{E}[|\varepsilon_i|^{1+q}] \leq M$ for constants $q \in (0, 1), M > 0$ independent of n, p . Let $\Omega = \{\|\mathbf{X}\Sigma^{-1/2}\|_{op} \leq 2\sqrt{n} + \sqrt{p}\} \cap \{\|\varepsilon\|^2 \leq n^{2/(1+q)}\}$. Then for any $\eta > 0$ independent of n, p there exists $C(\gamma, \mu, \eta, q, M) > 0$ depending only on $\{\gamma, \mu, \eta, q, M\}$ such that

$$\begin{aligned} \mathbb{P}\left(\|\Sigma^{1/2}(\widehat{\beta} - \beta^*)\|^2 - \|\Sigma^{1/2}(\widetilde{\beta} - \beta^*)\|^2 > \eta, \|\mathbf{r} + \text{tr}[\Sigma\widehat{\mathbf{A}}]\psi(\mathbf{r})\|^2 \leq \|\widetilde{\mathbf{r}} + \text{tr}[\Sigma\widetilde{\mathbf{A}}]\widetilde{\psi}(\widetilde{\mathbf{r}})\|^2\right) \\ \leq C(\gamma, \mu, \eta, q, M)n^{-q/(1+q)} + \mathbb{P}(\Omega^c) \rightarrow 0. \end{aligned}$$

Provided that the noise random variables ε_i have at least $1 + q$ moments, Theorem 7 implies that with probability approaching one given two M -estimators $\widehat{\beta}$ and $\widetilde{\beta}$, choosing the estimator corresponding to the smallest criteria among $\|\mathbf{r} + \text{tr}[\Sigma\widehat{\mathbf{A}}]\mathbf{r}\|^2$ and $\|\widetilde{\mathbf{r}} + \text{tr}[\Sigma\widetilde{\mathbf{A}}]\widetilde{\mathbf{r}}\|^2$ leads to the smallest out-of-sample error, up to any small constant $\eta > 0$. This allows noise random variables ε_i with infinite variance. A similar result can be obtained to select among K different M -estimators (1.1).

Corollary 8 As in Theorem 7, assume $\mathbb{E}[|\varepsilon_i|^{1+q}] \leq M$ and let $\widehat{\beta}_1, \dots, \widehat{\beta}_K$ be M -estimators of the form (1.1) with loss-penalty pair (ρ_k, g_k) satisfying Assumptions A and B. For each $k = 1, \dots, K$, let $\mathbf{r}_k = \mathbf{y} - \mathbf{X}\widehat{\beta}_k$ be the residuals and $\widehat{\mathbf{A}}_k$ be the corresponding matrix of size $p \times p$ from Theorem 1. Let $\hat{k} \in \text{argmin}_{k=1, \dots, K} \|\mathbf{r}_k + \text{tr}[\Sigma\widehat{\mathbf{A}}_k]\psi_k(\mathbf{r}_k)\|$ where $\psi_k = \rho'_k$. Then if $(\gamma, \mu, \eta, q, M)$ are constants independent of n, p

$$\mathbb{P}(\|\Sigma^{1/2}(\widehat{\beta}_{\hat{k}} - \beta^*)\|^2 > \min_{k=1, \dots, K} \|\Sigma^{1/2}(\widehat{\beta}_k - \beta^*)\|^2 + \eta) \rightarrow 0 \quad \text{if } K = o(n^{q/(1+q)}).$$

In other words, if $K = o(n^{q/(1+q)})$, the selector \hat{k} picks an optimal M -estimator in the sense

$$\|\Sigma^{1/2}(\widehat{\beta}_{\hat{k}} - \beta^*)\|^2 - \min_{k=1, \dots, K} \|\Sigma^{1/2}(\widehat{\beta}_k - \beta^*)\|^2 \xrightarrow{P} 0.$$

Given K different loss-penalty pairs and the corresponding M -estimators in (1.1), minimizing the criterion $\|\mathbf{r} + \text{tr}[\Sigma\widehat{\mathbf{A}}]\mathbf{r}\|$ thus provably selects a loss-penalty pair that leads to an optimal out-of-sample error, up to an arbitrary small constant $\eta > 0$ independent of n, p . The requirement $K = o(n^{q/(1+q)})$ means that the cardinality of the collection of M -estimators to select from should grow more slowly than a power of n . This is typically satisfied for default tuning parameter grids in popular libraries (e.g., `sklearn.linear_model.Lasso` from Pedregosa et al. (2011)) with tuning parameters evenly spaced in a log-scale that consequently have cardinality logarithmic in the parameter range. The major drawback of the criterion $\|\mathbf{r} + \text{tr}[\Sigma\widehat{\mathbf{A}}]\mathbf{r}\|$ is the dependence through $\text{tr}[\Sigma\widehat{\mathbf{A}}]$ on the covariance Σ of the design, which is typically unknown. The next section introduces an estimator of $\text{tr}[\Sigma\widehat{\mathbf{A}}]$ that does not require the knowledge of Σ .

5. Degrees of freedom and estimating $\text{tr}[\Sigma\widehat{\mathbf{A}}]$ without the knowledge of Σ

This section focuses on estimating $\text{tr}[\Sigma\widehat{\mathbf{A}}]$. The matrix $\widehat{\mathbf{A}}$ from Theorem 1 can be estimated from the data (\mathbf{y}, \mathbf{X}) in the sense that $\widehat{\mathbf{A}}$ is a measurable function of (\mathbf{y}, \mathbf{X}) (thanks to the observation

that derivatives are limits, and limits of measurable functions are again measurable). The difficulty is thus to estimate $\text{tr}[\Sigma\hat{\mathbf{A}}]$ without the knowledge of Σ . To illustrate this difficulty, consider Ridge regression with square loss $\rho(u) = u^2/2$ and penalty $g(\mathbf{b}) = \tau\|\mathbf{b}\|^2/2$. Then $\hat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) = (\mathbf{X}^\top \mathbf{X} + \tau n \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$ and $\hat{\mathbf{A}}$ in Theorem 1 is given explicitly by $\hat{\mathbf{A}} = (\mathbf{X}^\top \mathbf{X} + \tau n \mathbf{I}_p)^{-1}$ and

$$\text{tr}[\Sigma\hat{\mathbf{A}}] = \text{tr}[(\mathbf{G}^\top \mathbf{G} + n\tau\Sigma^{-1})^{-1}], \quad \text{where } \mathbf{G} = \mathbf{X}\Sigma^{-1/2}.$$

Above, \mathbf{G} is a random matrix with iid $N(0, 1)$ entries the value of $\text{tr}[\Sigma\hat{\mathbf{A}}]$ is highly dependent on the spectrum of Σ^{-1} . In this particular case, the limit of $\text{tr}[(\mathbf{G}^\top \mathbf{G} + n\tau\Sigma^{-1})^{-1}]$ can be obtained using random matrix theory (Marčenko and Pastur, 1967) as the limiting behavior of the Stieltjes transform of $\mathbf{G}^\top \mathbf{G}/n + \tau\Sigma^{-1}$ and its spectral distribution is known; however the limit of the spectral distribution depends on the spectrum of $\tau\Sigma^{-1}$. This is not desirable here as we wish to construct estimators that require no knowledge on Σ . For more involved loss-penalty pairs such as the Elastic-Net in Example 1, such random matrix theory results do not apply as $\text{tr}[\Sigma\hat{\mathbf{A}}]$ depends on the random support of $\hat{\boldsymbol{\beta}}$.

Instead, we do not rely on known random matrix theory results. With the matrix $\hat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ given by Theorem 1, our proposal to estimate $\text{tr}[\Sigma\hat{\mathbf{A}}]$ is the ratio $\hat{\text{d}}\text{f}/\text{tr}[\mathbf{V}]$ with $\hat{\text{d}}\text{f}$ and \mathbf{V} in (2.2)-(2.3). Both the scalar $\hat{\text{d}}\text{f}$ and the matrix $\mathbf{V} \in \mathbb{R}^{n \times n}$ are observable; in particular they do not depend on Σ .

Theorem 9 *Let Assumption A be fulfilled and $\hat{\mathbf{A}}$ be given by Theorem 1. Then*

$$\mathbb{E}[|\text{tr}[\Sigma\hat{\mathbf{A}}] \text{tr}[\mathbf{V}]/n - \hat{\text{d}}\text{f}/n|] \leq C_2(\gamma, \mu)n^{-1/2}. \quad (5.1)$$

Simulations in Figure 3 and table 1 confirm that the approximation $\text{tr}[\Sigma\hat{\mathbf{A}}] \approx \hat{\text{d}}\text{f}/\text{tr}[\mathbf{V}]$ is accurate for the Huber loss with Elastic-Net penalty. For the square loss, $\psi' = 1$ and $\text{tr}[\mathbf{V}] = n - \hat{\text{d}}\text{f}$ so that (5.1) becomes $\mathbb{E}[(1 - \hat{\text{d}}\text{f}/n)(1 + \text{tr}[\Sigma\hat{\mathbf{A}}]) - 1] \leq C_3(\gamma, \mu)n^{-1/2}$ and the following result holds.

Corollary 10 *Let Assumption A be fulfilled with $\rho(u) = u^2/2$ and $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. Then $(1 - \hat{\text{d}}\text{f}/n)(1 + \text{tr}[\Sigma\hat{\mathbf{A}}]) \xrightarrow{\mathbb{P}} 1$ and the normality (3.4) holds with $1 + \text{tr}[\Sigma\hat{\mathbf{A}}]$ replaced by $(1 - \hat{\text{d}}\text{f}/n)^{-1}$.*

For general loss ρ , the criterion (1.2) replaces $\text{tr}[\Sigma\hat{\mathbf{A}}]$ by $\hat{\text{d}}\text{f}/\text{tr}[\mathbf{V}]$ in the proxy of the out-of-sample error $\|\mathbf{r} + \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(\mathbf{r})\|^2$ studied in the previous section. Thanks to (5.1), this replacement preserves the good properties of $\|\mathbf{r} + \text{tr}[\Sigma\hat{\mathbf{A}}]\psi(\mathbf{r})\|^2$ proved in Theorems 7 and 8.

Theorem 11 *For $k = 1, \dots, K$, let (ρ_k, g_k) be a loss-penalty pair satisfying Assumptions A and B with $\psi_k = \rho'_k$, let $\hat{\boldsymbol{\beta}}_k, \mathbf{r}_k, \hat{\mathbf{A}}_k$ be the corresponding M -estimator residual vector and matrix of size $p \times p$ given by Theorem 1 as in Theorem 8 and let $\hat{\text{d}}\text{f}_k = \text{tr}[\mathbf{X}\hat{\mathbf{A}}_k\mathbf{X}^\top \text{diag}\{\psi'_k(\mathbf{r}_k)\}]$ and $\mathbf{V}_k = \text{diag}\{\psi'_k(\mathbf{r}_k)\}(\mathbf{I}_n - \mathbf{X}\hat{\mathbf{A}}_k\mathbf{X}^\top \text{diag}\{\psi'_k(\mathbf{r}_k)\})$. For a small constant $\eta > 0$ independent of n, p , say $\eta = 0.05$, define*

$$\hat{k} \in \underset{k=1, \dots, K}{\text{argmin}} \left\| \mathbf{r}_k + \frac{\hat{\text{d}}\text{f}_k}{\text{tr}[\mathbf{V}_k]} \psi_k(\mathbf{r}_k) \right\|^2 \quad \text{subject to} \quad \frac{1}{n} \sum_{i=1}^n \psi'_k(r_{ki}) \geq \eta.$$

If ε_i has $1 + q$ moments in the sense that $\mathbb{E}[|\varepsilon_i|^{1+q}] \leq M$ for constants $q \in (0, 1), M > 0$. If $(M, q, \eta, \mu, \gamma)$ and $\tilde{\eta} > 0$ are independent of n, p then

$$\mathbb{P}\left(\|\Sigma^{1/2}(\hat{\boldsymbol{\beta}}_{\hat{k}} - \boldsymbol{\beta}^*)\| > \min_{k=1, \dots, K; \frac{1}{n} \sum_{i=1}^n \psi'_k(r_{ki}) \geq \eta} \|\Sigma^{1/2}(\hat{\boldsymbol{\beta}}_k - \boldsymbol{\beta}^*)\| + \tilde{\eta}\right) \rightarrow 0 \quad \text{if } K = o(n^{q/(1+q)}).$$

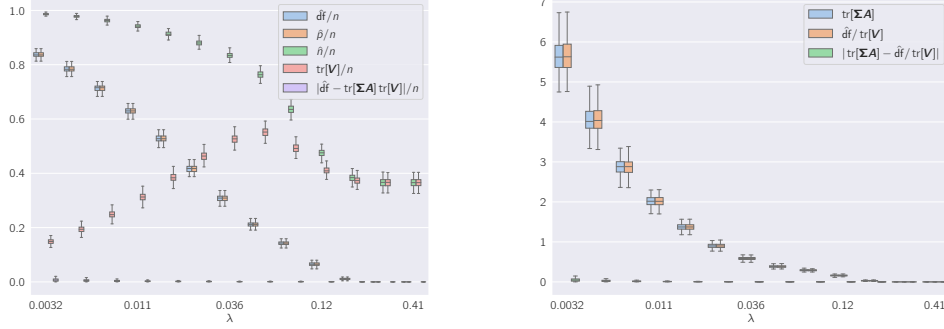


Figure 3: Above: Boxplots for $\hat{d}f$, $\hat{\rho}$, \hat{n} , $\text{tr}[\mathbf{V}]$, $\text{tr}[\Sigma\hat{\mathbf{A}}]$ and $|\text{tr}[\Sigma\hat{\mathbf{A}}] - \hat{d}f/\text{tr}[\mathbf{V}]|$ in Huber Elastic-Net regression with $\tau = 10^{-10}$ and $\lambda \in [0.0032, 0.41]$. Each box contains 200 data points. Below: heatmaps for $\hat{d}f/n$, $\text{tr}[\mathbf{V}]/n$ and $\hat{n}/n = \sum_{i=1}^n \psi'(r_i)/n$ under the simulation setup in Figure 1. The detailed simulation setup is given in Section 6.

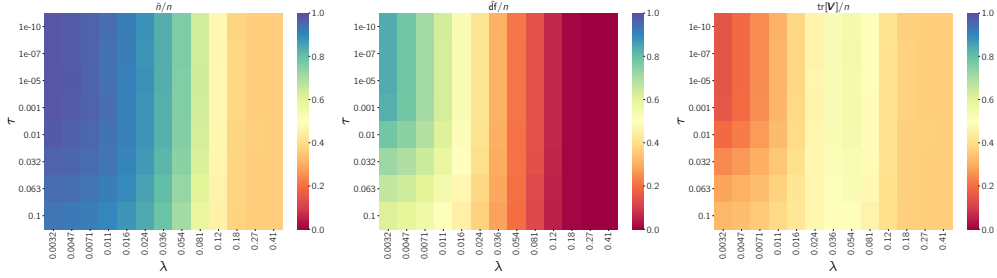


Figure 1 illustrates on simulations the success of the criterion (1.2) over a grid of tuning parameters for M -estimators with the Huber loss and Elastic-Net penalty. The criterion (1.2) is thus successful at selecting a M -estimator with smallest out-of-sample error up to an additive constant $\tilde{\eta}$, among those M -estimators indexed in $\{1, \dots, K\}$ that are such that $\frac{1}{n} \sum_{i=1}^n \psi'_k(r_{ki}) \geq \eta$. On the one hand it is unclear to us whether the restriction $\frac{1}{n} \sum_{i=1}^n \psi'_k(r_{ki}) \geq \eta$ can be omitted. On the other hand there is a practical meaning in excluding M -estimators with small $\frac{1}{n} \sum_{i=1}^n \psi'_k(r_{ki})$: For the Huber loss $H(u) := u^2/2$ for $|u| \leq 1$ and $|u| - 1/2$ for $|u| \geq 1$ the quantity $\frac{1}{n} \sum_{i=1}^n \psi'_k(r_{ki})$ is the number of data points in $\{1, \dots, n\}$ such that the residual $y_i - \mathbf{x}_i^\top \hat{\beta}_k$ fall within the quadratic regime of the loss function. Observations $i \in \{1, \dots, n\}$ that fall in the linear regime of the loss are excluded from the fit, in the sense that for some i with $r_{ki} = y_i - \mathbf{x}_i^\top \hat{\beta}_k > 1$, replacing y_i by $\tilde{y}_i = y_i + 1000$ (or any positive value) does not change the M -estimator solution $\hat{\beta}_k$ (this can be seen from the KKT conditions directly, or by integration the derivative with respect to y_i in (2.1)). Thus the constraint $\frac{1}{n} \sum_{i=1}^n \psi'_k(r_{ki}) \geq \eta$ requires that at most a constant fraction of the observations are excluded from the fit (or equivalently, at least a constant fraction of the n observations participate in the fit). For scaled versions of the Huber loss, $\rho_k(u) = a^2 H(a^{-1}u)$ for some $a > 0$, the value $\hat{n} = \frac{1}{n} \sum_{i=1}^n \psi'_k(r_{ki})$ again counts the number of residuals falling in the quadratic regime of the loss, i.e., the number of observations participating in the fit. The heatmaps of Figure 3 illustrate \hat{n} in a simulation for a wide range of parameters. Similarly, for smooth robust loss functions such as

$\rho_k(u) = \sqrt{1+u^2}$, the constraint $\frac{1}{n} \sum_{i=1}^n \psi'_k(r_{ki}) \geq \eta$ requires that at most a constant fraction of the n observations are such that $\psi'_k(r_{ki}) < \eta/2$, i.e., such that the second derivative ψ''_k is too small (and the loss ρ_k too flat).

Theorems 1, 5, 6 and 9 provide our general results applicable to a single regularized M -estimator (1.1) while corollaries such as Theorem 11 are obtained using the union bound. The next section specializes our results and notation to the Huber loss with Elastic-Net penalty and details the simulation setup used in the figures.

6. Example and simulation setting: Huber loss with Elastic-Net penalty

In simulations and in the example below, we focus on the loss-penalty pair

$$\rho(u; \Lambda) = \Lambda^2 H(\Lambda^{-1}u), \quad g(\mathbf{b}; \lambda, \tau) = \lambda \|\mathbf{b}\|_1 + (\tau/2) \|\mathbf{b}\|_2^2 \quad (6.1)$$

for tuning parameters $\Lambda, \lambda, \tau \geq 0$ where $H(u) := u^2/2$ for $|u| \leq 1$ and $|u| - 1/2$ for $|u| \geq 1$.

Example 1 With (ρ, g) in (6.1), matrix $\hat{\mathbf{A}}$ in (2.1) matrix \mathbf{V} in (2.3) and $\hat{\text{df}}$ in (2.2) we have

$$\begin{aligned} \hat{\mathbf{A}}_{\hat{S}, \hat{S}} &= (\mathbf{X}_{\hat{S}}^\top \text{diag}\{\psi'(\mathbf{r})\} \mathbf{X}_{\hat{S}} + n\tau \mathbf{I}_{\hat{p}})^{-1}, \quad A_{i,j} = 0 \text{ if } i \notin \hat{S} \text{ or } j \notin \hat{S}, \\ \mathbf{V} &= \text{diag}\{\psi'(\mathbf{r})\} - \text{diag}\{\psi'(\mathbf{r})\} \mathbf{X}_{\hat{S}} \hat{\mathbf{A}}_{\hat{S}, \hat{S}} \mathbf{X}_{\hat{S}}^\top \text{diag}\{\psi'(\mathbf{r})\}, \\ \hat{\text{df}} &= \text{tr}[\mathbf{X}_{\hat{S}} \hat{\mathbf{A}}_{\hat{S}, \hat{S}} \mathbf{X}_{\hat{S}}^\top \text{diag}\{\psi'(\mathbf{r})\}], \end{aligned} \quad (6.2)$$

where \hat{S} is the active set $\{j \in [p] : \hat{\beta}_j \neq 0\}$ and \hat{p} is the size of \hat{S} ; $\mathbf{X}_{\hat{S}}$ is the submatrix of \mathbf{X} selecting columns with index in \hat{S} and $\hat{\mathbf{A}}_{\hat{S}, \hat{S}}$ is the submatrix of $\hat{\mathbf{A}}$ with entries indexed in $\hat{S} \times \hat{S}$.

(λ, τ)	(0.036, 10^{-10})	(0.054, 0.01)	(0.036, 0.01)	(0.024, 0.1)
$\hat{\text{df}}/n$	0.31 ± 0.012	0.21 ± 0.0095	0.3 ± 0.011	0.37 ± 0.0093
\hat{p}/n	0.31 ± 0.012	0.22 ± 0.0098	0.31 ± 0.012	0.47 ± 0.014
\hat{n}/n	0.83 ± 0.011	0.76 ± 0.014	0.83 ± 0.012	0.84 ± 0.012
$\text{tr}[\Sigma \mathbf{A}]$	0.58 ± 0.039	0.39 ± 0.027	0.58 ± 0.038	0.8 ± 0.038
$ \text{tr}[\Sigma \mathbf{A}] - \hat{\text{df}}/\text{tr}[\mathbf{V}] $	0.0019 ± 0.0015	0.0015 ± 0.0012	0.0021 ± 0.0016	0.0023 ± 0.0017
$\ \Sigma^{1/2}(\hat{\beta} - \beta^*)\ ^2$	1.3 ± 0.18	1.7 ± 0.25	1.3 ± 0.19	1.9 ± 0.21
ζ_1	0.056 ± 1	0.021 ± 1	0.0044 ± 1	0.042 ± 0.97

Table 1: Simulation for Huber Elastic-Net regression under different choices of (λ, τ) . $(n, p) = (1001, 1000)$. For each choice of (λ, τ) , 600 data points are simulated with anisotropic design matrix and i.i.d. t -distributed noises with 2 degrees of freedom. A detailed setup is provided in Section 6.

The identities (6.2) are proved in (Bellec, 2020, §2.6). Simulations in Figures 1 to 3 and table 1 illustrate typical values for $\hat{\text{df}}$, $\text{tr}[\mathbf{V}]$, $\text{tr}[\Sigma \hat{\mathbf{A}}]$, the out-of-sample error and the criterion (1.2), $\hat{n} = \sum_{i=1}^n \psi'(r_i)$ and $\hat{p} = |\hat{S}|$ under anisotropic Gaussian design and heavy-tailed ε_i . The simulation setup is as follows.

Data Generation Process. Simulation data are generated from a linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\varepsilon}$ with anisotropic Gaussian design $\boldsymbol{\Sigma}$ and heavy-tail noise vector $\boldsymbol{\varepsilon}$. The design matrix \mathbf{X} has $n = 1001$ rows and $p = 1000$ columns. Each row of \mathbf{X} is i.i.d. $N(\mathbf{0}, \boldsymbol{\Sigma})$, with the same $\boldsymbol{\Sigma}$ across all repetitions, generated once by $\boldsymbol{\Sigma} = \mathbf{R}^\top \mathbf{R}/(2p)$ with $\mathbf{R} \in \mathbb{R}^{2p \times p}$ being a Rademacher matrix with i.i.d. entries $\mathbb{P}(\mathbf{R}_{ij} = \pm 1) = \frac{1}{2}$. The true signal vector $\boldsymbol{\beta}^* \in \mathbb{R}^p$ has its first 100 coordinates set to $p^{1/2}/100 = \sqrt{10}/10$ and the rest 900 coordinates set to 0. The noise vector $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ has i.i.d. entries from the t-distribution with 2 degrees of freedom (so that $\text{Var}[\varepsilon_i] = \infty$, i.e., ε_i is heavy-tailed).

Estimation Process. Each dataset (\mathbf{y}, \mathbf{X}) is fitted by a Huber Elastic-Net estimator with loss-penalty pair in (6.1). We focus on 2d heatmaps with respect to the two penalty parameters (λ, τ) of the penalty; to this end the Huber loss parameter Λ is set to $\Lambda = 0.054n^{1/2}$ and a grid for (λ, τ) is then set so that $\hat{\text{df}}/n$ varies on the grid from 0 to 1 (cf. the middle heatmap in Figure 3). The Elastic-Net penalty $g(\mathbf{b}; \lambda, \tau) = \lambda \|\mathbf{b}\|_1 + (\tau/2) \|\mathbf{b}\|_2^2$ is used with $(\lambda, \tau) \in \{(0.036, 10^{-10}), (0.054, 0.01), (0.036, 0.01), (0.024, 0.1)\}$ in Figure 2 and table 1, $(\lambda, \tau) \in [0.0032, 0.41] \times \{10^{-10}\}$ in Figure 3, and $(\lambda, \tau) \in [0.0032, 0.041] \times [10^{-10}, 0.1]$ in Figure 1. More simulation results are provided in the supplementary materials. In these simulations, the criterion (1.5) from Rad and Maleki (2020) was also computed and was not noticeably different from (1.2), cf. the lower half of Figures 6 and 7.

7. Relaxing the strong convexity assumption

While previous results rely heavily on the μ -strong convexity assumption (with respect to $\boldsymbol{\Sigma}$, as stated in the last part of Assumption A), the proof of the following proposition presents a device that lets us generalize the results under the following condition: For any ε , there exists an open set $U_\varepsilon \subset \mathbb{R}^{n \times p}$ such that the mapping

$$\Phi_\varepsilon : \begin{cases} U_\varepsilon \rightarrow \mathbb{R}^{n+p}, \\ \mathbf{X} \mapsto \frac{(n^{-1/2}\psi(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X})), \boldsymbol{\Sigma}^{1/2}(\hat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) - \boldsymbol{\beta}))}{(\frac{1}{n}\|\psi(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}))\|^2 + \|\boldsymbol{\Sigma}^{1/2}(\hat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) - \boldsymbol{\beta})\|^2)^{1/2}} \end{cases} \text{ is } \frac{L}{\sqrt{n}}\text{-Lipschitz.} \quad (7.1)$$

In this definition, ε is held fixed and $\hat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X})$ is the composition of $\hat{\boldsymbol{\beta}}$ with the function $\mathbf{X} \mapsto (\varepsilon + \mathbf{X}\boldsymbol{\beta}, \mathbf{X})$ so that Φ_ε is a function of \mathbf{X} only. The following proposition shows that strong convexity on the penalty function can be relaxed, provided that the above Lipschitz condition holds and the expectations are restricted to the event $\{\mathbf{X} \in U_\varepsilon\}$.

Proposition 12 *Let $L, \gamma > 0$ be constants and assume $p/n \leq \gamma$. Consider a convex differentiable loss $\rho : \mathbb{R} \rightarrow \mathbb{R}$ such that $\psi = \rho'$ is 1-Lipschitz and a convex penalty g , assume \mathbf{X} has iid $N(\mathbf{0}, \boldsymbol{\Sigma})$ rows with invertible $\boldsymbol{\Sigma}$ and the noise $\boldsymbol{\varepsilon}$ is independent of \mathbf{X} with continuous distribution. Assume that for some open $U_\varepsilon \subset \mathbb{R}^{n \times p}$, the Lipschitz condition (7.1) holds and that almost everywhere in U_ε , the derivative formulae (2.1) hold for some matrix $\hat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ satisfying $\|\boldsymbol{\Sigma}^{1/2} \hat{\mathbf{A}} \boldsymbol{\Sigma}^{1/2}\|_{op} \leq L/n$. Then*

$$\mathbb{E} \left[I\{\mathbf{X} \in U_\varepsilon\} \left| \frac{\text{tr}[\boldsymbol{\Sigma} \hat{\mathbf{A}}] \text{tr}[\mathbf{V}]}{n} - \frac{\text{df}}{n} \right| \right] \leq \frac{C_4(\gamma, L)}{\sqrt{n}} \quad (7.2)$$

$$\mathbb{E} \left[\frac{I\{\mathbf{X} \in U_\varepsilon\}}{\text{Rem}} \left\| \boldsymbol{\Sigma}^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) \right\|^2 + \frac{\|\boldsymbol{\varepsilon}\|^2}{n} - \frac{\|\mathbf{r} + \text{tr}[\boldsymbol{\Sigma} \hat{\mathbf{A}}] \psi(\mathbf{r})\|^2}{n} \right] \leq \frac{C_5(\gamma, L)}{\sqrt{n}} \quad (7.3)$$

where $I\{\mathbf{X} \in U_\varepsilon\}$ is the indicator function of the event $\{\mathbf{X} \in U_\varepsilon\}$ and Rem is defined in Theorem 6.

Proposition 12 is proved in Appendix E. Consequently, if the event $\{\mathbf{X} \in U_\varepsilon\}$ has high probability and the Lipschitz condition (7.1) holds in this event, the main results Theorems 6 and 9 still hold, with no strong convexity assumption on the penalty. The proof relies on an application of Kirszbraun’s theorem already presented in (Bellec, 2020).

The Lipschitz condition (7.1) and inequality $\|\Sigma^{1/2}\widehat{\mathbf{A}}\Sigma^{1/2}\|_{op} \leq L/n$ have been proved to hold in the regime $n \asymp p$ for covariance Σ such that $\Sigma_{jj} = 1, \forall j \in [p]$ in the following two cases:

- The Lasso (i.e., square loss and L1 penalty) under the assumption that $\|\beta^*\|_0 \leq s_*n$ for some enough small constant s_* ;
- The Huber Lasso (i.e., Huber Loss and L1 penalty) under the assumption that $\|\beta^*\|_0 \leq s_*n$ for some small enough constant s_* , and that at least $(1 - s_*)n$ components of the noise (the “inliers”) are iid standard normal;

cf. (Bellec, 2020, Assumption 2.3 and Proposition 12.1). In both cases, the constant s_* only depends on γ , the condition number of Σ and the multiplicative constant of the tuning parameters.

Our results can thus be extended on a case-by-case basis for loss-penalty pairs such that the Lipschitz condition (7.1) holds, for instance for the two above examples.

Acknowledgments

P.C.B.’s research was partially supported by the NSF Grant DMS-1811976 and DMS-1945428 and by Ecole Des Ponts ParisTech.

References

- Mohsen Bayati and Andrea Montanari. The lasso risk for gaussian matrices. *IEEE Transactions on Information Theory*, 58(4):1997–2017, 2012.
- Mohsen Bayati, Murat A Erdogdu, and Andrea Montanari. Estimating lasso risk and noise level. In *Advances in Neural Information Processing Systems*, pages 944–952, 2013.
- Pierre C Bellec. Out-of-sample error estimate for robust m-estimators with convex penalty. *arXiv:2008.11840*, 2020. URL <https://arxiv.org/pdf/2008.11840.pdf>.
- Pierre C Bellec and Cun-Hui Zhang. Second order stein: Sure for sure and other applications in high-dimensional inference. *Annals of Statistics*, *accepted, to appear*, 2018. URL <https://arxiv.org/pdf/1811.04121.pdf>.
- Pierre C Bellec and Cun-Hui Zhang. De-biasing convex regularized estimators and interval estimation in linear models. *arXiv:1912.11943*, 2019. URL <https://arxiv.org/pdf/1912.11943.pdf>.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 2013.

- Michael Celentano and Andrea Montanari. Fundamental barriers to high-dimensional regression with convex penalties. *arXiv preprint arXiv:1903.10603*, 2019.
- Michael Celentano, Andrea Montanari, and Yuting Wei. The lasso with general gaussian designs with applications to hypothesis testing. *arXiv preprint arXiv:2007.13716*, 2020.
- Sourav Chatterjee. Fluctuations of eigenvalues and second order poincaré inequalities. *Probability Theory and Related Fields*, 143(1):1–40, 2009.
- Kenneth R Davidson and Stanislaw J Szarek. Local operator theory, random matrices and banach spaces. *Handbook of the geometry of Banach spaces*, 1(317-366):131, 2001.
- Edgar Dobriban and Stefan Wager. High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- 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.
- Noureddine El Karoui, Derek Bean, Peter J Bickel, Chinghay Lim, and Bin Yu. On robust regression with high-dimensional predictors. *Proceedings of the National Academy of Sciences*, 110(36):14557–14562, 2013.
- Vladimir A Marčenko and Leonid Andreevich Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of the USSR-Sbornik*, 1(4):457, 1967.
- 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.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.
- Iosif Pinelis. Large deviations: Growth of empirical average of iid non-negative random variables with infinite expectations? MathOverflow, 2021. URL <https://mathoverflow.net/q/390939>. URL:<https://mathoverflow.net/q/390939> (version: 2021-05-24).
- Kamiar Rahnema Rad and Arian Maleki. A scalable estimate of the out-of-sample prediction error via approximate leave-one-out cross-validation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(4):965–996, 2020.
- Kamiar Rahnema Rad, Wenda Zhou, and Arian Maleki. Error bounds in estimating the out-of-sample prediction error using leave-one-out cross validation in high-dimensions. In *International Conference on Artificial Intelligence and Statistics*, pages 4067–4077. PMLR, 2020.
- Charles M Stein. Estimation of the mean of a multivariate normal distribution. *The annals of Statistics*, pages 1135–1151, 1981.
- Mihailo Stojnic. A framework to characterize performance of lasso algorithms. *arXiv preprint arXiv:1303.7291*, 2013.

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.

Ji Xu, Arian Maleki, Kamiar Rahnama Rad, and Daniel Hsu. Consistent risk estimation in moderately high-dimensional linear regression. *IEEE Transactions on Information Theory*, 67(9): 5997–6030, 2021.

William P Ziemer. *Weakly differentiable functions: Sobolev spaces and functions of bounded variation*, volume 120. Springer-Verlag New York, 1989. doi: 10.1007/978-1-4612-1015-3.

Appendix A. Proof of the main results

Notation. For vectors in \mathbb{R}^q or \mathbb{R}^n , the Euclidean norm is $\|\cdot\|$ and $\|\cdot\|_q$ is the ℓ_q -norm for $1 \leq q \leq +\infty$. For matrices, $\|\cdot\|_{op}$ is the operator norm (largest singular value), $\|\cdot\|_F$ the Frobenius norm. We use index i only to loop or sum over $[n] = \{1, \dots, n\}$ and j only to loop or sum over $[p] = \{1, \dots, p\}$, so that $e_i \in \mathbb{R}^n$ refers to the i -th canonical basis vector in \mathbb{R}^n and $e_j \in \mathbb{R}^p$ the j -th canonical basis vector in \mathbb{R}^p . Positive absolute constants are denoted C_0, C_1, C_2, \dots , constants that depend on γ only are denoted $C_0(\gamma), C_1(\gamma), \dots$ and constant that depend on γ, μ only are denoted by $C_0(\gamma, \mu), C_1(\gamma, \mu), \dots$. If $\mathbf{f} : \mathbb{R}^q \rightarrow \mathbb{R}^n$ is differentiable at $\mathbf{z} \in \mathbb{R}^q$, we denote the Jacobian matrix in $\mathbb{R}^{n \times q}$ by $\frac{\partial \mathbf{f}}{\partial \mathbf{z}}$ or $\partial \mathbf{f} / \partial \mathbf{z}$. For an event Ω , its indicator function is denoted by I_Ω or $I\{\Omega\}$.

Organization of the proofs. Appendix B provides the proof of the main results from the main text (Theorems 4 to 9 and 11) and the overall proof strategy. Appendix C gives the proof of the probabilistic tools used in Appendix B. Appendix D proves the differentiability formulae in Theorems 1 and 3.

Additional simulations. Additional simulations and figures are given in Appendix F for Gaussian designs and in Appendix G for non-Gaussian Rademacher design. The simulations for Rademacher design suggests that our results generalize to non-Gaussian design, although it is unclear at this point how to extend the proofs to non-Gaussian \mathbf{X} .

Appendix B. Proof of the main results

We perform the following change of variable to reduce the anisotropic design regression problem to an isotropic one, $\mathbf{G} = \mathbf{X}\Sigma^{-1/2} \in \mathbb{R}^{n \times p}$ a Gaussian matrix with iid $N(0, 1)$ entries and

$$\mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{G}) = \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \rho(\varepsilon_i - \mathbf{e}_i^\top \mathbf{G} \mathbf{u}) + g(\boldsymbol{\beta}^* + \Sigma^{-1/2} \mathbf{u}) \quad (\text{B.1})$$

and denote by $(h_j)_{j=1, \dots, p}$ the components of (B.1). Then $\Sigma^{1/2}(\widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) - \boldsymbol{\beta}^*) = \mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{X})$ with $\widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X})$ the M -estimator in (1.1). With $\mathbf{y} = \mathbf{G}\Sigma^{1/2}\boldsymbol{\beta}^* + \boldsymbol{\varepsilon}$, by the chain rule and (2.1),

$$\begin{aligned} & \Sigma^{-1/2}(\partial / \partial g_{ij}) \mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{G}) \\ &= (\partial / \partial g_{ij}) \widehat{\boldsymbol{\beta}}(\mathbf{G}\Sigma^{1/2}\boldsymbol{\beta}^* + \boldsymbol{\varepsilon}, \mathbf{G}\Sigma^{1/2}) \\ &= \widehat{\mathbf{A}}\mathbf{X}^\top \mathbf{e}_i \psi'(r_i)(\Sigma^{1/2}\boldsymbol{\beta}^*) \mathbf{e}_j + \widehat{\mathbf{A}}\Sigma^{1/2} \mathbf{e}_j \psi(r_i) - \widehat{\mathbf{A}}\mathbf{X}^\top \mathbf{e}_i \psi'(r_i)(\Sigma^{1/2}\widehat{\boldsymbol{\beta}}) \mathbf{e}_j \end{aligned}$$

where $\mathbf{e}_i \in \mathbb{R}^n$, $\mathbf{e}_j \in \mathbb{R}^p$ denote canonical basis vectors. Define $\boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{G}) = \boldsymbol{\psi}(\boldsymbol{\varepsilon} - \mathbf{G}\mathbf{h})$ and let

$$\mathbf{A} := \boldsymbol{\Sigma}^{1/2} \widehat{\mathbf{A}} \boldsymbol{\Sigma}^{1/2}.$$

Then we have

$$(\partial/\partial g_{ij})\mathbf{h}(\boldsymbol{\varepsilon}, \mathbf{G}) = \mathbf{A}\mathbf{e}_j\boldsymbol{\psi}(r_i) - \mathbf{A}\mathbf{G}^\top \mathbf{e}_i\boldsymbol{\psi}'(r_i)h_j \quad (\text{B.2})$$

$$(\partial/\partial g_{ij})\boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{G}) = -\text{diag}\{\boldsymbol{\psi}'(\mathbf{r})\}\mathbf{G}\mathbf{A}\mathbf{e}_j\boldsymbol{\psi}(r_i) - \mathbf{V}\mathbf{e}_i h_j \quad (\text{B.3})$$

where the second line follows by the chain rule for Lipschitz functions in in (Ziemer, 1989, Theorem 2.1.11). The crux of the argument is that the quantities of interest appearing in our results, $\|\mathbf{h}\|^2 = \|\boldsymbol{\Sigma}^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|^2$, $\|\boldsymbol{\psi}(\mathbf{r})\|^2$, $\text{tr}[\widehat{\mathbf{A}}\boldsymbol{\Sigma}] = \text{tr}[\mathbf{A}]$, $\text{tr}[\mathbf{V}]$ and $\widehat{\text{df}}$ naturally appear from tensor contractions involving the derivatives in (B.2)-(B.3). For instance, denoting $\mathbf{D} = \text{diag}\{\boldsymbol{\psi}'(\mathbf{r})\} \in \mathbb{R}^{n \times n}$ if h_j, ψ_i are the j -th and i -th component of (B.1) and $\boldsymbol{\psi}(\boldsymbol{\varepsilon}, \mathbf{G})$ and denoting $\sum_{i=1}^n \sum_{j=1}^p$ by \sum_{ij} for brevity,

$$\sum_{j=1}^p \frac{\partial h_j}{g_{ij}} = \text{tr}[\mathbf{A}]\psi_i - \mathbf{h}^\top \mathbf{A}\mathbf{G}^\top \mathbf{D}\mathbf{e}_i \quad \text{for a given } i = 1, \dots, n, \quad (\text{B.4})$$

$$\sum_{i=1}^n \frac{\partial \psi_i}{\partial g_{ij}} = -\boldsymbol{\psi}^\top \mathbf{D}\mathbf{G}\mathbf{A}\mathbf{e}_j - \text{tr}[\mathbf{V}]h_j \quad \text{for a given } j = 1, \dots, p, \quad (\text{B.5})$$

$$\sum_{ij} \frac{\partial (h_j \psi_i)}{g_{ij}} = \|\boldsymbol{\psi}\|^2 \text{tr}[\mathbf{A}] - \mathbf{h}^\top \mathbf{G}^\top \mathbf{D}\boldsymbol{\psi} - \boldsymbol{\psi}^\top \mathbf{D}\mathbf{G}\mathbf{A}\mathbf{h} - \|\mathbf{h}\|^2 \text{tr}[\mathbf{V}], \quad (\text{B.6})$$

$$\sum_{ij} \frac{\partial (h_j \mathbf{e}_i^\top \mathbf{G}\mathbf{h})}{g_{ij}} = \text{tr}[\mathbf{A}]\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h} - \mathbf{h}^\top \mathbf{A}\mathbf{G}^\top \mathbf{G}\mathbf{h} + n\|\mathbf{h}\|^2 + \boldsymbol{\psi}^\top \mathbf{G}\mathbf{A}\mathbf{h} - \|\mathbf{h}\|^2 \widehat{\text{df}}, \quad (\text{B.7})$$

$$\sum_{ij} \frac{\partial (\psi_i \mathbf{e}_j^\top \mathbf{G}^\top \boldsymbol{\psi})}{g_{ij}} = -\boldsymbol{\psi}^\top \mathbf{D}\mathbf{G}\mathbf{A}\mathbf{G}^\top \boldsymbol{\psi} - \text{tr}[\mathbf{V}]\boldsymbol{\psi}^\top \mathbf{G}\mathbf{h} - \mathbf{h}^\top \mathbf{G}^\top \mathbf{V}\boldsymbol{\psi} + (p - \widehat{\text{df}})\|\boldsymbol{\psi}\|^2 \quad (\text{B.8})$$

where we used that $\widehat{\text{df}} = \sum_{i=1}^n \mathbf{e}_i^\top \mathbf{G}\mathbf{A}\mathbf{G}^\top \mathbf{D}\mathbf{e}_i = \text{tr}[\mathbf{G}\mathbf{A}\mathbf{G}^\top \mathbf{D}]$ in the fourth line and $\widehat{\text{df}} = \sum_{j=1}^p \mathbf{e}_j^\top \mathbf{G}^\top \mathbf{D}\mathbf{G}\mathbf{A}\mathbf{e}_j = \text{tr}[\mathbf{G}^\top \mathbf{D}\mathbf{G}\mathbf{A}]$ in the fifth thanks to the commutation property of the trace. The terms in colored purple indicate terms that will be proved to be negligible later on. The probabilistic tool that leads to asymptotic normality of the residuals is the following.

Proposition 13 [Variant of Bellec and Zhang (2019)] *Let $\mathbf{z} \in N(\mathbf{0}, \mathbf{I}_q)$ and $\mathbf{f} := \mathbf{f}(\mathbf{z}) : \mathbb{R}^q \rightarrow \mathbb{R}^q \setminus \{\mathbf{0}\}$ be locally Lipschitz in \mathbf{z} with $\mathbb{E}[\|\mathbf{f}\|^{-2} \sum_{k=1}^q \|\frac{\partial \mathbf{f}}{\partial z_k}\|^2] < +\infty$. Then*

$$\mathbb{E}\left[\left(\frac{\mathbf{f}^\top \mathbf{z} - \sum_{k=1}^q (\partial/\partial z_k) f_k}{\|\mathbf{f}\|_2} - Z\right)^2\right] \leq (7 + 2\sqrt{6})\mathbb{E}\left[\|\mathbf{f}\|^{-2} \sum_{k=1}^q \|\frac{\partial \mathbf{f}}{\partial z_k}\|^2\right] < +\infty. \quad (\text{B.9})$$

Proposition 13 is proved in Appendix C. From here, asymptotic normality of the residuals in the square loss case is readily obtained using the explicit formulae for the derivatives and the contraction (B.4). We start with the square loss and the proof of Theorem 5.

Proof (Proof of Theorem 5) Apply Proposition 13 with $q = p + 1$ and $\mathbf{z} = (\mathbf{g}_i, \varepsilon_i/\sigma) \sim N(\mathbf{0}, \mathbf{I}_{p+1})$ conditionally on $(\mathbf{g}_l, \varepsilon_l)_{l \in [n] \setminus \{i\}}$, and with $\mathbf{f} = (\mathbf{h}, -\sigma) \in \mathbb{R}^{p+1}$. Note that the last

component of \mathbf{f} is constant and $\|\mathbf{f}\|^2 = \|\mathbf{h}\|^2 + \sigma^2$. By (B.4) and $\mathbf{D} = \mathbf{I}_n$ for the square loss, $\text{tr}[\partial\mathbf{f}/\partial\mathbf{z}] = \text{tr}[\mathbf{A}]\psi_i - \mathbf{h}^\top \mathbf{A} \mathbf{G}^\top \mathbf{e}_i$ and by symmetry in $i = 1, \dots, n$, $\mathbb{E}[\|\mathbf{h}^\top \mathbf{A} \mathbf{G}^\top \mathbf{e}_i\|^2 / \|\mathbf{f}\|^2] = \frac{1}{n} \sum_{l=1}^n \mathbb{E}[\|\mathbf{h}^\top \mathbf{A} \mathbf{G}^\top \mathbf{e}_l\|^2 / \|\mathbf{f}\|^2] = \frac{1}{n} \|\mathbf{G} \mathbf{A}^\top \mathbf{h}\|^2 / \|\mathbf{f}\|^2 \leq \frac{1}{n} \mathbb{E}[\|\mathbf{G}\|_{op}^2 \|\mathbf{A}\|_{op}^2] \leq n^{-2} C_6(\gamma, \mu)$ thanks to $\|\mathbf{A}\|_{op} \leq 1/(n\mu)$ and $\mathbb{E}[\|\mathbf{G}\|_{op}^2] \leq C_7(\gamma)n$. Similarly, for the square loss $r_i = \psi_i = \varepsilon_i - \mathbf{g}_i^\top \mathbf{h}$ and

$$\begin{aligned} \|\mathbf{f}\|^{-1} \|\partial\mathbf{f}/\partial\mathbf{z}\|_F &= (\|\mathbf{h}\|^2 + \sigma^2)^{-1/2} \|\mathbf{A}\psi_i - \mathbf{A} \mathbf{G}^\top \mathbf{e}_i \mathbf{h}^\top\|_F \\ &\leq \|\mathbf{A}\|_{op} [\sqrt{p}|\varepsilon_i|/\sigma + \sqrt{p}\|\mathbf{h}\|^{-1} |\mathbf{g}_i^\top \mathbf{h}|] + \|\mathbf{G}\|_{op}. \end{aligned}$$

By the triangle inequality, $\|\mathbf{A}\|_{op} \leq 1/(n\mu)$ and $p \leq \gamma n$,

$$\mathbb{E}[\|\mathbf{f}\|^{-2} \|\partial\mathbf{f}/\partial\mathbf{z}\|_F^2]^{1/2} \leq \frac{\sqrt{p}}{n\mu} (\mathbb{E}[\varepsilon_i^2/\sigma^2]^{1/2} + \mathbb{E}[(\mathbf{g}_i^\top \mathbf{h})^2 / \|\mathbf{h}\|^2]^{1/2}) + \frac{1}{n\mu} \mathbb{E}[\|\mathbf{G}\|_{op}^2]^{1/2}.$$

By symmetry in $i = 1, \dots, n$, $\mathbb{E}[(\mathbf{g}_i^\top \mathbf{h})^2 / \|\mathbf{h}\|^2] = \frac{1}{n} \sum_{l=1}^n \mathbb{E}[(\mathbf{g}_l^\top \mathbf{h})^2 / \|\mathbf{h}\|^2] \leq \frac{1}{n} \mathbb{E}[\|\mathbf{G}\|_{op}^2]$. Since $\frac{1}{n} \mathbb{E}[\|\mathbf{G}\|_{op}^2] \leq C_8(\gamma)$, the right-hand side in the previous display is bounded from above by $C_9(\gamma, \mu)n^{-1/2}$. Since $\mathbf{f}^\top \mathbf{z} = -r_i$ we obtain $-r_i - \text{tr}[\mathbf{A}]r_i = (\|\mathbf{h}\|^2 + \sigma^2)^{1/2}(Z + O_P(n^{-1/2}))$ which completes the proof of (3.4). \blacksquare

Proof (Proof of Theorem 4) Let $U \sim N(0, 1)$ be independent of everything else. We apply the previous proposition with $\mathbf{z} = (\mathbf{g}_i, U) \sim N(\mathbf{0}, \mathbf{I}_{p+1})$ conditionally on $(\varepsilon, \mathbf{g}_l, l \in [n] \setminus \{i\})$ to $\mathbf{f} = (\mathbf{h}, n^{-1/4}\psi(\varepsilon_i))$. Note that the last component of \mathbf{f} is constant. By (B.4), $\text{tr}[\partial\mathbf{f}/\partial\mathbf{z}] = \text{tr}[\mathbf{A}]\psi_i - \mathbf{h}^\top \mathbf{A} \mathbf{G}^\top \mathbf{D} \mathbf{e}_i$ and by (B.2),

$$\|\mathbf{f}\|^{-1} \|\partial\mathbf{f}/\partial\mathbf{z}\|_F = (\|\mathbf{h}\|^2 + n^{-1/2}\psi(\varepsilon_i)^2)^{-1/2} \|\mathbf{A}\psi_i - \mathbf{A} \mathbf{G}^\top \mathbf{D} \mathbf{e}_i \mathbf{h}^\top\|_F \quad (\text{B.10})$$

$$\leq \|\mathbf{A}\|_{op} [n^{1/4}\sqrt{p} + \sqrt{p}\|\mathbf{h}\|^{-1} |\mathbf{g}_i^\top \mathbf{h}|] + \|\mathbf{G}\|_{op} \quad (\text{B.11})$$

where we used $\|\mathbf{A}\|_F \leq \sqrt{p}\|\mathbf{A}\|_{op}$ and $|\psi_i| \leq \psi(\varepsilon_i) + |\mathbf{g}_i^\top \mathbf{h}|$ thanks to ψ being 1-Lipschitz. We have $\|\mathbf{A}\|_{op} \leq 1/(n\mu)$ and $\mathbb{E}[\|\mathbf{h}\|^{-2} |\mathbf{g}_i^\top \mathbf{h}|^2] = \frac{1}{n} \sum_{l=1}^n \mathbb{E}[\|\mathbf{h}\|^{-2} |\mathbf{g}_l^\top \mathbf{h}|^2] = \frac{1}{n} \mathbb{E}[\|\mathbf{h}\|^{-2} \|\mathbf{G} \mathbf{h}\|^2] \leq \frac{1}{n} \mathbb{E}[\|\mathbf{G}\|_{op}^2]$ by symmetry in $i = 1, \dots, n$, so that $\mathbb{E}[\|\mathbf{f}\|^{-2} \|\partial\mathbf{f}/\partial\mathbf{z}\|_F^2] \leq n^{-1/2} C_{10}(\gamma, \mu)$. Thus by Proposition 13,

$$\begin{aligned} (-r_i - \text{tr}[\mathbf{A}]\psi_i) + (\varepsilon_i - \|\mathbf{h}\|Z) &= \mathbf{g}_i^\top \mathbf{h} - \text{tr}[\mathbf{A}]\psi_i - \|\mathbf{h}\|Z \\ &= -Un^{-1/4}\psi(\varepsilon_i) + [\|\mathbf{f}\| - \|\mathbf{h}\|]Z + \|\mathbf{f}\| \text{Rem} - \mathbf{h}^\top \mathbf{A} \mathbf{G}^\top \mathbf{D} \mathbf{e}_i \end{aligned}$$

where $\mathbb{E}[\text{Rem}^2] \leq C_{11} \mathbb{E}[\|\mathbf{f}\|^{-2} \|\partial\mathbf{f}/\partial\mathbf{z}\|_F^2] \leq n^{-1/2} C_{12}(\gamma, \mu)$. By properties of the operator norm and symmetry in $i = 1, \dots, n$,

$$\mathbb{E}[\|\mathbf{h}\|^{-2} |\mathbf{h}^\top \mathbf{A} \mathbf{G}^\top \mathbf{D} \mathbf{e}_i|^2] = \frac{1}{n} \mathbb{E}[\|\mathbf{h}\|^{-2} \|\mathbf{D} \mathbf{G} \mathbf{A}^\top \mathbf{h}\|^2] \leq \frac{1}{n} \mathbb{E}[\|\mathbf{G}\|_{op}^2 \|\mathbf{A}\|_{op}^2] \leq \frac{C_{13}(\gamma, \mu)}{n^2}. \quad (\text{B.12})$$

By the triangle inequality, $|\|\mathbf{f}\| - \|\mathbf{h}\|| \leq n^{-1/4} |\psi(\varepsilon_i)|$ so that the right-hand side is of the form $O_P(n^{-1/4})(|\psi(\varepsilon_i)| + \|\mathbf{h}\|)$ as desired. The previous display can be rewritten as $r_i + \text{tr}[\mathbf{A}]\psi_i = \tilde{\varepsilon}_i^n + \|\mathbf{h}\| \tilde{Z}_i^n$ for

$$\tilde{\varepsilon}_i^n = \varepsilon_i + Un^{-1/4}\psi(\varepsilon_i) - [\|\mathbf{f}\| - \|\mathbf{h}\|](Z + \text{Rem}), \quad \tilde{Z}_i^n = -Z - \text{Rem} + \|\mathbf{h}\|^{-1} \mathbf{h}^\top \mathbf{A} \mathbf{G}^\top \mathbf{D} \mathbf{e}_i.$$

If ε_i has a fixed distribution F , then $|\psi(\varepsilon_i)| \leq |\psi(0)| + |\varepsilon_i| = |\varepsilon_i| = O_P(1)$ thanks to $\psi(0) = 0$ and ψ being 1-Lipschitz so that $(\tilde{\varepsilon}_i^n, \tilde{Z}_i^n) = (\varepsilon_i, -Z) + O_P(n^{-1/4})$. Since $(\varepsilon_i, -Z)$ are independent, by Slutsky's theorem this proves that $(\tilde{\varepsilon}_i^n, \tilde{Z}_i^n)$ converges weakly to the product measure $F \otimes N(0, 1)$. \blacksquare

Proposition 14 *Let $\mathbf{h} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^p$, $\boldsymbol{\psi} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ be locally Lipschitz functions. If $\mathbf{G} \in \mathbb{R}^{n \times p}$ has iid $N(0, 1)$ entries then*

$$\begin{aligned} & \mathbb{E} \left[\left(\frac{\boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - \sum_{ij} \frac{\partial(\psi_i h_j)}{g_{ij}}}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right)^2 + \left(\frac{\|\mathbf{G} \mathbf{h}\|^2 - \sum_{ij} \frac{\partial(h_j \mathbf{e}_i^\top \mathbf{G} \mathbf{h})}{g_{ij}}}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right)^2 + \left(\frac{\|\mathbf{G}^\top \boldsymbol{\psi}\|^2 - \sum_{ij} \frac{\partial(\psi_i \mathbf{e}_j^\top \mathbf{G}^\top \boldsymbol{\psi})}{g_{ij}}}{n\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2} \right)^2 \right] \\ & \leq C_{14} \mathbb{E} \left[n + p + \|\mathbf{G}\|_{op}^2 + (n+p) \sum_{i=1}^n \sum_{j=1}^p \frac{1 + \|\mathbf{G}\|_{op}^2/n}{(\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^2} \left(\left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \right) \right] \end{aligned} \quad (\text{B.13})$$

for some positive absolute constant in the second line.

Proposition 14 is proved in Appendix C; it is a consequence of (Bellec, 2020, Proposition 6.3). By Proposition 14 combined with the identities (B.6)-(B.7)-(B.8), and by showing that the purple-colored terms in (B.6)-(B.7)-(B.8) are negligible, we obtain the following.

Proposition 15 *Let Assumption A be fulfilled. Then*

$$\mathbb{E} \left[\left\{ n^{-\frac{1}{2}} (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{-1} (\boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - \text{tr}[\mathbf{A}] \|\boldsymbol{\psi}\|^2 + \text{tr}[\mathbf{V}] \|\mathbf{h}\|^2) \right\}^2 \right] \leq C_{15}(\gamma, \mu), \quad (\text{B.14})$$

$$\mathbb{E} \left[\left\{ n^{-\frac{1}{2}} (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{-1} \left(\frac{1}{n} \|\mathbf{G}^\top \boldsymbol{\psi}\|^2 - \frac{p-\text{df}}{n} \|\boldsymbol{\psi}\|^2 + \frac{\text{tr}[\mathbf{V}]}{n} \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} \right) \right\}^2 \right] \leq C_{16}(\gamma, \mu), \quad (\text{B.15})$$

$$\mathbb{E} \left[\left\{ n^{-\frac{1}{2}} (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{-1} (\|\mathbf{G} \mathbf{h}\|^2 - \text{tr}[\mathbf{A}] \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - (n - \hat{\text{d}}\text{f}) \|\mathbf{h}\|^2) \right\}^2 \right] \leq C_{17}(\gamma, \mu). \quad (\text{B.16})$$

Proof We bound from above the derivatives in (B.13). For the norm of $(\partial/\partial g_{ij})\mathbf{h}$ and $(\partial/\partial g_{ij})\boldsymbol{\psi}$, by (B.3)-(B.2) and $\frac{1}{2}(a+b)^2 \leq a^2 + b^2$,

$$\sum_{ij} \frac{1}{2} \left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 \leq \|\mathbf{A}\|_F^2 \|\boldsymbol{\psi}\|^2 + \|\mathbf{A} \mathbf{G}^\top \mathbf{D}\|_F^2 \|\mathbf{h}\|^2, \quad \sum_{ij} \frac{1}{2n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \leq \frac{\|\mathbf{D} \mathbf{G} \mathbf{A}\|_F^2 \|\boldsymbol{\psi}\|^2 + \|\mathbf{V}\|_F^2 \|\mathbf{h}\|^2}{n}.$$

Using $\|\mathbf{A}\|_{op} \leq 1/(n\mu)$, $\|\mathbf{D}\|_{op} \leq 1$, $p/n \leq \gamma$ and \mathbf{V} in (2.3), it follows that in (B.13) we have

$$\frac{1}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \sum_{i=1}^n \sum_{j=1}^p \left(\left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \right) \leq C_{18}(\gamma, \mu) \left(1 + \frac{\|\mathbf{G}\|_{op}^2}{n} \right). \quad (\text{B.17})$$

Since $\mathbb{E}[\|n^{-1/2} \mathbf{G}\|_{op}^4] \leq C_{19}(\gamma)$ (Davidson and Szarek, 2001, Theorem II.13), this shows that (B.13) is bounded from above by $C_{20}(\gamma, \mu)n$. The contractions appearing in the left-hand side of (B.13) are given in (B.6)-(B.7)-(B.8), so that it remains to bound from above the purple colored terms in these three equations. This is done by using the upper bounds on the operator norms $\|\mathbf{A}\|_{op} \leq 1/(n\mu)$, $\|\mathbf{D}\|_{op} \leq 1$ and again that $\mathbb{E}[\|n^{-1/2} \mathbf{G}\|_{op}^4] \leq C_{21}(\gamma)$, so that (B.13) yields the three inequalities in Proposition 15. \blacksquare

The next result is another probabilistic result where the contractions in (B.4)-(B.5) appear.

Proposition 16 *Let $\mathbf{h} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^p$, $\boldsymbol{\psi} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ be locally Lipschitz functions. If $\mathbf{G} \in \mathbb{R}^{n \times p}$ has iid $N(0, 1)$ entries then*

$$\begin{aligned} & \mathbb{E} \left[\left| \frac{\frac{p}{n} \|\boldsymbol{\psi}\|^2 - \frac{1}{n} \sum_{j=1}^p (\boldsymbol{\psi}^\top \mathbf{G} \mathbf{e}_j - \sum_{i=1}^n \frac{\partial \psi_i}{\partial g_{ij}})^2}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right| \right] + \mathbb{E} \left[\left| \frac{n \|\mathbf{h}\|^2 - \sum_{i=1}^n (\mathbf{g}_i^\top \mathbf{h} - \sum_{j=1}^p \frac{\partial h_j}{\partial g_{ij}})^2}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right| \right] \\ & \leq C_{22} \left(\sqrt{n+p} (1 + \Xi^{1/2}) + \Xi \right) \text{ where } \Xi = \mathbb{E} \left[\frac{1}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \sum_{i=1}^n \sum_{j=1}^p \left(\left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \right) \right]. \end{aligned}$$

Proposition 16 is proved in Appendix C; it is a consequence of (Bellec, 2020, Theorem 7.1). Using the contractions (B.4)-(B.5) in the left-hand side of Proposition 16, and by showing that the purple colored terms are negligible, we obtain the following two inequalities.

Proposition 17 *Let Assumption A be fulfilled. Then*

$$\mathbb{E} \left| n^{-\frac{1}{2}} (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{-1} \left(\frac{p}{n} \|\boldsymbol{\psi}\|^2 - \frac{1}{n} \|\mathbf{G}^\top \boldsymbol{\psi} + \text{tr}[\mathbf{V}]\mathbf{h}\|^2 \right) \right| \leq C_{23}(\gamma, \mu), \quad (\text{B.18})$$

$$\mathbb{E} \left| n^{-\frac{1}{2}} (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{-1} (n \|\mathbf{h}\|^2 - \|\mathbf{G}\mathbf{h} - \text{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2) \right| \leq C_{24}(\gamma, \mu). \quad (\text{B.19})$$

Proof For Ξ in Proposition 16, the fact that $\Xi \leq C_{25}(\gamma, \mu)$ is already proved in (B.17). For the first inequality we use Proposition 16 and the contraction (B.5). To control the purple terms in (B.5) inside the left-hand side of Proposition 17,

$$\begin{aligned} & \left| \sum_{j=1}^p (\boldsymbol{\psi}^\top \mathbf{G} \mathbf{e}_j - \sum_{i=1}^n \frac{\partial \psi_i}{\partial g_{ij}})^2 - \|\mathbf{G}^\top \boldsymbol{\psi} + \text{tr}[\mathbf{V}]\mathbf{h}\|^2 \right| = \left| \boldsymbol{\psi}^\top \mathbf{D} \mathbf{G} \mathbf{A} \left(2\mathbf{G}^\top \boldsymbol{\psi} + 2\text{tr}[\mathbf{V}]\mathbf{h} + \mathbf{A}^\top \mathbf{G}^\top \mathbf{D} \boldsymbol{\psi} \right) \right| \\ & \leq (\|\boldsymbol{\psi}\|^2/n + \|\mathbf{h}\|^2) (2n \|\mathbf{G}\|_{op}^2 \|\mathbf{A}\|_{op} + 2\sqrt{n} \|\mathbf{G}\|_{op} \|\mathbf{A}\|_{op} + n \|\mathbf{A}\|_{op}^2 \|\mathbf{G}\|_{op}^2) \end{aligned}$$

thanks to $|\text{tr} \mathbf{V}| \leq n$ in Theorem 1. With the bound obtained by multiplying the previous display by $n^{-3/2} (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{-1}$, and using the previous bounds on $\|\mathbf{A}\|_{op}$ and $\mathbb{E}[\|n^{-1/2} \mathbf{G}\|_{op}^2]$, we obtain (B.18) from Proposition 16 and (B.5). The second claim is obtained by Proposition 16, the contraction (B.4) and an argument similar to the previous display bound the purple term in (B.4). ■

We are now ready to prove Theorem 9.

Proof (Proof of Theorem 9) Define

$$\begin{aligned} \xi_I &= \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - \text{tr}[\mathbf{A}] \|\boldsymbol{\psi}\|^2 + \text{tr}[\mathbf{V}] \|\mathbf{h}\|^2 && \text{(bounded in (B.14))}, \\ \xi_{II} &= \frac{1}{n} \|\mathbf{G}^\top \boldsymbol{\psi}\|^2 - \frac{p-\hat{\text{d}}\text{f}}{n} \|\boldsymbol{\psi}\|^2 + \frac{\text{tr}[\mathbf{V}]}{n} \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} && \text{(bounded in (B.15))}, \\ \xi_{III} &= \|\mathbf{G}\mathbf{h}\|^2 - \text{tr}[\mathbf{A}] \boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - (n - \hat{\text{d}}\text{f}) \|\mathbf{h}\|^2 && \text{(bounded in (B.16))}, \\ \xi_{IV} &= \frac{p}{n} \|\boldsymbol{\psi}\|^2 - \frac{1}{n} \|\mathbf{G}^\top \boldsymbol{\psi} + \text{tr}[\mathbf{V}]\mathbf{h}\|^2 && \text{(bounded in (B.18))}, \\ \xi_V &= n \|\mathbf{h}\|^2 - \|\mathbf{G}\mathbf{h} - \text{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2 && \text{(bounded in (B.19))}. \end{aligned}$$

Then by expanding the square in ξ_{IV} and ξ_V and simple algebra (for instance by computing first $\xi_{II} + \xi_{IV}$ and $\xi_{III} + \xi_V$ separately),

$$(\text{tr}[\mathbf{V}]/n - \text{tr} \mathbf{A}) \xi_I + \xi_{II} + \xi_{III} + \xi_{IV} + \xi_V = (\|\boldsymbol{\psi}\|^2/n + \|\mathbf{h}\|^2) (\hat{\text{d}}\text{f} - \text{tr}[\mathbf{A}] \text{tr}[\mathbf{V}]).$$

Since $|\operatorname{tr}[\mathbf{V}]/n \leq 1$, $\operatorname{tr}[\mathbf{A}] \leq \gamma/\mu$ by Theorem 1, the previous display divided by $n^{1/2}(\|\boldsymbol{\psi}\|^2/n + \|\mathbf{h}\|^2)$ and the bounds (B.14), (B.15), (B.16), (B.18) and (B.19) complete the proof. \blacksquare

To prove Theorem 6, we need this extra proposition whose proof is closely related to Proposition 15.

Proposition 18 *Let Assumption A be fulfilled. Then*

$$\mathbb{E}[\{(\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{-\frac{1}{2}}\|\boldsymbol{\varepsilon}\|^{-1}\xi_{VI}\}^2] \leq C_{26}(\gamma, \mu) \quad \text{for} \quad \xi_{VI} = \boldsymbol{\varepsilon}^\top(\mathbf{G}\mathbf{h} - \operatorname{tr}[\mathbf{A}]\boldsymbol{\psi}). \quad (\text{B.20})$$

Proposition 18 is proved in Appendix C. We are now ready to prove Theorem 6.

Proof (Proof of Theorem 6) We have $n\|\mathbf{h}\|^2 + \|\boldsymbol{\varepsilon}\|^2 - \|\mathbf{r} + \operatorname{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2 = \xi_V + 2\xi_{VI}$ by simple algebra and the definitions of ξ_V and ξ_{VI} . Hence

$$\mathbb{E}\left[\frac{\|\mathbf{h}\|^2 + \|\boldsymbol{\varepsilon}\|^2/n - \|\mathbf{r} + \operatorname{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2/n}{\max\{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n, (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{1/2}(\|\boldsymbol{\varepsilon}\|^2/n)^{1/2}\}}\right] \leq n^{-1/2}C_{27}(\gamma, \mu) \quad (\text{B.21})$$

thanks to (B.20) and (B.19). \blacksquare

Proof (Proof of Theorem 7) We perform the change of variable (B.1) to $\tilde{\boldsymbol{\beta}}$ as well, giving $\tilde{\mathbf{h}}$ (the counterpart of \mathbf{h}), $\tilde{\boldsymbol{\psi}}$ (counterpart of $\boldsymbol{\psi}$) and $\tilde{\mathbf{A}}$ (counterpart of \mathbf{A}). Let Ω be the event defined in the theorem, i.e.,

$$\Omega = \{\|\mathbf{G}\|_{op} \leq 2\sqrt{n} + \sqrt{p}\} \cap \{\|\boldsymbol{\varepsilon}\|^2 \leq n^{2/(1+q)}\}. \quad (\text{B.22})$$

Then $\mathbb{P}(\Omega^c) \rightarrow 0$ by (Davidson and Szarek, 2001, Theorem 2.13) for the first event and Pinelis (2021) to show that $\|\boldsymbol{\varepsilon}\|^2/n^{2/(1+q)} \xrightarrow{\mathbb{P}} 0$ under the assumption that $\mathbb{E}[\|\boldsymbol{\varepsilon}_i\|^{1+q}]$ is bounded.

Under Assumption B, $I_\Omega(\|\boldsymbol{\psi}\|^2/n + \|\mathbf{h}\|^2)$ is bounded by a constant. Indeed, since the penalty g is minimized at $\mathbf{0}$, $(\hat{\boldsymbol{\beta}} - \mathbf{0})^\top \mathbf{X}^\top \boldsymbol{\psi} \in n(\hat{\boldsymbol{\beta}} - \mathbf{0})^\top (\partial g(\hat{\boldsymbol{\beta}}) - \partial g(\mathbf{0}))$ since $\mathbf{0} \in \partial g(\mathbf{0})$. By strong convexity of g in Assumption A, $(\hat{\boldsymbol{\beta}} - \mathbf{0})^\top \mathbf{X}^\top \boldsymbol{\psi} \geq \mu\|\boldsymbol{\Sigma}^{1/2}\hat{\boldsymbol{\beta}}\|^2$. In Ω , this implies $\|\boldsymbol{\Sigma}^{1/2}\hat{\boldsymbol{\beta}}\| \leq \frac{1}{\mu n}\|\mathbf{G}\|_{op}\|\boldsymbol{\psi}\| \leq C_{28}(\gamma, \mu)\|\boldsymbol{\psi}\|/\sqrt{n}$ and $\|\boldsymbol{\psi}\|/\sqrt{n} \leq M$ in Assumption B. Since $\|\boldsymbol{\Sigma}^{1/2}\boldsymbol{\beta}^*\|^2 \leq M$ in Assumption B, this yields $I_\Omega(\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n) \leq C_{29}(\gamma, \mu, M)$ and the same holds for $\tilde{\mathbf{h}}, \tilde{\boldsymbol{\psi}}$: $I_\Omega(\|\tilde{\mathbf{h}}\|^2 + \|\tilde{\boldsymbol{\psi}}\|^2/n) \leq C_{30}(\gamma, \mu, M)$.

Inequality (B.21) thus implies

$$\begin{aligned} \mathbb{E}[I_\Omega(\|\mathbf{h}\|^2 + \|\boldsymbol{\varepsilon}\|^2/n - \|\mathbf{r} + \operatorname{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2/n + \|\tilde{\mathbf{h}}\|^2 + \|\boldsymbol{\varepsilon}\|^2/n - \|\tilde{\mathbf{r}} + \operatorname{tr}[\tilde{\mathbf{A}}]\tilde{\boldsymbol{\psi}}\|^2/n)] \\ \leq C_{31}(\gamma, \mu, M)(n^{-1/2} \vee n^{-q/(1+q)}). \end{aligned}$$

Since $q \in (0, 1)$ we have $n^{-1/2} \vee n^{-q/(1+q)} = n^{-q/(1+q)}$ in the right-hand side. Let $\hat{\Omega} = \{\|\mathbf{h}\|^2 - \|\tilde{\mathbf{h}}\|^2 > \eta, \|\mathbf{r} + \operatorname{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2 \leq \|\tilde{\mathbf{r}} + \operatorname{tr}[\tilde{\mathbf{A}}]\tilde{\boldsymbol{\psi}}\|^2\}$ be the event for which we are trying to control the probability. By the triangle inequality,

$$\mathbb{E}[I_\Omega(\|\mathbf{h}\|^2 - \|\tilde{\mathbf{h}}\|^2 - \|\mathbf{r} + \operatorname{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2/n + \|\tilde{\mathbf{r}} + \operatorname{tr}[\tilde{\mathbf{A}}]\tilde{\boldsymbol{\psi}}\|^2/n)] \leq C_{32}(\gamma, \mu, M)n^{-q/(1+q)}.$$

In $\hat{\Omega}$, the random variable in the expectation sign is larger than ηI_Ω . Thus $\eta\mathbb{E}[I_\Omega I_{\hat{\Omega}}] \leq C_{33}(\gamma, \mu, M)n^{-q/(1+q)}$ and $\mathbb{P}(\hat{\Omega}) \leq \eta^{-1}C_{34}(\gamma, \mu, M)n^{-q/(1+q)} + \mathbb{P}(\Omega^c)$. \blacksquare

Proof (Proof of Theorem 8) We follow the same strategy. Let Ω be the same event as in the previous proof, so that $\mathbb{P}(\Omega^c) \rightarrow 0$ as before. We perform the change of variable (B.1) for each

$k = 1, \dots, K$ giving $\mathbf{h}_k, \boldsymbol{\psi}_k$ and \mathbf{A}_k . We have $I_\Omega \max_{k=1, \dots, K} (\|\mathbf{h}_k\|^2 + \|\boldsymbol{\psi}_k\|^2/n) \leq C_{35}(\gamma, \mu, M)$ as explained in the previous proof.

Summing over k the inequality (B.21) gives $\mathbb{E}[I_\Omega \sum_{k=1}^K \|\mathbf{h}_k\|^2 + \|\boldsymbol{\varepsilon}\|^2 - \|\mathbf{r}_k + \text{tr}[\mathbf{A}_k]\boldsymbol{\psi}_k\|^2] \leq KC_{36}(\gamma, \mu, M)n^{-q/(1+q)}$. Let \hat{k} be the minimizer of $\|\mathbf{r}_k + \text{tr}[\mathbf{A}_k]\boldsymbol{\psi}_k\|^2$ as defined in the statement of Theorem 8 and let $\tilde{k} \in \{1, \dots, K\}$ be such that $\|\mathbf{h}_{\hat{k}}\|^2 \geq \|\mathbf{h}_{\tilde{k}}\|^2 + \eta$ in the event $\tilde{\Omega}$ where such \tilde{k} exists. Then by the triangle inequality, $\eta \mathbb{E}[I_\Omega I_{\tilde{\Omega}}] \leq C_{37}(\gamma, \mu, M)n^{-q/(1+q)}$. It follows that $\mathbb{P}(\tilde{\Omega}) \leq \eta^{-1}C_{38}(\gamma, \mu, M)n^{-q/(1+q)} + \mathbb{P}(\Omega^c) \rightarrow 0$ as desired. \blacksquare

Proof (Proof of Theorem 11) Using $\|\mathbf{a}\|^2 - \|\mathbf{b}\|^2 = (\mathbf{a} - \mathbf{b})^\top (\mathbf{a} + \mathbf{b})$ we have

$$\|\mathbf{G}\mathbf{h} - \text{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2 - \|\mathbf{G}\mathbf{h} - (\hat{\text{d}}\mathbf{f}/\text{tr}[\mathbf{V}])\|^2 = (\hat{\text{d}}\mathbf{f}/\text{tr}[\mathbf{V}] - \text{tr}[\mathbf{A}])\boldsymbol{\psi}^\top (2\mathbf{G}\mathbf{h} - (\text{tr}[\mathbf{A}] + \hat{\text{d}}\mathbf{f}/\text{tr}[\mathbf{V}])\boldsymbol{\psi}).$$

Hence using $|\text{tr}[\mathbf{A}]| \leq \gamma/\mu, |\hat{\text{d}}\mathbf{f}| \leq n$ and the Cauchy-Schwarz inequality

$$\begin{aligned} & \left| \|\mathbf{G}\mathbf{h} - \text{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2 - \|\mathbf{G}\mathbf{h} - (\hat{\text{d}}\mathbf{f}/\text{tr}[\mathbf{V}])\|^2 \right| \\ & \leq C_{39}(\gamma, \mu) \left(\frac{n}{\text{tr}[\mathbf{V}]} \vee 1 \right) |\hat{\text{d}}\mathbf{f}/n - \text{tr}[\mathbf{V}] \text{tr}[\mathbf{A}]/n| (\|\boldsymbol{\psi}\|^2 + \|\mathbf{G}\|_{op} \|\mathbf{h}\|^2). \end{aligned}$$

Let Ω be the event in Theorem 7. Using the bound on the operator norm of \mathbf{G} in Ω , for any deterministic $\eta > 0$ we have proved

$$\mathbb{E} \left[I\{\Omega\} I\{\text{tr}[\mathbf{V}]n \geq \eta\} \frac{\left| \|\mathbf{G}\mathbf{h} - \text{tr}[\mathbf{A}]\boldsymbol{\psi}\|^2 - \|\mathbf{G}\mathbf{h} - (\hat{\text{d}}\mathbf{f}/\text{tr}[\mathbf{V}])\|^2 \right|}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right] \leq \frac{C_{40}(\gamma, \mu)}{\eta \wedge 1} n^{1/2}$$

thanks to Theorem 9. By (D.8), in the event Ω where the operator norm of $\|n^{-1/2}\mathbf{G}\|_{op}$ is bounded by a constant, $\text{tr}[\mathbf{V}] \geq \text{tr}[\text{diag}\{\boldsymbol{\psi}'(\mathbf{r})\}]/C_{41}(\gamma, \mu)$. Hence combining the previous display with (B.21), we have proved

$$\mathbb{E} \left[\frac{I\{\Omega\} I\{\sum_{i=1}^n \boldsymbol{\psi}'(r_i) \geq n\eta\} \left| \|\mathbf{h}\|^2 + \|\boldsymbol{\varepsilon}\|^2/n - \|\mathbf{r} + \frac{\hat{\text{d}}\mathbf{f}}{\text{tr}[\mathbf{V}]} \boldsymbol{\psi}\|^2/n \right|}{\max\{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n, (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{1/2} (\|\boldsymbol{\varepsilon}\|^2/n)^{1/2}\}} \right] \leq \frac{C_{42}(\gamma, \mu, \eta)}{\sqrt{n}}.$$

At this point the proof is similar to that of Theorem 8: We perform the change of variable (B.1) for each $k = 1, \dots, K$ giving $\mathbf{h}_k, \boldsymbol{\psi}_k, \hat{\text{d}}\mathbf{f}_k$ and \mathbf{V}_k . We have $I_\Omega \max_{k=1, \dots, K} (\|\mathbf{h}_k\|^2 + \|\boldsymbol{\psi}_k\|^2/n) \leq C_{43}(\gamma, \mu, M)$ as explained in the previous proofs. Summing over $k = 1, \dots, K$ the previous display, using $I_\Omega \max_{k=1, \dots, K} (\|\mathbf{h}_k\|^2 + \|\boldsymbol{\psi}_k\|^2/n) \leq C_{44}(\gamma, \mu, M)$ and $I_\Omega \|\boldsymbol{\varepsilon}\|^2 \leq n^{2/(1+q)}$ we find

$$\mathbb{E} \left[\sum_{k=1}^K I\{\Omega\} I\left\{ \sum_{i=1}^n \boldsymbol{\psi}'_k(r_{ki}) \geq n\eta \right\} \left| \|\mathbf{h}_k\|^2 + \|\boldsymbol{\varepsilon}\|^2/n - \|\mathbf{r}_k + \frac{\hat{\text{d}}\mathbf{f}_k}{\text{tr}[\mathbf{V}_k]} \boldsymbol{\psi}_k\|^2/n \right| \right] \leq \frac{KC_{45}(\gamma, \mu, \eta)}{n^{q/(1+q)}}.$$

Let $\tilde{\Omega}$ be the event that there exists \tilde{k} with $\frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}'_{\tilde{k}}(r_{\tilde{k}i}) \geq \eta$ satisfying $\|\mathbf{h}_{\tilde{k}}\|^2 + \tilde{\eta} \leq \|\mathbf{h}_{\hat{k}}\|^2$, then by the previous display and the triangle inequality, using $\|\mathbf{r}_{\hat{k}} + \frac{\hat{\text{d}}\mathbf{f}_{\tilde{k}}}{\text{tr}[\mathbf{V}_{\tilde{k}}]} \boldsymbol{\psi}_{\tilde{k}}\|^2 \leq \|\mathbf{r}_{\tilde{k}} + \frac{\hat{\text{d}}\mathbf{f}_{\tilde{k}}}{\text{tr}[\mathbf{V}_{\tilde{k}}]} \boldsymbol{\psi}_{\tilde{k}}\|^2$ by definition of \hat{k} , we obtain $\tilde{\eta} \mathbb{P}(I_\Omega I_{\tilde{\Omega}}) = O(K/n^{q/(1+q)})$. Since $\tilde{\eta}$ is a constant independent of n, p and $\mathbb{P}(\Omega) \rightarrow 1$, the probability $\mathbb{P}(\tilde{\Omega})$ converge to 0 if $K = o(n^{q/(1+q)})$. \blacksquare

Appendix C. Probabilistic results and their proofs

Proposition 13 [Variant of [Bellec and Zhang \(2019\)](#)] Let $\mathbf{z} \in N(\mathbf{0}, \mathbf{I}_q)$ and $\mathbf{f} := \mathbf{f}(\mathbf{z}) : \mathbb{R}^q \rightarrow \mathbb{R}^q \setminus \{\mathbf{0}\}$ be locally Lipschitz in \mathbf{z} with $\mathbb{E}[\|\mathbf{f}\|^{-2} \sum_{k=1}^q \|\frac{\partial \mathbf{f}}{\partial z_k}\|^2] < +\infty$. Then

$$\mathbb{E}\left[\left(\frac{\mathbf{f}^\top \mathbf{z} - \sum_{k=1}^q (\partial/\partial z_k) f_k}{\|\mathbf{f}\|_2} - Z\right)^2\right] \leq (7 + 2\sqrt{6})\mathbb{E}\left[\|\mathbf{f}\|^{-2} \sum_{k=1}^q \|\frac{\partial \mathbf{f}}{\partial z_k}\|^2\right] < +\infty. \quad (\text{B.9})$$

Proof Let $\mathbf{g} := \mathbf{g}(\mathbf{z}) = \frac{\mathbf{f}(\mathbf{z})}{\|\mathbf{f}(\mathbf{z})\|} - \mathbb{E}\left[\frac{\mathbf{f}(\mathbf{z})}{\|\mathbf{f}(\mathbf{z})\|}\right]$ and set

$$Z = \mathbf{z}^\top \mathbb{E}\left[\frac{\mathbf{f}(\mathbf{z})}{\|\mathbf{f}(\mathbf{z})\|}\right] / \sqrt{V}, \quad V = \|\mathbb{E}\left[\frac{\mathbf{f}(\mathbf{z})}{\|\mathbf{f}(\mathbf{z})\|}\right]\|^2$$

so that $Z \sim N(0, 1)$ and V is deterministic with $V \leq 1$ by Jensen's inequality. As a first step, we proceed to prove inequality

$$\mathbb{E}\left[\left(\frac{\mathbf{f}^\top \mathbf{z} - \sum_{k=1}^q (\partial/\partial z_k) f_k}{\|\mathbf{f}\|_2} - \sqrt{V}Z\right)^2\right] \leq 6 \mathbb{E}\left[\|\mathbf{f}\|^{-2} \sum_{k=1}^q \|\frac{\partial \mathbf{f}}{\partial z_k}\|^2\right]. \quad (\text{C.1})$$

Then at any point \mathbf{z} where \mathbf{f} is differentiable we have

$$\frac{\partial \mathbf{g}}{\partial z_k} = \|\mathbf{f}(\mathbf{z})\|^{-1} \hat{\mathbf{P}} \frac{\partial \mathbf{f}}{\partial z_k}, \quad \text{where} \quad \hat{\mathbf{P}} = \mathbf{I}_q - \frac{\mathbf{f}\mathbf{f}^\top}{\|\mathbf{f}\|^2}.$$

This implies that almost surely,

$$\frac{\mathbf{f}^\top \mathbf{z} - \sum_{k=1}^q (\partial/\partial z_k) f_k}{\|\mathbf{f}\|_2} - \sqrt{V}Z = \mathbf{g}^\top \mathbf{z} - \sum_{k=1}^q (\partial/\partial z_k) g_k - \frac{\mathbf{f}^\top (\partial \mathbf{f} / \partial \mathbf{z}) \mathbf{f}}{\|\mathbf{f}\|^3}$$

where $\partial \mathbf{f} / \partial \mathbf{z}$ is the matrix with entries (l, k) entry $(\partial/\partial z_k) f_l$ for all, $k, l = 1, \dots, q$.

By the triangle inequality and $(a+b)^2 \leq 2a^2 + 2b^2$, this implies that the left-hand side of (C.1) is bounded from above by $2\mathbb{E}[(\mathbf{z}^\top \mathbf{g} - \text{tr}[\partial \mathbf{g} / \partial \mathbf{z}])^2] + 2\mathbb{E}[\|\mathbf{f}\|^{-2} \|\partial \mathbf{f} / \partial \mathbf{z}\|_F^2]$. The first term can be bounded using the main result of [Bellec and Zhang \(2018\)](#) and the Gaussian Poincaré inequality ([Boucheron et al., 2013](#), Theorem 3.20)

$$\mathbb{E}[(\mathbf{z}^\top \mathbf{g} - \text{tr}[\partial \mathbf{g} / \partial \mathbf{z}])^2] = \mathbb{E}[\|\mathbf{g}\|^2] + \mathbb{E} \text{tr}[(\partial \mathbf{g} / \partial \mathbf{z})^2] \leq 2\mathbb{E}[\|\partial \mathbf{g} / \partial \mathbf{z}\|_F^2].$$

This proves (C.1). To bound $|\sqrt{V} - 1|$, we have by the triangle inequality

$$|\sqrt{V} - 1| = |\sqrt{V} - \|\frac{\mathbf{f}}{\|\mathbf{f}\|}\| | \leq \|\mathbb{E}[\frac{\mathbf{f}}{\|\mathbf{f}\|}] - \frac{\mathbf{f}}{\|\mathbf{f}\|}\| = \|\mathbf{g}\|.$$

By another application of the Gaussian Poincaré inequality,

$$|\sqrt{V} - 1|^2 \leq \mathbb{E}[\|\mathbf{g}\|_2^2] \leq \mathbb{E}[\|\partial \mathbf{g} / \partial \mathbf{z}\|_F^2] \leq \mathbb{E}[\|\mathbf{f}\|^{-2} \|\partial \mathbf{f} / \partial \mathbf{z}\|_F^2]. \quad (\text{C.2})$$

Combining Equations (C.1) and (C.2) using $(a+b)^2 = a^2 + 2ab + b^2 \leq a^2 + 1/\sqrt{6}a^2 + \sqrt{6}b^2 + b^2$, we obtain the constant $7 + 2\sqrt{6}$. ■

Proposition 14 *Let $\mathbf{h} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^p$, $\boldsymbol{\psi} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ be locally Lipschitz functions. If $\mathbf{G} \in \mathbb{R}^{n \times p}$ has iid $N(0, 1)$ entries then*

$$\begin{aligned} & \mathbb{E} \left[\left(\frac{\boldsymbol{\psi}^\top \mathbf{G} \mathbf{h} - \sum_{ij} \frac{\partial(\psi_i h_j)}{g_{ij}}}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right)^2 + \left(\frac{\|\mathbf{G} \mathbf{h}\|^2 - \sum_{ij} \frac{\partial(h_j \mathbf{e}_i^\top \mathbf{G} \mathbf{h})}{g_{ij}}}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right)^2 + \left(\frac{\|\mathbf{G}^\top \boldsymbol{\psi}\|^2 - \sum_{ij} \frac{\partial(\psi_i \mathbf{e}_j^\top \mathbf{G}^\top \boldsymbol{\psi})}{g_{ij}}}{n\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2} \right)^2 \right] \\ & \leq C_{46} \mathbb{E} \left[n + p + \|\mathbf{G}\|_{op}^2 + (n+p) \sum_{i=1}^n \sum_{j=1}^p \frac{1 + \|\mathbf{G}\|_{op}^2/n}{(\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^2} \left(\left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \right) \right] \end{aligned} \quad (\text{B.13})$$

for some positive absolute constant in the second line.

Proof (Proof of Proposition 14) We prove the claim separately for the three terms in the left-hand side of Proposition 14; we start with the first of the three terms. We will apply the probabilistic result given in Proposition 6.3 in Bellec (2020): if $\boldsymbol{\eta} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^p$ and $\boldsymbol{\rho} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ are locally Lipschitz and $\mathbf{G} \in \mathbb{R}^{n \times p}$ has iid $N(0, 1)$ entries,

$$\mathbb{E} \left[\left(\boldsymbol{\rho}^\top \mathbf{G} \boldsymbol{\eta} - \sum_{ij} \frac{\partial(\rho_i \eta_j)}{g_{ij}} \right)^2 \right] \leq \mathbb{E} \left[\|\boldsymbol{\rho}\|^2 \|\boldsymbol{\eta}\|^2 \right] + 2\mathbb{E} \left[\sum_{ij} \|\boldsymbol{\eta}\|^2 \left\| \frac{\partial \boldsymbol{\rho}}{\partial g_{ij}} \right\|^2 + \|\boldsymbol{\rho}\|^2 \left\| \frac{\partial \boldsymbol{\eta}}{\partial g_{ij}} \right\|^2 \right]. \quad (\text{C.3})$$

The proof only relies on Gaussian integration by parts to transform the left-hand side. Let $\mathbf{f} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n+p}$ be locally Lipschitz. For any i, j and at a point where both \mathbf{h} and $\boldsymbol{\psi}$ are differentiable and $\mathbf{f} \neq \mathbf{0}$,

$$\frac{\partial}{\partial g_{ij}} \left(\frac{\mathbf{f}}{\|\mathbf{f}\|} \right) = \frac{1}{\|\mathbf{f}\|} \left(\mathbf{I}_{n+p} - \frac{\mathbf{f} \mathbf{f}^\top}{\|\mathbf{f}\|^2} \right) \frac{\partial \mathbf{f}}{\partial g_{ij}} \quad \text{so that} \quad \left\| \frac{\partial}{\partial g_{ij}} \left(\frac{\mathbf{f}}{\|\mathbf{f}\|} \right) \right\|^2 \leq \frac{1}{\|\mathbf{f}\|^2} \left\| \frac{\partial \mathbf{f}}{\partial g_{ij}} \right\|^2.$$

We use this inequality applied with

$$\mathbf{f} = (\mathbf{h}, \frac{1}{\sqrt{n}} \boldsymbol{\psi}), \quad \boldsymbol{\rho} = \frac{1}{\sqrt{n}} \frac{\boldsymbol{\psi}}{\|\mathbf{f}\|}, \quad \boldsymbol{\eta} = \frac{\mathbf{h}}{\|\mathbf{f}\|}. \quad (\text{C.4})$$

To bound from above the right-hand side of (C.3), the inequality in the previous display can be rewritten

$$\left\| \frac{\partial \boldsymbol{\eta}}{\partial g_{ij}} \right\|^2 + \left\| \frac{\partial \boldsymbol{\rho}}{\partial g_{ij}} \right\|^2 \leq \frac{1}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \left(\left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \right). \quad (\text{C.5})$$

Since $\|\boldsymbol{\rho}\| \leq 1$ and $\|\boldsymbol{\eta}\| \leq 1$ by definition, the right-hand side of (C.3) is bounded from above by $1 + 2\mathbb{E} \left[\frac{1}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \left(\left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \right) \right]$. Thus the proof of Proposition 14 for the first term in the left-hand side is almost complete; it remains to control inside the parenthesis of the left-hand side,

$$\sum_{ij} \frac{1}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \frac{\partial(\psi_i n^{-1/2} h_j)}{\partial g_{ij}} - \frac{\partial}{\partial g_{ij}} \left(\frac{\psi_i n^{-1/2} h_j}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right) = 2 \sum_{ij} \psi_i n^{-1/2} h_j \frac{\mathbf{h}^\top \frac{\partial \mathbf{h}}{\partial g_{ij}} + \frac{1}{n} \boldsymbol{\psi}^\top \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}}}{(\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^2}.$$

By multiple applications of the Cauchy-Schwartz inequality, the absolute value of the previous display is bounded from above by $2(\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{-1/2} (\sum_{ij} \left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\| + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|)^{1/2}$. This completes the proof of Proposition 14 for the first term in the left-hand side.

For the second and third term in the left-hand side of Proposition 14, apply instead (C.3) to $\boldsymbol{\rho} = \mathbf{G}\boldsymbol{\eta}$ and $\boldsymbol{\eta} = \mathbf{G}^\top \boldsymbol{\rho}$ to obtain

$$\begin{aligned} \mathbb{E} \left[\left(\|\mathbf{G}\boldsymbol{\eta}\|^2 - \sum_{ij} \frac{\partial(\eta_j e_i^\top \mathbf{G}\boldsymbol{\eta})}{g_{ij}} \right)^2 \right] &\leq \mathbb{E} \left[\|\mathbf{G}\boldsymbol{\eta}\|^2 \|\boldsymbol{\eta}\|^2 \right] + 2\mathbb{E} \left[\sum_{ij} \|\boldsymbol{\eta}\|^2 \|e_i \eta_j + \mathbf{G} \frac{\partial \boldsymbol{\eta}}{\partial g_{ij}}\|^2 + \|\mathbf{G}\boldsymbol{\eta}\|^2 \left\| \frac{\partial \boldsymbol{\eta}}{\partial g_{ij}} \right\|^2 \right], \\ \mathbb{E} \left[\left(\|\mathbf{G}^\top \boldsymbol{\rho}\|^2 - \sum_{ij} \frac{\partial(\rho_i \boldsymbol{\rho}^\top \mathbf{G}e_j)}{g_{ij}} \right)^2 \right] &\leq \mathbb{E} \left[\|\mathbf{G}^\top \boldsymbol{\rho}\|^2 \|\boldsymbol{\rho}\|^2 \right] + 2\mathbb{E} \left[\sum_{ij} \|\mathbf{G}^\top \boldsymbol{\rho}\|^2 \left\| \frac{\partial \boldsymbol{\rho}}{\partial g_{ij}} \right\|^2 + \|\boldsymbol{\rho}\|^2 \|e_j \rho_i + \mathbf{G}^\top \frac{\partial \boldsymbol{\rho}}{\partial g_{ij}}\|^2 \right]. \end{aligned}$$

Setting $\boldsymbol{\rho} = \frac{1}{\sqrt{n}}\boldsymbol{\psi}/\|\mathbf{f}\|$, $\boldsymbol{\eta} = \mathbf{h}/\|\mathbf{f}\|$ we obtain the claim in Equation (C.3) by bounding the right-hand side of the previous displays using the operator norm of \mathbf{G} and arguments similar to (C.5). The term involving $\frac{\partial}{\partial g_{ij}} \left(\frac{1}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right)$ in the left-hand side is controlled similarly to the previous paragraph. ■

Proposition 16 *Let $\mathbf{h} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^p$, $\boldsymbol{\psi} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ be locally Lipschitz functions. If $\mathbf{G} \in \mathbb{R}^{n \times p}$ has iid $N(0, 1)$ entries then*

$$\begin{aligned} &\mathbb{E} \left[\frac{\left| \frac{p}{n} \|\boldsymbol{\psi}\|^2 - \frac{1}{n} \sum_{j=1}^p (\boldsymbol{\psi}^\top \mathbf{G}e_j - \sum_{i=1}^n \frac{\partial \psi_i}{\partial g_{ij}})^2 \right|}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right] + \mathbb{E} \left[\frac{\left| n \|\mathbf{h}\|^2 - \sum_{i=1}^n (\mathbf{g}_i^\top \mathbf{h} - \sum_{j=1}^p \frac{\partial h_j}{\partial g_{ij}})^2 \right|}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \right] \\ &\leq C_{47} \left(\sqrt{n+p} (1 + \Xi^{1/2}) + \Xi \right) \text{ where } \Xi = \mathbb{E} \left[\frac{1}{\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n} \sum_{i=1}^n \sum_{j=1}^p \left(\left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2 \right) \right]. \end{aligned}$$

Proof (Proof of Proposition 16) We first focus on the first term in the left-hand side. Theorem 7.1 in Bellec (2020) provides that if $\boldsymbol{\rho} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ is locally Lipschitz with $\|\boldsymbol{\rho}\| \leq 1$ then

$$\mathbb{E} \left| p \|\boldsymbol{\rho}\|^2 - \sum_{j=1}^p \left(\boldsymbol{\rho}^\top \mathbf{G}e_j - \sum_{i=1}^n \frac{\partial \rho_i}{\partial g_{ij}} \right)^2 \right| \leq C_{48} \sqrt{p} \left(1 + \mathbb{E} \sum_{ij} \left\| \frac{\partial \boldsymbol{\rho}}{\partial g_{ij}} \right\|^2 \right)^{1/2} + C_{49} \mathbb{E} \sum_{ij} \left\| \frac{\partial \boldsymbol{\rho}}{\partial g_{ij}} \right\|^2. \quad (\text{C.6})$$

Let $\boldsymbol{\rho} = n^{-1/2}\boldsymbol{\psi}/\|\mathbf{f}\|$ as in (C.4). Inequality (C.5) lets us bound from above the right-hand side of the previous display by the right-hand side of Proposition 16. In the left-hand side, $p\|\boldsymbol{\rho}\|^2 = \frac{p}{n}\|\boldsymbol{\psi}\|^2/(\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)$ as desired. For the left-hand side, using some algebra in (Bellec, 2020, Section 7), for any random vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ by the triangle and Cauchy-Schwarz inequalities we have

$$\begin{aligned} |p\|\boldsymbol{\rho}\|^2 - \|\mathbf{a}\|^2| - |p\|\boldsymbol{\rho}\|^2 - \|\mathbf{b}\|^2| &\leq \|\mathbf{a} - \mathbf{b}\| \|\mathbf{a} + \mathbf{b}\| \\ &\leq \|\mathbf{a} - \mathbf{b}\|^2 + 2\|\mathbf{a} - \mathbf{b}\| \|\mathbf{b}\| \\ &\leq \|\mathbf{a} - \mathbf{b}\|^2 + 2\|\mathbf{a} - \mathbf{b}\| (\sqrt{\|\mathbf{b}\|^2 - p\|\boldsymbol{\rho}\|^2} + \sqrt{p\|\boldsymbol{\rho}\|^2}) \\ &\leq 3\|\mathbf{a} - \mathbf{b}\|^2 + \frac{1}{2}\|\mathbf{b}\|^2 - p\|\boldsymbol{\rho}\|^2 + 2\|\mathbf{a} - \mathbf{b}\| \sqrt{p\|\boldsymbol{\rho}\|^2} \end{aligned}$$

so that $|p\|\boldsymbol{\rho}\|^2 - \|\mathbf{a}\|^2| \leq \frac{3}{2}|p\|\boldsymbol{\rho}\|^2 - \|\mathbf{b}\|^2| + 3\|\mathbf{a} - \mathbf{b}\|^2 + 2\|\mathbf{a} - \mathbf{b}\| \sqrt{p\|\boldsymbol{\rho}\|^2}$. Applying this to $b_j = \boldsymbol{\rho}^\top \mathbf{G}e_j - \sum_{i=1}^n \frac{\partial \rho_i}{\partial g_{ij}}$ we use (C.6) to bound $|p\|\boldsymbol{\rho}\|^2 - \|\mathbf{b}\|^2|$ and $\|\boldsymbol{\rho}\| \leq 1$ to bound

$\sqrt{p\|\boldsymbol{\rho}\|^2} \leq \sqrt{p}$. It remains to specify \mathbf{a} so that $|p\|\boldsymbol{\rho}\|^2 - \|\mathbf{a}\|^2|$ coincides with the first term in the left-hand side of Proposition 16 and bound $\|\mathbf{a} - \mathbf{b}\|$. Consequently, we set

$$a_j = \frac{\boldsymbol{\psi}^\top \mathbf{G} \mathbf{e}_j - \sum_{i=1}^n \frac{\partial \psi_i}{\partial g_{ij}}}{\sqrt{n}(\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{1/2}} = \boldsymbol{\rho}^\top \mathbf{G} \mathbf{e}_j - \frac{\sum_{i=1}^n \frac{\partial \psi_i}{\partial g_{ij}}}{\sqrt{n}(\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{1/2}} = b_j - \sum_{i=1}^n \frac{\psi_i}{\sqrt{n}} \frac{\partial(D^{-1})}{\partial g_{ij}}$$

where $D = (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{1/2}$ so that by the Cauchy-Schwarz inequality $\|\mathbf{a} - \mathbf{b}\|^2 \leq \frac{1}{n} \|\boldsymbol{\psi}\|^2 \sum_{ij} \left(\frac{\partial(D^{-1})}{\partial g_{ij}}\right)^2$ and

$$\sum_{ij} \left(\frac{\partial(D^{-1})}{\partial g_{ij}}\right)^2 = \frac{1}{D^6} \sum_{ij} \left(\mathbf{h}^\top \frac{\partial \mathbf{h}}{\partial g_{ij}} + \frac{\boldsymbol{\psi}^\top}{\sqrt{n}} \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}}\right)^2 \leq \frac{2}{D^4} \sum_{ij} \left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2. \quad (\text{C.7})$$

using again the Cauchy-Schwarz inequality and $\max\{\|\mathbf{h}\|^2, \|\boldsymbol{\psi}\|^2/n\} \leq D^2$. We obtain $\|\mathbf{a} - \mathbf{b}\|^2 \leq D^{-2} \sum_{ij} \left\| \frac{\partial \mathbf{h}}{\partial g_{ij}} \right\|^2 + \frac{1}{n} \left\| \frac{\partial \boldsymbol{\psi}}{\partial g_{ij}} \right\|^2$ which completes the proof for the first term in the left-hand side of Proposition 16. For the second term in the left-hand side, the proof is similar with by exchanging the role of n and p in (C.6) and applying (C.6) to \mathbf{h}/D instead of $\boldsymbol{\psi}/(\sqrt{n}D)$. ■

Proposition 18 *Let Assumption A be fulfilled. Then*

$$\mathbb{E} \left[\left\{ (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{-\frac{1}{2}} \|\boldsymbol{\varepsilon}\|^{-1} \xi_{VI} \right\}^2 \right] \leq C_{50}(\gamma, \mu) \quad \text{for} \quad \xi_{VI} = \boldsymbol{\varepsilon}^\top (\mathbf{G} \mathbf{h} - \text{tr}[\mathbf{A}] \boldsymbol{\psi}). \quad (\text{B.20})$$

Proof (Proof of Proposition 18) Apply (C.3) with $\boldsymbol{\rho} = \boldsymbol{\varepsilon}/\|\boldsymbol{\varepsilon}\|$ and $\boldsymbol{\eta} = \mathbf{h}/D$ where $D = (\|\mathbf{h}\|^2 + \|\boldsymbol{\psi}\|^2/n)^{1/2}$ as in the previous proof (this scalar D is not related to the diagonal matrix $\mathbf{D} = \text{diag}\{\psi'(\mathbf{r})\}$). Since $\boldsymbol{\varepsilon}$ has 0 derivative with respect to \mathbf{G} we find

$$\mathbb{E} \left[\left(\frac{\boldsymbol{\varepsilon}^\top \mathbf{G} \mathbf{h}}{\|\boldsymbol{\varepsilon}\| D} - \sum_{ij} \frac{\varepsilon_i}{\|\boldsymbol{\varepsilon}\|} \frac{\partial(h_j D^{-1})}{\partial g_{ij}} \right)^2 \right] \leq 1 + 2 \sum_{ij} \mathbb{E} \left[\left\| \frac{\partial \boldsymbol{\eta}}{\partial g_{ij}} \right\|^2 \right].$$

The right-hand side is bounded from above by $C_{51}(\gamma, \mu)$ thanks to (C.5) and (B.17). For the second term above we use product rule and (B.4),

$$\sum_{ij} \frac{\varepsilon_i}{\|\boldsymbol{\varepsilon}\|} \frac{\partial(h_j D^{-1})}{\partial g_{ij}} = \frac{\text{tr}[\mathbf{A}] \boldsymbol{\psi}^\top \boldsymbol{\varepsilon}}{D \|\boldsymbol{\varepsilon}\|} - \frac{\mathbf{h}^\top \mathbf{A} \mathbf{G}^\top \text{diag}(\boldsymbol{\psi}'(\mathbf{r})) \boldsymbol{\varepsilon}}{D \|\boldsymbol{\varepsilon}\|} + \sum_{ij} \frac{\varepsilon_i h_j}{\|\boldsymbol{\varepsilon}\|} \frac{\partial(D^{-1})}{\partial g_{ij}}.$$

To complete the proof we need to prove that the expectation of the square of the second and third terms colored in purple are bounded by $C_{52}(\gamma, \mu)$. Since $\|\mathbf{h}\| \leq D$, the second term is bounded from above by $\|\mathbf{A}\|_{op} \|\mathbf{G}\|_{op}$ since $|\psi'| \leq 1$ and $\mathbb{E}[\|\mathbf{A}\|_{op}^2 \|\mathbf{G}\|_{op}^2] \leq C_{53}(\gamma, \mu)$ thanks to $\|\mathbf{A}\|_{op} \leq 1/(n\mu)$ and (Davidson and Szarek, 2001, Theorem II.13). For the third term, we use the Cauchy-Schwarz inequality $(\sum_{ij} \frac{\varepsilon_i h_j}{\|\boldsymbol{\varepsilon}\|} \frac{\partial(D^{-1})}{\partial g_{ij}})^2 \leq \|\mathbf{h}\|^2 \sum_{ij} \left(\frac{\partial(D^{-1})}{\partial g_{ij}}\right)^2$, (C.7) and (B.17). ■

Appendix D. Proof of differentiability results

Theorem 1 *Let Assumption A be fulfilled. For almost every (\mathbf{y}, \mathbf{X}) the map $(\mathbf{y}, \mathbf{X}) \mapsto \widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X})$ is differentiable at (\mathbf{y}, \mathbf{X}) and there exists a matrix $\widehat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ depending on (\mathbf{X}, \mathbf{y}) with $\|\Sigma^{1/2} \widehat{\mathbf{A}} \Sigma^{1/2}\|_{op} \leq (n\mu)^{-1}$ s.t.*

$$\begin{aligned} (\partial/\partial y_i) \widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) &= \widehat{\mathbf{A}} \mathbf{X}^\top \mathbf{e}_i \psi'(r_i), \\ (\partial/\partial x_{ij}) \widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X}) &= \widehat{\mathbf{A}} \mathbf{e}_j \psi'(r_i) - \widehat{\mathbf{A}} \mathbf{X}^\top \mathbf{e}_i \psi'(r_i) \widehat{\boldsymbol{\beta}}_j, \end{aligned} \quad \text{where } r_i = y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}, \quad (2.1)$$

$\mathbf{e}_i \in \mathbb{R}^n$, $\mathbf{e}_j \in \mathbb{R}^p$ are canonical basis vectors, $\psi := \rho'$ and ψ' denote the derivatives. Furthermore,

$$\widehat{\text{df}} = \text{tr}[\mathbf{X}(\partial/\partial \mathbf{y}) \widehat{\boldsymbol{\beta}}] = \text{tr}[\mathbf{X} \widehat{\mathbf{A}} \mathbf{X}^\top \text{diag}\{\psi'(\mathbf{r})\}], \quad (2.2)$$

$$\mathbf{V} = \text{diag}\{\psi'(\mathbf{r})\}(\mathbf{I}_n - \mathbf{X}(\partial/\partial \mathbf{y}) \widehat{\boldsymbol{\beta}}) = \text{diag}\{\psi'(\mathbf{r})\} - \text{diag}\{\psi'(\mathbf{r})\} \mathbf{X} \widehat{\mathbf{A}} \mathbf{X}^\top \text{diag}\{\psi'(\mathbf{r})\} \quad (2.3)$$

satisfy $0 \leq \widehat{\text{df}} \leq n$ and $0 \leq \text{tr}[\mathbf{V}] \leq n$.

The first part of the following proof is similar to the argument using the KKT conditions in Bellec (2020). After (D.3), the argument is novel and lets us derive the convenient formula (2.1) and the existence of matrix $\widehat{\mathbf{A}}$ which plays a central role in the contractions (B.4)-(B.8).

Proof (Proof of Theorem 1) $\mathbf{X}_t = \mathbf{X} + t\mathbf{U}$ and $\mathbf{y}_t = \mathbf{y} + t\mathbf{v}$ with $t \in \mathbb{R}$ where $\mathbf{U} \in \mathbb{R}^{n \times p}$ and $\mathbf{v} \in \mathbb{R}^n$ are fixed. Let $\widehat{\boldsymbol{\beta}}_t = \widehat{\boldsymbol{\beta}}(\mathbf{y}_t, \mathbf{X}_t)$ and $\widehat{\mathbf{r}}_t = \mathbf{y}_t - \mathbf{X}_t \widehat{\boldsymbol{\beta}}_t$ and $\widehat{\boldsymbol{\psi}}_t(\mathbf{y}_t, \mathbf{X}_t) = \psi(\widehat{\mathbf{r}}_t)$. By convention, without arguments $\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\psi}}$ refer to (\mathbf{y}, \mathbf{X}) which is $(\mathbf{y}_t, \mathbf{X}_t)$ at $t = 0$. By the KKT conditions, $\mathbf{X}^\top \widehat{\boldsymbol{\psi}} \in n \text{dg}(\widehat{\boldsymbol{\beta}})$ and $\mathbf{X}_t^\top \widehat{\boldsymbol{\psi}}_t \in n \text{dg}(\widehat{\boldsymbol{\beta}}_t)$, by strong convexity of g , we have

$$n\mu \|\Sigma^{1/2}(\widehat{\boldsymbol{\beta}}_t - \widehat{\boldsymbol{\beta}})\|^2 \leq (\widehat{\boldsymbol{\beta}}_t - \widehat{\boldsymbol{\beta}})^\top (\mathbf{X}_t^\top \widehat{\boldsymbol{\psi}}_t - \mathbf{X}^\top \widehat{\boldsymbol{\psi}}). \quad (\text{D.1})$$

By the fact that ψ is non-decreasing and 1-Lipschitz, for any two real numbers $a < b$, $0 \leq \psi(b) - \psi(a) \leq b - a$. Multiplying $\psi(b) - \psi(a)$, we have $(\psi(b) - \psi(a))^2 \leq (\psi(b) - \psi(a))(b - a)$. Thus

$$\|\widehat{\boldsymbol{\psi}}_t - \widehat{\boldsymbol{\psi}}\|^2 \leq (\widehat{\boldsymbol{\psi}}_t - \widehat{\boldsymbol{\psi}})^\top (\widehat{\mathbf{r}}_t - \widehat{\mathbf{r}}).$$

Adding up the above two displays we have

$$n\mu \|\Sigma^{1/2}(\widehat{\boldsymbol{\beta}}_t - \widehat{\boldsymbol{\beta}})\|^2 + \|\widehat{\boldsymbol{\psi}}_t - \widehat{\boldsymbol{\psi}}\|^2 \leq (\widehat{\boldsymbol{\beta}}_t - \widehat{\boldsymbol{\beta}})^\top (\mathbf{X}_t^\top \widehat{\boldsymbol{\psi}}_t - \mathbf{X}^\top \widehat{\boldsymbol{\psi}}) + (\widehat{\boldsymbol{\psi}}_t - \widehat{\boldsymbol{\psi}})^\top (\widehat{\mathbf{r}}_t - \widehat{\mathbf{r}}). \quad (\text{D.2})$$

By $\mathbf{X}_t^\top \widehat{\boldsymbol{\psi}}_t - \mathbf{X}^\top \widehat{\boldsymbol{\psi}} = (\mathbf{X}_t - \mathbf{X})^\top \widehat{\boldsymbol{\psi}} + \mathbf{X}_t^\top (\widehat{\boldsymbol{\psi}}_t - \widehat{\boldsymbol{\psi}})$ and $\mathbf{X}_t(\widehat{\boldsymbol{\beta}}_t - \widehat{\boldsymbol{\beta}}) + \widehat{\mathbf{r}}_t - \widehat{\mathbf{r}} = \mathbf{y}_t - \mathbf{y} - (\mathbf{X}_t - \mathbf{X})^\top \widehat{\boldsymbol{\beta}}$, we have

$$n\mu \|\Sigma^{1/2}(\widehat{\boldsymbol{\beta}}_t - \widehat{\boldsymbol{\beta}})\|^2 + \|\widehat{\boldsymbol{\psi}}_t - \widehat{\boldsymbol{\psi}}\|^2 \leq (\widehat{\boldsymbol{\beta}}_t - \widehat{\boldsymbol{\beta}})^\top (\mathbf{X}_t - \mathbf{X})^\top \widehat{\boldsymbol{\psi}} + (\mathbf{y}_t - \mathbf{y} - (\mathbf{X}_t - \mathbf{X})^\top \widehat{\boldsymbol{\beta}})^\top (\widehat{\boldsymbol{\psi}}_t - \widehat{\boldsymbol{\psi}}).$$

By the Cauchy-Schwartz inequality, the above implies

$$(n\mu \|\Sigma^{1/2}(\widehat{\boldsymbol{\beta}}_t - \widehat{\boldsymbol{\beta}})\|^2 + \|\widehat{\boldsymbol{\psi}}_t - \widehat{\boldsymbol{\psi}}\|^2)^{1/2} \leq (n\mu)^{-1/2} \|\Sigma^{-1/2}(\mathbf{X}_t - \mathbf{X})^\top \widehat{\boldsymbol{\psi}}\|_2 + \|\mathbf{y}_t - \mathbf{y} - (\mathbf{X}_t - \mathbf{X})^\top \widehat{\boldsymbol{\beta}}\|_2,$$

Since $t, \mathbf{U}, \mathbf{v}$ are arbitrary, for $(\mathbf{y}_t, \mathbf{X}_t)$ and (\mathbf{y}, \mathbf{X}) both in a compact subset K of $\mathbb{R}^p \times \mathbb{R}^{n \times p}$, the above display also implies

$$(n\mu \|\Sigma^{1/2}(\widehat{\boldsymbol{\beta}}_t - \widehat{\boldsymbol{\beta}})\|^2 + \|\widehat{\boldsymbol{\psi}}_t - \widehat{\boldsymbol{\psi}}\|^2)^{1/2} \leq \text{const}(K) (\|\Sigma^{-1/2}(\mathbf{X}_t - \mathbf{X})\|_{op} + \|\mathbf{y}_t - \mathbf{y}\|_2),$$

where $\text{const}(K) := \sup_{(\mathbf{y}, \mathbf{X}) \in K} \{(n\mu)^{-1/2} \|\widehat{\boldsymbol{\psi}}\|_2 + 1 + \|\boldsymbol{\Sigma}^{1/2} \widehat{\boldsymbol{\beta}}\|_2\}$. This says that $\widehat{\boldsymbol{\beta}}(\mathbf{y}, \mathbf{X})$, $\widehat{\boldsymbol{\psi}}(\mathbf{y}, \mathbf{X})$ are locally Lipschitz in (\mathbf{y}, \mathbf{X}) . By Rademacher's Theorem, $\partial \widehat{\boldsymbol{\beta}}/\partial y_i$ and $\partial \widehat{\boldsymbol{\beta}}/\partial x_{ij}$ exist almost everywhere.

Taking the limit $t \rightarrow 0^+$ in (D.1) and using the chain rule, where the derivatives exist we have

$$\begin{aligned} & n\mu \left\| \boldsymbol{\Sigma}^{1/2} \left(\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \mathbf{v} + \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) \right) \right\|_2^2 \\ & \leq \left(\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \mathbf{v} + \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) \right)^\top \left(\mathbf{U}^\top \widehat{\boldsymbol{\psi}} + \mathbf{X}^\top \text{diag}(\widehat{\boldsymbol{\psi}}')(-\mathbf{U} \widehat{\boldsymbol{\beta}} - \mathbf{X} \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) + \left(I_n - \mathbf{X} \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \right) \mathbf{v}) \right) \\ & = \left(\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \mathbf{v} + \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) \right)^\top B(\mathbf{U}, \mathbf{v}) - \left\| \text{diag}(\widehat{\boldsymbol{\psi}}')^{\frac{1}{2}} \mathbf{X} \left(\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \mathbf{v} + \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) \right) \right\|_2^2 \end{aligned} \quad (\text{D.3})$$

where $(\partial \widehat{\boldsymbol{\beta}}/\partial \mathbf{y})\mathbf{v} := \sum_{i \in [n]} (\partial \widehat{\boldsymbol{\beta}}/\partial y_i) v_i$, the Jacobian with respect to \mathbf{X} and the linear map $B : \mathbb{R}^{n \times p} \times \mathbb{R}^n \rightarrow \mathbb{R}^p$ are defined as

$$\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) := \sum_{i,j \in [n] \times [p]} \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial x_{ij}} u_{ij} \in \mathbb{R}^p, \quad B(\mathbf{U}, \mathbf{v}) := \mathbf{U}^\top \widehat{\boldsymbol{\psi}} + \mathbf{X}^\top \text{diag}(\widehat{\boldsymbol{\psi}}')(-\mathbf{U} \widehat{\boldsymbol{\beta}} + \mathbf{v}) \in \mathbb{R}^p$$

where $(u_{ij})_{i=1, \dots, n, j=1, \dots, p}$ are the entries of \mathbf{U} . By the Cauchy-Schwartz inequality, (D.3) provides us the following two main ingredients:

$$\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \mathbf{v} + \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) = 0 \text{ for all } (\mathbf{U}, \mathbf{v}) \text{ such that } B(\mathbf{U}, \mathbf{v}) = 0, \quad (\text{D.4})$$

$$\left\| \boldsymbol{\Sigma}^{1/2} \left(\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \mathbf{v} + \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) \right) \right\|_2 \leq \mu^{-1} n^{-1} \|\boldsymbol{\Sigma}^{-1/2} B(\mathbf{U}, \mathbf{v})\|_2. \quad (\text{D.5})$$

Since both $\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \mathbf{v} + \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U})$ and $B(\mathbf{U}, \mathbf{v})$ are linear in $(\mathbf{U}, \mathbf{v}) \in \mathbb{R}^{n \times p} \times \mathbb{R}^n$ into \mathbb{R}^p , Proposition 19 implies that there exists a matrix $\widehat{\mathbf{A}} \in \mathbb{R}^{p \times p}$ such that $\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \mathbf{v} + \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) = \widehat{\mathbf{A}} B(\mathbf{U}, \mathbf{v})$ for all (\mathbf{U}, \mathbf{v}) , and by (D.5), $\widehat{\mathbf{A}}$ can be chosen such that $\|\boldsymbol{\Sigma}^{1/2} \widehat{\mathbf{A}} \boldsymbol{\Sigma}^{1/2}\|_{op} \leq (n\mu)^{-1}$ thanks to the operator norm identity in Proposition 19. With $(\mathbf{U}, \mathbf{v}) = (e_i e_j^\top, \mathbf{0})$ for $(i, j) \in [n] \times [p]$ and $(\mathbf{U}, \mathbf{v}) = (\mathbf{0}, e_k)$ for $k \in [n]$, we obtain the stated formulae for $(\partial/\partial x_{ij})\widehat{\boldsymbol{\beta}}$ and $(\partial/\partial y_k)\widehat{\boldsymbol{\beta}}$ in (2.1).

Now we show that both $\text{tr}[\mathbf{V}] := \text{tr}[\mathbf{D} - \mathbf{D}\mathbf{X}\widehat{\mathbf{A}}\mathbf{X}^\top \mathbf{D}]$ and $\text{df} := \text{tr}[\mathbf{X}\widehat{\mathbf{A}}\mathbf{X}^\top \mathbf{D}]$ are in $[0, n]$ where $\mathbf{D} := \text{diag}\{\psi'(\mathbf{r})\}$. Using the symmetric part of $\widehat{\mathbf{A}}$ defined as $\widetilde{\mathbf{A}} := (\widehat{\mathbf{A}} + \widehat{\mathbf{A}}^\top)/2$ we have $\text{tr}[\mathbf{V}] = \text{tr}[\mathbf{D} - \mathbf{D}\mathbf{X}\widetilde{\mathbf{A}}\mathbf{X}^\top \mathbf{D}]$ and $\text{df} = \text{tr}[\mathbf{D}^{1/2} \mathbf{X}\widetilde{\mathbf{A}}\mathbf{X}^\top \mathbf{D}^{1/2}]$ by property of the trace. In (D.3), take $\mathbf{U} = \mathbf{0}$ so that $\frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{y}} \mathbf{v} + \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{X}}(\mathbf{U}) = \widehat{\mathbf{A}} B(\mathbf{U}, \mathbf{v}) = \widehat{\mathbf{A}} \mathbf{X}^\top \mathbf{D} \mathbf{v}$ and we have with $\mathbf{G} = \mathbf{X} \boldsymbol{\Sigma}^{-1/2}$

$$\left(1 + \frac{n\mu}{\|\mathbf{D}^{1/2} \mathbf{G}\|_{op}^2} \right) \|\mathbf{D}^{1/2} \mathbf{X} \widehat{\mathbf{A}} \mathbf{X}^\top \mathbf{D} \mathbf{v}\|^2 \leq n\mu \|\widehat{\mathbf{A}} \mathbf{X}^\top \mathbf{D} \mathbf{v}\|^2 + \|\mathbf{D}^{1/2} \mathbf{X} \widehat{\mathbf{A}} \mathbf{X}^\top \mathbf{D} \mathbf{v}\|^2 \quad (\text{D.6})$$

$$\leq \mathbf{v}^\top \mathbf{D} \mathbf{X} \widetilde{\mathbf{A}} \mathbf{X}^\top \mathbf{D} \mathbf{v} = \mathbf{v}^\top \mathbf{D} \mathbf{X} \widetilde{\mathbf{A}} \mathbf{X}^\top \mathbf{D} \mathbf{v} \quad (\text{D.7})$$

for all \mathbf{v} . This implies the positive semi-definite property of the symmetric matrix $\mathbf{D}\mathbf{X}\widetilde{\mathbf{A}}\mathbf{X}^\top \mathbf{D}$, and thus $\text{df} \geq 0$ and $\text{tr}[\mathbf{V}] \leq \text{tr}[\mathbf{D}] \leq n$. With $\widetilde{\mathbf{v}} = \mathbf{D}^{1/2} \mathbf{v}$, it also implies

$$\left(1 + n\mu / \|\mathbf{D}^{1/2} \mathbf{G}\|_{op}^2 \right) \|\mathbf{D}^{1/2} \mathbf{X} \widehat{\mathbf{A}} \mathbf{X}^\top \mathbf{D}^{1/2} \widetilde{\mathbf{v}}\|^2 \leq \widetilde{\mathbf{v}}^\top \mathbf{D}^{1/2} \mathbf{X} \widehat{\mathbf{A}} \mathbf{X}^\top \mathbf{D}^{1/2} \widetilde{\mathbf{v}},$$

which, by the Cauchy-Schwartz inequality, yields $(1 + n\mu/\|\mathbf{D}^{1/2}\mathbf{G}\|_{op}^2)\|\mathbf{D}^{1/2}\mathbf{X}\widehat{\mathbf{A}}\mathbf{X}^\top\mathbf{D}^{1/2}\|_{op} \leq 1$. The same operator norm inequality with $\widehat{\mathbf{A}}$ replaced by its symmetric part $\widetilde{\mathbf{A}} = (\widehat{\mathbf{A}} + \widehat{\mathbf{A}}^\top)/2$ thanks to the triangle inequality. Thus $\widehat{\mathbf{d}} \leq \text{tr}[\mathbf{D}](1 + n\mu/\|\mathbf{D}^{1/2}\mathbf{G}\|_{op}^2)^{-1} \leq n$ as well as

$$\begin{aligned} \text{tr}[\mathbf{V}] &= \text{tr}[\mathbf{D}^{1/2}(\mathbf{I}_n - \mathbf{D}^{1/2}\mathbf{X}\widetilde{\mathbf{A}}\mathbf{X}^\top\mathbf{D}^{1/2})\mathbf{D}^{1/2}] \geq \text{tr}[\mathbf{D}](1 - (1 + n\mu/\|\mathbf{D}^{1/2}\mathbf{G}\|_{op}^2)^{-1}) \\ &= \text{tr}[\mathbf{D}]/(\|\mathbf{D}^{1/2}\mathbf{G}\|_{op}^2/(n\mu) + 1) \\ &\geq \text{tr}[\mathbf{D}]/(\|\mathbf{G}\|_{op}^2/(n\mu) + 1) \quad (\text{D.8}) \\ &\geq 0 \end{aligned}$$

thanks to $\psi' \in [0, 1]$. Inequality (D.7) with $\tilde{\mathbf{v}} = \mathbf{D}^{1/2}\mathbf{v}$ and $\mathbf{M} = \mathbf{I}_n - \mathbf{D}^{1/2}\mathbf{X}\widehat{\mathbf{A}}\mathbf{X}^\top\mathbf{D}^{1/2}$ implies $\|(\mathbf{M} - \mathbf{I}_n)\tilde{\mathbf{v}}\|^2 \leq \tilde{\mathbf{v}}^\top(\mathbf{I}_n - \mathbf{M})\tilde{\mathbf{v}}$. As the left-hand side is $\|\mathbf{M}\tilde{\mathbf{v}}\|^2 - 2\tilde{\mathbf{v}}^\top\mathbf{M}\tilde{\mathbf{v}} + \|\tilde{\mathbf{v}}\|^2$, this yields $\|\mathbf{M}\tilde{\mathbf{v}}\|^2 \leq \tilde{\mathbf{v}}^\top\mathbf{M}\tilde{\mathbf{v}} \leq \|\tilde{\mathbf{v}}\|\|\mathbf{M}\tilde{\mathbf{v}}\|$. If $\tilde{\mathbf{v}}$ has unit norm and is such that $\|\mathbf{M}\tilde{\mathbf{v}}\| = \|\mathbf{M}\|_{op}$ this gives $\|\mathbf{M}\|_{op} \leq 1$ so that $\|\mathbf{V}\|_{op} = \|\mathbf{D}^{1/2}\mathbf{M}\mathbf{D}^{1/2}\|_{op} \leq \|\mathbf{D}\|_{op} \leq 1$. This gives another proof of $\text{tr}[\mathbf{V}] \leq n$. \blacksquare

Proof (Proof of Theorem 3)

The proof for the intercept term included is the same to that of Theorem 1. The only difference is that when computing the derivatives,

$$\begin{aligned} \frac{d\widehat{\psi}_t}{dt}\Big|_{t=0} &= \mathbf{U}^\top\widehat{\psi} + \mathbf{X}^\top\left(\frac{\partial\widehat{\psi}}{\partial\mathbf{y}}\mathbf{v} + \frac{\partial\widehat{\psi}}{\partial\mathbf{X}}(\mathbf{U})\right), \quad \frac{\partial\widehat{\psi}}{\partial\mathbf{y}}\mathbf{v} = \text{diag}(\widehat{\psi}')(I_n - \mathbf{1}\frac{\partial\widehat{\beta}_0}{\partial\mathbf{y}} - \mathbf{X}\frac{\partial\widehat{\beta}}{\partial\mathbf{y}})\mathbf{v}, \\ \frac{\partial\widehat{\psi}}{\partial\mathbf{X}}(\mathbf{U}) &= \text{diag}(\widehat{\psi}')(-\mathbf{1}\frac{\partial\widehat{\beta}_0}{\partial\mathbf{X}}(\mathbf{U}) - \mathbf{U}\widehat{\beta} - \mathbf{X}\frac{\partial\widehat{\beta}}{\partial\mathbf{X}}(\mathbf{U})) \\ \implies \frac{d\widehat{\psi}_t}{dt}\Big|_{t=0} &= -\widehat{\psi}'\frac{d\widehat{\beta}_{0,t}}{dt}\Big|_{t=0} - \text{diag}(\widehat{\psi}')\mathbf{X}\frac{d\widehat{\beta}_t}{dt}\Big|_{t=0} + \text{diag}(\widehat{\psi}')\mathbf{v} - \text{diag}(\widehat{\psi}')\mathbf{U}\widehat{\beta}. \end{aligned}$$

We have an additional KKT conditions providing us $0 = \mathbf{1}^\top(d\widehat{\psi}_t/dt)|_{t=0}$. Multiplying $\mathbf{1}^\top$ on both sides of the above display, we have

$$\begin{aligned} \frac{d\widehat{\beta}_{0,t}}{dt}\Big|_{t=0} &= -\frac{\widehat{\psi}'^\top\mathbf{X}}{\mathbf{1}^\top\widehat{\psi}'}\frac{d\widehat{\beta}_t}{dt}\Big|_{t=0} + \frac{\widehat{\psi}'^\top\mathbf{v}}{\mathbf{1}^\top\widehat{\psi}'} - \frac{\widehat{\psi}'^\top\mathbf{U}\widehat{\beta}}{\mathbf{1}^\top\widehat{\psi}'}, \\ \implies \frac{d\widehat{\psi}_t}{dt}\Big|_{t=0} &= -\Psi'\mathbf{X}\frac{d\widehat{\beta}_t}{dt}\Big|_{t=0} + \Psi'\mathbf{v} - \Psi'\mathbf{U}\widehat{\beta}, \end{aligned}$$

where $\Psi' := \text{diag}(\widehat{\psi}') - \widehat{\psi}'\widehat{\psi}'^\top/\mathbf{1}^\top\widehat{\psi}'$. By taking limit of $t \rightarrow 0$ in Equation (D.2),

$$\begin{aligned} n\mu\left\|\frac{d\widehat{\beta}_t}{dt}\Big|_{t=0}\right\|_2^2 &\leq \frac{d\widehat{\beta}_t}{dt}\Big|_{t=0}^\top\frac{d(\mathbf{X}^\top\widehat{\psi})}{dt}\Big|_{t=0} = \frac{d\widehat{\beta}_t}{dt}\Big|_{t=0}^\top\left(\mathbf{U}^\top\widehat{\psi} + \mathbf{X}^\top\frac{d\widehat{\psi}_t}{dt}\Big|_{t=0}\right) \\ &= \frac{d\widehat{\beta}_t}{dt}\Big|_{t=0}^\top\left(\mathbf{U}^\top\widehat{\psi} + \mathbf{X}^\top\left(-\Psi'\mathbf{X}\frac{d\widehat{\beta}_t}{dt}\Big|_{t=0} + \Psi'\mathbf{v} - \Psi'\mathbf{U}\widehat{\beta}\right)\right) \\ &= \frac{d\widehat{\beta}_t}{dt}\Big|_{t=0}^\top\left(\mathbf{U}^\top\widehat{\psi} + \mathbf{X}^\top\Psi'\mathbf{v} - \mathbf{X}^\top\Psi'\mathbf{U}\widehat{\beta}\right) - \left\|\Psi'^{1/2}\mathbf{X}\frac{d\widehat{\beta}_t}{dt}\Big|_{t=0}\right\|_2^2. \end{aligned}$$

\blacksquare

Proposition 19 *Let \mathbf{A} and \mathbf{B} be two real matrices with shape n by p . Assume that $\mathbf{B}\mathbf{v} = \mathbf{0}$ for all $\mathbf{v} \in \mathbb{R}^p$ such that $\mathbf{A}\mathbf{v} = \mathbf{0}$. Then the matrix $\mathbf{C} := \mathbf{B}\mathbf{A}^+$ where \mathbf{A}^+ is the Moore-Penrose pseudoinverse of \mathbf{A} satisfies $\mathbf{B} = \mathbf{C}\mathbf{A}$ and $\|\mathbf{C}\|_{op} = \max_{\mathbf{u} \in \mathbb{R}^n: \mathbf{A}\mathbf{u} \neq \mathbf{0}} \{\|\mathbf{B}\mathbf{u}\|_2 / \|\mathbf{A}\mathbf{u}\|_2\}$.*

Proof Let r be the rank of \mathbf{A} . We let $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ be the SVD of \mathbf{A} with $\mathbf{U} \in \mathbb{R}^{n \times r}$, $\mathbf{D} \in \mathbb{R}^{r \times r}$ diagonal with positive entries, $\mathbf{V} \in \mathbb{R}^{p \times r}$ and \mathbf{U}, \mathbf{V} both with orthonormal columns. Then $\mathbf{A}^+ = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^\top$ and $\mathbf{C}\mathbf{A} = \mathbf{B}\mathbf{V}\mathbf{V}^\top = \mathbf{B} - \mathbf{B}(\mathbf{I}_p - \mathbf{V}\mathbf{V}^\top)$. Since $\ker \mathbf{A} \subset \ker \mathbf{B}$ and $(\mathbf{I}_p - \mathbf{V}\mathbf{V}^\top)$ is the orthogonal projection onto $\ker \mathbf{A}$ we have $\mathbf{B}(\mathbf{I}_p - \mathbf{V}\mathbf{V}^\top) = \mathbf{0}$. This proves $\mathbf{B} = \mathbf{C}\mathbf{A}$.

For $\|\mathbf{B}\mathbf{A}^+\|_{op}$, for any vector \mathbf{u} , $\|\mathbf{B}\mathbf{u}\| = \|\mathbf{C}\mathbf{A}\mathbf{u}\| \leq \|\mathbf{C}\|_{op}\|\mathbf{A}\mathbf{u}\|$ by definition of $\|\mathbf{C}\|_{op}$. This proves $\|\mathbf{C}\|_{op} \geq M$ for $M = \max_{\mathbf{u} \in \mathbb{R}^n: \mathbf{A}\mathbf{u} \neq \mathbf{0}} \{\|\mathbf{B}\mathbf{u}\|_2 / \|\mathbf{A}\mathbf{u}\|_2\}$. For the inequality $M \geq \|\mathbf{C}\|_{op}$, if $\|\mathbf{C}\|_{op} = \|\mathbf{C}\mathbf{v}\|$ for some unit vector \mathbf{v} then $\mathbf{A}\mathbf{A}^+\mathbf{v} = \mathbf{U}\mathbf{U}^\top\mathbf{v} \neq \mathbf{0}$ since \mathbf{v} is a right singular vector of $\mathbf{C} = \mathbf{B}\mathbf{V}\mathbf{D}\mathbf{U}^\top$ and cannot belong to $\ker(\mathbf{U}^\top)$. Next, $\|\mathbf{C}\mathbf{v}\| = \|\mathbf{B}\mathbf{A}^+\mathbf{v}\| \leq M\|\mathbf{A}\mathbf{A}^+\mathbf{v}\|$ and the conclusion follows since $\|\mathbf{A}\mathbf{A}^+\|_{op} \leq 1$ and $\|\mathbf{v}\| = 1$. \blacksquare

Appendix E. Relaxing strong convexity: Proof of Proposition 12

Consider the notation for $\mathbf{G}, \psi(\varepsilon, \mathbf{G}), \mathbf{h}(\varepsilon, \mathbf{G}), D$ defined around (B.1). Let $\Omega = \{\mathbf{X} \in U_\varepsilon\}$. Let us first rewrite the Lipschitz condition (7.1) using the change of variable $\mathbf{G} = \mathbf{X}\Sigma^{-1/2}$ explained around equation (B.1): the mapping

$$\tilde{\Phi}_\varepsilon \begin{cases} \{M\Sigma^{-1/2}, M \in U_\varepsilon\} \rightarrow \mathbb{R}^{n+p}, \\ \mathbf{G} \mapsto D^{-1}(n^{-1/2}\psi(\varepsilon, \mathbf{G}), \mathbf{h}(\varepsilon, \mathbf{G})) \end{cases} \text{ is } \frac{L}{\sqrt{n}}\text{-Lipschitz}$$

where we recall the notation $D = (\frac{1}{n}\|\psi(\varepsilon, \mathbf{G})\|^2 + \|\mathbf{h}(\varepsilon, \mathbf{G})\|^2)^{1/2}$. After the change of variable, the identities (B.2)-(B.8) all hold in the event Ω .

All previous calculations made in the strongly convex case hold here as well, but only in the event Ω . Outside of event Ω , the derivatives may not exist at all. As in Theorem 6, the device that lets us work around this is Kirszbraun's theorem: there exists an L/\sqrt{n} Lipschitz function $\bar{\Phi}_\varepsilon$ (the "extension") such that $\bar{\Phi}_\varepsilon(\mathbf{G}) = \tilde{\Phi}_\varepsilon(\mathbf{G})$ for all \mathbf{G} in the domain $\{M\Sigma^{-1/2}, M \in U_\varepsilon\}$ of $\tilde{\Phi}_\varepsilon$.

Now define $\rho : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ and $\eta : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^p$ by $\bar{\Phi}(\mathbf{G}) = (\rho(\mathbf{G}), \eta(\mathbf{G}))$. Using the Lipschitz condition and the fact that the Frobenius norm of a Jacobian is bounded from above by its rank times the square of its operator norm,

$$\sum_{ij} \left\| \frac{\partial \rho}{\partial g_{ij}} \right\|^2 \leq L^2, \quad \sum_{ij} \left\| \frac{\partial \eta}{\partial g_{ij}} \right\|^2 \leq \frac{p}{n} L^2$$

so that the right-hand side of (C.6) is bounded above by $C_{54}(L, \gamma)\sqrt{p}$ and the right-hand side of (C.3) is bounded from above by $C_{55}(L, \gamma)$. After we have bounded the right-hand side, we are allowed to add the indicator function $I\{\Omega\}$ in the left-hand sides of (C.6) and of (C.3) since adding an indicator function only makes it smaller. This device lets us obtain analogs of Proposition 14 and Proposition 16 where the right-hand sides are of the same order as in the strongly convex case, provided that we add the indicator function $I\{\Omega\}$ in the left-hand sides.

From here, in the event Ω we use the bound $\|\Sigma^{1/2}\hat{\mathbf{A}}\Sigma^{1/2}\|_{op} \leq L/n$ assumed in Proposition 12 instead of the bound $\|\Sigma^{1/2}\hat{\mathbf{A}}\Sigma^{1/2}\|_{op} \leq 1/(n\mu)$ from Theorem 1. This device provides the bounds

(B.14), (B.15), (B.16), (B.18), (B.19) and (B.20) on ξ_I, \dots, ξ_{V_I} , with the modification that the left-hand sides present the indicator function $I\{\Omega\}$ and the right-hand sides are $C_{56}(\gamma, L)$ instead of $C_{57}(\gamma, \mu)$ in the strongly convex case. The algebra is then the same as in the strongly convex case and Proposition 12 follows.

Appendix F. Additional Figures (anisotropic Gaussian design)

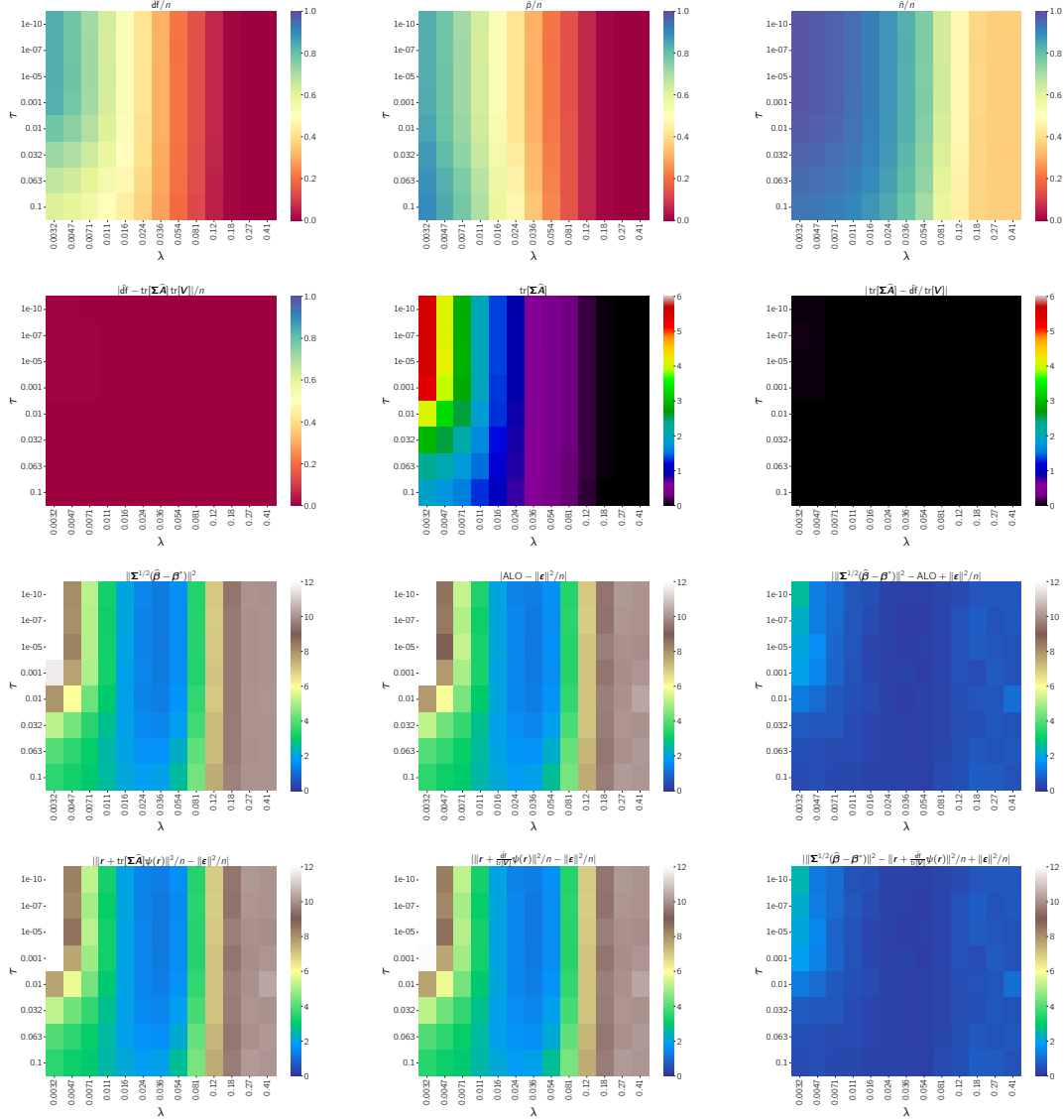


Figure 4: Heatmaps for the Huber loss and Elastic-Net penalty on a grid of tuning parameters with $\Lambda = 0.054n^{1/2}$ and (λ, τ) where $\lambda \in [0.0032, 0.41]$ and $\tau \in [10^{-10}, 0.1]$. Each cell is the average over 100 repetitions. See the simulation setup in Section 6 in the paper for more details.

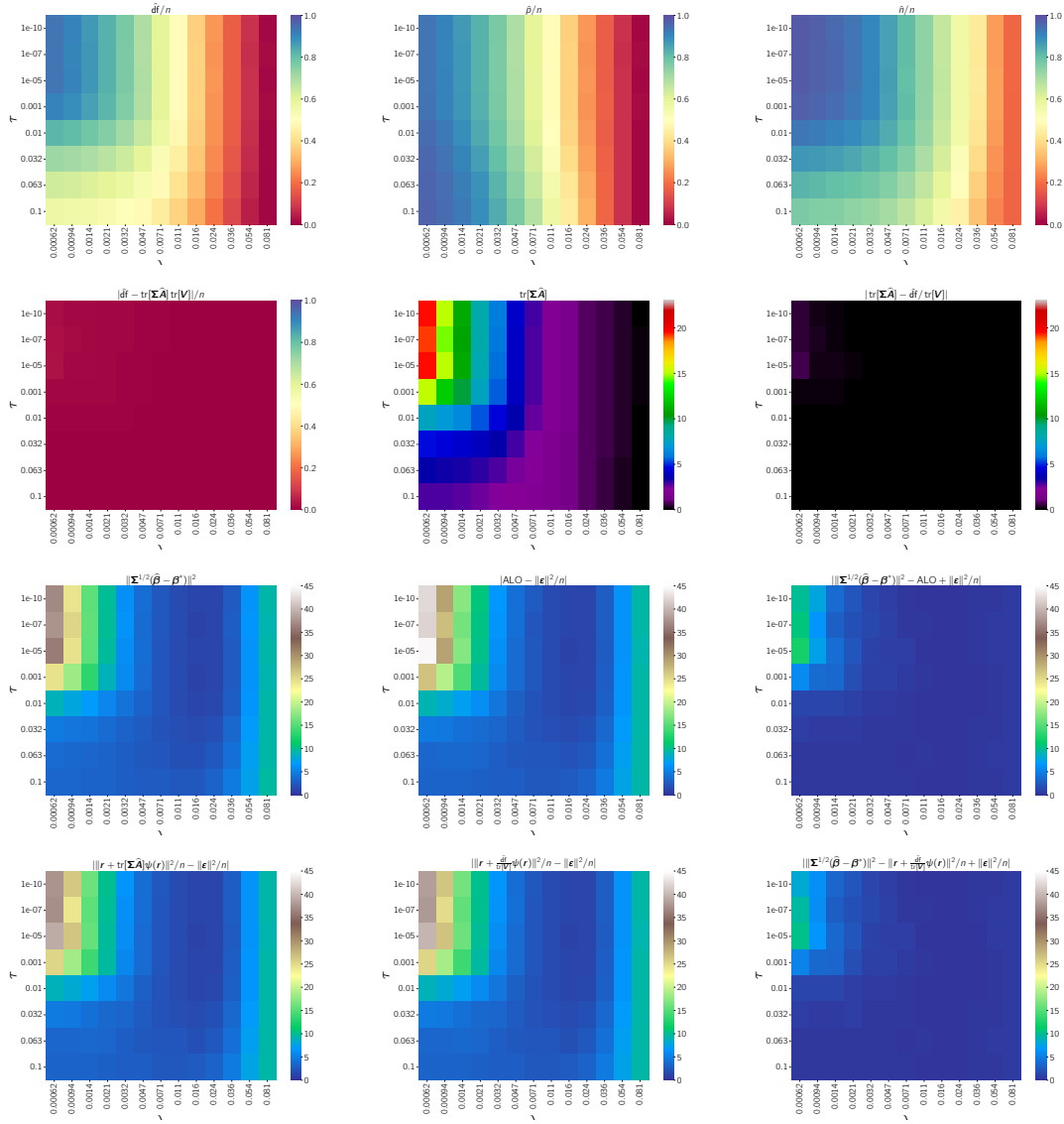


Figure 5: Heatmaps for the Huber loss and Elastic-Net penalty on a grid of tuning parameters with $\Lambda = 0.024n^{1/2}$ and (λ, τ) where $\lambda \in [0.00062, 0.081]$ and $\tau \in [10^{-10}, 0.1]$. Each cell is the average over 50 repetitions. See the simulation setup in Section 6 in the paper for more details.

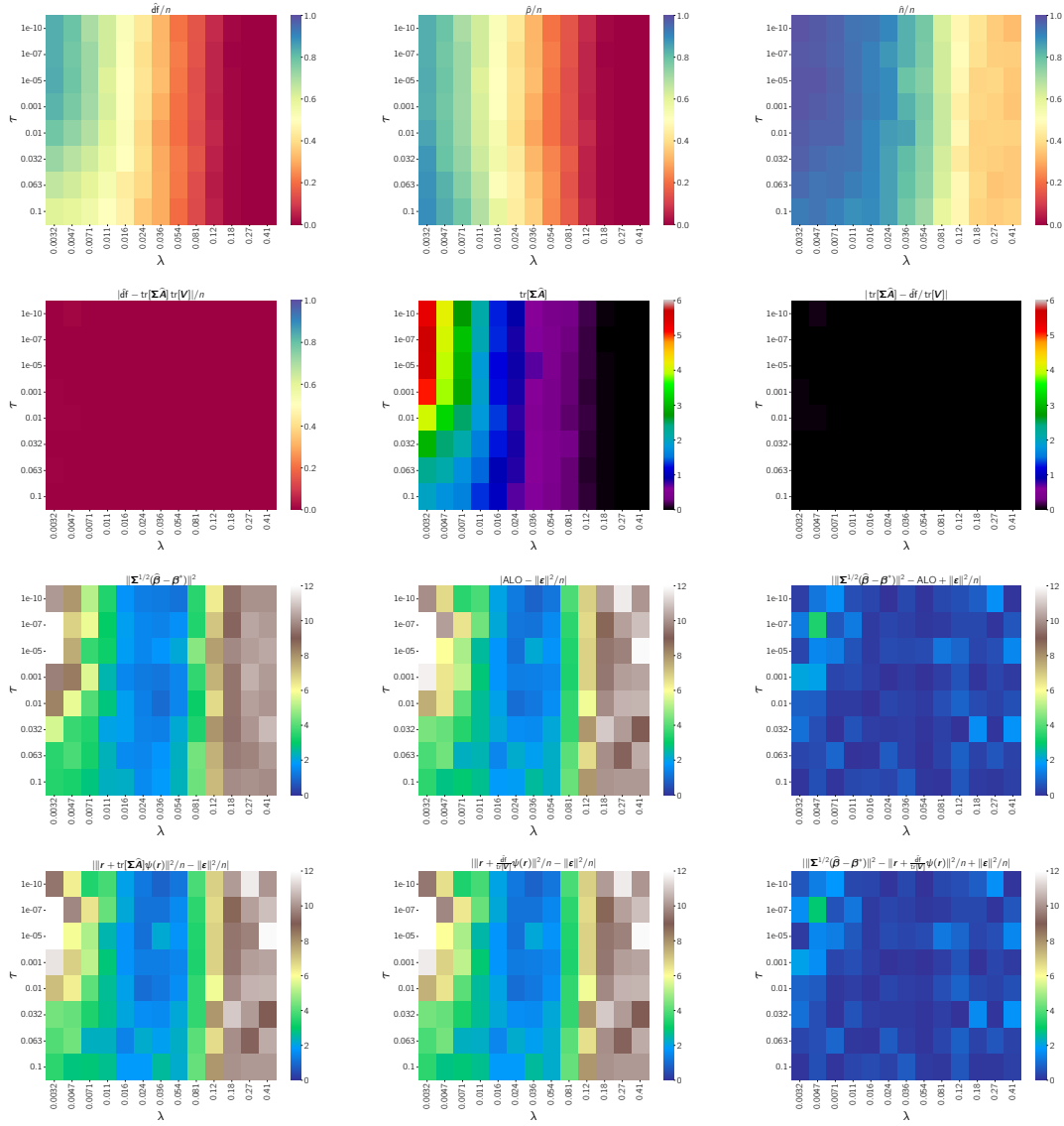


Figure 6: Heatmaps for the Huber loss and Elastic-Net penalty on a grid of tuning parameters with $\Lambda = 0.054n^{1/2}$ and (λ, τ) where $\lambda \in [0.0032, 0.41]$ and $\tau \in [10^{-10}, 0.1]$. Each cell is over 1 repetition. See the simulation setup in Section 6 in the paper for more details.

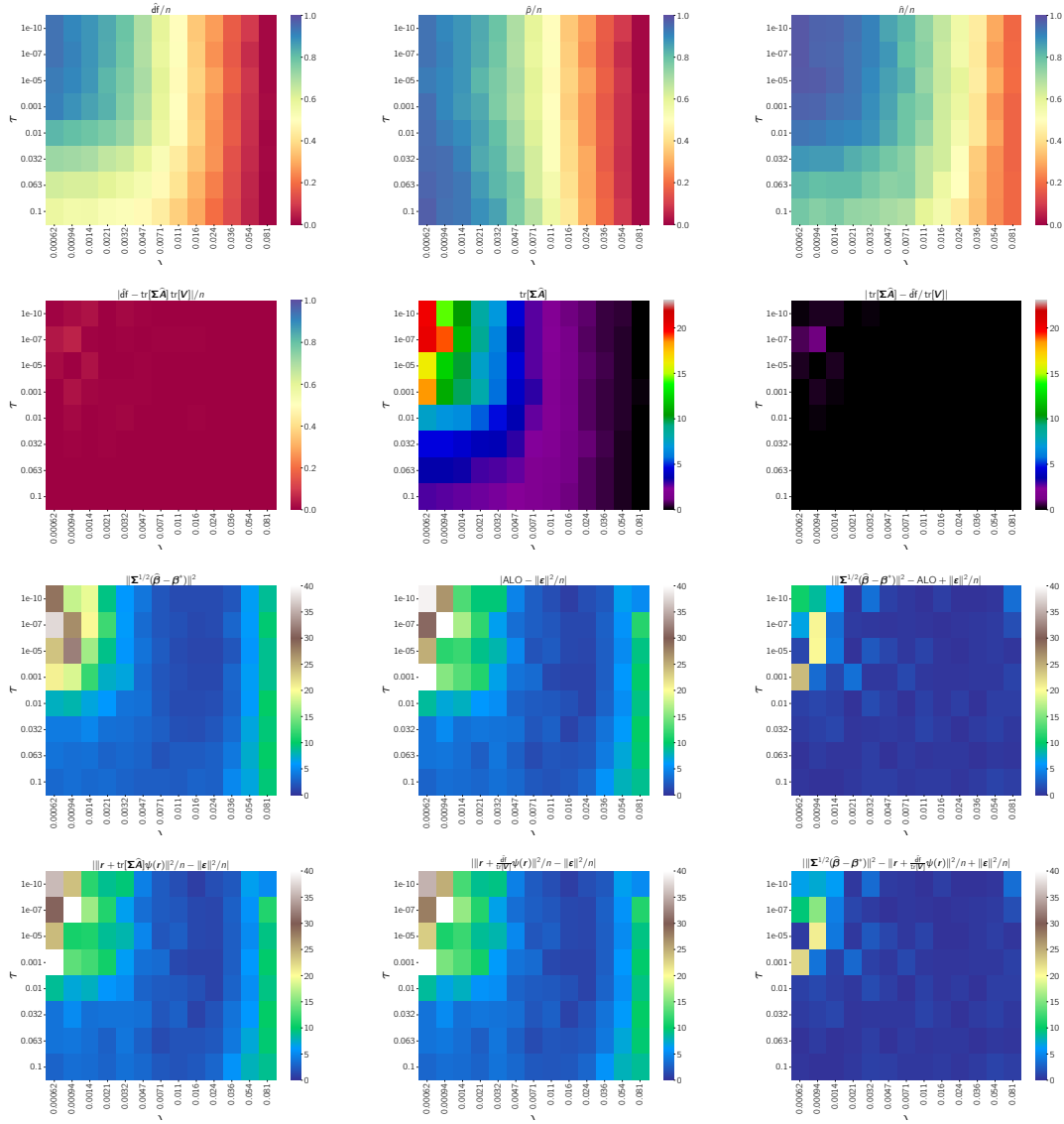


Figure 7: Heatmaps for the Huber loss and Elastic-Net penalty on a grid of tuning parameters with $\Lambda = 0.024n^{1/2}$ and (λ, τ) where $\lambda \in [0.00062, 0.081]$ and $\tau \in [10^{-10}, 0.1]$. Each cell over 1 repetition. See the simulation setup in Section 6 in the paper for more details.

Appendix G. Additional Figures (non-Gaussian, Rademacher design)

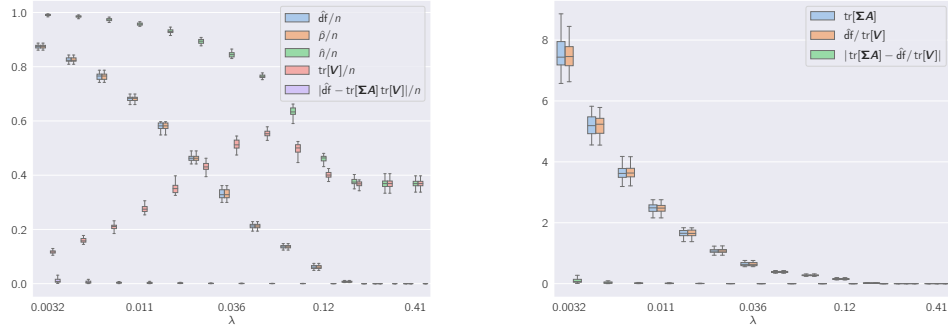


Figure 8: Boxplots for $\hat{d}\hat{f}$, $\hat{\rho}$, \hat{n} , $\text{tr}[\mathbf{V}]$, $\text{tr}[\Sigma\hat{\mathbf{A}}]$ and $|\text{tr}[\Sigma\hat{\mathbf{A}}] - \hat{d}\hat{f}/\text{tr}[\mathbf{V}]|$ in Huber Elastic-Net regression with $\tau = 10^{-10}$ and $\lambda \in [0.0032, 0.41]$. The data are generated with \mathbf{X} having iid entries taking value ± 1 each with probability 0.5 (so that $\Sigma = \mathbf{I}_p$). Each box contains 30 data points.

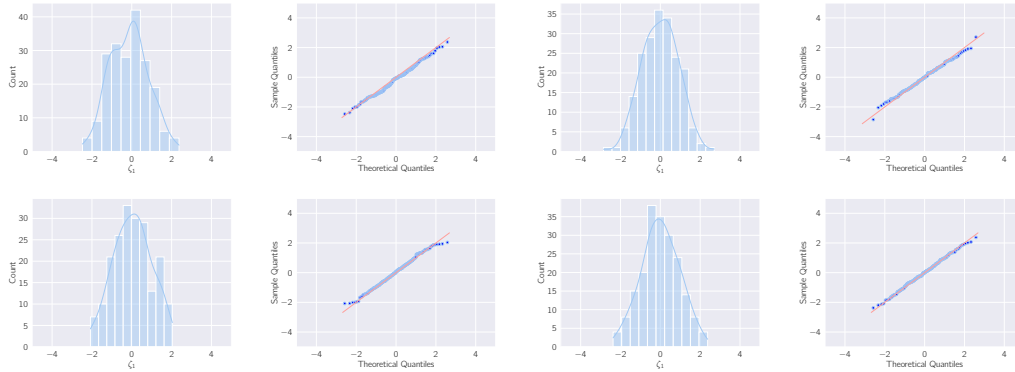


Figure 9: Histogram and QQ-plot for ζ_1 in (3.3) under Huber Elastic-Net regression for different choices of tuning parameters (λ, τ) . Left Top: $(0.036, 10^{-10})$, Right Top: $(0.054, 0.01)$, Left Bottom: $(0.036, 0.01)$, Right Bottom: $(0.024, 0.1)$. Each figure contains 100 data points generated with Rademacher design matrix (each entry has value ± 1 with probability 0.5) and iid ε_i from the t -distribution with 2 degrees of freedom.