

# The likelihood ratio test in high-dimensional logistic regression is asymptotically a *rescaled* Chi-square

Pragya Sur<sup>1</sup> · Yuxin Chen<sup>2</sup> · Emmanuel J. Candès<sup>1,3</sup>

Received: 17 June 2017 / Revised: 15 August 2018 / Published online: 23 January 2019 © Springer-Verlag GmbH Germany, part of Springer Nature 2019

# Abstract

Logistic regression is used thousands of times a day to fit data, predict future outcomes, and assess the statistical significance of explanatory variables. When used for the purpose of statistical inference, logistic models produce *p*-values for the regression coefficients by using an approximation to the distribution of the likelihood-ratio test (LRT). Indeed, Wilks' theorem asserts that whenever we have a fixed number p of variables, twice the log-likelihood ratio (LLR)  $2\Lambda$  is distributed as a  $\chi_k^2$  variable in the limit of large sample sizes n; here,  $\chi_k^2$  is a Chi-square with k degrees of freedom and k the number of variables being tested. In this paper, we prove that when p is not negligible compared to n, Wilks' theorem does not hold and that the Chi-square approximation is grossly incorrect; in fact, this approximation produces p-values that are far too small (under the null hypothesis). Assume that n and p grow large in such a way that  $p/n \to \kappa$  for some constant  $\kappa < 1/2$ . (For  $\kappa > 1/2, 2\Lambda \xrightarrow{\mathbb{P}} 0$  so that the LRT is not interesting in this regime.) We prove that for a class of logistic models, the LLR converges to a *rescaled* Chi-square, namely,  $2\Lambda \xrightarrow{d} \alpha(\kappa)\chi_k^2$ , where the scaling factor  $\alpha(\kappa)$  is greater than one as soon as the dimensionality ratio  $\kappa$  is positive. Hence, the LLR is larger than classically assumed. For instance, when  $\kappa = 0.3$ ,  $\alpha(\kappa) \approx 1.5$ . In general, we show how to compute the scaling factor by solving a nonlinear system of two equations with two unknowns. Our mathematical arguments are involved and use techniques from approximate message passing theory, from non-asymptotic random matrix theory and from convex geometry. We also complement our mathematical study by showing that the new limiting distribution is accurate for finite sample sizes. Finally, all the results from this paper extend to some other regression models such as the probit regression model.

☑ Pragya Sur pragya@stanford.edu

Extended author information available on the last page of the article

**Electronic supplementary material** The online version of this article (https://doi.org/10.1007/s00440-018-00896-9) contains supplementary material, which is available to authorized users.

**Keywords** Logistic regression  $\cdot$  Likelihood-ratio tests  $\cdot$  Wilks' theorem  $\cdot$  High-dimensionality  $\cdot$  Goodness of fit  $\cdot$  Approximate message passing  $\cdot$  Concentration inequalities  $\cdot$  Convex geometry  $\cdot$  Leave-one-out analysis

Mathematics Subject Classification 62Fxx

# 1 Introduction

Logistic regression is by far the most widely used tool for relating a binary response to a family of explanatory variables. This model is used to infer the importance of variables and nearly all standard statistical softwares have inbuilt packages for obtaining *p*-values for assessing the significance of their coefficients. For instance, one can use the snippet of R code below to fit a logistic regression model from a vector y of binary responses and a matrix X of covariates:

```
fitted <- glm(y ~ X+0, family = `binomial')
pvals <- summary(fitted)$coefficients[,4]</pre>
```

The vector pvals stores p-values for testing whether a variable belongs to a model or not, and it is well known that the underlying calculations used to produce these pvalues can also be used to construct confidence intervals for the regression coefficients. Since logistic models are used hundreds of times every day for inference purposes, it is important to know whether these calculations—e.g. these p-values—are accurate and can be trusted.

## 1.1 Binary regression

Imagine we have *n* samples of the form  $(y_i, X_i)$ , where  $y_i \in \{0, 1\}$  and  $X_i \in \mathbb{R}^p$ . In a generalized linear model, one postulates the existence of a link function  $\mu(\cdot)$  relating the conditional mean of the response variable to the linear predictor  $X_i^{\top} \beta$ ,

$$\mathbb{E}[y_i|X_i] = \mu(X_i^{\top}\boldsymbol{\beta}), \tag{1}$$

where  $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_p]^\top \in \mathbb{R}^p$  is an unknown vector of parameters. We focus here on the two most commonly used binary regression models, namely, the logistic and the probit models for which

$$\mu(t) := \begin{cases} e^t / (1 + e^t) & \text{in the logistic model,} \\ \Phi(t) & \text{in the probit model;} \end{cases}$$
(2)

here,  $\Phi$  is the cumulative distribution function (CDF) of a standard normal random variable. In both cases, the *Symmetry Condition* 

$$\mu(t) + \mu(-t) = 1 \tag{3}$$

holds, which says that the two types  $y_i = 0$  and  $y_i = 1$  are treated in a symmetric fashion. Assuming that the observations are independent, the negative log-likelihood function is given by [1, Section 4.1.2]

$$\ell(\boldsymbol{\beta}) := -\sum_{i=1}^{n} \left\{ y_i \log\left(\frac{\mu_i}{1-\mu_i}\right) + \log\left(1-\mu_i\right) \right\}, \qquad \mu_i := \mu(\boldsymbol{X}_i^{\top}\boldsymbol{\beta}).$$

Invoking the symmetry condition, a little algebra reveals an equivalent expression

$$\ell(\boldsymbol{\beta}) := \sum_{i=1}^{n} \rho(-\tilde{y}_i \boldsymbol{X}_i^{\top} \boldsymbol{\beta}), \qquad (4)$$

where

$$\tilde{y}_i := \begin{cases} 1 & \text{if } y_i = 1, \\ -1 & \text{if } y_i = 0, \end{cases} \quad \text{and} \quad \rho(t) := \begin{cases} \log(1 + e^t) & \text{in the logistic case,} \\ -\log \Phi(-t) & \text{in the probit case.} \end{cases}$$
(5)

Throughout we refer to this function  $\rho$  as the *effective link*.

## 1.2 The likelihood-ratio test and Wilks' phenomenon

Researchers often wish to determine which covariates are of importance, or more precisely, to test whether the *j*th variable belongs to the model or not: formally, we wish to test the hypothesis

$$H_i: \quad \beta_i = 0 \quad \text{versus} \quad \beta_i \neq 0.$$
 (6)

Arguably, one of the most commonly deployed techniques for testing  $H_j$  is the likelihood-ratio test (LRT), which is based on the log-likelihood ratio (LLR) statistic

$$\Lambda_j := \ell(\hat{\boldsymbol{\beta}}_{(-j)}) - \ell(\hat{\boldsymbol{\beta}}).$$
(7)

Here,  $\hat{\beta}$  and  $\hat{\beta}_{(-j)}$  denote respectively the maximum likelihood estimates (MLEs) under the full model and the reduced model on dropping the *j*th predictor; that is,

$$\hat{\boldsymbol{\beta}} = \arg\min_{\boldsymbol{\beta}\in\mathbb{R}^p} \ell(\boldsymbol{\beta})$$
 and  $\hat{\boldsymbol{\beta}}_{(-j)} = \arg\min_{\boldsymbol{\beta}\in\mathbb{R}^p, \beta_j=0} \ell(\boldsymbol{\beta}).$ 

Inference based on such log-likelihood ratio statistics has been studied extensively in prior literature [13,44,66]. Arguably, one of the most celebrated results in the large-sample regime is the Wilks' theorem.

To describe the Wilks' phenomenon, imagine we have a sequence of observations  $(y_i, X_i)$  where  $y_i \in \{0, 1\}, X_i \in \mathbb{R}^p$  with p fixed. Since we are interested in the limit of large samples, we may want to assume that the covariates are i.i.d. drawn from some population with non-degenerate covariance matrix so that the problem is fully p-dimensional. As before, we assume a conditional logistic model for the response. In this setting, Wilks' theorem [66] calculates the asymptotic distribution of  $\Lambda_j(n)$  when n grows to infinity:



**Fig. 1** Histogram of *p*-values for logistic regression under i.i.d. Gaussian design, when  $\beta = 0$ , n = 4000, p = 1200, and  $\kappa = 0.3$ : **a** classically computed *p*-values; **b** Bartlett-corrected *p*-values; **c** adjusted *p*-values by comparing the LLR to the rescaled chi square  $\alpha(\kappa)\chi_1^2$  (27)

(Wilks' phenomenon) Under suitable regularity conditions which, for instance, guarantee that the MLE exists and is unique,<sup>1</sup> the LLR statistic for testing  $H_j$ :  $\beta_j = 0$  vs.  $\beta_j \neq 0$  has asymptotic distribution under the null given by

$$2\Lambda_j(n) \xrightarrow{d} \chi_1^2, \quad \text{as } n \to \infty.$$
 (8)

This fixed-*p* large-*n* asymptotic result, which is a consequence of asymptotic normality properties of the MLE [64, Theorem 5.14], applies to a much broader class of testing problems in parametric models; for instance, it applies to the probit model as well. We refer the readers to [41, Chapter 12] and [64, Chapter 16] for a thorough exposition and details on the regularity conditions under which Wilks' theorem holds. Finally, there is a well-known extension which states that if we were to drop *k* variables from the model, then the LLR would converge to a Chi-square distribution with *k* degrees of freedom under the hypothesis that the reduced model is correct.

# 1.3 Inadequacy of Wilks' theorem in high dimensions

The Chi-square approximation to the distribution of the LLR statistic is used in standard statistical softwares to provide *p*-values for the single or multiple coefficient likelihood ratio tests. Here, we perform a simple experiment on synthetic data to study the accuracy of the Chi-square approximation when *p* and *n* are both decently large. Specifically, we set  $\beta = 0$  and test  $\beta_1 = 0$  versus  $\beta_1 \neq 0$  using the LRT in a setting where p = 1200. In each trial, n = 4000 observations are produced with  $y_i \stackrel{\text{i.i.d.}}{\sim}$  Bernoulli(1/2), and  $X := [X_1, \ldots, X_n]^\top \in \mathbb{R}^{n \times p}$  is obtained by generating a random matrix composed of i.i.d.  $\mathcal{N}(0, 1)$  entries. We fit a logistic regression of *y* on *X* using R, and extract the *p*-values for each coefficient. Figure 1 plots the pooled histogram that aggregates  $4.8 \times 10^5$  *p*-values in total (400 trials with 1200 *p*-values obtained in each trial).

<sup>&</sup>lt;sup>1</sup> Such conditions would also typically imply asymptotic normality of the MLE.

If the  $\chi_1^2$  approximation were true, then we would expect to observe uniformly distributed *p*-values. The histrogram from Fig. 1 is, however, far from uniform. This is an indication of the inadequacy of Wilks' theorem when *p* and *n* are both large. The same issue was also reported in [12], where the authors observed that this discrepancy is highly problematic since the distribution is skewed towards smaller values. Hence, such *p*-values cannot be trusted to construct level- $\alpha$  tests and the problem is increasingly severe when we turn attention to smaller *p*-values as in large-scale multiple testing applications.

#### 1.4 The Bartlett correction?

A natural question that arises immediately is whether the observed discrepancy could be an outcome of a finite-sample effect. It has been repeatedly observed that the Chisquare approximation does not yield accurate results with finite sample size. One correction to the LRT that is widely used in finite samples is the Bartlett correction, which dates back to Bartlett [5] and has been extensively studied over the past few decades (e.g. [9,11,14,16,40]). In the context of testing for a single coefficient in the logistic model, this correction can be described as follows [45]: compute the expectation of the LLR statistic up to terms of order 1/n; that is, compute a parameter  $\alpha$  such that

$$\mathbb{E}[2\Lambda_j] = 1 + \frac{\alpha}{n} + O\left(\frac{1}{n^2}\right),$$

which suggests a corrected LLR statistic

$$\frac{2\Lambda_j}{1+\frac{\alpha_n}{n}}\tag{9}$$

with  $\alpha_n$  being an estimator of  $\alpha$ . With a proper choice of  $\alpha_n$ , one can ensure

$$\mathbb{E}\left[\frac{2\Lambda_j}{1+\frac{\alpha_n}{n}}\right] = 1 + O\left(\frac{1}{n^2}\right)$$

in the classical setting where *p* is fixed and *n* diverges. In expectation, this corrected statistic is closer to a  $\chi_1^2$  distribution than the original LLR for finite samples. Notably, the correction factor may in general be a function of the unknown  $\beta$  and, in that case, must be estimated from the null model via maximum likelihood estimation.

In the context of GLMs, Cordeiro [14] derived a general formula for the Bartlett corrected LLR statistic, see [15,20] for a detailed survey. In the case where there is no signal ( $\beta = 0$ ), one can compute  $\alpha_n$  for the logistic regression model following [14,45], which yields

$$\alpha_n = \frac{n}{2} \left[ \operatorname{Tr} \left( \boldsymbol{D}_p^2 \right) - \operatorname{Tr} \left( \boldsymbol{D}_{p-1}^2 \right) \right].$$
(10)

🖄 Springer

Here,  $D_p$  is the diagonal part of  $X(X^{\top}X)^{-1}X^{\top}$  and  $D_{p-1}$  is that of  $X_{(-j)}(X_{(-j)}^{\top}X_{(-j)})^{-1}X_{(-j)}^{\top}$  in which  $X_{(-j)}$  is the design matrix X with the *j*th column removed. Comparing the adjusted LLRs to a  $\chi_1^2$  distribution yields adjusted *p*-values. In the setting of Fig. 1a, the histogram of Bartlett corrected *p*-values is shown in Fig. 1b. As we see, these *p*-values are still far from uniform.

If the mismatch is not due to finite sample-size effects, what is the distribution of the LLR in high dimensions? Our main contribution is to provide a very precise answer to this question; below, we derive the high-dimensional asymptotic distribution of the log-likelihood ratios, i.e. in situations where the dimension p is not necessarily negligible compared to the sample size n.

Under suitable assumptions, we establish in Theorem 1 that the asymptotic distribution of the LLR is a rescaled  $\chi_1^2$  and identify the rescaling factor precisely. Figure 1c shows the histogram of the computed *p*-values by comparing the LLR to this distribution. Clearly, these adjusted *p*-values are much closer to uniform random variables.

# 2 Main results

#### 2.1 Modelling assumptions

In this paper, we focus on the high-dimensional regime where the sample size is not much larger than the number of parameters to be estimated—a setting which has attracted a flurry of activity in recent years. In particular, we assume that the number p(n) of covariates grows proportionally with the number n of observations; that is,

$$\lim_{n \to \infty} \frac{p(n)}{n} = \kappa, \tag{11}$$

where  $\kappa > 0$  is a fixed constant independent of *n* and *p*(*n*). In fact, we shall also assume  $\kappa < 1/2$  for both the logistic and the probit models, as the MLE does not exist otherwise; see Sect. 2.2.

To formalize the notion of high-dimensional asymptotics when both *n* and p(n) diverge, we consider a sequence of instances  $\{X(n), y(n)\}_{n>0}$  such that for any *n*,

- $X(n) \in \mathbb{R}^{n \times p(n)}$  has i.i.d. rows  $X_i(n) \sim \mathcal{N}(0, \Sigma)$ , where  $\Sigma \in \mathbb{R}^{p(n) \times p(n)}$  is positive definite;
- $y_i(n) | X(n) \sim y_i(n) | X_i(n) \stackrel{\text{ind.}}{\sim} \text{Bernoulli} (\mu(X_i(n)^\top \beta(n)))$ , where  $\mu$  satisfies the Symmetry Condition;
- we further assume  $\beta(n) = 0$ . From the Symmetry Condition it follows that  $\mu(0) = 1/2$ , which directly implies that y(n) is a vector with i.i.d Bernoulli(1/2) entries.

The MLE is denoted by  $\beta(n)$  and there are p(n) LLR statistics  $\Lambda_j(n)$   $(1 \le j \le p(n))$ , one for each of the p(n) regression coefficients. In the sequel, the dependency on n shall be suppressed whenever it is clear from the context.

## 2.2 When does the MLE exist?

Even though we are operating in the regime where n > p, the existence of the MLE cannot be guaranteed for all p and n. Interestingly, the norm of the MLE undergoes a sharp phase transition in the sense that with high probability,

$$\|\hat{\boldsymbol{\beta}}\| = \infty \qquad \text{if } \kappa > 1/2; \tag{12}$$

$$\|\boldsymbol{\Sigma}^{1/2}\hat{\boldsymbol{\beta}}\| = O(1) \quad \text{if } \kappa < 1/2.$$
(13)

The first result (12) concerns the *separating capacity* of linear inequalities, which dates back to Cover's Ph. D. thesis [17]. Specifically, given that  $\rho(t) \ge \rho(-\infty) = 0$  for both the logistic and probit models, each summand in (4) is minimized if  $\tilde{y}_i X_i^{\top} \beta = \infty$ , which occurs when  $\operatorname{sign}(X_i^{\top}\beta) = \operatorname{sign}(\tilde{y}_i)$  and  $\|\beta\| = \infty$ . As a result, if there exists a nontrivial ray  $\beta$  such that

$$X_i^{\top} \boldsymbol{\beta} > 0 \quad \text{if } \tilde{y}_i = 1 \quad \text{and} \quad X_i^{\top} \boldsymbol{\beta} < 0 \quad \text{if } \tilde{y}_i = -1$$
 (14)

for any  $1 \le i \le n$ , then pushing  $\|\boldsymbol{\beta}\|$  to infinity leads to an optimizer of (4). In other words, the solution to (4) becomes unbounded (the MLE is at  $\infty$ ) whenever there is a hyperplane perfectly separating the two sets of samples  $\{i \mid \tilde{y}_i = 1\}$  and  $\{i \mid \tilde{y}_i = -1\}$ . According to [17,18], the probability of separability tends to one when  $\kappa > 1/2$ , in which case the MLE does not exist. This separability result was originally derived using classical combinatorial geometry.

The current paper follows another route by resorting to convex geometry, which, as we will demonstrate later, also tells us how to control the norm of the MLE when  $\kappa < 1/2$ .<sup>2</sup> To begin with, we observe that  $\tilde{y}_i$  is independent of X and the distribution of X is symmetric under the assumptions from Sect. 2.1. Hence, to calculate the chance that there exists a separating hyperplane, we can assume  $\tilde{y}_i = 1$  ( $1 \le i \le n$ ) without loss of generality. In this case, the event (14) becomes

$$\left\{ \boldsymbol{X}\boldsymbol{\beta} \mid \boldsymbol{\beta} \in \mathbb{R}^p \right\} \cap \mathbb{R}^n_{++} \neq \emptyset, \tag{15}$$

where  $\mathbb{R}^{n}_{++}$  is the positive orthant. Write  $X = Z\Sigma^{1/2}$  so that Z is an  $n \times p$  matrix with i.i.d. standard Gaussian entries, and  $\theta = \Sigma^{1/2}\beta$ . Then the event (15) is equivalent to

$$\left\{ \boldsymbol{Z}\boldsymbol{\theta} \mid \boldsymbol{\theta} \in \mathbb{R}^p \right\} \cap \mathbb{R}^n_{++} \neq \emptyset.$$
(16)

Now the probability that (16) occurs is the same as that

$$\left\{ \boldsymbol{Z}\boldsymbol{\theta} \mid \boldsymbol{\theta} \in \mathbb{R}^p \right\} \cap \mathbb{R}^n_+ \neq \{ \boldsymbol{0} \}$$
(17)

occurs, where  $\mathbb{R}^n_+$  denotes the non-negative orthant.

<sup>&</sup>lt;sup>2</sup> The separability results in [17,18] do not imply the control on the norm of the MLE when  $\kappa < 1/2$ .

From the approximate kinematic formula [3, Theorem I] in the literature on convex geometry, the event (17) happens with high probability if and only if the total statistical dimension of the two closed convex cones exceeds the ambient dimension, i.e.

$$\delta\left(\left\{\boldsymbol{Z\boldsymbol{\theta}} \mid \boldsymbol{\theta} \in \mathbb{R}^p\right\}\right) + \delta\left(\mathbb{R}^n_+\right) > n + o(n).$$
(18)

Here, the statistical dimension of a closed convex cone  $\mathcal{K}$  is defined as

$$\delta(\mathcal{K}) := \mathbb{E}_{\boldsymbol{g} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})} \left[ \| \Pi_{\mathcal{K}} \left( \boldsymbol{g} \right) \|^2 \right]$$
(19)

with  $\Pi_{\mathcal{K}}(g) := \arg \min_{z \in \mathcal{K}} \|g - z\|$  the Euclidean projection. Recognizing that [3, Proposition 2.4]

$$\delta\left(\left\{\mathbf{Z}\boldsymbol{\theta} \mid \boldsymbol{\theta} \in \mathbb{R}^p\right\}\right) = p \text{ and } \delta(\mathbb{R}^n_+) = n/2,$$

we reduce the condition (18) to

$$p + n/2 > n + o(n)$$
 or  $p/n > 1/2 + o(1)$ ,

thus indicating that  $\|\hat{\beta}\| = \infty$  with high probability when  $\kappa = \lim p/n > 1/2$ . (Hence, in this regime the LLR converges in probability to 0.)

The preceding argument only reveals that the MLE does not exist with high probability when  $\kappa > 1/2$ . If  $\kappa = p/n < 1/2$ , we establish precise control on the norm of the MLE, properly scaled. In fact, in Theorem 4 we prove that in this regime,  $\|\mathbf{\Sigma}^{1/2}\hat{\boldsymbol{\beta}}\| = O(1)$  with high probability. In light of these observations we work with the additional condition

$$\kappa < 1/2. \tag{20}$$

#### 2.3 The high-dimensional limiting distribution of the LLR

In contrast to the classical Wilks' result, our findings reveal that the LLR statistic follows a *rescaled* Chi-square distribution with a rescaling factor that can be explicitly pinned down through the solution to a system of equations.

## 2.3.1 A system of equations

We start by setting up the crucial system of equations. Before proceeding, we first recall the proximal operator

$$\operatorname{prox}_{b\rho}(z) := \arg\min_{x \in \mathbb{R}} \left\{ b\rho(x) + \frac{1}{2}(x-z)^2 \right\}$$
(21)

defined for any b > 0 and convex function  $\rho(\cdot)$ . As in [22], we introduce the operator

$$\Psi(z;b) := b\rho'(\operatorname{prox}_{b\rho}(z)), \tag{22}$$

which is simply the proximal operator of the conjugate  $(b\rho)^*$  of  $b\rho$ .<sup>3</sup> To see this, we note that  $\Psi$  satisfies the relation [22, Proposition 6.4]

$$\Psi(z;b) + \operatorname{prox}_{bo}(z) = z. \tag{23}$$

The claim that  $\Psi(\cdot; b) = \text{prox}_{(b\rho)^*}(\cdot)$  then follows from the Moreau decomposition

$$\operatorname{prox}_{f}(z) + \operatorname{prox}_{f^{*}}(z) = z, \quad \forall z,$$
(24)

which holds for a closed convex function f [47, Section 2.5]. Interested readers are referred to [22, Appendix 1] for more properties of prox<sub>bo</sub> and  $\Psi$ .

We are now in position to present the system of equations that plays a crucial role in determining the distribution of the LLR statistic in high dimensions:

$$\tau^{2} = \frac{1}{\kappa} \mathbb{E}\left[ \left( \Psi \left( \tau Z; \ b \right) \right)^{2} \right], \tag{25}$$

$$\kappa = \mathbb{E}\left[\Psi'\left(\tau Z; b\right)\right],\tag{26}$$

where  $Z \sim \mathcal{N}(0, 1)$ , and  $\Psi'(\cdot, \cdot)$  denotes differentiation with respect to the first variable. The fact that this system of equations would admit a unique solution in  $\mathbb{R}^2_+$  is not obvious *a priori*. We shall establish this for the logistic and the probit models later in Sect. 6.

# 2.3.2 Main result

**Theorem 1** Consider a logistic or probit regression model under the assumptions from Sect. 2.1. If  $\kappa \in (0, 1/2)$ , then for any  $1 \le j \le p$ , the log-likelihood ratio statistic  $\Lambda_j$  as defined in (7) obeys

$$2\Lambda_j \stackrel{\mathrm{d}}{\to} \alpha(\kappa) \chi_1^2, \quad \alpha(\kappa) = \tau_*^2/b_*, \quad as \ n \to \infty,$$
 (27)

where  $(\tau_*, b_*) \in \mathbb{R}^2_+$  is the unique solution to the system of Eqs. (25) and (26). Furthermore, the LLR statistic obtained by dropping k variables for any fixed k converges to  $\alpha(\kappa) \chi_k^2$ . Finally, these results extend to all binary regression models with links obeying the assumptions listed in Sect. 2.3.3.

Hence, the limiting distribution is a rescaled Chi-square with a rescaling factor  $\alpha(\kappa)$  that only depends on the aspect ratio  $\kappa$ . Figure 2 illustrates the dependence of the rescaling factor on the limiting aspect ratio  $\kappa$  for logistic regression. The figures for the probit model are similar as the rescaling constants actually differ by very small values.

To study the quality of approximation for finite samples, we repeat the same numerical experiments as before but now obtain the p-values by comparing the LLR statistic with the rescaled Chi-square suggested by Theorem 1. For a particular run of the

<sup>&</sup>lt;sup>3</sup> The conjugate  $f^*$  of a function f is defined as  $f^*(x) = \sup_{u \in \text{dom}(f)} \{ \langle u, x \rangle - f(u) \}.$ 



**Fig. 2** Rescaling constant  $\alpha(\kappa)$  as a function of  $\kappa$  for the logistic model. Note the logarithmic scale in the right panel. The curves for the probit model are nearly identical

experiment (n = 4000, p = 1200,  $\kappa = 0.3$ ), we compute the adjusted LLR statistic  $2\Lambda_j/\alpha(\kappa)$  for each coefficient and obtain the *p*-values based on the  $\chi_1^2$  distribution. The pooled histogram that aggregates  $4.8 \times 10^5 p$ -values in total is shown in Fig. 1c.

As we clearly see, the *p*-values are much closer to a uniform distribution now. One can compute the Chi-square goodness of fit statistic to test the closeness of the above distribution to uniformity. To this end, we divide the interval [0, 1] into 20 equally spaced bins of width 0.05 each. For each bin we compute the observed number of times a *p*-value falls in the bin out of the  $4.8 \times 10^5$  values. Then a Chi-square goodness of fit statistic is computed, noting that the expected frequency is 24000 for each bin. The Chi-square statistic in this case is 16.049, which gives a *p*-value of 0.654 in comparison with a  $\chi^2_{19}$  variable. The same test when performed with the Bartlett corrected *p*-values (Fig. 1b) yields a Chi-square statistic 5599 with a *p*-value of 0.<sup>4</sup> Thus, our correction gives the desired uniformity in the *p*-values when the true signal  $\beta = 0$ .

Practitioners would be concerned about the validity of *p*-values when they are small—again, think about multiple testing applications. In order to study whether our correction yields valid results for small *p*-values, we compute the proportion of times the *p*-values (in all the three cases) lie below 5%, 1%, 0.5%, 0.1%, 0.05%, 0.01% out of the  $4.8 \times 10^5$  times. The results are summarized in Table 1. This further illustrates the deviation from uniformity for the classical and Bartlett corrected *p*-values, whereas the "adjusted" *p*-values obtained invoking Theorem 1 are still valid.

Last but not least, it is seen from Fig. 2 that the rescaling factor  $\alpha(\kappa) \rightarrow 1$  as  $\kappa \rightarrow 0$ . This reveals an important message: *the classical Wilks phenomenon remains valid as long as p/n*  $\rightarrow 0$ . This fact arises even though the classical asymptotic normality of the MLE may fail to hold [30], see Sect. 2.7 for a more detailed discussion.

<sup>&</sup>lt;sup>4</sup> Note that the *p*-values obtained at each trial are not exactly independent. However, they are exchangeable, and weakly dependent (see the proof of Corollary 1 for a formal justification of this fact). Therefore, we expect the goodness of fit test to be an approximately valid procedure in this setting.

	Classical	Bartlett-corrected	Adjusted
$\mathbb{P}\{p\text{-values} \le 5\%\}$	11.1044% (0.0668%)	6.9592% (0.0534%)	5.0110% (0.0453%)
$\mathbb{P}\{p\text{-values} \le 1\%\}$	3.6383% (0.038%)	1.6975% (0.0261%)	0.9944% (0.0186%)
$\mathbb{P}\{p\text{-values} \le 0.5\%\}$	2.2477% (0.0292%)	0.9242% (0.0178%)	0.4952% (0.0116%)
$\mathbb{P}\{p\text{-values} \le 0.1\%\}$	0.7519% (0.0155%)	0.2306% (0.0078%)	0.1008% (0.0051%)
$\mathbb{P}\{p\text{-values} \le 0.05\%\}$	0.4669% (0.0112%)	0.124% (0.0056%)	0.0542% (0.0036%)
$\mathbb{P}\{p\text{-values} \le 0.01\%\}$	0.1575% (0.0064%)	0.0342% (0.0027%)	0.0104% (0.0014%)

 Table 1
 Estimates of p-value probabilities with estimated Monte Carlo standard errors in parentheses under

 i.i.d. Gaussian design
 1

# 2.3.3 Extensions

As noted in Sect. 1.1, the Symmetry Condition (3) allows to express the negative log-likelihood in the form (4), which makes use of the effective link  $\rho(\cdot)$ . Theorem 1 applies to any  $\rho(\cdot)$  obeying the following properties:

- 1.  $\rho$  is non-negative, has up to three derivatives, and obeys  $\rho(t) \ge t$ .
- 2.  $\rho'$  may be unbounded but it should grow sufficiently slowly, in particular, we assume  $|\rho'(t)| = O(|t|)$  and  $\rho'(\operatorname{prox}_{c\rho}(Z))$  is a sub-Gaussian random variable for any constant c > 0 and any  $Z \sim \mathcal{N}(0, \sigma^2)$  for some finite  $\sigma > 0$ .
- 3.  $\rho''(t) > 0$  for any t which implies that  $\rho$  is convex, and  $\sup_t \rho''(t) < \infty$ .
- 4.  $\sup_{t} |\rho'''(t)| < \infty$ .
- 5. Given any  $\tau > 0$ , the Eq. (26) has a unique solution in *b*.
- 6. The map  $\mathcal{V}(\tau^2)$  as defined in (59) has a fixed point.

It can be checked that the effective links for both the logistic and the probit models (5) obey *all* of the above. The last two conditions are assumed to ensure existence of a unique solution to the system of Eqs. (25) and (26) as will be seen in Sect. 6; we shall justify these two conditions for the logistic and the probit models in Sect. 6.1.

## 2.4 Reduction to independent covariates

In order to derive the asymptotic distribution of the LLR statistics, it in fact suffices to consider the special case  $\Sigma = I_p$ .

**Lemma 1** Let  $\Lambda_j(X)$  be the LLR statistic based on the design matrix X, where the rows of X are i.i.d.  $\mathcal{N}(\mathbf{0}, \Sigma)$  and  $\Lambda_j(\mathbf{Z})$  that where the rows are i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . Then

$$\Lambda_j(X) \stackrel{\mathrm{d}}{=} \Lambda_j(Z).$$

**Proof** Recall from (4) that the LLR statistic for testing the jth coefficient can be expressed as

$$\Lambda_j(X) = \min_{\boldsymbol{\beta}} \sum_{i=1}^n \rho(-\tilde{y}_i \boldsymbol{e}_i^\top X \boldsymbol{\beta}) - \min_{\boldsymbol{\beta}: \beta_j = 0} \sum_{i=1}^n \rho(-\tilde{y}_i \boldsymbol{e}_i^\top X \boldsymbol{\beta}).$$

Write  $\mathbf{Z}' = \mathbf{X} \mathbf{\Sigma}^{-1/2}$  so that the rows of  $\mathbf{Z}'$  are i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  and set  $\boldsymbol{\theta}' = \mathbf{\Sigma}^{1/2} \boldsymbol{\beta}$ . With this reparameterization, we observe that the constraint  $\beta_j = 0$  is equivalent to  $\mathbf{a}_j^{\top} \boldsymbol{\theta}' = 0$  for some non-zero vector  $\mathbf{a}_j \in \mathbb{R}^p$ . This gives

$$\Lambda_j(X) = \min_{\boldsymbol{\theta}'} \sum_{i=1}^n \rho(-\tilde{y}_i \boldsymbol{e}_i^\top \boldsymbol{Z}' \boldsymbol{\theta}') - \min_{\boldsymbol{\theta}': \boldsymbol{a}_j^\top \boldsymbol{\theta}' = 0} \sum_{i=1}^n \rho(-\tilde{y}_i \boldsymbol{e}_i^\top \boldsymbol{Z}' \boldsymbol{\theta}').$$

Now let Q be an orthogonal matrix mapping  $a_j \in \mathbb{R}^p$  into the vector  $||a_j||e_j \in \mathbb{R}^p$ , i.e.  $Qa_j = ||a_j||e_j$ . Additionally, set Z = Z'Q (the rows of Z are still i.i.d.  $\mathcal{N}(\mathbf{0}, I_p)$ ) and  $\theta = Q\theta'$ . Since  $a_j^\top \theta' = 0$  occurs if and only if  $\theta_j = 0$ , we obtain

$$\Lambda_j(\boldsymbol{X}) = \min_{\boldsymbol{\theta}} \sum_{i=1}^n \rho(-\tilde{y}_i \boldsymbol{e}_i^\top \boldsymbol{Z} \boldsymbol{\theta}) - \min_{\boldsymbol{\theta}: \boldsymbol{\theta}_j = 0} \sum_{i=1}^n \rho(-\tilde{y}_i \boldsymbol{e}_i^\top \boldsymbol{Z} \boldsymbol{\theta}) = \Lambda_j(\boldsymbol{Z}),$$

which proves the lemma.

In the remainder of the paper we, therefore, assume  $\Sigma = I_p$ .

## 2.5 Proof architecture

This section presents the main steps for proving Theorem 1. We will only prove the theorem for  $\{\Lambda_i\}$ , the LLR statistic obtained by dropping a single variable. The analysis for the LLR statistic obtained on dropping *k* variables (for some fixed *k*) follows very similar steps and is hence omitted for the sake of conciseness. As discussed before, we are free to work with any configuration of the  $y_i$ 's. For the two steps below, we will adopt two different configurations for convenience of presentation.

# 2.5.1 Step 1: Characterizing the asymptotic distributions of $\hat{\beta}_i$

Without loss of generality, we assume here that  $y_i = 1$  (and hence  $\tilde{y}_i = 1$ ) for all  $1 \le i \le n$  and, therefore, the MLE problem reduces to

We would first like to characterize the marginal distribution of  $\hat{\beta}$ , which is crucial in understanding the LLR statistic. To this end, our analysis follows by a reduction to the setup of [22,24–26], with certain modifications that are called for due to the specific choices of  $\rho(\cdot)$  we deal with here. Specifically, consider the linear model

$$y = X\boldsymbol{\beta} + \boldsymbol{w},\tag{28}$$

and prior work [22,24-26] investigating the associated M-estimator

Our problem reduces to (29) on setting y = w = 0 in (29). When  $\rho(\cdot)$  satisfies certain assumptions (e.g. strong convexity), the asymptotic distribution of  $\|\hat{\beta}\|$  has been studied in a series of works [24–26] using a leave-one-out analysis and independently in [22] using approximate message passing (AMP) machinery. An outline of their main results is described in Sect. 2.7. However, the function  $\rho(\cdot)$  in our cases has vanishing curvature and, therefore, lacks the essential strong convexity assumption that was utilized in both the aforementioned lines of work. To circumvent this issue, we propose to invoke the AMP machinery as in [22], in conjunction with the following critical additional ingredients:

• (*Norm bound condition*) We utilize results from the conic geometry literature (e.g. [3]) to establish that

$$\|\hat{\boldsymbol{\beta}}\| = O(1)$$

with high probability as long as  $\kappa < 1/2$ . This will be elaborated in Theorem 4.

- (*Likelihood curvature condition*) We establish some regularity conditions on the Hessian of the log-likelihood function, generalizing the strong convexity condition, which will be detailed in Lemma 4.
- (*Uniqueness of the solution to* (25) *and* (26)) We establish that for both the logistic and the probit case, the system of Eqs. (25) and (26) admits a unique solution.

We emphasize that these elements are not straightforward, require significant effort and a number of novel ideas, which form our primary technical contributions for this step.

These ingredients enable the use of the AMP machinery even in the absence of strong convexity on  $\rho(\cdot)$ , finally leading to the following theorem:

**Theorem 2** Under the conditions of Theorem 1,

$$\lim_{n \to \infty} \|\hat{\boldsymbol{\beta}}\|^2 =_{\text{a.s.}} \tau_*^2.$$
(30)

This theorem immediately implies that the marginal distribution of  $\hat{\beta}_i$  is normal.

**Corollary 1** Under the conditions of Theorem 1, for every  $1 \le j \le p$ , it holds that

$$\sqrt{p}\hat{\beta}_j \stackrel{\mathrm{d}}{\to} \mathcal{N}(0, \tau_*^2), \quad as \ n \to \infty.$$
 (31)

**Proof** From the rotational invariance of our i.i.d. Gaussian design, it can be easily verified that  $\hat{\boldsymbol{\beta}}/\|\hat{\boldsymbol{\beta}}\|$  is uniformly distributed on the unit sphere  $\mathbb{S}^{p-1}$  and is independent of  $\|\hat{\boldsymbol{\beta}}\|$ . Therefore,  $\hat{\beta}_j$  has the same distribution as  $\|\hat{\boldsymbol{\beta}}\|Z_j/\|Z\|$ , where

 $\mathbf{Z} = (Z_1, \ldots, Z_p) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  independent of  $\|\hat{\boldsymbol{\beta}}\|$ . Since  $\sqrt{p} \|\hat{\boldsymbol{\beta}}\| / \|\mathbf{Z}\|$  converges in probability to  $\tau_*$ , we have, by Slutsky's theorem, that  $\sqrt{p}\hat{\beta}_j$  converges to  $\mathcal{N}(0, \tau_*^2)$  in distribution.

# 2.5.2 Step 2: Connecting $\Lambda_i$ with $\hat{\beta}_i$

Now that we have derived the asymptotic distribution of  $\hat{\beta}_j$ , the next step involves a reduction of the LLR statistic to a function of the relevant coordinate of the MLE. Before continuing, we note that the distribution of  $\Lambda_j$  is the same for all  $1 \le j \le p$ due to exchangeability. As a result, going forward we will only analyze  $\Lambda_1$  without loss of generality. In addition, we introduce the following convenient notations and assumptions:

- the design matrix on dropping the first column is written as X
   *X* and the MLE in the corresponding reduced model as β
   *β*;
- write  $X = [X_1, \ldots, X_n]^\top \in \mathbb{R}^{n \times p}$  and  $\tilde{X} = [\tilde{X}_1, \ldots, \tilde{X}_n]^\top \in \mathbb{R}^{n \times (p-1)}$ ;
- without loss of generality, assume that  $\tilde{y}_i = -1$  for all *i* in this subsection, and hence the MLEs under the full and the reduced models reduce to

$$\hat{\boldsymbol{\beta}} = \arg\min_{\boldsymbol{\beta}\in\mathbb{R}^p} \ell(\boldsymbol{\beta}) := \sum_{i=1}^n \rho(\boldsymbol{X}_i^{\top}\boldsymbol{\beta}),$$
(32)

$$\tilde{\boldsymbol{\beta}} = \arg\min_{\boldsymbol{\beta} \in \mathbb{R}^{p-1}} \tilde{\ell}(\boldsymbol{\beta}) := \sum_{i=1}^{n} \rho(\tilde{\boldsymbol{X}}_{i}^{\top} \boldsymbol{\beta}).$$
(33)

With the above notations in place, the LLR statistic for testing  $\beta_1 = 0$  vs.  $\beta_1 \neq 0$  can be expressed as

$$\Lambda_1 := \tilde{\ell}(\tilde{\boldsymbol{\beta}}) - \ell(\hat{\boldsymbol{\beta}}) = \sum_{i=1}^n \left\{ \rho(\tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}}) - \rho(\boldsymbol{X}_i^{\top} \hat{\boldsymbol{\beta}}) \right\}.$$
(34)

To analyze  $\Lambda_1$ , we invoke Taylor expansion to reach

$$\Lambda_{1} = \underbrace{\sum_{i=1}^{n} \rho'\left(\boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}}\right)\left(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}} - \boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}}\right)}_{:=\mathcal{Q}_{\text{lin}}} + \frac{1}{2} \sum_{i=1}^{n} \rho''\left(\boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}}\right)\left(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}} - \boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}}\right)^{2} + \frac{1}{6} \sum_{i=1}^{n} \rho'''(\gamma_{i})\left(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}} - \boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}}\right)^{3},$$
(35)

where  $\gamma_i$  lies between  $\tilde{X}_i^{\top} \hat{\beta}$  and  $X_i^{\top} \hat{\beta}$ . A key observation is that the linear term  $Q_{\text{lin}}$  in the above equation vanishes. To see this, note that the first-order optimality conditions for the MLE  $\hat{\beta}$  is given by

$$\sum_{i=1}^{n} \rho'(\boldsymbol{X}_{i}^{\top} \hat{\boldsymbol{\beta}}) \boldsymbol{X}_{i} = \boldsymbol{0}.$$
(36)

Replacing  $\tilde{X}_i^{\top} \tilde{\beta}$  with  $X_i^{\top} \begin{bmatrix} 0 \\ \tilde{\beta} \end{bmatrix}$  in  $Q_{\text{lin}}$  and using the optimality condition, we obtain

$$Q_{\rm lin} = \left(\sum_{i=1}^{n} \rho'\left(X_i^{\top} \hat{\boldsymbol{\beta}}\right) X_i\right)^{\top} \left(\begin{bmatrix}0\\\tilde{\boldsymbol{\beta}}\end{bmatrix} - \hat{\boldsymbol{\beta}}\right) = 0.$$

Consequently,  $\Lambda_1$  simplifies to the following form

$$\Lambda_1 = \frac{1}{2} \sum_{i=1}^n \rho''(\boldsymbol{X}_i^{\top} \hat{\boldsymbol{\beta}}) \left( \tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}} - \boldsymbol{X}_i^{\top} \hat{\boldsymbol{\beta}} \right)^2 + \frac{1}{6} \sum_{i=1}^n \rho'''(\gamma_i) \left( \tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}} - \boldsymbol{X}_i^{\top} \hat{\boldsymbol{\beta}} \right)^3.$$
(37)

Thus, computing the asymptotic distribution of  $\Lambda_1$  boils down to analyzing  $X_i^{\top} \hat{\beta} - \tilde{X}_i^{\top} \tilde{\beta}$ . Our argument is inspired by the leave-one-predictor-out approach developed in [24,25].

We re-emphasize that our setting is not covered by that of [24,25], due to the violation of strong convexity and some other technical assumptions. We sidestep this issue by utilizing the *Norm Bound Condition* and the *Likelihood Curvature Condition*. In the end, our analysis establishes the equivalence of  $\Lambda_1$  and  $\hat{\beta}_1$  up to some explicit multiplicative factors modulo negligible error terms. This is summarized as follows.

**Theorem 3** Under the assumptions of Theorem 1,

$$2\Lambda_1 - \frac{p}{b_*}\hat{\beta}_1^2 \xrightarrow{\mathbb{P}} 0, \quad as \ n \to \infty.$$
(38)

Theorem 3 reveals a simple yet surprising connection between the LLR statistic  $\Lambda_1$ and the MLE  $\hat{\beta}$ . As we shall see in the proof of the theorem, the quadratic term in (37) is  $\frac{1}{2} \frac{p}{b_*} \hat{\beta}_1^2 + o(1)$ , while the remaining third-order term of (37) is vanishingly small. Finally, putting Corollary 1 and Theorem 3 together directly establishes Theorem 1.

#### 2.6 Comparisons with the classical regime

We pause to shed some light on the interpretation of the correction factor  $\tau_*^2/b_*$  in Theorem 1 and understand the differences from classical results. Classical theory (e.g. [35,36]) asserts that when *p* is fixed and *n* diverges, the MLE for a fixed design *X* is asymptotically normal, namely,

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{\mathrm{d}} \mathcal{N}(0, \mathcal{I}_{\boldsymbol{\beta}}^{-1}),$$
 (39)

where

$$\mathcal{I}_{\boldsymbol{\beta}} = \frac{1}{n} \boldsymbol{X}^{\top} \boldsymbol{D}_{\boldsymbol{\beta}} \boldsymbol{X} \quad \text{with} \quad \boldsymbol{D}_{\boldsymbol{\beta}} := \begin{bmatrix} \rho^{\prime\prime} \left( \boldsymbol{X}_{1}^{\top} \boldsymbol{\beta} \right) & & \\ & \ddots & \\ & & \rho^{\prime\prime} \left( \boldsymbol{X}_{n}^{\top} \boldsymbol{\beta} \right) \end{bmatrix}$$
(40)



Fig. 3 Ratio of asymptotic variance and dimensionality factor  $\kappa$  as a function of  $\kappa$ 

is the normalized Fisher information at the true value  $\beta$ . In particular, under the global null and i.i.d. Gaussian design, this converges to

$$\mathbb{E}_{X}[\mathcal{I}_{\beta}] = \begin{cases} \frac{1}{4}I, & \text{for the logistic model} \\ \frac{2}{\pi}I, & \text{for the probit model} \end{cases}$$

as *n* tends to infinity [64, Example 5.40].

The behavior in high dimensions is different. In particular, Corollary 1 states that under the global null, we have

$$\sqrt{p}(\hat{\beta}_j - \beta_j) \stackrel{\mathrm{d}}{\to} \mathcal{N}(0, \tau_*^2).$$
 (41)

Comparing the variances in the logistic model, we have that

$$\lim_{n \to \infty} \operatorname{Var}\left(\sqrt{p}\hat{\beta}_{j}\right) = \begin{cases} 4\kappa, & \text{in classical large-sample theory;} \\ \tau_{*}^{2}, & \text{in high dimensions.} \end{cases}$$

Figure 3 illustrates the behavior of the ratio  $\tau_*^2/\kappa$  as a function of  $\kappa$ . Two observations are immediate:

- First, in Fig. 3a we have  $\tau_*^2 \ge 4\kappa$  for all  $\kappa \ge 0$ . This indicates an inflation in variance or an "extra Gaussian noise" component that appears in high dimensions, as discussed in [22]. The variance of the "extra Gaussian noise" component increases as  $\kappa$  grows.
- Second, as  $\kappa \to 0$ , we have  $\tau_*^2/4\kappa \to 1$  in the logistic model, which indicates that classical theory becomes accurate in this case. In other words, our theory recovers the classical prediction in the regime where p = o(n).

Further, for the testing problem considered here, the LLR statistic in the classical setup can be expressed, through Taylor expansion, as

$$2\Lambda_1 = n(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})^\top \mathcal{I}_{\boldsymbol{\beta}}(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}) + o_P(1), \qquad (42)$$

where  $\hat{\beta}$  is defined in (33). In the high-dimensional setting, we will also establish a quadratic approximation of the form

$$2\Lambda_1 = n(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})^\top \boldsymbol{G}(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}) + o_P(1), \qquad \boldsymbol{G} = \frac{1}{n} \boldsymbol{X}^\top \boldsymbol{D}_{\hat{\boldsymbol{\beta}}} \boldsymbol{X}.$$

In Theorem 7, we shall see that  $b_*$  is the limit of  $\frac{1}{n}$ Tr( $G^{-1}$ ), the Stieltjes transform of the empirical spectral distribution of G evaluated at 0. Thus, this quantity in some sense captures the spread in the eigenvalues of G one would expect to happen in high dimensions.

## 2.7 Prior art

Wilks' type of phenomenon in the presence of a diverging dimension p has received much attention in the past. For instance, Portnoy [51] investigated simple hypotheses in regular exponential families, and established the asymptotic Chi-square approximation for the LLR test statistic as long as  $p^{3/2}/n \rightarrow 0$ . This phenomenon was later extended in [54] to accommodate the MLE with a quadratic penalization, and in [67] to account for parametric models underlying several random graph models. Going beyond parametric inference, Fan et al. [27,29] explored extensions to infinitedimensional non-parametric inference problems, for which the MLE might not even exist or might be difficult to derive. While the classical Wilks' phenomenon fails to hold in such settings, Fan et al. [27,29] proposed a generalization of the likelihood ratio statistics based on suitable non-parametric estimators and characterized the asymptotic distributions. Such results have further motivated Boucheron and Massart [10] to investigate the non-asymptotic Wilks' phenomenon or, more precisely, the concentration behavior of the difference between the excess empirical risk and the true risk, from a statistical learning theory perspective. The Wilks' phenomenon for penalized empirical likelihood has also been established [59]. However, the precise asymptotic behavior of the LLR statistic in the regime that permits p to grow proportional to n is still beyond reach.

On the other hand, as demonstrated in Sect. 2.5.1, the MLE here under the global null can be viewed as an M-estimator for a linear regression problem. Questions regarding the behavior of robust linear regression estimators in high dimensions—where p is allowed to grow with n—were raised in Huber [35], and have been extensively studied in subsequent works, e.g. [43,48–50]. When it comes to logistic regression, the behavior of the MLE was studied for a diverging number of parameters by [33], which characterized the squared estimation error of the MLE if  $(p \log p)/n \rightarrow 0$ . In addition, the asymptotic normality properties of the MLE and the penalized MLE for logistic regression have been established by [28,42], respectively. A very recent paper

by Fan et al. [30] studied the logistic model under the global null  $\beta = 0$ , and investigated the classical asymptotic normality as given in (39). It was discovered in [30] that the convergence property breaks down even in terms of the marginal distribution, namely,

$$\frac{\sqrt{n}\hat{\beta}_i}{\left(\mathcal{I}_{\boldsymbol{\beta}}\right)_{i\,i}^{-1/2}} \stackrel{\mathrm{d}}{\nrightarrow} \mathcal{N}\left(0,1\right), \qquad \mathcal{I}_{\boldsymbol{\beta}} = \frac{1}{4n} \boldsymbol{X}^\top \boldsymbol{X},$$

as soon as p grows at a rate exceeding  $n^{2/3}$ . In other words, classical theory breaks down. Having said this, when  $\kappa = p/n \rightarrow 0$ , Wilks' theorem still holds since Theorem 1 and Fig. 2 demonstrate that the LLR statistic  $2\Lambda_j$  converges in distribution to a Chi-square.<sup>5</sup>

The line of work that is most relevant to the present paper was initially started by El Karoui et al. [26]. Focusing on the regime where p is comparable to n, the authors uncovered, via a non-rigorous argument, that the asymptotic  $\ell_2$  error of the MLE could be characterized by a system of nonlinear equations. This seminal result was later made rigorous independently by Donoho et al. [22,23] under i.i.d. Gaussian design and by El Karoui [24,25] under more general i.i.d. random design as well as certain assumptions on the error distribution. Both approaches rely on strong convexity on the function  $\rho(\cdot)$  that defines the M-estimator, which does not hold in the models considered herein. The case of ridge regularized M-estimators were also studied in [24,25]. Thrampoulidis et al. [61] studied the asymptotic behavior of the squared error for regularized M-estimators for a broad class of regularizers. They also examined the unregularized case under more general assumptions on the loss function. Our setup is not covered by this work; several conditions are violated including some pertaining to the growth rate of  $\rho'$ . We also have completely different error distributions since we are not working under the linear model as they are. Simultaneously, penalized likelihood procedures have been studied extensively under high-dimensional setups with sparsity constraints imposed on the underlying signal; see, for instance, [37,39,62,63] and the references cited therein.

Finally, we remark that the AMP machinery has already been successfully applied to study other statistical problems, including but not limited to the mean square estimation error of the Lasso [8], the tradeoff between the type I and type II errors along the Lasso path [55], and the hidden clique problem [21].

## 2.8 Notations

We adopt the standard notation f(n) = O(g(n)) or  $f(n) \leq g(n)$  which means that there exists a constant c > 0 such that  $|f(n)| \leq c|g(n)|$ . Likewise,  $f(n) = \Omega(g(n))$ or  $f(n) \geq g(n)$  means that there exists a constant c > 0 such that  $|f(n)| \geq c|g(n)|$ ,

<sup>&</sup>lt;sup>5</sup> Mathematically, the convex geometry and the leave-one-out analyses employed in our proof naturally extend to the case where p = o(n). It remains to develop a formal AMP theory for the regime where p = o(n). Alternatively, we note that the AMP theory has been mainly invoked to characterize  $\|\hat{\beta}\|$ , which can also be accomplished via the leave-one-out argument (cf. [25]). This alternative proof strategy can easily extend to the regime p = o(n).



**Fig. 4** Histogram of *p*-values for logistic regression under i.i.d. Bernoulli design (this is, therefore, not the setup of Fig. 1), when  $\beta = 0$ , n = 4000, p = 1200, and  $\kappa = 0.3$ : **a** classically computed *p*-values; **b** Bartlett corrected *p*-values; **c** adjusted *p*-values

 $f(n) \approx g(n)$  means that there exist constants  $c_1, c_2 > 0$  such that  $c_1|g(n)| \leq |f(n)| \leq c_2|g(n)|$ , and f(n) = o(g(n)) means that  $\lim_{n\to\infty} \frac{f(n)}{g(n)} = 0$ . Any mention of  $C, C_i, c$ ,  $c_i$  for  $i \in \mathbb{N}$  refers to some positive universal constants whose value may change from line to line. For a square symmetric matrix M, the minimum eigenvalue is denoted by  $\lambda_{\min}(M)$ . Logarithms are base e.

# **3 Numerics**

# 3.1 Non-Gaussian covariates

In this section we first study the sensitivity of our result to the Gaussianity assumption on the design matrix. To this end, we consider a high dimensional binary regression set up with a Bernoulli design matrix. We simulate n = 4000 i.i.d. observations  $(y_i, X_i)$ with  $y_i \stackrel{\text{i.i.d.}}{\sim}$  Bernoulli(1/2), and  $X_i$  generated independent of  $y_i$ , such that each entry takes on values in  $\{1, -1\}$  w.p. 1/2. At each trial, we fit a logistic regression model to the data and obtain the classical, Bartlett corrected and adjusted *p*-values (using the rescaling factor  $\alpha(\kappa)$ ). Figure 4 plots the histograms for the pooled *p*-values, obtained across 400 trials.

It is instructive to compare the histograms to that obtained in the Gaussian case (Fig. 1). The classical and Bartlett corrected p-values exhibit similar deviations from uniformity as in the Gaussian design case, whereas our adjusted p-values continue to have an approximate uniform distribution. We test for deviations from uniformity using a formal Chi-squared goodness of fit test as in Sect. 2.3.2. For the Bartlett corrected p-values, the Chi-squared statistic turns out to be 5885, with a p-value 0. For the adjusted p-values, the Chi-squared statistic is 24.1024, with a p-value 0.1922.<sup>6</sup>

Once again, the Bartlett correction fails to provide valid p-values whereas the adjusted p-values are consistent with a uniform distribution. These findings indicate that the distribution of the LLR statistic under the i.i.d. Bernoulli design is in agreement

<sup>&</sup>lt;sup>6</sup> Recall our earlier footnote about the use of a  $\chi^2$  test.

Table 2Estimates of p-valueprobabilities with estimatedMonte Carlo standard errors inparentheses underi.i.d. Bernoulli design		Adjusted	
	$\mathbb{P}\{p\text{-values} \le 5\%\}$	5.0222% (0.0412%)	
	$\mathbb{P}\{p\text{-values} \le 1\%\}$	1.0048% (0.0174%)	
	$\mathbb{P}\{p\text{-values} \le 0.5\%\}$	0.5123% (0.0119%)	
	$\mathbb{P}\{p\text{-values} \le 0.1\%\}$	0.1108% (0.005%)	
	$\mathbb{P}\{p\text{-values} \le 0.05\%\}$	0.0521% (0.0033%)	
	$\mathbb{P}\{p\text{-values} \le 0.01\%\}$	0.0102% (0.0015%)	

to the rescaled  $\chi_1^2$  derived under the Gaussian design in Theorem 1, suggesting that the distribution is not too sensitive to the Gaussianity assumption. Estimates of *p*-value probabilities for our method are provided in Table 2.

# 3.2 Quality of approximations for finite sample sizes

In the rest of this section, we report some numerical experiments which study the applicability of our theory in finite sample setups.

Validity of tail approximation The first experiment explores the efficacy of our correction for extremely small *p*-values. This is particularly important in the context of multiple comparisons, where practitioners care about the validity of exceedingly small *p*-values. To this end, the empirical cumulative distribution of the adjusted *p*-values is estimated under a standard Gaussian design with n = 4000, p = 1200 and  $4.8 \times 10^5$  *p*-values. The range [0.1/p, 12/p] is divided into points which are equispaced with a distance of 1/p between any two consecutive points. The estimated empirical CDF at each of these points is represented in Fig. 5. The estimated CDF is in near-perfect agreement with the diagonal, suggesting that the adjusted *p*-values computed using the rescaled Chi-square distribution are remarkably close to a uniform,





even when we zoom in at very small resolutions as would be the case when applying Bonferroni-style corrections.

**Moderate sample sizes** The final experiment studies the accuracy of our asymptotic result for moderately large samples. This is especially relevant for applications where the sample sizes are not too large. We repeat our numerical experiments with n = 200, p = 60 for i.i.d. Gaussian design, and  $4.8 \times 10^5 p$ -values. The empirical CDF for these *p*-values are estimated and Fig. 6 shows that the adjusted *p*-values are nearly uniformly distributed even for moderate sample sizes such as n = 200.

# **4** Preliminaries

This section gathers a few preliminary results that will be useful throughout the paper. We start by collecting some facts regarding i.i.d. Gaussian random matrices.

**Lemma 2** Let  $X = [X_1, X_2, ..., X_n]^\top$  be an  $n \times p$  matrix with i.i.d. standard Gaussian entries. Then

$$\mathbb{P}\left(\|\boldsymbol{X}^{\top}\boldsymbol{X}\| \le 9n\right) \ge 1 - 2\exp(-n/2); \tag{43}$$

$$\mathbb{P}\left(\sup_{1 \le i \le n} \|X_i\| \le 2\sqrt{p}\right) \ge 1 - 2n \exp(-(\sqrt{p} - 1)^2/2).$$
(44)

**Proof** This is a straighforward application of [65, Corollary 5.35] and the union bound.  $\Box$ 

**Lemma 3** Suppose X is an  $n \times p$  matrix with entries i.i.d  $\mathcal{N}(0, 1)$ , then there exists a constant  $\epsilon_0$  such that whenever  $0 \le \epsilon \le \epsilon_0$  and  $0 \le t \le \sqrt{1-\epsilon} - \sqrt{p/n}$ ,

$$\lambda_{\min}\left(\frac{1}{n}\sum_{i\in S} X_i X_i^{\top}\right) \ge \left(\sqrt{1-\epsilon} - \sqrt{\frac{p}{n}} - t\right)^2, \quad \forall S \subseteq [n] \text{ with } |S| = (1-\epsilon)n \quad (45)$$

🖉 Springer

The message of Lemma 4 is this: take 
$$\epsilon > 0$$
 to Then

**Proof** See Appendix A.2. be a sufficiently small constant.

holds simultaneously for all  $\boldsymbol{\beta} \in \mathbb{R}^p$ .

**Proof** See Appendix A.1.

and the result is this:

 $-\epsilon \log \epsilon - (1 - \epsilon) \log(1 - \epsilon).$ 

The above facts are useful in establishing an eigenvalue lower bound on the Hessian of the log-likelihood function. Specifically, recall that

with probability exceeding  $1 - 2 \exp\left(-\left(\frac{(1-\epsilon)t^2}{2} - H(\epsilon)\right)n\right)$ . Here,  $H(\epsilon) =$ 

$$\nabla^2 \ell(\boldsymbol{\beta}) = \sum_{i=1}^n \rho'' \left( \boldsymbol{X}_i^\top \boldsymbol{\beta} \right) \boldsymbol{X}_i \boldsymbol{X}_i^\top, \tag{46}$$

**Lemma 4** (Likelihood curvature condition) Suppose that 
$$p/n < 1$$
 and that  $\rho''(\cdot) \ge 0$ .  
Then there exists a constant  $\epsilon_0$  such that whenever  $0 \le \epsilon \le \epsilon_0$ , with probability at least  $1 - 2 \exp(-nH(\epsilon)) - 2 \exp(-n/2)$ , the matrix inequality

$$\frac{1}{n} \nabla^2 \ell(\boldsymbol{\beta}) \succeq \left( \inf_{z: |z| \le \frac{3\|\boldsymbol{\beta}\|}{\sqrt{\epsilon}}} \rho''(z) \right) \left( \sqrt{1-\epsilon} - \sqrt{\frac{p}{n}} - 2\sqrt{\frac{H(\epsilon)}{1-\epsilon}} \right)^2 \boldsymbol{I}$$
(47)

 $\frac{1}{n} \nabla^2 \ell(\boldsymbol{\beta}) \succeq \omega(\|\boldsymbol{\beta}\|) \boldsymbol{I}$ 

for some non-increasing, continuous and positive function 
$$\omega(\cdot)$$
 independent of *n*. This is a generalization of the strong convexity condition.

# 5 When is the MLE bounded?

# 5.1 Phase transition

In Sect. 2.2, we argued that the MLE is at infinity if we have less than two observations per dimension or  $\kappa > 1/2$ . In fact, a stronger version of the phase transition phenemonon occurs in the sense that

$$\|\hat{\boldsymbol{\beta}}\| = O(1)$$

as soon as  $\kappa < 1/2$ . This is formalized in the following theorem.

**Theorem 4** (Norm Bound Condition) Fix any small constant  $\epsilon > 0$ , and let  $\hat{\beta}$  be the MLE for a model with effective link satisfying the conditions from Sect. 2.3.3.

- (i) If  $p/n \ge 1/2 + \epsilon$ , then the MLE does not exist with probability exceeding  $1 4 \exp(-\epsilon^2 n/8)$ .
- (ii) There exist universal constants  $c_1, c_2, C_2 > 0$  such that if  $p/n < 1/2 c_1 \epsilon^{3/4}$ , then<sup>7</sup>

$$\|\hat{\boldsymbol{\beta}}\| < \frac{4\log 2}{\epsilon^2}$$

with probability at least  $1 - C_2 \exp(-c_2 \epsilon^2 n)$ .

These conclusions clearly continue to hold if  $\hat{\beta}$  is replaced by  $\tilde{\beta}$  (the MLE under the restricted model obtained on dropping the first predictor).

The rest of this section is devoted to proving this theorem. As we will see later, the fact that  $\|\hat{\boldsymbol{\beta}}\| = O(1)$  is crucial for utilizing the AMP machinery in the absence of strong convexity.

## 5.2 Proof of Theorem 4

As in Sect. 2.5.1, we assume  $\tilde{y}_i \equiv 1$  throughout this section, and hence the MLE reduces to

#### 5.2.1 Proof of Part (i)

Invoking [3, Theorem I] yields that if

$$\delta\left(\left\{\boldsymbol{X\boldsymbol{\beta}} \mid \boldsymbol{\beta} \in \mathbb{R}^p\right\}\right) + \delta\left(\mathbb{R}^n_+\right) \ge (1+\epsilon)\,n,$$

or equivalently, if  $p/n \ge 1/2 + \epsilon$ , then

$$\mathbb{P}\left\{\left\{\boldsymbol{X\boldsymbol{\beta}}\mid\boldsymbol{\beta}\in\mathbb{R}^{p}\right\}\cap\mathbb{R}^{n}_{+}\neq\{\boldsymbol{0}\}\right\}\geq1-4\exp\left(-\epsilon^{2}n/8\right).$$

As is seen in Sect. 2.2,  $\|\hat{\boldsymbol{\beta}}\| = \infty$  when  $\{\boldsymbol{X}\boldsymbol{\beta} \mid \boldsymbol{\beta} \in \mathbb{R}^p\} \cap \mathbb{R}^n_+ \neq \{\mathbf{0}\}$ , establishing Part (i) of Theorem 4.

 $<sup>\</sup>overline{{}^7 \text{ When } X_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}) \text{ for a general } \boldsymbol{\Sigma} \succ \mathbf{0}, \text{ one has } \|\boldsymbol{\Sigma}^{1/2} \hat{\boldsymbol{\beta}}\| \lesssim 1/\epsilon^2 \text{ with high probability.}}$ 

## 5.2.2 Proof of Part (ii)

We now turn to the regime in which  $p/n \le 1/2 - O(\epsilon^{3/4})$ , where  $0 < \epsilon < 1$  is any fixed constant. Begin by observing that the least singular value of *X* obeys

$$\sigma_{\min}\left(X\right) \ge \sqrt{n}/4 \tag{49}$$

with probability at least  $1 - 2 \exp\left(-\frac{1}{2}\left(\frac{3}{4} - \frac{1}{\sqrt{2}}\right)^2 n\right)$  (this follows from Lemma 3 using  $\epsilon = 0$ ). Then for any  $\beta \in \mathbb{R}^p$  obeying

$$\ell_0(\boldsymbol{\beta}) = \sum_{j=1}^n \rho\left(-\boldsymbol{X}_j^{\top} \boldsymbol{\beta}\right) \le n \log 2 = \ell_0(\boldsymbol{0})$$
(50)

and 
$$\|\boldsymbol{\beta}\| \ge \frac{4\log 2}{\epsilon^2},$$
 (51)

we must have

$$\sum_{j=1}^{n} \max\left\{-X_{j}^{\top}\boldsymbol{\beta}, 0\right\} = \sum_{j: X_{j}^{\top}\boldsymbol{\beta}<0} \left(-X_{j}^{\top}\boldsymbol{\beta}\right) \stackrel{(a)}{\leq} \sum_{j: X_{j}^{\top}\boldsymbol{\beta}<0} \rho\left(-X_{j}^{\top}\boldsymbol{\beta}\right) \stackrel{(b)}{\leq} n \log 2;$$

(a) follows since  $t \le \rho(t)$  and (b) is a consequence of (50). Continuing, (49) and (51) give

$$n \log 2 \le 4\sqrt{n} \frac{\|\boldsymbol{X}\boldsymbol{\beta}\|}{\|\boldsymbol{\beta}\|} \log 2 \le \epsilon^2 \sqrt{n} \|\boldsymbol{X}\boldsymbol{\beta}\|.$$

This implies the following proposition: if the solution  $\hat{\beta}$ —which necessarily satisfies  $\ell_0(\hat{\beta}) \leq \ell_0(\mathbf{0})$ —has norm exceeding  $\|\hat{\beta}\| \geq \frac{4\log 2}{\epsilon^2}$ , then  $X\hat{\beta}$  must fall within the cone

$$\mathcal{A} := \left\{ \boldsymbol{u} \in \mathbb{R}^n \, \middle| \, \sum_{j=1}^n \max\left\{ -u_j, 0 \right\} \le \epsilon^2 \sqrt{n} \|\boldsymbol{u}\| \right\}.$$
(52)

Therefore, if one wishes to rule out the possibility of having  $\|\hat{\boldsymbol{\beta}}\| \ge \frac{4\log 2}{\epsilon^2}$ , it suffices to show that with high probability,

$$\left\{ \boldsymbol{X}\boldsymbol{\beta} \mid \boldsymbol{\beta} \in \mathbb{R}^p \right\} \cap \mathcal{A} = \left\{ \boldsymbol{0} \right\}.$$
(53)

This is the content of the remaining proof.

We would like to utilize tools from conic geometry [3] to analyze the probability of the event (53). Note, however, that A is not convex, while the theory developed in [3] applies only to convex cones. To bypass the non-convexity issue, we proceed in the following three steps:

1. Generate a set of  $N = \exp(2\epsilon^2 p)$  closed *convex* cones  $\{\mathcal{B}_i \mid 1 \le i \le N\}$  such that it forms a cover of  $\mathcal{A}$  with probability exceeding  $1 - \exp(-\Omega(\epsilon^2 p))$ .

2. Show that if  $p < \left(\frac{1}{2} - 2\sqrt{2}\epsilon^{\frac{3}{4}} - 2H(2\sqrt{\epsilon})\right)n$  and if *n* is sufficiently large, then

$$\mathbb{P}\left\{\left\{\boldsymbol{X\boldsymbol{\beta}} \mid \boldsymbol{\beta} \in \mathbb{R}^{p}\right\} \cap \mathcal{B}_{i} \neq \{\mathbf{0}\}\right\}$$
  
$$\leq 4 \exp\left\{-\frac{1}{8}\left(\frac{1}{2} - 2\sqrt{2}\epsilon^{\frac{3}{4}} - 10H(2\sqrt{\epsilon}) - \frac{p}{n}\right)^{2}n\right\}$$

for each  $1 \le i \le N$ .

3. Invoke the union bound to reach

$$\mathbb{P}\left\{\left\{\boldsymbol{X\boldsymbol{\beta}}\mid\boldsymbol{\beta}\in\mathbb{R}^{p}\right\}\cap\mathcal{A}\neq\{\boldsymbol{0}\}\right\}\leq\mathbb{P}\left\{\left\{\mathcal{B}_{i}\mid1\leq i\leq N\right\}\text{ does not form a cover of }\mathcal{A}\right\}$$
$$+\sum_{i=1}^{N}\mathbb{P}\left\{\left\{\boldsymbol{X\boldsymbol{\beta}\mid\boldsymbol{\beta}\in\mathbb{R}^{p}\right\}\cap\mathcal{B}_{i}\neq\{\boldsymbol{0}\}\right\}$$
$$\leq\exp\left(-\Omega(\epsilon^{2}p)\right),$$

where we have used the fact that

$$\begin{split} &\sum_{i=1}^{N} \mathbb{P}\left\{\left\{\boldsymbol{X}\boldsymbol{\beta} \mid \boldsymbol{\beta} \in \mathbb{R}^{p}\right\} \cap \mathcal{B}_{i} \neq \{\mathbf{0}\}\right\} \\ &\leq 4N \exp\left\{-\frac{1}{8}\left(\frac{1}{2} - 2\sqrt{2}\epsilon^{\frac{3}{4}} - 10H(2\sqrt{\epsilon}) - \frac{p}{n}\right)^{2}n\right\} \\ &< 4\exp\left\{-\left(\frac{1}{8}\left(\frac{1}{2} - 2\sqrt{2}\epsilon^{\frac{3}{4}} - 10H(2\sqrt{\epsilon}) - \frac{p}{n}\right)^{2} - 2\epsilon^{2}\right)n\right\} \\ &< 4\exp\left\{-\epsilon^{2}n\right\}. \end{split}$$

Here, the last inequality holds if  $\left(\frac{1}{2} - 2\sqrt{2}\epsilon^{\frac{3}{4}} - 10H(2\sqrt{\epsilon}) - \frac{p}{n}\right)^2 > 24\epsilon^2$ , or equivalently,  $\frac{p}{n} < \frac{1}{2} - 2\sqrt{2}\epsilon^{\frac{3}{4}} - 10H(2\sqrt{\epsilon}) - \sqrt{24}\epsilon$ .

Taken collectively, these steps establish the following claim: if  $\frac{p}{n} < \frac{1}{2} - 2\sqrt{2}\epsilon^{\frac{3}{4}} - 10H(2\sqrt{\epsilon}) - \sqrt{24}\epsilon$ , then

$$\mathbb{P}\left\{\|\hat{\boldsymbol{\beta}}\| > \frac{4\log 2}{\epsilon^2}\right\} < \exp\left\{-\Omega(\epsilon^2 n)\right\},\,$$

thus establishing Part (ii) of Theorem 4. We defer the complete details of the preceding steps to Appendix D.

# 6 Asymptotic $\ell_2$ error of the MLE

This section aims to establish Theorem 2, which characterizes precisely the asymptotic squared error of the MLE  $\hat{\beta}$  under the global null  $\beta = 0$ . As described in Sect. 2.5.1, it suffices to assume that  $\hat{\beta}$  is the solution to the following problem

In what follows, we derive the asymptotic convergence of  $\|\hat{\beta}\|$  under the assumptions from our main theorem.

**Theorem 5** Under the assumptions of Theorem 1, the solution  $\hat{\beta}$  to (54) obeys

$$\lim_{n \to \infty} \|\hat{\boldsymbol{\beta}}\|^2 =_{\text{a.s.}} \tau_*^2.$$
(55)

Theorem 5 is derived by invoking the AMP machinery [7,8,38]. The high-level idea is the following: in order to study  $\hat{\beta}$ , one introduces an iterative algorithm (called AMP) where a sequence of iterates  $\hat{\beta}^t$  is formed at each time *t*. The algorithm is constructed so that the iterates asymptotically converge to the MLE in the sense that

$$\lim_{t \to \infty} \lim_{n \to \infty} \|\hat{\boldsymbol{\beta}}^t - \hat{\boldsymbol{\beta}}\|^2 =_{\text{a.s.}} 0.$$
(56)

On the other hand, the asymptotic behavior (asymptotic in *n*) of  $\hat{\boldsymbol{\beta}}^{l}$  for each *t* can be described accurately by a scalar sequence  $\{\tau_t\}$ —called *state evolution* (SE)—following certain update equations [7]. This, in turn, provides a characterization of the  $\ell_2$  loss of  $\hat{\boldsymbol{\beta}}$ .

Further, in order to prove Theorem 2, one still needs to justify

- (a) the existence of a solution to the system of Eqs. (25) and (26),
- (b) and the existence of a fixed point for the iterative map governing the SE sequence updates.

We will elaborate on these steps in the rest of this section.

#### 6.1 State evolution

We begin with the SE sequence  $\{\tau_t\}$  introduced in [22]. Starting from some initial point  $\tau_0$ , we produce two sequences  $\{b_t\}$  and  $\{\tau_t\}$  following a two-step procedure.

- For t = 0, 1, ...:
  - Set  $b_t$  to be the solution in b to

$$\kappa = \mathbb{E}\left[\Psi'(\tau_t Z; b)\right];\tag{57}$$

- Set  $\tau_{t+1}$  to be

$$\tau_{t+1}^2 = \frac{1}{\kappa} \mathbb{E}\left[ (\Psi^2(\tau_t Z; b_t)) \right].$$
(58)

Suppose that for given any  $\tau > 0$ , the solution in *b* to (57) with  $\tau_t = \tau$  exists and is unique, then one can denote the solution as  $b(\tau)$ , which in turn allows one to write the sequence  $\{\tau_t\}$  as

$$\tau_{t+1}^2 = \mathcal{V}(\tau_t^2)$$

with the variance map

$$\mathcal{V}(\tau^2) = \frac{1}{\kappa} \mathbb{E}\Big[\Psi^2(\tau Z; b(\tau))\Big].$$
(59)

As a result, if there exists a fixed point  $\tau_*$  obeying  $\mathcal{V}(\tau_*^2) = \tau_*^2$  and if we start with  $\tau_0 = \tau_*$ , then by induction,

$$\tau_t \equiv \tau_*$$
 and  $b_t \equiv b_* := b(\tau_*), \quad t = 0, 1, \dots$ 

Notably,  $(\tau_*, b_*)$  solves the system of Eqs. (25) and (26). We shall work with this choice of initial condition throughout our proof.

The preceding arguments hold under two conditions: (i) the solution to (57) exists and is unique for any  $\tau_t > 0$ ; (ii) the variance map (59) admits a fixed point. To verify these two conditions, we make two observations.

• Condition (i) holds if one can show that the function

$$G(b) := \mathbb{E}\left[\Psi'(\tau Z; b)\right], \qquad b > 0 \tag{60}$$

is strictly monotone for any given  $\tau > 0$ , and that  $\lim_{b\to 0} G(b) < \kappa < \lim_{b\to\infty} G(b)$ .

• Since  $\mathcal{V}(\cdot)$  is a continuous function, Condition (ii) becomes self-evident once we show that  $\mathcal{V}(0) > 0$  and that there exists  $\tau > 0$  obeying  $\mathcal{V}(\tau^2) < \tau^2$ . The behavior of the variance map is illustrated in Fig. 7 for the logistic and probit regression when  $\kappa = 0.3$ . One can in fact observe that the fixed point is unique. For other values of  $\kappa$ , the variance map shows the same behavior.

In fact, the aforementioned properties can be proved for a certain class of effective links, as summarized in the following lemmas. In particular, they can be shown for the logistic and the probit models.

**Lemma 5** Suppose the effective link  $\rho$  satisfies the following two properties:

- (a)  $\rho'$  is log-concave.
- (b) For any fixed  $\tau > 0$  and any fixed z,  $b\rho''(\operatorname{prox}_{b\rho}(\tau z)) \to \infty$  when  $b \to \infty$ .



**Fig. 7** The variance map for both the logistic and the probit models when  $\kappa = 0.3$ : (solid line) variance map  $\mathcal{V}(\tau^2)$  as a function of  $\tau^2$ ; (dotted line) diagonal

Then for any  $\tau > 0$ , the function G(b) defined in (60) is an increasing function in b (b > 0), and the equation

$$G(b) = \kappa$$

has a unique positive solution.

**Proof** See Appendix **B**.

**Lemma 6** Suppose that  $0 < \kappa < 1/2$  and that  $\rho = \log(1 + e^t)$  or  $\rho = -\log \Phi(-t)$ . Then

(i) V(0) > 0;
(ii) V(τ<sup>2</sup>) < τ<sup>2</sup> for some sufficiently large τ<sup>2</sup>.

**Proof** See Appendix C and the supplemental material [58].

*Remark 1* A byproduct of the proof is that the following relations hold for any constant  $0 < \kappa < 1/2$ :

• In the logistic case,

$$\begin{cases} \lim_{\tau \to \infty} \frac{\mathcal{V}(\tau^2)}{\tau^2} &= \frac{x^2 \mathbb{P}\{Z > x\} + \mathbb{E}[Z^2 \mathbf{1}_{\{0 < Z < x\}}]}{\mathbb{P}\{0 < Z < x\}} \bigg|_{x = \Phi^{-1}(\kappa + 0.5)};\\ \lim_{\tau \to \infty} \frac{b(\tau)}{\tau} &= \Phi^{-1}(\kappa + 0.5). \end{cases}$$

• In the probit case,

$$\lim_{\tau \to \infty} b(\tau) = \frac{2\kappa}{1 - 2\kappa} \quad \text{and} \quad \lim_{\tau \to \infty} \frac{\mathcal{V}(\tau^2)}{\tau^2} = 2\kappa.$$
(61)

**Remark 2** Lemma 6 is proved for the two special effective link functions, the logistic and the probit cases. However, the proof sheds light on general conditions on the effective link that suffice for the lemma to hold. Such general sufficient conditions are also discussed in the supplemental material [58].

#### 6.2 AMP recursion

In this section, we construct the AMP trajectory tracked by two sequences  $\{\hat{\boldsymbol{\beta}}^t(n) \in \mathbb{R}^p\}$  and  $\{\boldsymbol{\eta}^t(n) \in \mathbb{R}^n\}$  for  $t \ge 0$ . Going forward we suppress the dependence on *n* to simplify presentation. Picking  $\hat{\boldsymbol{\beta}}^0$  such that

$$\lim_{n \to \infty} \|\hat{\boldsymbol{\beta}}^{0}\|^{2} = \tau_{0}^{2} = \tau_{*}^{2}$$

and taking  $\eta^{-1} = 0$  and  $b_{-1} = 0$ , the AMP path is obtained via Algorithm 1, which is adapted from the algorithm in [22, Section 2.2].

A B	-		•		•
Algorithm		Λn	nrovimata	maccona	naccing
AISOLUUU		$-\Delta U$	טווואות	IIICSSARC	Dassing
					P

For  $t = 0, 1, \cdots$ :

1. Set

$$\boldsymbol{\eta}^{t} = \boldsymbol{X} \hat{\boldsymbol{\beta}}^{t} + \Psi \left( \boldsymbol{\eta}^{t-1}; b_{t-1} \right);$$
(62)

2. Let  $b_t$  be the solution to

$$\kappa = \mathbb{E}\left[\Psi'(\tau_t Z; b)\right],\tag{63}$$

where  $\tau_t$  is the SE sequence value at that time. 3. Set

$$\hat{\boldsymbol{\beta}}^{t+1} = \hat{\boldsymbol{\beta}}^t - \frac{1}{p} \boldsymbol{X}^\top \Psi \left( \boldsymbol{\eta}^t; \boldsymbol{b}_t \right).$$
(64)

Here,  $\Psi(\cdot)$  is applied in an entrywise manner, and  $\Psi'(.,.)$  denotes derivative w.r.t the first variable.

As asserted by [22], the SE sequence  $\{\tau_t\}$  introduced in Sect. 6.1 proves useful as it offers a formal procedure for predicting operating characteristics of the AMP iterates at any fixed iteration. In particular it assigns predictions to two types of observables: observables which are functions of the  $\hat{\beta}^t$  sequence and those which are functions of  $\eta^t$ . Repeating identical argument as in [22, Theorem 3.4], we obtain

$$\lim_{n \to \infty} \|\hat{\boldsymbol{\beta}}^t\|^2 =_{\text{a.s.}} \tau_t^2 \equiv \tau_*^2, \quad t = 0, 1, \dots.$$
(65)

## 6.3 AMP converges to the MLE

We are now in position to show that the AMP iterates  $\{\hat{\beta}^t\}$  converge to the MLE in the large *n* and *t* limit. Before continuing, we state below two properties that are satisfied under our assumptions.

• The MLE  $\hat{\boldsymbol{\beta}}$  obeys

$$\lim_{n \to \infty} \|\hat{\boldsymbol{\beta}}\| < \infty \tag{66}$$

almost surely.

• And there exists some non-increasing continuous function  $0 < \omega(\cdot) < 1$  independent of *n* such that

$$\mathbb{P}\left\{\frac{1}{n}\nabla^{2}\ell\left(\boldsymbol{\beta}\right) \succeq \omega\left(\|\boldsymbol{\beta}\|\right) \cdot \boldsymbol{I}, \; \forall \boldsymbol{\beta}\right\} \ge 1 - c_{1}e^{-c_{2}n}.$$
(67)

In fact, the norm bound (66) follows from Theorem 4 together with Borel-Cantelli, while the likelihood curvature condition (67) is an immediate consequence of Lemma 4. With this in place, we have:

**Theorem 6** Suppose (66) and (67) hold. Let  $(\tau_*, b_*)$  be a solution to the system (25) and (26), and assume that  $\lim_{n\to\infty} \|\hat{\boldsymbol{\beta}}^0\|^2 = \tau_*^2$ . Then the AMP trajectory as defined in Algorithm 1 obeys

$$\lim_{t\to\infty}\lim_{n\to\infty}\|\hat{\boldsymbol{\beta}}^t-\hat{\boldsymbol{\beta}}\|=_{\text{a.s.}}0$$

Taken collectively, Theorem 6 and Eq. (65) imply that

$$\lim_{n \to \infty} \|\hat{\boldsymbol{\beta}}\| =_{\text{a.s.}} \lim_{t \to \infty} \lim_{n \to \infty} \|\hat{\boldsymbol{\beta}}^{T}\| =_{\text{a.s.}} \tau_{*}, \tag{68}$$

thus establishing Theorem 5. In addition, an upshot of these theorems is a uniqueness result:

# **Corollary 2** The solution to the system of Eqs. (25) and (26) is unique.

**Proof** When the AMP trajectory  $\hat{\boldsymbol{\beta}}^t$  is started with the initial condition from Theorem 6,  $\lim_{n\to\infty} \|\hat{\boldsymbol{\beta}}\|^2 =_{\text{a.s.}} \tau_*^2$ . This holds for any  $\tau_*$  such that  $(\tau_*, b_*)$  is a solution to (25) and (26). However, since the MLE problem is strongly convex and hence admits a unique solution  $\hat{\boldsymbol{\beta}}$ , this implies that  $\tau_*$  must be unique, which together with the monotonicity of  $G(\cdot)$  (cf. (60)) implies that  $b_*$  is unique as well.

**Proof of Theorem 6** To begin with, repeating the arguments in [22, Lemma 6.9] we reach

$$\lim_{t \to \infty} \lim_{n \to \infty} \|\hat{\boldsymbol{\beta}}^{t+1} - \hat{\boldsymbol{\beta}}^t\|^2 =_{\text{a.s.}} 0;$$
(69)

$$\lim_{t \to \infty} \lim_{n \to \infty} \frac{1}{n} \| \boldsymbol{\eta}^{t+1} - \boldsymbol{\eta}^t \|^2 =_{\text{a.s.}} 0.$$
(70)

To show that the AMP iterates converge to the MLE, we shall analyze the loglikelihood function. Recall from Taylor's theorem that

$$\ell(\hat{\boldsymbol{\beta}}) = \ell(\hat{\boldsymbol{\beta}}^{t}) + \left\langle \nabla \ell(\hat{\boldsymbol{\beta}}^{t}), \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{t} \right\rangle + \frac{1}{2} \left( \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{t} \right)^{\top} \nabla^{2} \ell\left( \hat{\boldsymbol{\beta}}^{t} + \lambda(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{t}) \right) \left( \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{t} \right)$$

holds for some  $0 < \lambda < 1$ . To deal with the quadratic term, we would like to control the Hessian of the likelihood at a point between  $\hat{\beta}$  and  $\hat{\beta}^{t}$ . Invoking the likelihood curvature condition (67), one has

$$\ell(\hat{\boldsymbol{\beta}}^{t}) \geq \ell(\hat{\boldsymbol{\beta}}) \geq \ell(\hat{\boldsymbol{\beta}}^{t}) + \left\langle \nabla \ell(\hat{\boldsymbol{\beta}}^{t}), \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{t} \right\rangle + \frac{1}{2} n \omega \left( \max\left\{ \|\hat{\boldsymbol{\beta}}\|, \|\hat{\boldsymbol{\beta}}^{t}\| \right\} \right) \|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{t}\|^{2}$$

$$(71)$$

with high probability. Apply Cauchy-Schwarz to yield that with exponentially high probability,

$$\|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{t}\| \leq \frac{2}{\omega \left(\max\left\{\|\hat{\boldsymbol{\beta}}\|, \|\hat{\boldsymbol{\beta}}^{t}\|\right\}\right)} \left\|\frac{1}{n} \nabla \ell(\hat{\boldsymbol{\beta}}^{t})\right\| \leq \frac{2}{\omega \left(\|\hat{\boldsymbol{\beta}}\|\right) \omega \left(\|\hat{\boldsymbol{\beta}}^{t}\|\right)} \left\|\frac{1}{n} \nabla \ell(\hat{\boldsymbol{\beta}}^{t})\right\|,$$

where the last inequality follows since  $0 < \omega(\cdot) < 1$  and  $\omega(\cdot)$  is non-increasing.

It remains to control  $\|\nabla \ell(\hat{\boldsymbol{\beta}}^t)\|$ . The identity  $\Psi(z; b_*) = z - \operatorname{prox}_{b_*\rho}(z)$  and (62) give

$$\operatorname{prox}_{b_*\rho}\left(\boldsymbol{\eta}^{t-1}\right) = X\hat{\boldsymbol{\beta}}^t + \boldsymbol{\eta}^{t-1} - \boldsymbol{\eta}^t.$$
(72)

In addition, substituting  $\Psi(z; b) = b\rho'(\operatorname{prox}_{\rho b}(z))$  into (64) yields

$$\frac{p}{b_*}(\hat{\boldsymbol{\beta}}^t - \hat{\boldsymbol{\beta}}^{t-1}) = -X^\top \rho' \left( \operatorname{prox}_{b_*\rho}(\boldsymbol{\eta}^{t-1}) \right) = -X^\top \rho' \left( X \hat{\boldsymbol{\beta}}^t + \boldsymbol{\eta}^{t-1} - \boldsymbol{\eta}^t \right).$$

We are now ready to bound  $\|\nabla \ell(\hat{\boldsymbol{\beta}}^t)\|$ . Recalling that

$$\begin{aligned} \nabla \ell(\hat{\boldsymbol{\beta}}^{t}) &= \boldsymbol{X}^{\top} \rho'(\boldsymbol{X}^{\top} \hat{\boldsymbol{\beta}}^{t}) = \boldsymbol{X}^{\top} \rho'\left(\boldsymbol{X} \hat{\boldsymbol{\beta}}^{t} + \boldsymbol{\eta}^{t-1} - \boldsymbol{\eta}^{t}\right) \\ &+ \boldsymbol{X}^{\top} \left(\rho'(\boldsymbol{X}^{\top} \hat{\boldsymbol{\beta}}^{t}) - \rho'\left(\boldsymbol{X} \hat{\boldsymbol{\beta}}^{t} + \boldsymbol{\eta}^{t-1} - \boldsymbol{\eta}^{t}\right)\right) \end{aligned}$$

and that  $\sup_{z} \rho''(z) < \infty$ , we have

$$\begin{split} \left\| \nabla \ell(\hat{\boldsymbol{\beta}}^{t}) \right\| &\leq \left\| - \boldsymbol{X}^{\top} \rho' \left( \boldsymbol{X} \hat{\boldsymbol{\beta}}^{t} + \boldsymbol{\eta}^{t-1} - \boldsymbol{\eta}^{t} \right) \right\| \\ &+ \left\| \boldsymbol{X} \right\| \left| \rho' \left( \boldsymbol{X} \hat{\boldsymbol{\beta}}^{t} + \boldsymbol{\eta}^{t-1} - \boldsymbol{\eta}^{t} \right) - \rho'(\boldsymbol{X} \hat{\boldsymbol{\beta}}^{t}) \right| \\ &\leq \frac{p}{b_{*}} \| \hat{\boldsymbol{\beta}}^{t} - \hat{\boldsymbol{\beta}}^{t-1} \| + \| \boldsymbol{X} \| \left( \sup_{z} \rho''(z) \right) \| \boldsymbol{\eta}^{t-1} - \boldsymbol{\eta}^{t} \|. \end{split}$$

🖄 Springer

This establishes that with probability at least  $1 - c_1 e^{-c_2 n}$ ,

$$\|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{t}\| \leq \frac{2}{\omega\left(\|\hat{\boldsymbol{\beta}}\|\right)\omega\left(\|\hat{\boldsymbol{\beta}}^{t}\|\right)} \left\{\frac{p}{b_{*}n}\|\hat{\boldsymbol{\beta}}^{t} - \hat{\boldsymbol{\beta}}^{t-1}\| + \frac{1}{n}\left(\sup_{z}\rho''(z)\right)\|\boldsymbol{X}\|\|\boldsymbol{\eta}^{t-1} - \boldsymbol{\eta}^{t}\|\right\}.$$
(73)

Using (44) together with Borel-Cantelli yields  $\lim_{n\to\infty} ||X||/\sqrt{n} < \infty$  almost surely. Further, it follows from (65) that  $\lim_{n\to\infty} ||\hat{\boldsymbol{\beta}}^t||$  is finite almost surely as  $\tau_* < \infty$ . These taken together with (66), (69) and (70) yield

$$\lim_{t \to \infty} \lim_{n \to \infty} \|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^t\| =_{\text{a.s.}} 0$$
(74)

as claimed.

## 7 Likelihood ratio analysis

This section presents the analytical details for Sect. 2.5.2, which relates the loglikelihood ratio statistic  $\Lambda_i$  with  $\hat{\beta}_i$ . Recall from (37) that the LLR statistic for testing  $\beta_1 = 0$  vs.  $\beta_1 \neq 0$  is given by

$$\Lambda_{1} = \frac{1}{2} \left( \tilde{X} \tilde{\boldsymbol{\beta}} - X \hat{\boldsymbol{\beta}} \right)^{\top} \boldsymbol{D}_{\hat{\boldsymbol{\beta}}} \left( \tilde{X} \tilde{\boldsymbol{\beta}} - X \hat{\boldsymbol{\beta}} \right) + \frac{1}{6} \sum_{i=1}^{n} \rho^{\prime\prime\prime}(\gamma_{i}) \left( \tilde{X}_{i}^{\top} \tilde{\boldsymbol{\beta}} - X_{i}^{\top} \hat{\boldsymbol{\beta}} \right)^{3}, \quad (75)$$

where

$$\boldsymbol{D}_{\hat{\boldsymbol{\beta}}} := \begin{bmatrix} \rho^{\prime\prime} \left( \boldsymbol{X}_{1}^{\top} \hat{\boldsymbol{\beta}} \right) & & \\ & \ddots & \\ & & \rho^{\prime\prime} \left( \boldsymbol{X}_{n}^{\top} \hat{\boldsymbol{\beta}} \right) \end{bmatrix}$$
(76)

and  $\gamma_i$  lies between  $X_i^{\top} \hat{\beta}$  and  $\tilde{X}_i^{\top} \tilde{\beta}$ . The asymptotic distribution of  $\Lambda_1$  claimed in Theorem 3 immediately follows from the result below, whose proof is the subject of the rest of this section.

**Theorem 7** Let  $(\tau_*, b_*)$  be the unique solution to the system of Eqs. (25) and (26), and define

$$\tilde{\boldsymbol{G}} = \frac{1}{n} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}} \quad and \quad \tilde{\boldsymbol{\alpha}} = \frac{1}{n} \operatorname{Tr}(\tilde{\boldsymbol{G}}^{-1}).$$
(77)

Suppose  $p/n \to \kappa \in (0, 1/2)$ . Then (a) the log-likelihood ratio statistic obeys

$$2\Lambda_1 - p\hat{\beta}_1^2 / \tilde{\alpha} \stackrel{\mathbb{P}}{\to} 0; \tag{78}$$

(b) and the scalar  $\tilde{\alpha}$  converges,

$$\tilde{\alpha} \stackrel{\mathbb{P}}{\to} b_*. \tag{79}$$

# 7.1 More notations and preliminaries

Before proceeding, we introduce some notations that will be used throughout. For any matrix X, denote by  $X_{ij}$  and  $X_{.j}$  its (i, j)-th entry and jth column, respectively.

We denote an analogue  $\mathbf{r} = \{r_i\}_{1 \le i \le n}$  (resp.  $\tilde{\mathbf{r}} = \{\tilde{r}_i\}_{1 \le i \le n}$ ) of residuals in the full (resp. reduced) model by

$$r_i := -\rho' (\boldsymbol{X}_i^{\top} \hat{\boldsymbol{\beta}}) \quad \text{and} \quad \tilde{r}_i := -\rho' (\tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}}).$$
 (80)

As in (76), set

$$\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} := \begin{bmatrix} \rho^{\prime\prime}(\tilde{\boldsymbol{X}}_{1}^{\top}\tilde{\boldsymbol{\beta}}) & & \\ & \ddots & \\ & & \rho^{\prime\prime}(\tilde{\boldsymbol{X}}_{n}^{\top}\tilde{\boldsymbol{\beta}}) \end{bmatrix} \text{ and } \boldsymbol{D}_{\hat{\boldsymbol{\beta}},\tilde{\boldsymbol{b}}} := \begin{bmatrix} \rho^{\prime\prime}(\boldsymbol{\gamma}_{1}^{*}) & & \\ & \ddots & \\ & & \rho^{\prime\prime}(\boldsymbol{\gamma}_{n}^{*}), \end{bmatrix},$$
(81)

where  $\gamma_i^*$  is between  $X_i^{\top} \hat{\beta}$  and  $X_i^{\top} \tilde{b}$ , and  $\tilde{b}$  is to be defined later in Sect. 7.2. Further, as in (77), introduce the Gram matrices

$$\boldsymbol{G} := \frac{1}{n} \boldsymbol{X}^{\top} \boldsymbol{D}_{\hat{\boldsymbol{\beta}}} \boldsymbol{X} \text{ and } \boldsymbol{G}_{\hat{\boldsymbol{\beta}}, \tilde{\boldsymbol{b}}} = \frac{1}{n} \boldsymbol{X}^{\top} \boldsymbol{D}_{\hat{\boldsymbol{\beta}}, \tilde{\boldsymbol{b}}} \boldsymbol{X}.$$
(82)

Let  $\tilde{G}_{(i)}$  denote the version of  $\tilde{G}$  without the term corresponding to the *i*th observation, that is,

$$\tilde{\boldsymbol{G}}_{(i)} = \frac{1}{n} \sum_{j: j \neq i} \rho''(\tilde{\boldsymbol{X}}_{j}^{\top} \tilde{\boldsymbol{\beta}}) \tilde{\boldsymbol{X}}_{j} \tilde{\boldsymbol{X}}_{j}^{\top}.$$
(83)

Additionally, let  $\hat{\beta}_{[-i]}$  be the MLE when the *i*th observation is dropped and let  $G_{[-i]}$  be the corresponding Gram matrix,

$$\boldsymbol{G}_{[-i]} = \frac{1}{n} \sum_{j:j \neq i} \rho''(\boldsymbol{X}_j^{\top} \hat{\boldsymbol{\beta}}_{[-i]}) \boldsymbol{X}_j \boldsymbol{X}_j^{\top}.$$
(84)

Further, let  $\tilde{\boldsymbol{\beta}}_{[-i]}$  be the MLE when the first predictor and *i* th observation are removed, i.e.

$$\tilde{\boldsymbol{\beta}}_{[-i]} := \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{p-1}} \sum_{j: j \neq i} \rho(\tilde{\boldsymbol{X}}_j^{\top} \boldsymbol{\beta}).$$

Deringer

Below  $\tilde{G}_{[-i]}$  is the corresponding version of  $\tilde{G}$ ,

$$\tilde{\boldsymbol{G}}_{[-i]} = \frac{1}{n} \sum_{j:j \neq i} \rho''(\tilde{\boldsymbol{X}}_{j}^{\top} \tilde{\boldsymbol{\beta}}_{[-i]}) \tilde{\boldsymbol{X}}_{j} \tilde{\boldsymbol{X}}_{j}^{\top}.$$
(85)

For these different versions of G, their least eigenvalues are all bounded away from 0, as asserted by the following lemma.

**Lemma 7** There exist some absolute constants  $\lambda_{lb}$ , C, c > 0 such that

$$\mathbb{P}(\lambda_{\min}(\boldsymbol{G}) > \lambda_{\mathrm{lb}}) \geq 1 - Ce^{-cn}$$

Moreover, the same result holds for  $\tilde{G}$ ,  $G_{\hat{\beta},\tilde{b}}$ ,  $\tilde{G}_{(i)}$ ,  $G_{[-i]}$  and  $\tilde{G}_{[-i]}$  for all  $i \in [n]$ .

**Proof** This result follows directly from Lemmas 2, 4, and Theorem 4.

Throughout the rest of this section, we restrict ourselves (for any given n) to the following event:

$$\mathcal{A}_{n} := \{\lambda_{\min}(\boldsymbol{G}) > \lambda_{\mathrm{lb}}\} \cap \{\lambda_{\min}(\boldsymbol{G}) > \lambda_{\mathrm{lb}}\} \cap \{\lambda_{\min}(\boldsymbol{G}_{\hat{\boldsymbol{\beta}}, \tilde{\boldsymbol{b}}}) > \lambda_{\mathrm{lb}}\}$$
$$\cap \{\bigcap_{i=1}^{n} \lambda_{\min}(\tilde{\boldsymbol{G}}_{(i)}) > \lambda_{\mathrm{lb}}\} \cap \{\bigcap_{i=1}^{n} \lambda_{\min}(\tilde{\boldsymbol{G}}_{[-i]}) > \lambda_{\mathrm{lb}}\}$$
$$\cap \{\bigcap_{i=1}^{n} \lambda_{\min}(\boldsymbol{G}_{[-i]}) > \lambda_{\mathrm{lb}}\}.$$
(86)

By Lemma 7,  $A_n$  arises with exponentially high probability, i.e.

$$\mathbb{P}(\mathcal{A}_n) \ge 1 - \exp(-\Omega(n)). \tag{87}$$

#### 7.2 A surrogate for the MLE

In view of (75), the main step in controlling  $\Lambda_1$  consists of characterizing the differences  $X\hat{\beta} - \tilde{X}\tilde{\beta}$  or  $\hat{\beta} - \begin{bmatrix} 0\\ \tilde{\beta} \end{bmatrix}$ . Since the definition of  $\hat{\beta}$  is implicit and not amenable to direct analysis, we approximate  $\hat{\beta}$  by a more amenable surrogate  $\tilde{b}$ , an idea introduced in [24–26]. We collect some properties of the surrogate which will prove valuable in the subsequent analysis.

To begin with, our surrogate is

$$\tilde{\boldsymbol{b}} = \begin{bmatrix} 0\\ \tilde{\boldsymbol{\beta}} \end{bmatrix} + \tilde{b}_1 \begin{bmatrix} 1\\ -\tilde{\boldsymbol{G}}^{-1}\boldsymbol{w} \end{bmatrix}, \qquad (88)$$

where  $\tilde{G}$  is defined in (82),

$$\boldsymbol{w} := \frac{1}{n} \sum_{i=1}^{n} \rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) \boldsymbol{X}_{i1} \tilde{\boldsymbol{X}}_{i} = \frac{1}{n} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \boldsymbol{X}_{\cdot 1},$$
(89)

and  $\tilde{b}_1$  is a scalar to be specified later. This vector is constructed in the hope that

$$\hat{\boldsymbol{\beta}} \approx \tilde{\boldsymbol{b}}, \quad \text{or equivalently,} \quad \begin{cases} \hat{\beta}_1 \approx \tilde{b}_1, \\ \hat{\boldsymbol{\beta}}_{2:p} - \tilde{\boldsymbol{\beta}} \approx - \tilde{b}_1 \tilde{\boldsymbol{G}}^{-1} \boldsymbol{w}, \end{cases}$$
(90)

where  $\hat{\beta}_{2:p}$  contains the 2nd through *p*th components of  $\hat{\beta}$ .

Before specifying  $\tilde{b}_1$ , we shall first shed some insights into the remaining terms in  $\tilde{b}$ . By definition,

$$\nabla^{2} \ell \left( \begin{bmatrix} 0 \\ \tilde{\boldsymbol{\beta}} \end{bmatrix} \right) = \boldsymbol{X}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \boldsymbol{X} = \begin{bmatrix} \boldsymbol{X}_{\cdot 1}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \boldsymbol{X}_{\cdot 1} & \boldsymbol{X}_{\cdot 1}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}} \\ \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \boldsymbol{X}_{\cdot 1} & \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}} \end{bmatrix} = \begin{bmatrix} \boldsymbol{X}_{\cdot 1}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \boldsymbol{X}_{\cdot 1} & n \boldsymbol{w}^{\top} \\ n \boldsymbol{w} & n \tilde{\boldsymbol{G}} \end{bmatrix}.$$

Employing the first-order approximation of  $\nabla \ell(\cdot)$  gives

$$\nabla^{2}\ell\left(\begin{bmatrix}0\\\tilde{\boldsymbol{\beta}}\end{bmatrix}\right)\left(\hat{\boldsymbol{\beta}}-\begin{bmatrix}0\\\tilde{\boldsymbol{\beta}}\end{bmatrix}\right)\approx\nabla\ell(\hat{\boldsymbol{\beta}})-\nabla\ell\left(\begin{bmatrix}0\\\tilde{\boldsymbol{\beta}}\end{bmatrix}\right).$$
(91)

Suppose  $\hat{\boldsymbol{\beta}}_{2:p}$  is well approximated by  $\tilde{\boldsymbol{\beta}}$ . Then all but the first coordinates of  $\nabla \ell(\tilde{\boldsymbol{\beta}})$  and  $\nabla \ell \left( \begin{bmatrix} 0 \\ \hat{\boldsymbol{\beta}} \end{bmatrix} \right)$  are also very close to each other. Therefore, taking the 2nd through *p*th components of (91) and approximating them by zero give

$$\begin{bmatrix} \boldsymbol{w}, \, \tilde{\boldsymbol{G}} \end{bmatrix} \left( \hat{\boldsymbol{\beta}} - \begin{bmatrix} 0 \\ \tilde{\boldsymbol{\beta}} \end{bmatrix} \right) \approx \boldsymbol{0}$$

This together with a little algebra yields

$$\hat{\boldsymbol{\beta}}_{2:p} - \tilde{\boldsymbol{\beta}} \approx -\hat{\beta}_1 \tilde{\boldsymbol{G}}^{-1} \boldsymbol{w} \approx -\tilde{b}_1 \tilde{\boldsymbol{G}}^{-1} \boldsymbol{w},$$

which coincides with (90). In fact, for all but the 1st entries,  $\tilde{b}$  is constructed by moving  $\tilde{\beta}$  one-step in the direction which takes it closest to  $\hat{\beta}$ .

Next, we come to discussing the scalar  $\tilde{b}_1$ . Introduce the projection matrix

$$\boldsymbol{H} := \boldsymbol{I} - \frac{1}{n} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{G}}^{-1} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}, \qquad (92)$$

and define  $\tilde{b}_1$  as

$$\tilde{b}_1 := \frac{X_{\cdot 1}^{\top} \tilde{\boldsymbol{r}}}{X_{\cdot 1}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \boldsymbol{H} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} X_{\cdot 1}},$$
(93)

where  $\tilde{r}$  comes from (80). In fact, the expression  $\tilde{b}_1$  is obtained through similar (but slightly more complicated) first-order approximation as for  $\tilde{b}_{2:p}$ , in order to make sure that  $b_1 \approx \hat{\beta}_1$ ; see [26, Pages 14560–14561] for a detailed description.

521

🖄 Springer

We now formally justify that the surrogate  $\tilde{b}$  and the MLE  $\hat{\beta}$  are close to each other.

**Theorem 8** The MLE  $\hat{\boldsymbol{\beta}}$  and the surrogate  $\tilde{\boldsymbol{b}}$  (88) obey

$$\|\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{b}}\| \lesssim n^{-1+o(1)},\tag{94}$$

$$|\tilde{b}_1| \lesssim n^{-1/2 + o(1)},$$
 (95)

and

$$\sup_{1 \le i \le n} |\boldsymbol{X}_i^{\top} \tilde{\boldsymbol{b}} - \tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}}| \lesssim n^{-1/2 + o(1)}$$
(96)

with probability tending to one as  $n \to \infty$ .

## Proof See Sect. 7.4.

The global accuracy (94) immediately leads to a coordinate-wise approximation result between  $\hat{\beta}_1$  and  $\tilde{b}_1$ .

**Corollary 3** With probability tending to one as  $n \to \infty$ ,

$$\sqrt{n}|\tilde{b}_1 - \hat{\beta}_1| \lesssim n^{-1/2 + o(1)}.$$
 (97)

Another consequence from Theorem 8 is that the value  $X_i^{\top} \hat{\beta}$  in the full model and its counterpart  $\tilde{X}_i^{\top} \tilde{\beta}$  in the reduced model are uniformly close.

**Corollary 4** The values  $X_i^{\top}\hat{\beta}$  and  $\tilde{X}_i^{\top}\tilde{\beta}$  are uniformly close in the sense that

$$\sup_{1 \le i \le n} \left| \boldsymbol{X}_i^{\top} \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}} \right| \lesssim n^{-1/2 + o(1)}$$
(98)

*holds with probability approaching one as*  $n \to \infty$ *.* 

Proof Note that

$$\sup_{1 \le i \le n} \left| X_i^\top \hat{\boldsymbol{\beta}} - \tilde{X}_i^\top \tilde{\boldsymbol{\beta}} \right| \le \sup_{1 \le i \le n} \left| X_i^\top (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{b}}) \right| + \sup_{1 \le i \le n} \left| X_i^\top \tilde{\boldsymbol{b}} - \tilde{X}_i^\top \tilde{\boldsymbol{\beta}} \right|.$$

The second term in the right-hand side is upper bounded by  $n^{-1/2+o(1)}$  with probability 1 - o(1) according to Theorem 8. Invoking Lemma 2 and Theorem 8 and applying Cauchy-Schwarz inequality yield that the first term is  $O(n^{-1/2+o(1)})$  with probability 1 - o(1). This establishes the claim.
#### 7.3 Analysis of the likelihood-ratio statistic

We are now positioned to use our surrogate  $\hat{b}$  to analyze the likelihood-ratio statistic. In this subsection we establish Theorem 7(a). The proof for Theorem 7(b) is deferred to Appendix I.

Recall from (37) that

$$2\Lambda_1 = (\tilde{X}\tilde{\boldsymbol{\beta}} - X\hat{\boldsymbol{\beta}})^{\top} \boldsymbol{D}_{\hat{\boldsymbol{\beta}}} (\tilde{X}\tilde{\boldsymbol{\beta}} - X\hat{\boldsymbol{\beta}}) + \underbrace{\frac{1}{3}\sum_{i=1}^n \rho^{\prime\prime\prime\prime}(\gamma_i)(\tilde{X}_i^{\top}\tilde{\boldsymbol{\beta}} - X_i^{\top}\hat{\boldsymbol{\beta}})^3}_{:=I_3}.$$

To begin with, Corollary 4 together with the assumption  $\sup_{z} \rho'''(z) < \infty$  implies that

$$I_3 \lesssim n^{-1/2+o(1)}$$

with probability 1 - o(1). Hence,  $I_3$  converges to zero in probability.

Reorganize the quadratic term as follows:

$$(\tilde{X}\tilde{\beta} - X\hat{\beta})^{\top} D_{\hat{\beta}}(\tilde{X}\tilde{\beta} - X\hat{\beta}) = \sum_{i} \rho''(X_{i}^{\top}\hat{\beta}) \left(X_{i}^{\top}\hat{\beta} - \tilde{X}_{i}^{\top}\tilde{\beta}\right)^{2}$$

$$= \sum_{i} \rho''(X_{i}^{\top}\hat{\beta}) \left[X_{i}^{\top}(\hat{\beta} - \tilde{b}) + (X_{i}^{\top}\tilde{b} - \tilde{X}_{i}^{\top}\tilde{\beta})\right]^{2}$$

$$= \sum_{i} \rho''(X_{i}^{\top}\hat{\beta})(X_{i}^{\top}(\hat{\beta} - \tilde{b}))^{2} + 2\sum_{i} \rho''(X_{i}^{\top}\hat{\beta})X_{i}^{\top}(\hat{\beta} - \tilde{b})(X_{i}^{\top}\tilde{b} - \tilde{X}_{i}^{\top}\tilde{\beta})$$

$$+ \sum_{i} \rho''(X_{i}^{\top}\hat{\beta}) \left(X_{i}^{\top}\tilde{b} - \tilde{X}_{i}^{\top}\tilde{\beta}\right)^{2}.$$
(99)

We control each of the three terms in the right-hand side of (99).

• Since  $\sup_{z} \rho''(z) < \infty$ , the first term in (99) is bounded by

$$\sum_{i} \rho''(\boldsymbol{X}_{i}^{\top} \hat{\boldsymbol{\beta}}) (\boldsymbol{X}_{i}^{\top} (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{b}}))^{2} \lesssim ||\tilde{\boldsymbol{\beta}} - \tilde{\boldsymbol{b}}||^{2} \left\| \sum_{i} \boldsymbol{X}_{i} \boldsymbol{X}_{i}^{\top} \right\| \lesssim n^{-1+o(1)}$$

with probability 1 - o(1), by an application of Theorem 8 and Lemma 2.

• From the definition of  $\mathbf{b}$ , the second term can be upper bounded by

$$2\sum_{i} \rho''(\boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}})(\hat{\boldsymbol{\beta}}-\tilde{\boldsymbol{b}})^{\top}\boldsymbol{X}_{i}\boldsymbol{X}_{i}^{\top}\tilde{\boldsymbol{b}}_{1}\begin{bmatrix}1\\-\tilde{\boldsymbol{G}}^{-1}\boldsymbol{w}\end{bmatrix}$$
$$\leq |\tilde{\boldsymbol{b}}_{1}| \cdot \|\hat{\boldsymbol{\beta}}-\tilde{\boldsymbol{b}}\| \cdot \left\|\sum_{i}\boldsymbol{X}_{i}\boldsymbol{X}_{i}^{\top}\right\| \cdot \sqrt{1+\boldsymbol{w}^{\top}\tilde{\boldsymbol{G}}^{-2}\boldsymbol{w}}$$
$$\lesssim n^{-\frac{1}{2}+o(1)}$$

Deringer

with probability 1 - o(1), where the last line follows from a combination of Theorem 8, Lemma 2 and the following lemma.

**Lemma 8** Let  $\tilde{G}$  and w be as defined in (82) and (89), respectively. Then

$$\mathbb{P}\left(\boldsymbol{w}^{\top}\tilde{\boldsymbol{G}}^{-2}\boldsymbol{w} \lesssim 1\right) \geq 1 - \exp(-\Omega(n)).$$
(100)

Proof See Appendix E.

• The third term in (99) can be decomposed as

$$\sum_{i} \rho''(\boldsymbol{X}_{i}^{\top}\boldsymbol{\hat{\beta}})(\boldsymbol{X}_{i}^{\top}\boldsymbol{\tilde{b}} - \boldsymbol{\tilde{X}}_{i}^{\top}\boldsymbol{\tilde{\beta}}))^{2}$$

$$= \sum_{i} \left( \rho''(\boldsymbol{X}_{i}^{\top}\boldsymbol{\hat{\beta}}) - \rho''(\boldsymbol{\tilde{X}}_{i}^{\top}\boldsymbol{\tilde{\beta}}) \right) (\boldsymbol{X}_{i}^{\top}\boldsymbol{\tilde{b}} - \boldsymbol{\tilde{X}}_{i}^{\top}\boldsymbol{\tilde{\beta}}))^{2}$$

$$+ \sum_{i} \rho''(\boldsymbol{\tilde{X}}_{i}^{\top}\boldsymbol{\tilde{\beta}})(\boldsymbol{X}_{i}^{\top}\boldsymbol{\tilde{b}} - \boldsymbol{\tilde{X}}_{i}^{\top}\boldsymbol{\tilde{\beta}})^{2}$$

$$= \sum_{i} \rho'''(\boldsymbol{\tilde{\gamma}}_{i})(\boldsymbol{X}_{i}^{\top}\boldsymbol{\hat{\beta}} - \boldsymbol{\tilde{X}}_{i}^{\top}\boldsymbol{\tilde{\beta}}) \left(\boldsymbol{X}_{i}^{\top}\boldsymbol{\tilde{b}} - \boldsymbol{\tilde{X}}_{i}^{\top}\boldsymbol{\tilde{\beta}}\right)^{2}$$

$$+ \sum_{i} \rho''(\boldsymbol{\tilde{X}}_{i}^{\top}\boldsymbol{\tilde{\beta}}) \left(\boldsymbol{X}_{i}^{\top}\boldsymbol{\tilde{b}} - \boldsymbol{\tilde{X}}_{i}^{\top}\boldsymbol{\tilde{\beta}}\right)^{2}$$
(101)

for some  $\tilde{\gamma}_i$  between  $X_i^{\top}\hat{\beta}$  and  $\tilde{X}_i^{\top}\tilde{\beta}$ . From Theorem 8 and Corollary 4, the first term in (101) is  $O(n^{-1/2+o(1)})$  with probability 1-o(1). Hence, the only remaining term is the second.

In summary, we have

$$2\Lambda_{1} - \underbrace{\sum_{i} \rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) \left(\boldsymbol{X}_{i}^{\top} \tilde{\boldsymbol{b}} - \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}\right)^{2}}_{=\boldsymbol{v}^{\top} \boldsymbol{X}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \boldsymbol{X} \boldsymbol{v}} \overset{\mathbb{P}}{\to} 0, \qquad (102)$$

where  $\boldsymbol{v} := \tilde{b}_1 \begin{bmatrix} 1 \\ -\tilde{\boldsymbol{G}}^{-1} \boldsymbol{w} \end{bmatrix}$  according to (88). On simplification, the quadratic form reduces to

$$\boldsymbol{v}^{\top}\boldsymbol{X}^{\top}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}\boldsymbol{X}\boldsymbol{v} = \tilde{b}_{1}^{2}\left(\boldsymbol{X}_{\cdot 1} - \tilde{\boldsymbol{X}}\tilde{\boldsymbol{G}}^{-1}\boldsymbol{w}\right)^{\top}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}\left(\boldsymbol{X}_{\cdot 1} - \tilde{\boldsymbol{X}}\tilde{\boldsymbol{G}}^{-1}\boldsymbol{w}\right)$$
$$= \tilde{b}_{1}^{2}\left(\boldsymbol{X}_{\cdot 1}^{\top}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}\boldsymbol{X}_{\cdot 1} - 2\boldsymbol{X}_{\cdot 1}^{\top}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}\tilde{\boldsymbol{X}}\tilde{\boldsymbol{G}}^{-1}\boldsymbol{w} + \boldsymbol{w}^{\top}\tilde{\boldsymbol{G}}^{-1}\tilde{\boldsymbol{X}}^{\top}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}\tilde{\boldsymbol{X}}\tilde{\boldsymbol{G}}^{-1}\boldsymbol{w}\right)$$

Deringer

$$= \tilde{b}_{1}^{2} \left( \boldsymbol{X}_{\cdot 1}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \boldsymbol{X}_{\cdot 1} - n \boldsymbol{w}^{\top} \tilde{\boldsymbol{G}}^{-1} \boldsymbol{w} \right)$$
$$= n \tilde{b}_{1}^{2} \left( \underbrace{\frac{1}{n} \boldsymbol{X}_{\cdot 1}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \boldsymbol{H} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \boldsymbol{X}_{\cdot 1}}_{:=\xi} \right),$$

recalling the definitions (82), (89), and (92). Hence, the log-likelihood ratio  $2\Lambda_1$  simplifies to  $n\tilde{b}_1^2\xi + o_P(1)$  on  $\mathcal{A}_n$ .

Finally, rewrite  $\boldsymbol{v}^{\top} \boldsymbol{X}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \boldsymbol{X} \boldsymbol{v}$  as  $n(\tilde{b}_1^2 - \hat{\beta}_1^2)\boldsymbol{\xi} + n\hat{\beta}_1^2\boldsymbol{\xi}$ . To analyze the first term, note that

$$n|\tilde{b}_{1}^{2} - \hat{\beta}_{1}^{2}| = n|\tilde{b}_{1} - \hat{\beta}_{1}| \cdot |\tilde{b}_{1} + \hat{\beta}_{1}| \le n|\tilde{b}_{1} - \hat{\beta}_{1}|^{2} + 2n|\tilde{b}_{1}| \cdot |\tilde{b}_{1} - \hat{\beta}_{1}| \lesssim n^{-\frac{1}{2} + o(1)}$$
(103)

with probability 1 - o(1) in view of Theorem 8 and Corollary 3. It remains to analyze  $\xi$ . Recognize that  $X_{\cdot 1}$  is independent of  $D_{\tilde{\beta}}^{1/2} H D_{\tilde{\beta}}^{1/2}$ . Applying the Hanson-Wright inequality [32,52] and the Sherman-Morrison-Woodbury formula (e.g. [31]) leads to the following lemma:

Lemma 9 Let  $\tilde{\alpha} = \frac{1}{n} \operatorname{Tr}(\tilde{\boldsymbol{G}}^{-1})$ , where  $\tilde{\boldsymbol{G}} = \frac{1}{n} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}}$ . Then one has

$$\left|\frac{p-1}{n} - \tilde{\alpha} \frac{1}{n} \boldsymbol{X}_{\cdot 1}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \boldsymbol{H} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \boldsymbol{X}_{\cdot 1}\right| \lesssim n^{-1/2 + o(1)}$$
(104)

with probability approaching one as  $n \to \infty$ .

#### **Proof** See Appendix **F**.

In addition, if one can show that  $\tilde{\alpha}$  is bounded away from zero with probability 1-o(1), then it is seen from Lemma 9 that

$$\xi - \frac{p}{n\tilde{\alpha}} \stackrel{\mathbb{P}}{\to} 0. \tag{105}$$

To justify the above claim, we observe that since  $\rho''$  is bounded,  $\lambda_{\max}(\tilde{G}) \lesssim \lambda_{\max}(\tilde{X}^{\top}\tilde{X})/n \lesssim 1$  with exponentially high probability (Lemma 2). This yields

$$\tilde{\alpha} = \mathrm{Tr}(\tilde{\boldsymbol{G}}^{-1})/n \gtrsim p/n$$

with probability 1 - o(1). On the other hand, on  $A_n$  one has

$$\tilde{\alpha} \leq p/(n\lambda_{\min}(\tilde{G})) \lesssim p/n.$$

Hence, it follows that  $\xi = \Omega(1)$  with probability 1 - o(1). Putting this together with (103) gives the approximation

$$\boldsymbol{v}^{\top}\boldsymbol{X}^{\top}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}\boldsymbol{X}\boldsymbol{v} = n\hat{\beta}_{1}^{2}\boldsymbol{\xi} + o(1).$$
(106)

🖄 Springer

Taken collectively (102), (105) and (106) yields the desired result

$$2\Lambda_1 - p\hat{\beta}_1^2/\tilde{\alpha} \stackrel{\mathbb{P}}{\to} 0.$$

#### 7.4 Proof of Theorem 8

This subsection outlines the main steps for the proof of Theorem 8. To begin with, we shall express the difference  $\hat{\beta} - \tilde{b}$  in terms of the gradient of the negative log-likelihood function. Note that  $\nabla \ell(\hat{\beta}) = 0$ , and hence

$$\nabla \ell(\tilde{\boldsymbol{b}}) = \nabla \ell(\tilde{\boldsymbol{b}}) - \nabla \ell(\hat{\boldsymbol{\beta}}) = \sum_{i=1}^{n} X_{i} [\rho'(\boldsymbol{X}_{i}^{\top} \tilde{\boldsymbol{b}}) - \rho'(\boldsymbol{X}_{i}' \hat{\boldsymbol{\beta}})]$$
$$= \sum_{i=1}^{n} \rho''(\gamma_{i}^{*}) X_{i} X_{i}^{\top} (\tilde{\boldsymbol{b}} - \hat{\boldsymbol{\beta}}),$$

where  $\gamma_i^*$  is between  $X_i^{\top} \hat{\beta}$  and  $X_i^{\top} \tilde{b}$ . Recalling the notation introduced in (82), this can be rearranged as

$$\tilde{\boldsymbol{b}} - \hat{\boldsymbol{\beta}} = \frac{1}{n} \boldsymbol{G}_{\hat{\boldsymbol{\beta}}, \tilde{\boldsymbol{b}}}^{-1} \nabla \ell(\tilde{\boldsymbol{b}}).$$

Hence, on  $\mathcal{A}_n$ , this yields

$$\|\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{b}}\| \le \frac{\|\nabla \ell(\boldsymbol{b})\|}{\lambda_{\text{lb}} n}.$$
(107)

The next step involves expressing  $\nabla \ell(\tilde{\boldsymbol{b}})$  in terms of the difference  $\tilde{\boldsymbol{b}} - \begin{vmatrix} 0\\ \tilde{\boldsymbol{\beta}} \end{vmatrix}$ .

**Lemma 10** On the event  $A_n$  (86), the negative log-likelihood evaluated at the surrogate  $\tilde{b}$  obeys

$$\nabla \ell(\tilde{\boldsymbol{b}}) = \sum_{i=1}^{n} \left[ \rho''(\gamma_{i}^{*}) - \rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}}) \right] \boldsymbol{X}_{i} \boldsymbol{X}_{i}^{\top} \left( \tilde{\boldsymbol{b}} - \begin{bmatrix} 0\\ \tilde{\boldsymbol{\beta}} \end{bmatrix} \right),$$

where  $\gamma_i^*$  is some quantity between  $X_i^{\top} \tilde{\boldsymbol{b}}$  and  $\tilde{X}_i^{\top} \tilde{\boldsymbol{\beta}}$ .

**Proof** The proof follows exactly the same argument as in the proof of [25, Proposition 3.11], and is thus omitted.  $\Box$ 

The point of expressing  $\nabla \ell(\tilde{b})$  in this way is that the difference  $\tilde{b} - \begin{bmatrix} 0\\ \tilde{\beta} \end{bmatrix}$  is known explicitly from the definition of  $\tilde{b}$ . Invoking Lemma 10 and the definition (88) allows one to further upper bound (107) as

$$\|\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{b}}\| \lesssim \frac{1}{n} \|\nabla \ell(\tilde{\boldsymbol{b}})\| \lesssim \sup_{i} \left|\rho''(\gamma_{i}^{*}) - \rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})\right| \left\|\frac{1}{n}\sum_{i=1}^{n}\boldsymbol{X}_{i}\boldsymbol{X}_{i}^{\top}\right\| \left\|\tilde{\boldsymbol{b}} - \begin{bmatrix}\boldsymbol{0}\\\tilde{\boldsymbol{\beta}}\end{bmatrix}\right\|$$
$$\lesssim \sup_{i} \left|\boldsymbol{X}_{i}^{\top}\tilde{\boldsymbol{b}} - \tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}}\right| |\rho'''|_{\infty}$$
$$\times \left\|\frac{1}{n}\sum_{i=1}^{n}\boldsymbol{X}_{i}\boldsymbol{X}_{i}^{\top}\right\| \cdot |\tilde{\boldsymbol{b}}_{1}|\sqrt{1 + \boldsymbol{w}^{\top}\tilde{\boldsymbol{G}}^{-2}\boldsymbol{w}}$$
$$\lesssim |\tilde{\boldsymbol{b}}_{1}|\sup_{i} \left|\boldsymbol{X}_{i}^{\top}\tilde{\boldsymbol{b}} - \tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}}\right|$$
(108)

with probability at least  $1 - \exp(-\Omega(n))$ . The last inequality here comes from our assumption that  $\sup_{z} |\rho'''(z)| < \infty$  together with Lemmas 2 and 8.

In order to bound (108), we first make use of the definition of  $\hat{b}$  to reach

$$\sup_{i} \left| \boldsymbol{X}_{i}^{\top} \tilde{\boldsymbol{b}} - \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}} \right| = |\tilde{b}_{1}| \sup_{i} |\boldsymbol{X}_{i1} - \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}^{-1} \boldsymbol{w}|.$$
(109)

The following lemma provides an upper bound on  $\sup_i |X_{i1} - \tilde{X}_i^{\top} \tilde{\boldsymbol{G}}^{-1} \boldsymbol{w}|$ .

Lemma 11 With  $\tilde{G}$  and w as defined in (82) and (89),

$$\mathbb{P}\left(\sup_{1\leq i\leq n} \left|X_{i1} - \tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}^{-1}\boldsymbol{w}\right| \leq n^{o(1)}\right) \geq 1 - o(1).$$
(110)

**Proof** See Appendix G.

In view of Lemma 11, the second term in the right-hand side of (109) is bounded above by  $n^{o(1)}$  with high probability. Thus, in both the bounds (108) and (109), it only remains to analyze the term  $\tilde{b}_1$ . To this end, we control the numerator and the denominator of  $\tilde{b}_1$  separately.

• Recall from the definition (93) that the numerator of  $\tilde{b}_1$  is given by  $X_{.1}^{\top}\tilde{r}$  and that  $\tilde{r}$  is independent of  $X_{.1}$ . Thus, conditional on  $\tilde{X}$ , the quantity  $X_{.1}^{\top}\tilde{r}$  is distributed as a Gaussian with mean zero and variance

$$\sigma^2 = \sum_{i=1}^n \left( \rho'(\tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}}) \right)^2.$$

Since  $|\rho'(x)| = O(|x|)$ , the variance is bounded by

$$\sigma^{2} \lesssim \tilde{\boldsymbol{\beta}}^{\top} \left( \sum_{i=1}^{n} \tilde{\boldsymbol{X}}_{i} \tilde{\boldsymbol{X}}_{i}^{\top} \right) \tilde{\boldsymbol{\beta}} \lesssim n \| \tilde{\boldsymbol{\beta}} \|^{2} \lesssim n$$
(111)

with probability at least  $1 - \exp(-\Omega(n))$ , a consequence from Theorem 4 and Lemma 2. Therefore, with probability 1 - o(1), we have

$$\frac{1}{\sqrt{n}} X_{\cdot 1}^{\top} \tilde{\boldsymbol{r}} \lesssim n^{o(1)}.$$
(112)

Deringer

• We now move on to the denominator of  $\tilde{b}_1$  in (93). In the discussion following Lemma 9 we showed  $\frac{1}{n} X_{\cdot 1}^{\top} D_{\tilde{\beta}}^{1/2} H D_{\tilde{\beta}}^{1/2} X_{\cdot 1} = \Omega(1)$  with probability 1 - o(1).

Putting the above bounds together, we conclude

$$\mathbb{P}\left(|\tilde{b}_1| \lesssim n^{-\frac{1}{2} + o(1)}\right) = 1 - o(1).$$
(113)

Substitution into (108) and (109) yields

 $\|\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{b}}\| \lesssim n^{-1+o(1)}$  and  $\sup_{i} \left| \boldsymbol{X}_{i}^{\top} \tilde{\boldsymbol{b}} - \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}} \right| \lesssim n^{-1/2+o(1)}$ 

with probability 1 - o(1) as claimed.

## 8 Discussion

In this paper, we derived the high-dimensional asymptotic distribution of the LLR under our modelling assumptions. In particular, we showed that the LLR is inflated vis a vis the classical Wilks' approximation and that this inflation grows as the dimensionality  $\kappa$  increases. This inflation is typical of high-dimensional problems, and one immediate practical consequence is that it explains why classically computed *p*-values are completely off since they tend to be far too small under the null hypothesis. In contrast, we have shown in our simulations that our new limiting distribution yields reasonably accurate *p*-values in finite samples. Having said this, our work raises a few important questions that we have not answered and we conclude this paper with a couple of them.

- We expect that our results continue to hold when the covariates are not normally distributed; see Sect. 3 for some numerical evidence in this direction. To be more precise, we expect the same limiting distribution to hold when the variables are simply sub-Gaussian so that our approximation would enjoy a form of universal validity. Notably, all three sets of analysis techniques we employ in this work have already been extended to accommodate non-Gaussian designs for other problems: the universality law of several convex geometry results has been established via the Lindeberg principle [46]; the AMP machinery has been extended to non-Gaussian designs via the moment method (e.g. [6,21]); and the leave-one-out analysis developed can readily cope with other types of i.i.d. random designs as well [25]. Taken collectively, these suggest a possible path to establish the universality of our results.
- A major limitation of our work is the fact that our limiting distribution holds under the global null; that is, under the assumption that all the regression coefficients vanish. A natural question is, therefore, *how would the distribution change in the case where the coefficients are not all zero?* Consider for instance the setup where the empirical distribution of the regression coefficients satisfies

$$\frac{1}{p}\sum_{i=1}^{p}\delta_{\beta_{i}}\stackrel{\mathrm{d}}{\to}\Pi,$$

where  $\Pi$  is some limiting distribution. How is the distribution of LLR affected by  $\Pi$ ? While this manuscript was under review, the first and third authors have established, among other results, an explicit characterization for the existence of the MLE and the limiting distribution of the LRT in this setup [57]. The proofs in this manuscript are the foundation of this follow-up work, although several new ideas are needed. We refer the interested reader to [56] for an in-depth explanation. In sum, the insights from this paper are essential for establishing a complete likelihood theory.

Acknowledgements E. C. was partially supported by the Office of Naval Research under grant N00014-16-1-2712, and by the Math + X Award from the Simons Foundation. P. S. was partially supported by the Ric Weiland Graduate Fellowship in the School of Humanities and Sciences, Stanford University. Y. C. is supported in part by the AFOSR YIP award FA9550-19-1-0030, by the ARO grant W911NF-18-1-0303, and by the Princeton SEAS innovation award. P. S. and Y. C. are grateful to Andrea Montanari for his help in understanding AMP and [22]. Y. C. thanks Kaizheng Wang and Cong Ma for helpful discussion about [25], and P. S. thanks Subhabrata Sen for several helpful discussions regarding this project. E. C. would like to thank Iain Johnstone for a helpful discussion as well.

#### A Proofs for eigenvalue bounds

#### A.1 Proof of Lemma 3

Fix  $\epsilon \ge 0$  sufficiently small. For any given  $S \subseteq [n]$  obeying  $|S| = (1 - \epsilon)n$  and  $0 \le t \le \sqrt{1 - \epsilon} - \sqrt{p/n}$  it follows from [65, Corollary 5.35] that

$$\lambda_{\min}\left(\frac{1}{n}\sum_{i\in S}X_{i}X_{i}^{\top}\right) < \frac{1}{n}\left(\sqrt{|S|} - \sqrt{p} - t\sqrt{n}\right)^{2} = \left(\sqrt{1-\epsilon} - \sqrt{\frac{p}{n}} - t\right)^{2}$$

holds with probability at most  $2 \exp\left(-\frac{t^2|S|}{2}\right) = 2 \exp\left(-\frac{(1-\epsilon)t^2n}{2}\right)$ . Taking the union bound over all possible subsets *S* of size  $(1-\epsilon)n$  gives

$$\mathbb{P}\left\{\exists S \subseteq [n] \text{ with } |S| = (1 - \epsilon)n \quad \text{s.t.} \quad \frac{1}{n}\lambda_{\min}\left(\sum_{i \in S} X_i X_i^{\top}\right) \\ < \left(\sqrt{1 - \epsilon} - \sqrt{\frac{p}{n}} - t\right)^2\right\} \\ \leq \left(\binom{n}{(1 - \epsilon)n} 2 \exp\left(-\frac{(1 - \epsilon)t^2n}{2}\right) \\ \leq 2 \exp\left(nH\left(\epsilon\right) - \frac{(1 - \epsilon)t^2}{2}n\right),$$

Deringer

where the last line is a consequence of the inequality  $\binom{n}{(1-\epsilon)n} \leq e^{nH(\epsilon)}$  [19, Example 11.1.3].

#### A.2 Proof of Lemma 4

Define

$$S_B(\boldsymbol{\beta}) := \left\{ i : |\boldsymbol{X}_i^{\top} \boldsymbol{\beta}| \le B \| \boldsymbol{\beta} \| \right\}$$

for any B > 0 and any  $\beta$ . Then

$$\sum_{i=1}^{n} \rho'' \left( X_{i}^{\top} \boldsymbol{\beta} \right) X_{i} X_{i}^{\top} \succeq \sum_{i \in S_{B}(\boldsymbol{\beta})} \rho'' \left( X_{i}^{\top} \boldsymbol{\beta} \right) X_{i} X_{i}^{\top} \succeq \inf_{z: |z| \le B \|\boldsymbol{\beta}\|} \rho''(z) \sum_{i \in S_{B}(\boldsymbol{\beta})} X_{i} X_{i}^{\top}.$$

If one also has  $|S_B(\beta)| \ge (1 - \epsilon)n$  (for  $\epsilon \ge 0$  sufficiently small), then this together with Lemma 3 implies that

$$\frac{1}{n}\sum_{i=1}^{n}\rho''\left(\boldsymbol{X}_{i}^{\top}\boldsymbol{\beta}\right)\boldsymbol{X}_{i}\boldsymbol{X}_{i}^{\top}\succeq\inf_{\boldsymbol{z}:|\boldsymbol{z}|\leq\boldsymbol{B}\|\boldsymbol{\beta}\|}\rho''(\boldsymbol{z})\left(\sqrt{1-\epsilon}-\sqrt{\frac{p}{n}}-t\right)^{2}\boldsymbol{I}$$

with probability at least  $1 - 2 \exp\left(-\left(\frac{(1-\epsilon)t^2}{2} - H(\epsilon)\right)n\right)$ .

Thus if we can ensure that with high probability,  $|S_B(\beta)| \ge (1 - \epsilon)n$  holds simultaneously for all  $\beta$ , then we are done. From Lemma 2 we see that  $\frac{1}{n} ||X^{\top}X|| \le 9$  with probability exceeding  $1 - 2 \exp(-n/2)$ . On this event,

$$\|\boldsymbol{X}\boldsymbol{\beta}\|^2 \le 9n\|\boldsymbol{\beta}\|^2, \quad \forall \boldsymbol{\beta}.$$
(114)

On the other hand, the definition of  $S_B(\beta)$  gives

$$\|\boldsymbol{X}\boldsymbol{\beta}\|^{2} \geq \sum_{i \notin S_{B}(\boldsymbol{\beta})} \left|\boldsymbol{X}_{i}^{\top}\boldsymbol{\beta}\right|^{2} \geq \left(n - |S_{B}(\boldsymbol{\beta})|\right) (B\|\boldsymbol{\beta}\|)^{2} = n \left(1 - \frac{|S_{B}(\boldsymbol{\beta})|}{n}\right) B^{2}\|\boldsymbol{\beta}\|^{2}.$$
(115)

Taken together, (114) and (115) yield

$$|S_B(\boldsymbol{\beta})| \ge \left(1 - \frac{9}{B^2}\right)n, \quad \forall \boldsymbol{\beta}$$

with probability at least  $1-2 \exp(-n/2)$ . Therefore, with probability  $1-2 \exp(-n/2)$ ,  $\left|S_{3/\sqrt{\epsilon}}(\boldsymbol{\beta})\right| \geq (1-\epsilon)n$  holds simultaneously for all  $\boldsymbol{\beta}$ . Putting the above results together and setting  $t = 2\sqrt{\frac{H(\epsilon)}{1-\epsilon}}$  give

$$\sum_{i=1}^{n} \rho'' \left( \boldsymbol{X}_{i}^{\top} \boldsymbol{\beta} \right) \boldsymbol{X}_{i} \boldsymbol{X}_{i}^{\top} \succeq \inf_{\substack{z: |z| \leq \frac{3\|\boldsymbol{\beta}\|}{\sqrt{\epsilon}}}} \rho''(z) \left( \sqrt{1-\epsilon} - \sqrt{\frac{p}{n}} - 2\sqrt{\frac{H(\epsilon)}{1-\epsilon}} \right)^{2} \boldsymbol{A}_{i}$$

simultaneously for all  $\beta$  with probability at least  $1 - 2 \exp(-nH(\epsilon)) - 2 \exp(-n/2)$ .

# B Proof of Lemma 5

Applying an integration by parts leads to

$$\mathbb{E}\left[\Psi'(\tau Z; b)\right] = \int_{-\infty}^{\infty} \Psi'(\tau z; b)\phi(z)dz = \frac{1}{\tau}\Psi(\tau z; b)\phi(z)\Big|_{-\infty}^{\infty}$$
$$-\frac{1}{\tau}\int_{-\infty}^{\infty}\Psi(\tau z; b)\phi'(z)dz$$
$$= -\frac{1}{\tau}\int_{-\infty}^{\infty}\Psi(\tau z; b)\phi'(z)dz$$

with  $\phi(z) = \frac{1}{\sqrt{2\pi}} \exp(-z^2/2)$ . This reveals that

$$G'(b) = -\frac{1}{\tau} \int_{-\infty}^{\infty} \frac{\partial \Psi(\tau z; b)}{\partial b} \phi'(z) dz = -\frac{1}{\tau} \int_{-\infty}^{\infty} \frac{\rho'\left(\operatorname{prox}_{b\rho}(\tau z)\right)}{1 + b\rho''\left(\operatorname{prox}_{b\rho}(\tau z)\right)} \phi'(z) dz$$
$$= \frac{1}{\tau} \int_{0}^{\infty} \left( \frac{\rho'\left(\operatorname{prox}_{b\rho}(-\tau z)\right)}{1 + x\rho''\left(\operatorname{prox}_{b\rho}(-\tau z)\right)} - \frac{\rho'\left(\operatorname{prox}_{b\rho}(\tau z)\right)}{1 + x\rho''\left(\operatorname{prox}_{b\rho}(\tau z)\right)} \right) \phi'(z) dz,$$
(116)

where the second identity comes from [22, Proposition 6.4], and the last identity holds since  $\phi'(z) = -\phi'(-z)$ .

Next, we claim that

(a) The function h (z) := μ'(z)/(1+bρ"(z)) is increasing in z;
(b) prox<sub>bρ</sub>(z) is increasing in z.

These two claims imply that

$$\frac{\rho'\left(\operatorname{prox}_{b\rho}(-\tau z)\right)}{1+b\rho''\left(\operatorname{prox}_{b\rho}(-\tau z)\right)} - \frac{\rho'\left(\operatorname{prox}_{b\rho}(\tau z)\right)}{1+b\rho''\left(\operatorname{prox}_{b\rho}(\tau z)\right)} < 0, \quad \forall z > 0.$$

which combined with the fact  $\phi'(z) < 0$  for z > 0 reveals

$$\operatorname{sign}\left(\left(\frac{\rho'\left(\operatorname{prox}_{b\rho}(-\tau z)\right)}{1+b\rho''\left(\operatorname{prox}_{b\rho}(-\tau z)\right)}-\frac{\rho'\left(\operatorname{prox}_{b\rho}(\tau z)\right)}{1+b\rho''\left(\operatorname{prox}_{b\rho}(\tau z)\right)}\right)\phi'(z)\right)=1,\quad\forall z>0.$$

In other words, the integrand in (116) is positive, which allows one to conclude that G'(b) > 0.

🖄 Springer

We then move on to justify (a) and (b). For the first, the derivative of h is given by

$$h'(z) = \frac{\rho''(z) + b(\rho''(z))^2 - b\rho'(z)\rho'''(z)}{(1 + b\rho''(z))^2}.$$

Since  $\rho'$  is log concave, this directly yields  $(\rho'')^2 - \rho' \rho''' > 0$ . As  $\rho'' > 0$  and  $b \ge 0$ , the above implies h'(z) > 0 for all z.

The second claim follows from  $\frac{\partial \operatorname{prox}_{b\rho}(z)}{\partial z} \ge \frac{1}{1+b\|\rho''\|_{\infty}} > 0$  (cf. [22, Equation (56)]). It remains to analyze the behavior of G in the limits when  $b \to 0$  and  $b \to \infty$ .

It remains to analyze the behavior of G in the limits when  $b \to 0$  and  $b \to \infty$ . From [22, Proposition 6.4], G(b) can also be expressed as

$$G(b) = 1 - \mathbb{E}\left[\frac{1}{1 + b\rho''(\mathsf{prox}_{b\rho}(\tau Z))}\right]$$

Since  $\rho''$  is bounded and the integrand is at most 1, the dominated convergence theorem gives

$$\lim_{b \to 0} G(b) = 0.$$

When  $b \to \infty$ ,  $b\rho''(\operatorname{prox}_{b\rho}(\tau z)) \to \infty$  for a fixed z. Again by applying the dominated convergence theorem,

$$\lim_{b \to \infty} G(b) = 1$$

It follows that  $\lim_{b\to 0} G(b) < \kappa < \lim_{b\to\infty} G(b)$  and, therefore,  $G(b) = \kappa$  has a unique positive solution.

**Remark 3** Finally, we show that the logistic and the probit effective links obey the assumptions of Lemma 5. We work with a fixed  $\tau > 0$ .

- A direct computation shows that ρ' is log-concave for the logistic model. For the probit, it is well-known that the reciprocal of the hazard function (also known as Mills' ratio) is strictly log-convex [4].
- To check the other condition, recall that the proximal mapping operator satisfies

$$b\rho'(\operatorname{prox}_{b\rho}(\tau z)) + \operatorname{prox}_{b\rho}(\tau z) = \tau z.$$
(117)

For a fixed z, we claim that if  $b \to \infty$ ,  $\operatorname{prox}_{b\rho}(\tau z) \to -\infty$ . To prove this claim, we start by assuming that this is not true. Then either  $\operatorname{prox}_{b\rho}(\tau z)$  is bounded or diverges to  $\infty$ . If it is bounded, the LHS above diverges to  $\infty$  while the RHS is fixed, which is a contradiction. Similarly if  $\operatorname{prox}_{b\rho}(\tau z)$  diverges to  $\infty$ , the left-hand side of (117) diverges to  $\infty$  while the right-hand side is fixed, which cannot be true as well. Further, when  $b \to \infty$ , we must have  $\operatorname{prox}_{b\rho}(\tau z) \to -\infty$ ,  $b\rho'(\operatorname{prox}_{b\rho}(\tau z)) \to \infty$ , such that the difference of these two is  $\tau z$ . Observe that for the logistic,  $\rho''(x) = \rho'(x)(1 - \rho'(x))$  and for the probit,  $\rho''(x) = \rho'(x)(\rho'(x) - x)$  [53]. Hence,

combining the asymptotic behavior of  $\operatorname{prox}_{b\rho}(\tau z)$  and  $b\rho'(\operatorname{prox}_{b\rho}(\tau z))$ , we obtain that  $b\rho''(\operatorname{prox}_{b\rho}(\tau z))$  diverges to  $\infty$  in both models when  $b \to \infty$ .

### C Proof of Lemma 6

#### C.1 Proof of Part (i)

Recall from [22, Proposition 6.4] that

$$\kappa = \mathbb{E}\left[\Psi'(\tau Z; b(\tau))\right] = 1 - \mathbb{E}\left[\frac{1}{1 + b(\tau)\rho''(\operatorname{prox}_{b(\tau)\rho}(\tau Z))}\right].$$
 (118)

If we denote  $c := \text{prox}_{bo}(0)$ , then b(0) is given by the following relation:

$$1 - \kappa = \frac{1}{1 + b(0)\rho''(c)} \implies b(0) = \frac{\kappa}{\rho''(c)(1 - \kappa)} > 0$$

as  $\rho''(c) > 0$  for any given c > 0. In addition, since  $\rho'(c) > 0$ , we have

$$\mathcal{V}(0) = \frac{\Psi(0, b(0))^2}{\kappa} \stackrel{(a)}{=} \frac{b(0)^2 \rho'(c)^2}{\kappa} > 0,$$

where (a) comes from (22).

#### C.2 Proof of Part (ii)

We defer the proof of this part to the supplemental materials [58].

### D Proof of Part (ii) of Theorem 4

As discussed in Sect. 5.2.2, it suffices to (1) construct a set  $\{\mathcal{B}_i \mid 1 \le i \le N\}$  that forms a cover of the cone  $\mathcal{A}$  defined in (52), and (2) upper bound  $\mathbb{P}\{\{X\boldsymbol{\beta} \mid \boldsymbol{\beta} \in \mathbb{R}^p\} \cap \mathcal{B}_i \ne \{0\}\}$ . In what follows, we elaborate on these two steps.

• Step 1. Generate  $N = \exp(2\epsilon^2 p)$  i.i.d. points  $z^{(i)} \sim \mathcal{N}(\mathbf{0}, \frac{1}{p} \mathbf{I}_p), 1 \le i \le N$ , and construct a collection of convex cones

$$C_i := \left\{ \boldsymbol{u} \in \mathbb{R}^p \left| \left\langle \boldsymbol{u}, \frac{\boldsymbol{z}^{(i)}}{\|\boldsymbol{z}^{(i)}\|} \right\rangle \ge \epsilon \|\boldsymbol{u}\| \right\}, \quad 1 \le i \le N.$$

In words,  $C_i$  consists of all directions that have nontrivial positive correlation with  $z^{(i)}$ . With high probability, this collection  $\{C_i \mid 1 \le i \le N\}$  forms a cover of  $\mathbb{R}^p$ , a fact which is an immediate consequence of the following lemma.

**Lemma 12** Consider any given constant  $0 < \epsilon < 1$ , and let  $N = \exp(2\epsilon^2 p)$ . Then there exist some positive universal constants  $c_5$ ,  $C_5 > 0$  such that with probability exceeding  $1 - C_5 \exp(-c_5\epsilon^2 p)$ ,

$$\sum_{i=1}^{N} \mathbf{1}_{\{\langle x, z^{(i)} \rangle \ge \epsilon \|x\| \| z^{(i)} \|\}} \ge 1$$

holds simultaneously for all  $\mathbf{x} \in \mathbb{R}^p$ .

With our family  $\{C_i \mid 1 \le i \le N\}$  we can introduce

$$\mathcal{B}_i := \mathcal{C}_i \cap \left\{ \boldsymbol{u} \in \mathbb{R}^n \mid \sum_{j=1}^n \max\left\{ -u_j, 0 \right\} \le \epsilon \sqrt{n} \left\langle \boldsymbol{u}, \frac{\boldsymbol{z}^{(i)}}{\|\boldsymbol{z}^{(i)}\|} \right\rangle \right\}, \quad 1 \le i \le N, \quad (119)$$

which in turn forms a cover of the nonconvex cone  $\mathcal{A}$  defined in (52). To justify this, note that for any  $\boldsymbol{u} \in \mathcal{A}$ , one can find  $i \in \{1, ..., N\}$  obeying  $\boldsymbol{u} \in \mathcal{C}_i$ , or equivalently,  $\left\langle \boldsymbol{u}, \frac{\boldsymbol{z}^{(i)}}{\|\boldsymbol{z}^{(i)}\|} \right\rangle \geq \epsilon \|\boldsymbol{u}\|$ , with high probability. Combined with the membership to  $\mathcal{A}$  this gives

$$\sum_{j=1}^{n} \max\left\{-u_{j}, 0\right\} \leq \epsilon^{2} \sqrt{n} \|\boldsymbol{u}\| \leq \epsilon \sqrt{n} \left\langle \boldsymbol{u}, \frac{\boldsymbol{z}^{(i)}}{\|\boldsymbol{z}^{(i)}\|} \right\rangle,$$

indicating that u is contained within some  $\mathcal{B}_i$ .

• Step 2. We now move on to control  $\mathbb{P} \{ \{ X \beta \mid \beta \in \mathbb{R}^p \} \cap \mathcal{B}_i \neq \{ 0 \} \}$ . If the statistical dimensions of the two cones obey  $\delta(\mathcal{B}_i) < n - \delta(\{ X \beta \mid \beta \in \mathbb{R}^p \}) = n - p$ , then an application of [3, Theorem I] gives

$$\mathbb{P}\left\{\left\{\boldsymbol{X\boldsymbol{\beta}}\mid\boldsymbol{\beta}\in\mathbb{R}^{p}\right\}\cap\mathcal{B}_{i}\neq\{\mathbf{0}\}\right\} \\
\leq 4\exp\left\{-\frac{1}{8}\left(\frac{n-\delta\left(\left\{\boldsymbol{X\boldsymbol{\beta}}\mid\boldsymbol{\beta}\in\mathbb{R}^{p}\right\}\right)-\delta\left(\mathcal{B}_{i}\right)}{\sqrt{n}}\right)^{2}\right\} \\
\leq 4\exp\left\{-\frac{(n-p-\delta(\mathcal{B}_{i}))^{2}}{8n}\right\}.$$
(120)

It then comes down to upper bounding  $\delta(\mathcal{B}_i)$ , which is the content of the following lemma.

**Lemma 13** Fix  $\epsilon > 0$ . When n is sufficiently large, the statistical dimension of the convex cone  $\mathcal{B}_i$  defined in (119) obeys

$$\delta(\mathcal{B}_i) \le \left(\frac{1}{2} + 2\sqrt{2}\epsilon^{\frac{3}{4}} + 10H(2\sqrt{\epsilon})\right)n,\tag{121}$$

where  $H(x) := -x \log x - (1 - x) \log(1 - x)$ .

Substitution into (120) gives

$$\mathbb{P}\left\{\left\{\boldsymbol{X\boldsymbol{\beta}}\mid\boldsymbol{\beta}\in\mathbb{R}^{p}\right\}\cap\mathcal{B}_{i}\neq\{\boldsymbol{0}\}\right\}$$

$$\leq4\exp\left\{-\frac{\left(\left(\frac{1}{2}-2\sqrt{2}\epsilon^{\frac{3}{4}}-10H(2\sqrt{\epsilon})\right)n-p\right)^{2}}{8n}\right\}$$

$$=4\exp\left\{-\frac{1}{8}\left(\frac{1}{2}-2\sqrt{2}\epsilon^{\frac{3}{4}}-10H(2\sqrt{\epsilon})-\frac{p}{n}\right)^{2}n\right\}.$$
(122)

Finally, we prove Lemmas 12 and 13 in the next subsections. These are the only remaining parts for the proof of Theorem 4.

#### D.1 Proof of Lemma 12

To begin with, it is seen that all  $||z^{(i)}||$  concentrates around 1. Specifically, apply [34, Proposition 1] to get

$$\mathbb{P}\left\{\|z^{(i)}\|^2 > 1 + 2\sqrt{\frac{t}{p}} + \frac{2t}{p}\right\} \le e^{-t},$$

and set  $t = 3\epsilon^2 p$  to reach

$$\mathbb{P}\left\{\|z^{(i)}\|^2 > 1 + 10\epsilon\right\} \leq \mathbb{P}\left\{\|z^{(i)}\|^2 > 1 + 2\sqrt{3}\epsilon + 6\epsilon^2\right\} \leq e^{-3\epsilon^2 p}.$$

Taking the union bound we obtain

$$\mathbb{P}\left\{\exists 1 \le i \le N \text{ s.t. } \|z^{(i)}\|^2 > 1 + 10\epsilon\right\} \le Ne^{-3\epsilon^2 p} = e^{-\epsilon^2 p}.$$
 (123)

Next, we note that it suffices to prove Lemma 12 for all unit vectors  $\mathbf{x}$ . The following lemma provides a bound on  $\langle \mathbf{z}^{(i)}, \mathbf{x} \rangle$  for any fixed unit vector  $\mathbf{x} \in \mathbb{R}^{p}$ .

**Lemma 14** Consider any fixed unit vector  $\mathbf{x} \in \mathbb{R}^p$  and any given constant  $0 < \epsilon < 1$ , and set  $N = \exp(2\epsilon^2 p)$ . There exist positive universal constants  $c_5, c_6, C_6 > 0$  such that

$$\mathbb{P}\left\{\sum_{i=1}^{N} \mathbf{1}_{\left\{\left\langle z^{(i)}, \mathbf{x} \right\rangle \geq \frac{1}{2}\epsilon \right\}} \leq \exp\left(\left(1 - o\left(1\right)\right) \frac{7}{4}\epsilon^{2}p\right)\right\} \\
\leq \exp\left\{-2\exp\left(\left(1 - o\left(1\right)\right) \frac{7}{4}\epsilon^{2}p\right)\right\}.$$
(124)

Recognizing that Lemma 12 is a uniform result, we need to extend Lemma 14 to all x simultaneously, which we achieve via the standard covering argument. Specifically,

one can find a set  $C := \{ \mathbf{x}^{(j)} \in \mathbb{R}^p \mid 1 \le j \le K \}$  of unit vectors with cardinality  $K = (1 + 2p^2)^p$  to form a cover of the unit ball of resolution  $p^{-2}$  [65, Lemma 5.2]; that is, for any unit vector  $\mathbf{x} \in \mathbb{R}^p$ , there exists a  $\mathbf{x}^{(j)} \in C$  such that

$$\|\boldsymbol{x}^{(j)} - \boldsymbol{x}\| \le p^{-2}.$$

Apply Lemma 14 and take the union bound to arrive at

$$\sum_{i=1}^{N} \mathbf{1}_{\left\{ \langle z^{(i)}, \mathbf{x}^{(j)} \rangle \geq \frac{1}{2}\epsilon \right\}} \ge \exp\left( (1 - o(1)) \frac{7}{4} \epsilon^2 p \right) > 1, \quad 1 \le j \le K$$
(125)

with probability exceeding  $1 - K \exp\left\{-2 \exp\left((1 - o(1))\frac{7}{4}\epsilon^2 p\right)\right\} \ge 1 - \exp\left\{-2(1 - o(1))\exp\left((1 - o(1))\frac{7}{4}\epsilon^2 p\right)\right\}$ . This guarantees that for each  $\mathbf{x}^{(j)}$ , one can find at least one  $\mathbf{z}^{(i)}$  obeying

$$\left\langle \boldsymbol{z}^{(i)}, \boldsymbol{x}^{(j)} \right\rangle \geq \frac{1}{2} \epsilon.$$

This result together with (123) yields that with probability exceeding  $1 - C \exp(-c\epsilon^2 p)$ , for some universal constants C, c > 0.

$$\begin{split} \left\langle \boldsymbol{z}^{(i)}, \boldsymbol{x} \right\rangle &\geq \left\langle \boldsymbol{z}^{(i)}, \boldsymbol{x}^{(j)} \right\rangle - \left\langle \boldsymbol{z}^{(i)}, \boldsymbol{x}^{(j)} - \boldsymbol{x} \right\rangle \geq \left\langle \boldsymbol{z}^{(i)}, \boldsymbol{x}^{(j)} \right\rangle - \|\boldsymbol{z}^{(i)}\| \cdot \|\boldsymbol{x}^{(j)} - \boldsymbol{x}\| \\ &\geq \frac{1}{2}\epsilon - \frac{1}{p^2} \|\boldsymbol{z}^{(i)}\| \geq \frac{\frac{1}{2}\epsilon}{\sqrt{1 + 10\epsilon}} \|\boldsymbol{z}^{(i)}\| \\ &\quad - \frac{1}{p^2} \|\boldsymbol{z}^{(i)}\| \\ &\geq \frac{1}{30}\epsilon \|\boldsymbol{z}^{(i)}\| \end{split}$$

holds simultaneously for all unit vectors  $x \in \mathbb{R}^p$ . Since  $\epsilon > 0$  can be an arbitrary constant, this concludes the proof.

**Proof of Lemma 14** Without loss of generality, it suffices to consider  $x = e_1 = [1, 0, ..., 0]^{\top}$ . For any t > 0 and any constant  $\zeta > 0$ , it comes from [2, Theorem A.1.4] that

$$\mathbb{P}\left\{\frac{1}{N}\sum_{i=1}^{N}\mathbf{1}_{\left\{\langle \boldsymbol{z}^{(i)},\boldsymbol{e}_{1}\rangle<\zeta\right\}}>(1+t)\Phi\left(\zeta\sqrt{p}\right)\right\}\leq\exp\left(-2t^{2}\Phi^{2}\left(\zeta\sqrt{p}\right)N\right).$$

Setting  $t = 1 - \Phi\left(\zeta \sqrt{p}\right)$  gives

$$\mathbb{P}\left\{\frac{1}{N}\sum_{i=1}^{N}\mathbf{1}_{\left\{\langle z^{(i)}, e_{1} \rangle < \zeta\right\}} > \left(2 - \Phi\left(\zeta\sqrt{p}\right)\right)\Phi\left(\zeta\sqrt{p}\right)\right\}$$
$$\leq \exp\left(-2\left(1 - \Phi\left(\zeta\sqrt{p}\right)\right)^{2}\Phi^{2}\left(\zeta\sqrt{p}\right)N\right).$$

D Springer

Recall that for any t > 1, one has  $(t^{-1} - t^{-3})\phi(t) \le 1 - \Phi(t) \le t^{-1}\phi(t)$  which implies that

$$1 - \Phi\left(\zeta\sqrt{p}\right) = \exp\left(-\frac{(1+o(1))\,\zeta^2 p}{2}\right).$$

Taking  $\zeta = \frac{1}{2}\epsilon$ , we arrive at

$$(2 - \Phi(\zeta\sqrt{p})) \Phi(\zeta\sqrt{p}) = 1 - \exp\left(-(1 + o(1))\zeta^2 p\right)$$
$$= 1 - \exp\left(-(1 + o(1))\frac{1}{4}\epsilon^2 p\right),$$
$$(1 - \Phi(\zeta\sqrt{p}))^2 \Phi^2(\zeta\sqrt{p}) = \exp\left(-(1 + o(1))\zeta^2 p\right)$$
$$= \exp\left(-(1 + o(1))\frac{1}{4}\epsilon^2 p\right) \gg \frac{1}{N}.$$

This justifies that

$$\mathbb{P}\left\{\sum_{i=1}^{N} \mathbf{1}_{\left\{\langle z^{(i)}, \boldsymbol{e}_{1} \rangle \geq \frac{1}{2} \epsilon\right\}} \leq N \exp\left(-(1+o(1))\frac{1}{4}\epsilon^{2}p\right)\right\}$$
$$= \mathbb{P}\left\{\frac{1}{N}\sum_{i=1}^{N} \mathbf{1}_{\left\{\langle z^{(i)}, \boldsymbol{e}_{1} \rangle < \zeta\right\}} > \left(2-\Phi\left(\zeta\sqrt{p}\right)\right)\Phi\left(\zeta\sqrt{p}\right)\right\}$$
$$\leq \exp\left\{-2\exp\left(-(1+o(1))\frac{1}{4}\epsilon^{2}p\right)N\right\}$$
$$= \exp\left\{-2\exp\left((1-o(1))\frac{7}{4}\epsilon^{2}p\right)\right\}$$

as claimed.

### D.2 Proof of Lemma 13

First of all, recall from the definition (19) that

$$\delta(\mathcal{B}_i) = \mathbb{E}\left[\left\|\Pi_{\mathcal{B}_i}\left(\boldsymbol{g}\right)\right\|^2\right] = \mathbb{E}\left[\left\|\boldsymbol{g}\right\|^2 - \min_{\boldsymbol{u}\in\mathcal{B}_i}\|\boldsymbol{g} - \boldsymbol{u}\|^2\right] = n - \mathbb{E}\left[\min_{\boldsymbol{u}\in\mathcal{B}_i}\|\boldsymbol{g} - \boldsymbol{u}\|^2\right]$$
$$\leq n - \mathbb{E}\left[\min_{\boldsymbol{u}\in\mathcal{D}_i}\|\boldsymbol{g} - \boldsymbol{u}\|^2\right],$$

where  $\boldsymbol{g} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_n)$ , and  $\mathcal{D}_i$  is a superset of  $\mathcal{B}_i$  defined by

$$\mathcal{D}_i := \left\{ \boldsymbol{u} \in \mathbb{R}^n \mid \sum_{j=1}^n \max\left\{ -u_j, 0 \right\} \le \epsilon \sqrt{n} \|\boldsymbol{u}\| \right\}.$$
(126)

Recall from the triangle inequality that

$$\|g - u\| \ge \|u\| - \|g\| > \|g\| = \|g - 0\|, \quad \forall u : \|u\| > 2\|g\|.$$

Since  $\mathbf{0} \in \mathcal{D}_i$ , this implies that

$$\left\| \arg\min_{\boldsymbol{u}\in\mathcal{D}_i} \|\boldsymbol{g}-\boldsymbol{u}\| \right\| \leq 2\|\boldsymbol{g}\|,$$

revealing that

$$\mathbb{E}\left[\min_{\boldsymbol{u}\in\mathcal{D}_i}\|\boldsymbol{g}-\boldsymbol{u}\|^2\right] = \mathbb{E}\left[\min_{\boldsymbol{u}\in\mathcal{D}_i,\|\boldsymbol{u}\|\leq 2\|\boldsymbol{g}\|}\|\boldsymbol{g}-\boldsymbol{u}\|^2\right].$$

In what follows, it suffices to look at the set of u's within  $\mathcal{D}_i$  obeying  $||u|| \le 2||g||$ , which verify

$$\sum_{j=1}^{n} \max\left\{-u_{j}, 0\right\} \le \epsilon \sqrt{n} \|\boldsymbol{u}\| \le 2\epsilon \sqrt{n} \|\boldsymbol{g}\|.$$
(127)

It is seen that

$$\|\boldsymbol{g} - \boldsymbol{u}\|^{2} \geq \sum_{i:g_{i}<0} (g_{i} - u_{i})^{2} = \left\{ \sum_{i:g_{i}<0, u_{i}\geq0} + \sum_{i:g_{i}<0, -\sqrt{\frac{\varepsilon}{n}} \|\boldsymbol{g}\| < u_{i}<0} + \sum_{i:g_{i}<0, u_{i}\leq-\sqrt{\frac{\varepsilon}{n}} \|\boldsymbol{g}\|} \right\} (g_{i} - u_{i})^{2}$$

$$\geq \sum_{i:g_{i}<0, u_{i}\geq0} g_{i}^{2} + \sum_{i:g_{i}<0, -\sqrt{\frac{\varepsilon}{n}} \|\boldsymbol{g}\| < u_{i}<0} (g_{i} - u_{i})^{2}$$

$$\geq \sum_{i:g_{i}<0, u_{i}\geq0} g_{i}^{2} + \sum_{i:g_{i}<0, -\sqrt{\frac{\varepsilon}{n}} \|\boldsymbol{g}\| < u_{i}<0} (g_{i}^{2} - 2u_{i}g_{i})$$

$$\geq \sum_{i:g_{i}<0, u_{i}>-\sqrt{\frac{\varepsilon}{n}} \|\boldsymbol{g}\|} g_{i}^{2} - \sum_{i:g_{i}<0, -\sqrt{\frac{\varepsilon}{n}} \|\boldsymbol{g}\| < u_{i}<0} 2u_{i}g_{i}. \quad (128)$$

1. Regarding the first term of (128), we first recognize that

$$\left\{i \mid u_i \leq -\sqrt{\frac{\epsilon}{n}} \|\boldsymbol{g}\|\right\} \leq \frac{\sum_{i:u_i < 0} |u_i|}{\sqrt{\frac{\epsilon}{n}} \|\boldsymbol{g}\|} = \frac{\sum_{i=1}^n \max\left\{-u_i, 0\right\}}{\sqrt{\frac{\epsilon}{n}} \|\boldsymbol{g}\|} \leq 2\sqrt{\epsilon}n,$$

where the last inequality follows from the constraint (127). As a consequence,

$$\sum_{\substack{i:g_i<0, \ u_i>-\sqrt{\frac{\epsilon}{n}}\|g\|}} g_i^2 \ge \sum_{\substack{i:g_i<0}} g_i^2 - \sum_{\substack{i:u_i\leq -\sqrt{\frac{\epsilon}{n}}\|g\|}} g_i^2$$
$$\ge \sum_{\substack{i:g_i<0}} g_i^2 - \max_{\substack{S\subseteq [n]: \ |S|=2\sqrt{\epsilon}n}} \sum_{i\in S} g_i^2.$$

2. Next, we turn to the second term of (128), which can be bounded by

$$\sum_{i:g_i<0, -\sqrt{\frac{\epsilon}{n}} \|\boldsymbol{g}\| < u_i < 0} u_i g_i$$

$$\leq \sqrt{\left(\sum_{i:g_i<0, -\sqrt{\frac{\epsilon}{n}} \|\boldsymbol{g}\| < u_i < 0} u_i^2\right) \left(\sum_{i:g_i<0, -\sqrt{\frac{\epsilon}{n}} \|\boldsymbol{g}\| < u_i < 0} g_i^2\right)}$$

$$\leq \sqrt{\left(\max_{i:-\sqrt{\frac{\epsilon}{n}} \|\boldsymbol{g}\| < u_i < 0} |u_i|\right) \left(\sum_{i:u_i<0} |u_i|\right) \cdot \|\boldsymbol{g}\|^2}$$

$$\leq \sqrt{\sqrt{\frac{\epsilon}{n}} \|\boldsymbol{g}\| \left(\sum_{i:u_i<0} |u_i|\right) \cdot \|\boldsymbol{g}\|^2} \leq \sqrt{2\epsilon^{\frac{3}{4}} \|\boldsymbol{g}\|^2},$$

where the last inequality follows from the constraint (127).

Putting the above results together, we have

$$\|\boldsymbol{g} - \boldsymbol{u}\|^{2} \ge \sum_{i:g_{i} < 0} g_{i}^{2} - \max_{S \subseteq [n]: |S| = 2\sqrt{\epsilon}n} \sum_{i \in S} g_{i}^{2} - 2\sqrt{2\epsilon^{\frac{3}{4}}} \|\boldsymbol{g}\|^{2}$$

for any  $\boldsymbol{u} \in \mathcal{D}_i$  obeying  $\|\boldsymbol{u}\| \le 2\|\boldsymbol{g}\|$ , whence

$$\mathbb{E}\left[\min_{\boldsymbol{u}\in\mathcal{D}_{i}}\|\boldsymbol{g}-\boldsymbol{u}\|^{2}\right] \geq \mathbb{E}\left[\sum_{i:g_{i}<0}g_{i}^{2}-\max_{S\subseteq[n]:|S|=2\sqrt{\epsilon}n}\sum_{i\in S}g_{i}^{2}-2\sqrt{2}\epsilon^{\frac{3}{4}}\|\boldsymbol{g}\|^{2}\right]$$
$$=\left(\frac{1}{2}-2\sqrt{2}\epsilon^{\frac{3}{4}}\right)n-\mathbb{E}\left[\max_{S\subseteq[n]:|S|=2\sqrt{\epsilon}n}\sum_{i\in S}g_{i}^{2}\right].$$
(129)

Finally, it follows from [34, Proposition 1] that for any  $t > 2\sqrt{\epsilon n}$ ,

$$\mathbb{P}\left\{\sum_{i\in S}g_i^2 \ge 5t\right\} \le \mathbb{P}\left\{\sum_{i\in S}g_i^2 \ge |S| + 2\sqrt{|S|t} + 2t\right\} \le e^{-t},$$

which together with the union bound gives

$$\mathbb{P}\left\{\max_{\substack{S\subseteq[n]: |S|=2\sqrt{\epsilon}n}}\sum_{i\in S}g_i^2 \ge 5t\right\} \le \sum_{\substack{S\subseteq[n]: |S|=2\sqrt{\epsilon}n}}\mathbb{P}\left\{\sum_{i\in S}g_i^2 \ge 5t\right\}$$
$$\le \exp\left\{H\left(2\sqrt{\epsilon}\right)n-t\right\}.$$

This gives

$$\mathbb{E}\left[\max_{S\subseteq[n]:\ |S|=2\sqrt{\epsilon}n}\sum_{i\in S}g_i^2\right] = \int_0^\infty \mathbb{P}\left\{\max_{S\subseteq[n]:\ |S|=2\sqrt{\epsilon}n}\sum_{i\in S}g_i^2 \ge t\right\}dt$$
$$\le 5H\left(2\sqrt{\epsilon}\right)n + \int_{5H(2\sqrt{\epsilon})n}^\infty \exp\left\{H\left(2\sqrt{\epsilon}\right)n - \frac{1}{5}t\right\}dt$$
$$< 10H\left(2\sqrt{\epsilon}\right)n,$$

for any given  $\epsilon > 0$  with the proviso that *n* is sufficiently large. This combined with (129) yields

$$\mathbb{E}\left[\min_{\boldsymbol{u}\in\mathcal{D}_{i}}\|\boldsymbol{g}-\boldsymbol{u}\|^{2}\right] \geq \left(\frac{1}{2}-2\sqrt{2}\epsilon^{\frac{3}{4}}-10H(2\sqrt{\epsilon})\right)n \tag{130}$$

as claimed.

# E Proof of Lemma 8

Throughout, we shall restrict ourselves on the event  $A_n$  as defined in (86), on which  $\tilde{G} \succeq \lambda_{\text{lb}} I$ . Recalling the definitions of  $\tilde{G}$  and w from (82) and (89), we see that

$$\boldsymbol{w}^{\top} \tilde{\boldsymbol{G}}^{-2} \boldsymbol{w} = \frac{1}{n^2} \boldsymbol{X}_{\cdot 1}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}} \left( \frac{1}{n} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}} \right)^{-2} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \boldsymbol{X}_{\cdot 1}$$
$$\leq \frac{\|\boldsymbol{X}_{\cdot 1}^{\top}\|^2}{n} \left\| \frac{1}{n} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}} \left( \frac{1}{n} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}} \right)^{-2} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \right\|.$$
(131)

If we let the singular value decomposition of  $\frac{1}{\sqrt{n}} D_{\hat{\beta}}^{1/2} \tilde{X}$  be  $U \Sigma V^{\top}$ , then a little algebra gives  $\Sigma \geq \sqrt{\lambda_{\text{lb}}} I$  and

$$\frac{1}{n} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \tilde{\boldsymbol{X}} \left( \frac{1}{n} \tilde{\boldsymbol{X}}' \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}} \right)^{-2} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} = \boldsymbol{U} \boldsymbol{\Sigma}^{-2} \boldsymbol{U}^{\top} \preceq \lambda_{\text{lb}}^{-1} \boldsymbol{I}.$$

Deringer

Substituting this into (131) and using the fact  $||X_{.1}||^2 \leq n$  with high probability (by Lemma 2), we obtain

$$\boldsymbol{w}^{ op} \tilde{\boldsymbol{G}}^{-2} \boldsymbol{w} \lesssim \frac{1}{n\lambda_{ ext{lb}}} \|\boldsymbol{X}_{\cdot 1}\|^2 \lesssim 1$$

with probability at least  $1 - \exp(-\Omega(n))$ .

# F Proof of Lemma 9

Throughout this and the subsequent sections, we consider  $H_n$  and  $K_n$  to be two diverging sequences with the following properties:

$$H_n = o\left(n^{\epsilon}\right), \quad K_n = o\left(n^{\epsilon}\right), \quad n^2 \exp\left(-c_1 H_n^2\right) = o(1), \quad n \exp\left(-c_2 K_n^2\right) = o(1),$$
(132)

for any constants  $c_i > 0$ , i = 1, 2 and any  $\epsilon > 0$ . This lemma is an analogue of [25, Proposition 3.18]. We modify and adapt the proof ideas to establish the result in our setup. Throughout we shall restrict ourselves to the event  $A_n$ , on which  $\tilde{G} \succeq \lambda_{lb} I$ .

Due to independence between  $X_{.1}$  and  $\{D_{\tilde{\beta}}, H\}$ , one can invoke the Hanson-Wright inequality [52, Theorem 1.1] to yield

$$\mathbb{P}\left(\left|\frac{1}{n}X_{\cdot 1}^{\top}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\boldsymbol{H}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}X_{\cdot 1}-\frac{1}{n}\operatorname{Tr}\left(\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\boldsymbol{H}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\right)\right| > t \mid \boldsymbol{H}, \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}\right) \\
\leq 2\exp\left(-c\min\left\{\frac{t^{2}}{\frac{K^{4}}{n^{2}}\|\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\boldsymbol{H}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\|_{\mathrm{F}}^{2}, \frac{t}{\frac{K^{2}}{n}}\|\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\boldsymbol{H}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\|_{\mathrm{F}}^{2}\right) \\
\leq 2\exp\left(-c\min\left\{\frac{t^{2}}{\frac{K^{4}}{n}\|\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\boldsymbol{H}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\|^{2}, \frac{t}{\frac{K^{2}}{n}}\|\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\boldsymbol{H}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\|_{\mathrm{F}}^{2}\right)\right),$$

where  $\|.\|_{\rm F}$  denotes the Frobenius norm. Choose  $t = C^2 \| D_{\hat{\beta}}^{1/2} H D_{\hat{\beta}}^{1/2} \| H_n / \sqrt{n}$  with C > 0 a sufficiently large constant, and take  $H_n$  to be as in (132). Substitution into the above inequality and unconditioning give

$$\mathbb{P}\left(\left|\frac{1}{n}X_{\cdot 1}^{\top}D_{\tilde{\beta}}^{1/2}HD_{\tilde{\beta}}^{1/2}X_{\cdot 1}-\frac{1}{n}\operatorname{Tr}\left(D_{\tilde{\beta}}^{1/2}HD_{\tilde{\beta}}^{1/2}\right)\right| > \frac{1}{\sqrt{n}}C^{2}H_{n}\|D_{\tilde{\beta}}^{1/2}HD_{\tilde{\beta}}^{1/2}\|\right) \\
\leq 2\exp\left(-c\min\left\{\frac{C^{4}H_{n}^{2}}{K^{4}},\frac{C^{2}\sqrt{n}H_{n}}{K^{2}}\right\}\right) = C\exp\left(-cH_{n}^{2}\right) = o(1), \quad (133)$$

for some universal constants C, c > 0.

Deringer

We are left to analyzing  $\operatorname{Tr}(D^{1/2}_{\tilde{\beta}}HD^{1/2}_{\tilde{\beta}})$ . Recall from the definition (92) of H that

$$\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \boldsymbol{H} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} = \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} - \frac{1}{n} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{G}}^{-1} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}},$$

and, hence,

$$\operatorname{Tr}\left(\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\boldsymbol{H}\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\right) = \sum_{i=1}^{n} \left(\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}}) - \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})^{2}}{n}\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}^{-1}\tilde{\boldsymbol{X}}_{i}\right).$$
(134)

This requires us to analyze  $\tilde{G}^{-1}$  carefully. To this end, recall that the matrix  $\tilde{G}_{(i)}$  defined in (83) obeys

$$\tilde{\boldsymbol{G}}_{(i)} = \tilde{\boldsymbol{G}} - \frac{1}{n} \rho''(\tilde{\boldsymbol{X}}^{\top} \tilde{\boldsymbol{\beta}}) \tilde{\boldsymbol{X}}_i \tilde{\boldsymbol{X}}_i^{\top}.$$

Invoking Sherman-Morrison-Woodbury formula (e.g. [31]), we have

$$\tilde{\boldsymbol{G}}^{-1} = \tilde{\boldsymbol{G}}_{(i)}^{-1} - \frac{\frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})}{n}\tilde{\boldsymbol{G}}_{(i)}^{-1}\tilde{\boldsymbol{X}}_{i}\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}_{(i)}^{-1}}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})}{n}\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}_{(i)}^{-1}\tilde{\boldsymbol{X}}_{i}}.$$
(135)

It follows that

$$\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}^{-1} \tilde{\boldsymbol{X}}_{i} = \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i} - \frac{\frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}})}{n} (\boldsymbol{X}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i})^{2}}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}})}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i}},$$

which implies that

$$\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}^{-1} \tilde{\boldsymbol{X}}_{i} = \frac{\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i}}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}})}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i}}.$$
(136)

The relations (134) and (136) taken collectively reveal that

$$\frac{1}{n} \operatorname{Tr} \left( \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \boldsymbol{H} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \right) = \frac{1}{n} \sum_{i=1}^{n} \frac{\rho''(\tilde{X}_{i} \tilde{\boldsymbol{\beta}})}{1 + \frac{\rho''(\tilde{X}_{i}^{\top} \tilde{\boldsymbol{\beta}})}{n} \tilde{X}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{X}_{i}}.$$
(137)

We shall show that the trace above is close to Tr(I - H) up to some factors. For this purpose we analyze the latter quantity in two different ways. To begin with, observe that

$$\operatorname{Tr}(\boldsymbol{I} - \boldsymbol{H}) = \operatorname{Tr}\left(\frac{\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{G}}^{-1} \tilde{\boldsymbol{X}}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}}{n}\right) = \operatorname{Tr}(\tilde{\boldsymbol{G}} \tilde{\boldsymbol{G}}^{-1}) = p - 1.$$
(138)

On the other hand, it directly follows from the definition of H and (136) that the *i*th diagonal entry of H is given by

$$H_{i,i} = \frac{1}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}})}{n} \tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_i}.$$

Applying this relation, we can compute Tr(I - H) analytically as follows:

$$\operatorname{Tr}(\boldsymbol{I} - \boldsymbol{H}) = \sum_{i} \frac{\frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i}}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i}}$$

$$= \sum_{i} \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) \tilde{\alpha} + \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}})}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i} - \rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) \tilde{\alpha}}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}})} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i}}$$

$$= \sum_{i} \rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) \tilde{\alpha} H_{i,i} + \sum_{i} \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) \left(\frac{1}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i} - \tilde{\alpha}\right)}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}})}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i}}, \quad (140)$$

where  $\tilde{\alpha} := \frac{1}{n} \operatorname{Tr} \left( \tilde{\boldsymbol{G}}^{-1} \right)$ .

Observe that the first quantity in the right-hand side above is simply  $\tilde{\alpha} \operatorname{Tr}(D_{\tilde{\beta}}^{1/2})$  $HD_{\tilde{\beta}}^{1/2}$ ). For simplicity, denote

$$\eta_i = \frac{1}{n} \tilde{\boldsymbol{X}}_i^\top \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_i - \tilde{\alpha}.$$
(141)

Note that  $\tilde{G}_{(i)} \succ 0$  on  $\mathcal{A}_n$  and that  $\rho'' > 0$ . Hence the denominator in the second term in (140) is greater than 1 for all *i*. Comparing (138) and (140), we deduce that

$$\left|\frac{p-1}{n} - \frac{1}{n} \operatorname{Tr}\left(\boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2} \boldsymbol{H} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^{1/2}\right) \tilde{\boldsymbol{\alpha}}\right| \leq \sup_{i} |\eta_{i}| \cdot \frac{1}{n} \sum_{i} |\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}})| \lesssim \sup_{i} |\eta_{i}|$$
(142)

on  $A_n$ . It thus suffices to control  $\sup_i |\eta_i|$ . The above bounds together with Lemma (87) and the proposition below complete the proof.

**Proposition 1** Let  $\eta_i$  be as defined in (141). Then there exist universal constants  $C_1, C_2, C_3 > 0$  such that

$$\mathbb{P}\left(\sup_{i} |\eta_{i}| \leq \frac{C_{1}K_{n}^{2}H_{n}}{\sqrt{n}}\right) \geq 1 - C_{2}n^{2}\exp\left(-c_{2}H_{n}^{2}\right) - C_{3}n\exp\left(-c_{3}K_{n}^{2}\right) - \exp\left(-C_{4}n\left(1+o(1)\right)\right) = 1 - o(1),$$

where  $K_n$ ,  $H_n$  are diverging sequences as specified in (132).

**Proof of Proposition 1** Fix any index *i*. Recall that  $\tilde{\boldsymbol{\beta}}_{[-i]}$  is the MLE when the 1st predictor and *i*th observation are removed. Also recall the definition of  $\tilde{\boldsymbol{G}}_{[-i]}$  in (85). The proof essentially follows three steps. First, note that  $\tilde{\boldsymbol{X}}_i$  and  $\tilde{\boldsymbol{G}}_{[-i]}$  are independent. Hence, an application of the Hanson-Wright inequality [52] yields that

$$\mathbb{P}\left(\left|\frac{1}{n}\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}_{[-i]}^{-1}\tilde{\boldsymbol{X}}_{i}-\frac{1}{n}\mathrm{Tr}\left(\tilde{\boldsymbol{G}}_{[-i]}^{-1}\right)\right| > t \mid \tilde{\boldsymbol{G}}_{[-i]}\right)$$

$$\leq 2\exp\left(-c\min\left\{\frac{t^{2}}{\frac{K^{4}}{n^{2}}\|\tilde{\boldsymbol{G}}_{[-i]}^{-1}\|_{\mathrm{F}}^{2}}, \frac{t}{\frac{K^{2}}{n}}\|\tilde{\boldsymbol{G}}_{[-i]}^{-1}\|\right\}\right)$$

$$\leq 2\exp\left(-c\min\left\{\frac{t^{2}}{\frac{K^{4}}{n}\|\tilde{\boldsymbol{G}}_{[-i]}^{-1}\|^{2}}, \frac{t}{\frac{K^{2}}{n}}\|\tilde{\boldsymbol{G}}_{[-i]}^{-1}\|\right\}\right)$$

We choose  $t = C^2 \|\tilde{\boldsymbol{G}}_{[-i]}^{-1}\| H_n / \sqrt{n}$ , where C > 0 is a sufficiently large constant. Now marginalizing gives

$$\begin{split} & \mathbb{P}\left(\left|\frac{1}{n}\tilde{X}_{i}^{\top}\tilde{\boldsymbol{G}}_{[-i]}^{-1}\tilde{X}_{i}-\frac{1}{n}\mathrm{Tr}\left(\tilde{\boldsymbol{G}}_{[-i]}^{-1}\right)\right| > C^{2}\left\|\tilde{\boldsymbol{G}}_{[-i]}^{-1}\right\|\frac{H_{n}}{\sqrt{n}}\right) \\ & \leq 2\exp\left(-c\min\left\{\frac{C^{4}H_{n}^{2}}{K^{4}},\frac{C^{2}\sqrt{n}H_{n}}{K^{2}}\right\}\right) \\ & \leq 2\exp\left(-C'H_{n}^{2}\right), \end{split}$$

where C' > 0 is a sufficiently large constant. On  $\mathcal{A}_n$ , the spectral norm  $\|\tilde{\boldsymbol{G}}_{(i)}^{-1}\|$  is bounded above by  $\lambda_{\text{Ib}}$  for all *i*. Invoking (87) we obtain that there exist universal constants  $C_1$ ,  $C_2$ ,  $C_3 > 0$  such that

$$\mathbb{P}\left(\sup_{i}\left|\frac{1}{n}\tilde{X}_{i}^{\top}\tilde{\boldsymbol{G}}_{[-i]}^{-1}\tilde{X}_{i}-\frac{1}{n}\mathrm{Tr}\left(\tilde{\boldsymbol{G}}_{[-i]}^{-1}\right)\right|>C_{1}\frac{H_{n}}{\sqrt{n}}\right)\leq C_{2}n\exp\left(-C_{3}H_{n}^{2}\right).$$
 (143)

The next step consists of showing that  $\operatorname{Tr}(\tilde{\boldsymbol{G}}_{[-i]}^{-1})$  (resp.  $\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}_{[-i]}^{-1}\tilde{\boldsymbol{X}}_{i}$ ) and  $\operatorname{Tr}(\tilde{\boldsymbol{G}}_{(i)}^{-1})$  (resp.  $\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}_{(i)}^{-1}\tilde{\boldsymbol{X}}_{i}$ ) are uniformly close across all *i*. This is established in the following lemma.

D Springer

**Lemma 15** Let  $\tilde{G}_{(i)}$  and  $\tilde{G}_{[-i]}$  be defined as in (83) and (85), respectively. Then there exist universal constants  $C_1, C_2, C_3, C_4, c_2, c_3 > 0$  such that

$$\mathbb{P}\left(\sup_{i}\left|\frac{1}{n}\tilde{X}_{i}^{\top}\tilde{G}_{(i)}^{-1}\tilde{X}_{i}-\frac{1}{n}\tilde{X}_{i}^{\top}\tilde{G}_{[-i]}^{-1}\tilde{X}_{i}\right| \leq C_{1}\frac{K_{n}^{2}H_{n}}{\sqrt{n}}\right) \\
= 1 - C_{2}n^{2}\exp\left(-c_{2}H_{n}^{2}\right) - C_{3}n\exp\left(-c_{3}K_{n}^{2}\right) \\
-\exp\left(-C_{4}n\left(1+o(1)\right)\right) = 1 - o(1), \quad (144) \\
\mathbb{P}\left(\sup_{i}\left|\frac{1}{n}\operatorname{Tr}\left(\tilde{G}_{(i)}^{-1}\right)-\frac{1}{n}\operatorname{Tr}\left(\tilde{G}_{[-i]}^{-1}\right)\right| \leq C_{1}\frac{K_{n}^{2}H_{n}}{\sqrt{n}}\right) \\
= 1 - C_{2}n^{2}\exp\left(-c_{2}H_{n}^{2}\right) - C_{3}n\exp\left(-c_{3}K_{n}^{2}\right) \\
-\exp\left(-C_{4}n\left(1+o(1)\right)\right) = 1 - o(1), \quad (145)$$

where  $K_n$ ,  $H_n$  are diverging sequences as defined in (132).

This together with (143) yields that

$$\mathbb{P}\left(\sup_{i}\left|\frac{1}{n}\tilde{X}_{i}^{\top}\tilde{\mathbf{G}}_{(i)}^{-1}\tilde{X}_{i}-\frac{1}{n}\operatorname{Tr}(\tilde{\mathbf{G}}_{(i)}^{-1})\right|>C_{1}\frac{K_{n}^{2}H_{n}}{\sqrt{n}}\right) \leq C_{2}n^{2}\exp\left(-c_{2}H_{n}^{2}\right)+C_{3}n\exp\left(-c_{3}K_{n}^{2}\right)+\exp\left(-C_{4}n\left(1+o(1)\right)\right).$$
(146)

The final ingredient is to establish that  $\frac{1}{n} \operatorname{Tr}(\tilde{\boldsymbol{G}}_{(i)}^{-1})$  and  $\frac{1}{n} \operatorname{Tr}(\tilde{\boldsymbol{G}}^{-1})$  are uniformly close across *i*.

**Lemma 16** Let  $\tilde{G}$  and  $\tilde{G}_{(i)}$  be as defined in (82) and (83), respectively. Then one has

$$\mathbb{P}\left(\left|\operatorname{Tr}\left(\tilde{\boldsymbol{G}}_{(i)}^{-1}\right) - \operatorname{Tr}\left(\tilde{\boldsymbol{G}}^{-1}\right)\right| \le \frac{1}{\lambda_{\mathrm{lb}}}\right) \ge 1 - \exp\left(-\Omega(n)\right).$$
(147)

This completes the proof.

**Proof of Lemma 15** For two invertible matrices **A** and **B** of the same dimensions, the difference of their inverses can be written as

$$A^{-1} - B^{-1} = A^{-1}(B - A)B^{-1}.$$

Applying this identity, we have

$$\tilde{\boldsymbol{G}}_{(i)}^{-1} - \tilde{\boldsymbol{G}}_{[-i]}^{-1} = \tilde{\boldsymbol{G}}_{(i)}^{-1} \left( \tilde{\boldsymbol{G}}_{[-i]} - \tilde{\boldsymbol{G}}_{(i)} \right) \tilde{\boldsymbol{G}}_{[-i]}^{-1}.$$

From the definition of these matrices, it follows directly that

$$\tilde{\boldsymbol{G}}_{[-i]} - \tilde{\boldsymbol{G}}_{(i)} = \frac{1}{n} \sum_{j:j \neq i} \left( \rho'' \left( \tilde{\boldsymbol{X}}_{j}^{\top} \tilde{\boldsymbol{\beta}}_{[-i]} \right) - \rho'' \left( \tilde{\boldsymbol{X}}_{j}^{\top} \tilde{\boldsymbol{\beta}} \right) \right) \tilde{\boldsymbol{X}}_{j} \tilde{\boldsymbol{X}}_{j}^{\top}.$$
(148)

Deringer

As  $\rho'''$  is bounded, by the mean-value theorem, it suffices to control the differences  $X_j^{\top} \tilde{\beta}_{[-i]} - \tilde{X}_j^{\top} \tilde{\beta}$  uniformly across all *j*. This is established in the following lemma, the proof of which is deferred to Appendix H.

**Lemma 17** Let  $\hat{\boldsymbol{\beta}}$  be the full model MLE and  $\hat{\boldsymbol{\beta}}_{[-i]}$  be the MLE when the *i*th observation is dropped. Let  $q_i$  be as described in Lemma 18 and  $K_n$ ,  $H_n$  be as in (132). Then there exist universal constants  $C_1$ ,  $C_2$ ,  $C_3$ ,  $C_4$ ,  $c_2$ ,  $c_3 > 0$  such that

$$\mathbb{P}\left(\sup_{j\neq i} \left| X_{j}^{\top} \hat{\boldsymbol{\beta}}_{[-i]} - X_{j}^{\top} \hat{\boldsymbol{\beta}} \right| \leq C_{1} \frac{K_{n}^{2} H_{n}}{\sqrt{n}} \right) \\
\geq 1 - C_{2}n \exp\left(-c_{2} H_{n}^{2}\right) - C_{3} \exp\left(-c_{3} K_{n}^{2}\right) \\
- \exp\left(-C_{4}n \left(1 + o(1)\right)\right) = 1 - o(1), \quad (149)$$

$$\mathbb{P}\left(\sup_{i} |X_{i}^{\top} \hat{\boldsymbol{\beta}} - prox_{q_{i}\rho}(X_{i}^{\top} \hat{\boldsymbol{\beta}}_{[-i]})| \leq C_{1} \frac{K_{n}^{2} H_{n}}{\sqrt{n}} \right) \\
\geq 1 - C_{2}n \exp\left(-c_{2} H_{n}^{2}\right) - C_{3} \exp\left(-c_{3} K_{n}^{2}\right) \\
- \exp\left(-C_{4}n \left(1 + o(1)\right)\right) = 1 - o(1). \quad (150)$$

Invoking this lemma, we see that the spectral norm of (148) is bounded above by some constant times

$$\frac{K_n^2 H_n}{\sqrt{n}} \Big\| \sum_{j: j \neq i} \tilde{X}_j \tilde{X}_j^\top / n \Big\|$$

with high probability as specified in (149). From Lemma 2, the spectral norm here is bounded by some constant with probability at least  $1 - c_1 \exp(-c_2 n)$ . These observations together with (87) and the fact that on  $A_n$  the minimum eigenvalues of  $\tilde{G}_{(i)}$  and  $\tilde{G}_{[-i]}$  are bounded by  $\lambda_{\rm lb}$  yield that

$$\mathbb{P}\left(\left\|\tilde{\boldsymbol{G}}_{(i)}^{-1} - \tilde{\boldsymbol{G}}_{[-i]}^{-1}\right\| \le C_1 \frac{K_n^2 H_n}{\sqrt{n}}\right) \ge 1 - C_2 n \exp\left(-c_2 H_n^2\right) - C_3 \exp\left(-c_3 K_n^2\right) - \exp\left(-c_4 n \left(1 + o(1)\right)\right).$$

This is true for any *i*. Hence, taking the union bound we obtain

$$\mathbb{P}\left(\sup_{i} \|\tilde{\boldsymbol{G}}_{(i)}^{-1} - \tilde{\boldsymbol{G}}_{[-i]}^{-1}\| \le C_1 \frac{K_n^2 H_n}{\sqrt{n}}\right)$$
  

$$\ge 1 - C_2 n^2 \exp\left(-c_2 H_n^2\right) - C_3 n \exp\left(-c_3 K_n^2\right) - \exp\left(-C_4 n \left(1 + o(1)\right)\right).$$
(151)

D Springer

In order to establish the first result, note that

$$\sup_{i} \frac{1}{n} \left| \tilde{X}_{i}^{\top} \tilde{G}_{(i)}^{-1} \tilde{X}_{i} - \tilde{X}_{i}^{\top} \tilde{G}_{[-i]}^{-1} \tilde{X}_{i} \right| \leq \sup_{i} \frac{\|X_{i}\|^{2}}{n} \sup_{i} \|\tilde{G}_{(i)}^{-1} - \tilde{G}_{[-i]}^{-1}\|$$

To obtain the second result, note that

$$\sup_{i} \left| \frac{1}{n} \operatorname{Tr}(\tilde{\boldsymbol{G}}_{(i)}^{-1}) - \frac{1}{n} \operatorname{Tr}(\tilde{\boldsymbol{G}}_{[-i]}^{-1}) \right| \le \frac{p-1}{n} \sup_{i} \|\tilde{\boldsymbol{G}}_{(i)}^{-1} - \tilde{\boldsymbol{G}}_{[-i]}^{-1} \|.$$

Therefore, combining (151) and Lemma 2 gives the desired result.

**Proof of Lemma 16** We restrict ourselves to the event  $A_n$  throughout. Recalling (135), one has

$$\operatorname{Tr}(\tilde{\boldsymbol{G}}_{(i)}^{-1}) - \operatorname{Tr}(\tilde{\boldsymbol{G}}^{-1}) = \frac{\rho''(\tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}})}{n} \frac{\tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-2} \tilde{\boldsymbol{X}}_i}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{\beta}})}{n} \tilde{\boldsymbol{X}}_i^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_i}$$

In addition, on  $A_n$  we have

$$\frac{1}{\lambda_{\mathrm{lb}}}\tilde{X}_{i}^{\top}\tilde{G}_{(i)}^{-1}\tilde{X}_{i}-\tilde{X}_{i}^{\top}\tilde{G}_{(i)}^{-2}\tilde{X}_{i}=\frac{1}{\lambda_{\mathrm{lb}}}\tilde{X}_{i}^{\top}\tilde{G}_{(i)}^{-1}\left(\tilde{G}_{(i)}-\lambda_{\mathrm{lb}}I\right)\tilde{G}_{(i)}^{-1}\tilde{X}_{i}\geq0.$$

Combining these results and recognizing that  $\rho'' > 0$ , we get

$$\left| \operatorname{Tr}(\tilde{\boldsymbol{G}}_{(i)}^{-1}) - \operatorname{Tr}(\tilde{\boldsymbol{G}}^{-1}) \right| \leq \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})}{n} \frac{\frac{1}{\lambda_{\mathrm{lb}}}\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}_{(i)}^{-1}\tilde{\boldsymbol{X}}_{i}}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})}{n}\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}_{(i)}^{-1}\tilde{\boldsymbol{X}}_{i}} \leq \frac{1}{\lambda_{\mathrm{lb}}} \quad (152)$$

as claimed.

## G Proof of Lemma 11

Again, we restrict ourselves to the event  $A_n$  on which  $\tilde{G} \succeq \lambda_{lb} I$ . Note that

$$\tilde{X}_i^{\top} \tilde{\boldsymbol{G}}^{-1} \boldsymbol{w} = \frac{1}{n} \tilde{X}_i^{\top} \tilde{\boldsymbol{G}}^{-1} \tilde{X}^{\top} \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} X_{\cdot 1}.$$

Note that  $\{\tilde{G}, \tilde{X}\}$  and  $X_{.1}$  are independent. Conditional on  $\tilde{X}$ , the left-hand side is Gaussian with mean zero and variance  $\frac{1}{n^2}\tilde{X}_i^{\top}\tilde{G}^{-1}\tilde{X}^{\top}D_{\tilde{\beta}}^2\tilde{X}\tilde{G}^{-1}\tilde{X}_i$ . The variance is bounded above by

$$\sigma_X^2 := \frac{1}{n^2} \tilde{X}_i^\top \tilde{\boldsymbol{G}}^{-1} \tilde{X}^\top \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}}^2 \tilde{X} \tilde{\boldsymbol{G}}^{-1} \tilde{X}_i \le \sup_i \left| \rho''(\tilde{X}_i^\top \tilde{\boldsymbol{\beta}}) \right| \cdot \frac{1}{n^2} \tilde{X}_i^\top \tilde{\boldsymbol{G}}^{-1} \tilde{X}^\top \boldsymbol{D}_{\tilde{\boldsymbol{\beta}}} \tilde{X} \tilde{\boldsymbol{G}}^{-1} \tilde{X}_i$$
$$= \frac{1}{n} \sup_i \left| \rho''(\tilde{X}_i^\top \tilde{\boldsymbol{\beta}}) \right| \cdot \tilde{X}_i^\top \tilde{\boldsymbol{G}}^{-1} \tilde{X}_i \lesssim \frac{1}{n} \| \tilde{X}_i \|^2$$
(153)

In turn, Lemma 2 asserts that  $n^{-1} \| \tilde{X}_i \|^2$  is bounded by a constant with high probability. As a result, applying Gaussian concentration results [60, Theorem 2.1.12] gives

$$| ilde{m{X}}_i^{ op} ilde{m{G}}^{-1} m{w}| \lesssim H_n$$

with probability exceeding  $1 - C \exp(-cH_n^2)$ , where C, c > 0 are universal constants.

In addition,  $\sup_i |X_{i1}| \leq H_n$  holds with probability exceeding  $1 - C \exp(-cH_n^2)$ . Putting the above results together, applying the triangle inequality  $|X_{i1} - \tilde{X}_i^{\top} \tilde{G}^{-1} w| \leq |X_{i1}| + |\tilde{X}_i^{\top} \tilde{G}^{-1} w|$ , and taking the union bound, we obtain

$$\mathbb{P}\left(\sup_{1\leq i\leq n}|X_{i1}-\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{G}}^{-1}\boldsymbol{w}| \lesssim H_{n}\right) \geq 1-Cn\exp\left(-cH_{n}^{2}\right) = 1-o(1).$$

## H Proof of Lemma 17

The goal of this section is to prove Lemma 17, which relates the full-model MLE  $\hat{\beta}$  and the MLE  $\hat{\beta}_{[-i]}$ . To this end, we establish the key lemma below.

**Lemma 18** Suppose  $\hat{\boldsymbol{\beta}}_{[-i]}$  denote the MLE when the *i*th observation is dropped. Further let  $\boldsymbol{G}_{[-i]}$  be as in (84), and define  $q_i$  and  $\hat{\boldsymbol{b}}$  as follows:

$$q_{i} = \frac{1}{n} X_{i}^{\top} G_{[-i]}^{-1} X_{i};$$
  
$$\hat{b} = \hat{\beta}_{[-i]} - \frac{1}{n} G_{[-i]}^{-1} X_{i} \left( \rho' \left( prox_{q_{i}\rho} \left( X_{i}^{\top} \hat{\beta}_{[-i]} \right) \right) \right).$$
(154)

Suppose  $K_n$ ,  $H_n$  are diverging sequences as in (132). Then there exist universal constants  $C_1$ ,  $C_2$ ,  $C_3 > 0$  such that

$$\mathbb{P}\left(\|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{b}}\| \le C_{1} \frac{K_{n}^{2} H_{n}}{n}\right) \ge 1 - C_{2} n \exp(-c_{2} H_{n}^{2}) 
- C_{3} \exp(-c_{3} K_{n}^{2}) - \exp(-C_{4} n (1 + o(1)));$$

$$\mathbb{P}\left(\sup_{j \ne i} |X_{j}^{\top} \hat{\boldsymbol{\beta}}_{[-i]} - X_{j}^{\top} \hat{\boldsymbol{b}}| \le C_{1} \frac{K_{n} H_{n}}{\sqrt{n}}\right) 
\ge 1 - C_{2} n \exp\left(-c_{2} H_{n}^{2}\right) - C_{3} \exp\left(-c_{3} K_{n}^{2}\right) 
- \exp\left(-C_{4} n (1 + o(1))\right).$$
(156)

Deringer

The proof ideas are inspired by the leave-one-observation-out approach of [25]. We however emphasize once more that the adaptation of these ideas to our setup is not straightforward and crucially hinges on Theorem 4, Lemma 7 and properties of the effective link function.

**Proof of Lemma 18** Invoking techniques similar to that for establishing Lemma 7, it can be shown that

$$\frac{1}{n}\sum_{i=1}^{n}\rho''(\gamma_{i}^{*})X_{i}X_{i}^{\top} \succeq \lambda_{\rm lb}I$$
(157)

with probability at least  $1 - \exp(\Omega(n))$ , where  $\gamma_i^*$  is between  $X_i^{\top} \hat{\boldsymbol{b}}$  and  $X_i^{\top} \hat{\boldsymbol{\beta}}$ . Denote by  $\mathcal{B}_n$  the event where (157) holds. Throughout this proof, we work on the event  $\mathcal{C}_n := \mathcal{A}_n \cap \mathcal{B}_n$ , which has probability  $1 - \exp(-\Omega(n))$ . As in (107) then,

$$\|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{b}}\| \le \frac{1}{n\lambda_{\rm lb}} \|\nabla \ell(\hat{\boldsymbol{b}})\|.$$
(158)

Next, we simplify (158). To this end, recall the defining relation of the proximal operator

$$b\rho'(\operatorname{prox}_{b\rho}(z)) + \operatorname{prox}_{b\rho}(z) = z,$$

which together with the definitions of  $\hat{b}$  and  $q_i$  gives

$$\boldsymbol{X}_{i}^{\top}\boldsymbol{\hat{b}} = \operatorname{prox}_{q_{i}\rho}\left(\boldsymbol{X}_{i}^{\top}\boldsymbol{\hat{\beta}}_{[-i]}\right).$$
(159)

Now, let  $\ell_{[-i]}$  denote the negative log-likelihood function when the *i*th observation is dropped, and hence  $\nabla \ell_{[-i]}(\hat{\beta}_{[-i]}) = 0$ . Expressing  $\nabla \ell(\hat{b})$  as  $\nabla \ell(\hat{b}) - \nabla \ell_{[-i]}(\hat{\beta}_{[-i]})$ , applying the mean value theorem, and using the analysis similar to that in [25, Proposition 3.4], we obtain

$$\frac{1}{n}\nabla\ell(\hat{\boldsymbol{b}}) = \frac{1}{n}\sum_{j:j\neq i} \left[\rho''(\boldsymbol{\gamma}_j^*) - \rho''(\boldsymbol{X}_j^{\top}\hat{\boldsymbol{\beta}}_{[-i]})\right]\boldsymbol{X}_j\boldsymbol{X}_j^{\top}\left(\hat{\boldsymbol{b}} - \hat{\boldsymbol{\beta}}_{[-i]}\right), \quad (160)$$

where  $\gamma_j^*$  is between  $X_j^{\top} \hat{\boldsymbol{b}}$  and  $X_j^{\top} \hat{\boldsymbol{\beta}}_{[-i]}$ . Combining (158) and (160) leads to the upper bound

$$\|\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{b}}\| \leq \frac{1}{\lambda_{\mathrm{lb}}} \left\| \frac{1}{n} \sum_{j:j \neq i} \boldsymbol{X}_{j} \boldsymbol{X}_{j}^{\top} \right\| \cdot \sup_{j \neq i} \left| \rho''(\boldsymbol{\gamma}_{j}^{*}) - \rho''(\boldsymbol{X}_{j}^{\top} \hat{\boldsymbol{\beta}}_{[-i]}) \right| \cdot \left\| \frac{1}{n} \boldsymbol{G}_{[-i]}^{-1} \boldsymbol{X}_{i} \right\|$$
$$\cdot \left| \rho'\left( \mathsf{prox}_{q_{i}\rho}(\boldsymbol{X}_{i}^{\top} \hat{\boldsymbol{\beta}}_{[-i]}) \right) \right|. \tag{161}$$

🖄 Springer

We need to control each term in the right-hand side. To start with, the first term is bounded by a universal constant with probability  $1 - \exp(-\Omega(n))$  (Lemma 2). For the second term, since  $\gamma_i^*$  is between  $X_i^{\top} \hat{\boldsymbol{b}}$  and  $X_i^{\top} \hat{\boldsymbol{\beta}}_{[-i]}$  and  $\|\rho^{\prime\prime\prime}\|_{\infty} < \infty$ , we get

$$\sup_{j\neq i} \left| \rho''(\gamma_j^*) - \rho''(\boldsymbol{X}_j^{\top} \hat{\boldsymbol{\beta}}_{[-i]}) \right| \le \|\rho'''\|_{\infty} \|\boldsymbol{X}_j^{\top} \hat{\boldsymbol{b}} - \boldsymbol{X}_j^{\top} \hat{\boldsymbol{\beta}}_{[-i]}\|$$
(162)

$$\leq \|\rho^{\prime\prime\prime}\|_{\infty} \left| \frac{1}{n} \boldsymbol{X}_{j}^{\top} \boldsymbol{G}_{[-i]}^{-1} \boldsymbol{X}_{i} \rho^{\prime} \left( \operatorname{prox}_{q_{i}\rho} \left( \boldsymbol{X}_{i}^{\top} \hat{\boldsymbol{\beta}}_{[-i]} \right) \right) \right|$$
(163)

$$\leq \|\rho^{\prime\prime\prime}\|_{\infty} \frac{1}{n} \sup_{j \neq i} \left| \boldsymbol{X}_{j}^{\top} \boldsymbol{G}_{[-i]}^{-1} \boldsymbol{X}_{i} \right| \cdot \left| \rho^{\prime} \left( \operatorname{prox}_{q_{i}\rho}(\boldsymbol{X}_{i}^{\top} \hat{\boldsymbol{\beta}}_{[-i]}) \right) \right|.$$
(164)

Given that  $\{X_j, G_{[-i]}\}$  and  $X_i$  are independent for all  $j \neq i$ , conditional on  $\{X_j, G_{[-i]}\}$  one has

$$\boldsymbol{X}_{j}^{\top}\boldsymbol{G}_{[-i]}^{-1}\boldsymbol{X}_{i}\sim\mathcal{N}\left(0,\boldsymbol{X}_{j}^{\top}\boldsymbol{G}_{[-i]}^{-2}\boldsymbol{X}_{j}
ight).$$

In addition, the variance satisfies

$$|X_{j}^{\top}G_{[-i]}^{-2}X_{j}| \leq \frac{\|X_{j}\|^{2}}{\lambda_{\rm lb}^{2}} \lesssim n$$
(165)

with probability at least  $1 - \exp(-\Omega(n))$ . Applying standard Gaussian concentration results [60, Theorem 2.1.12], we obtain

$$\mathbb{P}\left(\frac{1}{\sqrt{p}}\left|\boldsymbol{X}_{j}^{\top}\boldsymbol{G}_{[-i]}^{-1}\boldsymbol{X}_{i}\right| \geq C_{1}H_{n}\right) \leq C_{2}\exp\left(-c_{2}H_{n}^{2}\right) + \exp\left(-C_{3}n\left(1+o(1)\right)\right).$$
(166)

By the union bound

$$\mathbb{P}\left(\frac{1}{\sqrt{p}}\sup_{j\neq i} |X_{j}^{\top}G_{[-i]}^{-1}X_{i}| \le C_{1}H_{n}\right) \ge 1 - nC_{2}\exp\left(-c_{2}H_{n}^{2}\right) - \exp\left(-C_{3}n\left(1 + o(1)\right)\right). \quad (167)$$

Consequently,

$$\sup_{j \neq i} \left| \rho''(\gamma_j^*) - \rho''(\boldsymbol{X}_j^{\top} \hat{\boldsymbol{\beta}}_{[-i]}) \right| \lesssim \sup_{j \neq i} \|\boldsymbol{X}_j^{\top} \hat{\boldsymbol{b}} - \boldsymbol{X}_j^{\top} \hat{\boldsymbol{\beta}}_{[-i]} \|$$
$$\lesssim \frac{1}{\sqrt{n}} H_n \left| \rho' \left( \mathsf{prox}_{q_i \rho}(\boldsymbol{X}_i^{\top} \hat{\boldsymbol{\beta}}_{[-i]}) \right) \right|.$$
(168)

In addition, the third term in the right-hand side of (161) can be upper bounded as well since

$$\frac{1}{n} \|\boldsymbol{G}_{[-i]}^{-1} \boldsymbol{X}_i\| = \frac{1}{n} \sqrt{|\boldsymbol{X}_i^\top \boldsymbol{G}_{[-i]}^{-2} \boldsymbol{X}_i|} \lesssim \frac{1}{\sqrt{n}}$$
(169)

with high probability.

It remains to bound  $\left|\rho'\left(\operatorname{prox}_{q_i\rho}(X_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]})\right)\right|$ . To do this, we begin by considering  $\rho'(\operatorname{prox}_{c\rho}(Z))$  for any constant c > 0 (rather than a random variable  $q_i$ ). Recall that for any constant c > 0 and any  $Z \sim \mathcal{N}(0, \sigma^2)$  with finite variance, the random variable  $\rho'(\operatorname{prox}_{c\rho}(Z))$  is sub-Gaussian. Conditional on  $\hat{\boldsymbol{\beta}}_{[-i]}$ , one has  $X_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]} \sim \mathcal{N}(0, \|\hat{\boldsymbol{\beta}}_{[-i]}\|^2)$ . This yields

$$\mathbb{P}\left(\rho'\left(\mathsf{prox}_{c\rho}(\boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}}_{[-i]})\right) \geq C_{1}K_{n}\right) \leq C_{2}\mathbb{E}\left[\exp\left(-\frac{C_{3}^{2}K_{n}^{2}}{\|\hat{\boldsymbol{\beta}}_{[-i]}\|^{2}}\right)\right] \\ \leq C_{2}\exp\left(-C_{3}K_{n}^{2}\right) + C_{4}\exp\left(-C_{5}n\right)$$
(170)

for some constants  $C_1, C_2, C_3, C_4, C_5 > 0$  since  $\|\hat{\boldsymbol{\beta}}_{[-i]}\|$  is bounded with high probability (see Theorem 4).

Note that  $\frac{\partial \text{prox}_{b\rho}(z)}{\partial b} \leq 0$  by [22, Proposition 6.3]. Hence, in order to move over from the above concentration result established for a fixed constant *c* to the random variables  $q_i$ , it suffices to establish a uniform lower bound for  $q_i$  with high probability. Observe that for each *i*,

$$q_i \ge \frac{\|X_i\|^2}{n} \frac{1}{\|G_{[-i]}\|} \ge C^*$$

with probability  $1 - \exp(-\Omega(n))$ , where  $C^*$  is some universal constant. On this event, one has

$$\rho'\left(\operatorname{prox}_{q_i\rho}\left(\boldsymbol{X}_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]}\right)\right) \leq \rho'\left(\operatorname{prox}_{C^*\rho}\left(\boldsymbol{X}_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]}\right)\right).$$

This taken collectively with (170) yields

$$\mathbb{P}\left(\rho'(\operatorname{prox}_{q_i\rho}(\boldsymbol{X}_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]})) \leq C_1 K_n\right) \geq \mathbb{P}\left(\rho'(\operatorname{prox}_{C^*\rho}(\boldsymbol{X}_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]})) \leq C_1 K_n\right)$$

$$(171)$$

$$\geq 1 - C_2 \exp\left(-C_3 K_n^2\right) - C_4 \exp\left(-C_5 n\right).$$

$$(172)$$

This controls the last term.

To summarize, if  $\{K_n\}$  and  $\{H_n\}$  are diverging sequences satisfying the assumptions in (132), combining (161) and the bounds for each term in the right-hand side finally gives (155). On the other hand, combining (167) and (172) yields (156).

🖉 Springer

With the help of Lemma 18 we are ready to prove Lemma 17. Indeed, observe that

$$X_j^{ op}(\hat{\boldsymbol{\beta}}_{[-i]} - \hat{\boldsymbol{\beta}}) \Big| \leq \Big| X_j^{ op}(\hat{\boldsymbol{b}} - \hat{\boldsymbol{\beta}}) \Big| + \Big| X_j^{ op}(\hat{\boldsymbol{\beta}}_{[-i]} - \hat{\boldsymbol{b}}) \Big|,$$

and hence by combining Lemmas 2 and 18, we establish the first claim (149). The second claim (150) follows directly from Lemmas 2, 18 and (159).

## I Proof of Theorem 7(b)

This section proves that the random sequence  $\tilde{\alpha} = \text{Tr}(\tilde{\boldsymbol{G}}^{-1})/n$  converges in probability to the constant  $b_*$  defined by the system of Eqs. (25) and (26). To begin with, we claim that  $\tilde{\alpha}$  is close to a set of auxiliary random variables  $\{\tilde{q}_i\}$  defined below.

**Lemma 19** Define  $\tilde{q}_i$  to be

$$\tilde{q}_i = \frac{1}{n} \tilde{\boldsymbol{X}}_i^\top \tilde{\boldsymbol{G}}_{[-i]}^{-1} \tilde{\boldsymbol{X}}_i,$$

where  $\tilde{G}_{[-i]}$  is defined in 85.

Then there exist universal constants  $C_1$ ,  $C_2$ ,  $C_3$ ,  $C_4$ ,  $c_2$ ,  $c_3 > 0$  such that

$$\mathbb{P}\left(\sup_{i} |\tilde{q}_{i} - \tilde{\alpha}| \leq C_{1} \frac{K_{n}^{2} H_{n}}{\sqrt{n}}\right)$$
  
$$\geq 1 - C_{2} n^{2} \exp\left(c_{2} H_{n}^{2}\right) - C_{3} n \exp\left(-c_{3} K_{n}^{2}\right)$$
  
$$- \exp\left(-C_{4} n \left(1 + o(1)\right)\right) = 1 - o(1),$$

where  $K_n$ ,  $H_n$  are as in (132).

**Proof** This result follows directly from Proposition 1 and Eq. (144). A consequence is that  $\operatorname{prox}_{\tilde{q}_i\rho}\left(\boldsymbol{X}_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]}\right)$  becomes close to  $\operatorname{prox}_{\tilde{\alpha}\rho}\left(\boldsymbol{X}_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]}\right)$ .

**Lemma 20** Let  $\tilde{q}_i$  and  $\tilde{\alpha}$  be as defined earlier. Then one has

$$\mathbb{P}\left(\sup_{i}\left|\operatorname{prox}_{\tilde{q}_{i}\rho}\left(\boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}}_{\left[-i\right]}\right)-\operatorname{prox}_{\tilde{\alpha}\rho}\left(\boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}}_{\left[-i\right]}\right)\right|\leq C_{1}\frac{K_{n}^{3}H_{n}}{\sqrt{n}}\right)\\ \geq 1-C_{2}n^{2}\exp\left(-c_{2}H_{n}^{2}\right)-C_{3}n\exp\left(-c_{3}K_{n}^{2}\right)\\ -\exp\left(-C_{4}n\left(1+o(1)\right)\right)=1-o(1), \tag{173}$$

where  $K_n$ ,  $H_n$  are as in (132).

The key idea behind studying  $\operatorname{prox}_{\tilde{\alpha}\rho}\left(X_{i}^{\top}\hat{\boldsymbol{\beta}}_{[-i]}\right)$  is that it is connected to a random function  $\delta_{n}(\cdot)$  defined below, which happens to be closely related to the Eq. (26). In

fact, we will show that  $\delta_n(\tilde{\alpha})$  converges in probability to 0; the proof relies on the connection between  $\operatorname{prox}_{\tilde{\alpha}\rho}\left(\boldsymbol{X}_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]}\right)$  and the auxiliary quantity  $\operatorname{prox}_{\tilde{q}_i\rho}\left(\boldsymbol{X}_i^{\top}\hat{\boldsymbol{\beta}}_{[-i]}\right)$ . The formal results is this:

**Proposition 2** For any index *i*, let  $\hat{\boldsymbol{\beta}}_{[-i]}$  be the MLE obtained on dropping the *i*th observation. Define  $\delta_n(x)$  to be the random function

$$\delta_n(x) := \frac{p}{n} - 1 + \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + x\rho'' \left( \text{prox}_{x\rho} \left( X_i^\top \hat{\beta}_{[-i]} \right) \right)}.$$
 (174)

Then one has  $\delta_n(\tilde{\alpha}) \xrightarrow{\mathbb{P}} 0$ .

Furthermore, the random function  $\delta_n(x)$  converges to a deterministic function  $\Delta(x)$  defined by

$$\Delta(x) = \kappa - 1 + \mathbb{E}_Z \left[ \frac{1}{1 + x \rho''(\mathsf{prox}_{x\rho}(\tau_* Z))} \right],\tag{175}$$

where  $Z \sim \mathcal{N}(0, 1)$ , and  $\tau_*$  is such that  $(\tau_*, b_*)$  is the unique solution to (25) and (26).

**Proposition 3** With  $\Delta(x)$  as in (175),  $\Delta(\tilde{\alpha}) \xrightarrow{\mathbb{P}} 0$ .

In fact, one can easily verify that

$$\Delta(x) = \kappa - \mathbb{E}\left[\Psi'\left(\tau_* Z; x\right)\right],\tag{176}$$

and hence by Lemma 5, the solution to  $\Delta(x) = 0$  is exactly  $b_*$ . As a result, putting the above claims together, we show that  $\tilde{\alpha}$  converges in probability to  $b_*$ .

It remains to formally prove the preceding lemmas and propositions, which is the goal of the rest of this section.

Proof of Lemma 20 By [22, Proposition 6.3], one has

$$\frac{\partial \operatorname{prox}_{b\rho}(z)}{\partial b} = -\left. \frac{\rho'(x)}{1 + b\rho''(x)} \right|_{x = \operatorname{prox}_{b\rho}(z)}$$

which yields

$$\sup_{i} \left| \operatorname{prox}_{\tilde{q}_{i}\rho} \left( \boldsymbol{X}_{i}^{\top} \hat{\boldsymbol{\beta}}_{[-i]} \right) - \operatorname{prox}_{\tilde{\alpha}\rho} \left( \boldsymbol{X}_{i}^{\top} \hat{\boldsymbol{\beta}}_{[-i]} \right) \right|$$

$$= \sup_{i} \left[ \left| \frac{\rho'(x)}{1 + q_{\tilde{\alpha},i}\rho''(x)} \right|_{x = \operatorname{prox}_{q_{\tilde{\alpha},i}\rho} \left( \boldsymbol{X}_{i}^{\top} \hat{\boldsymbol{\beta}}_{[-i]} \right)} \right| \cdot |\tilde{q}_{i} - \tilde{\alpha}| \right]$$

$$\leq \sup_{i} \left| \rho' \left( \operatorname{prox}_{q_{\tilde{\alpha},i}} (\boldsymbol{X}_{i}^{\top} \hat{\boldsymbol{\beta}}_{[-i]}) \right) \right| \cdot \sup_{i} |\tilde{q}_{i} - \tilde{\alpha}|, \qquad (177)$$

where  $q_{\tilde{\alpha},i}$  is between  $\tilde{q}_i$  and  $\tilde{\alpha}$ . Here, the last inequality holds since  $q_{\tilde{\alpha},i}$ ,  $\rho'' \ge 0$ .

🖄 Springer

In addition, just as in the proof of Lemma 18, one can show that  $q_i$  is bounded below by some constant  $C^* > 0$  with probability  $1 - \exp(-\Omega(n))$ . Since  $q_{\tilde{\alpha},i} \ge \min\{\tilde{q}_i, \tilde{\alpha}\}$ , on the event  $\sup_i |\tilde{q}_i - \tilde{\alpha}| \le C_1 K_n^2 H_n / \sqrt{n}$ , which happens with high probability (Lemma 19),  $q_{\tilde{\alpha},i} \ge C_{\alpha}$  for some universal constant  $C_{\alpha} > 0$ . Hence, by an argument similar to that establishing (172), we have

$$\mathbb{P}\left(\sup_{i} \left|\rho'\left(\operatorname{prox}_{q_{\tilde{\alpha},i}}\left(\boldsymbol{X}_{i}^{\top}\hat{\boldsymbol{\beta}}_{\left[-i\right]}\right)\right)\right| \geq C_{1}K_{n}\right)$$
  
$$\leq C_{2}n^{2}\exp\left(-c_{2}H_{n}^{2}\right) + C_{3}n\exp\left(-c_{3}K_{n}^{2}\right) + \exp\left(-C_{4}n\left(1+o(1)\right)\right).$$

This together with (177) and Lemma 19 concludes the proof.

**Proof of Proposition 2** To begin with, recall from (138) and (139) that on  $A_n$ ,

$$\frac{p-1}{n} = \sum_{i=1}^{n} \frac{\frac{\rho''(\tilde{X}_{i}^{\top}\tilde{\boldsymbol{\beta}})}{n} \tilde{X}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{X}_{i}}{1 + \frac{\rho''(\tilde{X}_{i}^{\top}\tilde{\boldsymbol{\beta}})}{n} \tilde{X}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{X}_{i}} = 1 - \frac{1}{n} \sum_{i=1}^{n} \frac{1}{1 + \frac{\rho''(\tilde{X}_{i}^{\top}\tilde{\boldsymbol{\beta}})}{n} \tilde{X}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{X}_{i}}.$$
(178)

Using the fact that  $\left|\frac{1}{1+x} - \frac{1}{1+y}\right| \le |x - y|$  for  $x, y \ge 0$ , we obtain

$$\begin{aligned} \left| \frac{1}{n} \sum_{i=1}^{n} \frac{1}{1 + \frac{\rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}})}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i}} - \frac{1}{n} \sum_{i=1}^{n} \frac{1}{1 + \rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) \tilde{\alpha}} \right| \\ & \leq \frac{1}{n} \sum_{i=1}^{n} \rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) \left| \frac{1}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i} - \tilde{\alpha} \right| \\ & \leq \|\rho''\|_{\infty} \sup_{i} \left| \frac{1}{n} \tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{G}}_{(i)}^{-1} \tilde{\boldsymbol{X}}_{i} - \tilde{\alpha} \right| \\ & = \|\rho''\|_{\infty} \sup_{i} |\eta_{i}| \leq C_{1} \frac{K_{n}^{2} H_{n}}{\sqrt{n}}, \end{aligned}$$

with high probability (Proposition 1). This combined with (178) yields

$$\mathbb{P}\left(\left|\frac{p-1}{n}-1+\frac{1}{n}\sum_{i=1}^{n}\frac{1}{1+\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})\tilde{\boldsymbol{\alpha}}}\right| \geq C_{1}\frac{K_{n}^{2}H_{n}}{\sqrt{n}}\right)$$
$$\leq C_{2}n^{2}\exp\left(-c_{2}H_{n}^{2}\right)+C_{3}n\exp\left(-c_{3}K_{n}^{2}\right)+\exp\left(-C_{4}n\left(1+o(1)\right)\right).$$

The above bound concerns  $\frac{1}{n} \sum_{i=1}^{n} \frac{1}{1+\rho''(\tilde{X}_{i}^{\top}\tilde{\boldsymbol{\beta}})\tilde{\alpha}}$ , and it remains to relate it to  $\frac{1}{n} \sum_{i=1}^{n} \frac{1}{1+\rho''(\operatorname{prox}_{\tilde{\alpha}\rho}(\tilde{X}_{i}^{\top}\tilde{\boldsymbol{\beta}}))\tilde{\alpha}}$ . To this end, we first get from the uniform boundedness of  $\rho'''$  and Lemma 17 that

$$\mathbb{P}\left(\sup_{i}\left|\rho''(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}})-\rho''\left(\operatorname{prox}_{\tilde{q}_{i}\rho}(\tilde{\boldsymbol{X}}_{i}^{\top}\tilde{\boldsymbol{\beta}}_{[-i]})\right)\right|\geq C_{1}\frac{K_{n}^{2}H_{n}}{\sqrt{n}}\right)$$
  
$$\leq C_{2}n\exp(-c_{2}H_{n}^{2})+C_{3}\exp(-c_{3}K_{n}^{2})+\exp(-C_{4}n(1+o(1))). \quad (179)$$

Note that

$$\begin{aligned} \left| \frac{1}{n} \sum_{i=1}^{n} \frac{1}{1 + \rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) \tilde{\alpha}} - \frac{1}{n} \sum_{i=1}^{n} \frac{1}{1 + \rho''(\operatorname{prox}_{\tilde{\alpha}\rho}(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}_{[-i]})) \tilde{\alpha}} \right| \\ &\leq |\tilde{\alpha}| \sup_{i} \left| \rho''(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}) - \rho''\left(\operatorname{prox}_{\tilde{\alpha}\rho}(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}_{[-i]})\right) \right| \\ &\leq |\tilde{\alpha}| \sup_{i} \left\{ \left| \rho''\left(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}\right) - \rho''\left(\operatorname{prox}_{\tilde{q}_{i}\rho}(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}_{[-i]})\right) \right| \\ &+ \left| \rho''\left(\operatorname{prox}_{\tilde{q}_{i}\rho}(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}_{[-i]})\right) - \rho''\left(\operatorname{prox}_{\tilde{\alpha}\rho}(\tilde{\boldsymbol{X}}_{i}^{\top} \tilde{\boldsymbol{\beta}}_{[-i]})\right) \right| \right\}. \end{aligned}$$

By the bound (179), an application of Lemma 20, and the fact that  $\tilde{\alpha} \leq p/(n\lambda_{\rm lb})$  (on  $\mathcal{A}_n$ ), we obtain

$$\mathbb{P}\left(\left|\frac{p}{n}-1+\frac{1}{n}\sum_{i=1}^{n}\frac{1}{1+\rho''(\operatorname{prox}_{\tilde{\alpha}\rho}(X_{i}^{\top}\hat{\beta}_{[-i]}))\tilde{\alpha}}\right| \geq C_{1}\frac{K_{n}^{3}H_{n}}{\sqrt{n}}\right) \leq C_{2}n^{2}\exp\left(-c_{2}H_{n}^{2}\right)+C_{3}n\exp\left(-c_{3}K_{n}^{2}\right)+\exp\left(-C_{4}n(1+o(1))\right).$$

This establishes that  $\delta_n(\tilde{\alpha}) \xrightarrow{\mathbb{P}} 0$ .

**Proof of Proposition 3** Here we only provide the main steps of the proof. Note that since  $0 < \alpha \le p/(n\lambda_{lb}) := B$  on  $\mathcal{A}_n$ , it suffices to show that

$$\sup_{x\in[0,B]}|\delta_n(x)-\Delta(x)| \stackrel{\mathbb{P}}{\to} 0.$$

We do this by following three steps. Below, M > 0 is some sufficiently large constant.

1. First we truncate the random function  $\delta_n(x)$  and define

$$\tilde{\delta}_n(x) = \frac{p}{n} - 1 + \sum_{i=1}^n \frac{1}{1 + x\rho'' \left( \mathsf{prox}_{x\rho} \left( X_i^\top \hat{\beta}_{[-i]} \mathbf{1}_{\{\| \hat{\beta}_{[-i]} \| \le M\}} \right) \right)}.$$

The first step is to show that  $\sup_{x \in [0,B]} \left| \tilde{\delta}_n(x) - \delta_n(x) \right| \xrightarrow{\mathbb{P}} 0$ . This step can be established using Theorem 4 and some straightforward analysis. We stress that this truncation does not arise in [25], and it is required to keep track of the truncation throughout the rest of the proof.

- 2. Show that  $\sup_{x \in [0,B]} \left| \tilde{\delta}_n(x) \mathbb{E} \left[ \tilde{\delta}_n(x) \right] \right| \stackrel{\mathbb{P}}{\to} 0.$
- 3. Show that  $\sup_{x \in [0,B]} \left| \mathbb{E} \left[ \tilde{\delta}_n(x) \right] \Delta(x) \right| \xrightarrow{\mathbb{P}} 0.$

Steps 2 and 3 can be established by arguments similar to that in [25, Lemma 3.24,3.25], with neceassary modifications for our setup. We skip the detailed arguments here and refer the reader to [58].

## References

- 1. Agresti, A., Kateri, M.: Categorical Data Analysis. Springer, Berlin (2011)
- 2. Alon, N., Spencer, J.H.: The Probabilistic Method, 3rd edn. Wiley, Hoboken (2008)
- Amelunxen, D., Lotz, M., McCoy, M.B., Tropp, J.A.: Living on the edge: phase transitions in convex programs with random data. Inf. Inference 3, 224–294 (2014)
- Baricz, Á.: Mills' ratio: monotonicity patterns and functional inequalities. J. Math. Anal. Appl. 340(2), 1362–1370 (2008)
- Bartlett, M.S.: Properties of sufficiency and statistical tests. Proc. R. Soc. Lond. Ser. A Math. Phys. Sci. 160, 268–282 (1937)
- Bayati, M., Lelarge, M., Montanari, A., et al.: Universality in polytope phase transitions and message passing algorithms. Ann. Appl. Probab. 25(2), 753–822 (2015)
- Bayati, M., Montanari, A.: The dynamics of message passing on dense graphs, with applications to compressed sensing. IEEE Trans. Inf. Theory 57(2), 764–785 (2011)
- Bayati, M., Montanari, A.: The LASSO risk for Gaussian matrices. IEEE Trans. Inf. Theory 58(4), 1997–2017 (2012)
- Bickel, P.J., Ghosh, J.K.: A decomposition for the likelihood ratio statistic and the Bartlett correction—a Bayesian argument. Ann. Stat. 18, 1070–1090 (1990)
- Boucheron, S., Massart, P.: A high-dimensional Wilks phenomenon. Probab. Theory Relat. Fields 150(3–4), 405–433 (2011)
- Box, G.: A general distribution theory for a class of likelihood criteria. Biometrika 36(3/4), 317–346 (1949)
- Candès, E., Fan, Y., Janson, L., Lv, J.: Panning for gold: model-free knockoffs for high-dimensional controlled variable selection (2016). ArXiv preprint arXiv:1610.02351
- 13. Chernoff, H.: On the distribution of the likelihood ratio. Ann. Math. Stat. 25, 573–578 (1954)
- Cordeiro, G.M.: Improved likelihood ratio statistics for generalized linear models. J. R. Stat. Soc. Ser. B (Methodol.) 25, 404–413 (1983)
- Cordeiro, G.M., Cribari-Neto, F.: An Introduction to Bartlett Correction and Bias Reduction. Springer, New York (2014)
- Cordeiro, G.M., Cribari-Neto, F., Aubin, E.C.Q., Ferrari, S.L.P.: Bartlett corrections for one-parameter exponential family models. J. Stat. Comput. Simul. 53(3–4), 211–231 (1995)
- 17. Cover, T.M.: Geometrical and statistical properties of linear threshold devices. Ph.D. thesis (1964)
- Cover, T.M.: Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. IEEE Trans. Electron. Comput. 3, 326–334 (1965)
- 19. Cover, T.M., Thomas, J.A.: Elements of Information Theory. Wiley, Hoboken (2012)
- Cribari-Neto, F., Cordeiro, G.M.: On Bartlett and Bartlett-type corrections Francisco Cribari-Neto. Econom. Rev. 15(4), 339–367 (1996)
- 21. Deshpande, Y., Montanari, A.: Finding hidden cliques of size  $\sqrt{N/e}$  in nearly linear time. Found. Comput. Math. **15**(4), 1069–1128 (2015)
- Donoho, D., Montanari, A.: High dimensional robust M-estimation: asymptotic variance via approximate message passing. Probab. Theory Relat. Fields 3, 935–969 (2013)
- 23. Donoho, D., Montanari, A.: Variance breakdown of Huber (M)-estimators:  $n/p \rightarrow m \in (1, \infty)$ . Technical report (2015)
- El Karoui, N.: Asymptotic behavior of unregularized and ridge-regularized high-dimensional robust regression estimators: rigorous results (2013). ArXiv preprint arXiv:1311.2445
- El Karoui, N.: On the impact of predictor geometry on the performance on high-dimensional ridgeregularized generalized robust regression estimators. Probab. Theory Relat. Fields 170, 95–175 (2017)

- El Karoui, N., Bean, D., Bickel, P.J., Lim, C., Yu, B.: On robust regression with high-dimensional predictors. Proc. Natl. Acad. Sci. 110(36), 14557–14562 (2013)
- Fan, J., Jiang, J.: Nonparametric inference with generalized likelihood ratio tests. Test 16(3), 409–444 (2007)
- Fan, J., Lv, J.: Nonconcave penalized likelihood with NP-dimensionality. IEEE Trans. Inf. Theory 57(8), 5467–5484 (2011)
- Fan, J., Zhang, C., Zhang, J.: Generalized likelihood ratio statistics and Wilks phenomenon. Ann. Stat. 29, 153–193 (2001)
- Fan, Y., Demirkaya, E., Lv, J.: Nonuniformity of p-values can occur early in diverging dimensions (2017). arXiv:1705.03604
- 31. Hager, W.W.: Updating the inverse of a matrix. SIAM Rev. 31(2), 221-239 (1989)
- Hanson, D.L., Wright, F.T.: A bound on tail probabilities for quadratic forms in independent random variables. Ann. Math. Stat. 42(3), 1079–1083 (1971)
- 33. He, X., Shao, Q.-M.: On parameters of increasing dimensions. J. Multivar. Anal. 73(1), 120–135 (2000)
- Hsu, D., Kakade, S., Zhang, T.: A tail inequality for quadratic forms of subgaussian random vectors. Electron. Commun. Probab. 17(52), 1–6 (2012)
- 35. Huber, P.J.: Robust regression: asymptotics, conjectures and Monte Carlo. Ann. Stat. 1, 799-821 (1973)
- 36. Huber, P.J.: Robust Statistics. Springer, Berlin (2011)
- Janková, J., Van De Geer, S.: Confidence regions for high-dimensional generalized linear models under sparsity (2016). ArXiv preprint arXiv:1610.01353
- Javanmard, A., Montanari, A.: State evolution for general approximate message passing algorithms, with applications to spatial coupling. Inf. Inference 2, 115–144 (2013)
- Javanmard, A., Montanari, A.: De-biasing the lasso: optimal sample size for Gaussian designs (2015). ArXiv preprint arXiv:1508.02757
- Lawley, D.N.: A general method for approximating to the distribution of likelihood ratio criteria. Biometrika 43(3/4), 295–303 (1956)
- 41. Lehmann, E.L., Romano, J.P.: Testing Statistical Hypotheses. Springer, Berlin (2006)
- Liang, H., Pang, D., et al.: Maximum likelihood estimation in logistic regression models with a diverging number of covariates. Electron. J. Stat. 6, 1838–1846 (2012)
- Mammen, E.: Asymptotics with increasing dimension for robust regression with applications to the bootstrap. Ann. Stat. 17, 382–400 (1989)
- McCullagh, P., Nelder, J.A.: Generalized Linear Models. Monograph on Statistics and Applied Probability. Chapman & Hall, London (1989)
- Moulton, L.H., Weissfeld, L.A., Laurent, R.T.S.: Bartlett correction factors in logistic regression models. Comput. Stat. Data Anal. 15(1), 1–11 (1993)
- Oymak, S., Tropp, J.A.: Universality laws for randomized dimension reduction, with applications. Inf. Inference J. IMA 7, 337–446 (2015)
- 47. Parikh, N., Boyd, S.: Proximal algorithms. Found. Trends Optim. 1(3), 127-239 (2014)
- 48. Portnoy, S.: Asymptotic behavior of M-estimators of p regression parameters when  $p^2/n$  is large. I. Consistency. Ann. Stat. **12**, 1298–1309 (1984)
- Portnoy, S.: Asymptotic behavior of M-estimators of p regression parameters when p<sup>2</sup>/n is large; II. Normal approximation. Ann. Stat. 13, 1403–1417 (1985)
- Portnoy, S.: Asymptotic behavior of the empiric distribution of m-estimated residuals from a regression model with many parameters. Ann. Stat. 14, 1152–1170 (1986)
- 51. Portnoy, S., et al.: Asymptotic behavior of likelihood methods for exponential families when the number of parameters tends to infinity. Ann. Stat. **16**(1), 356–366 (1988)
- Rudelson, M., Vershynin, R., et al.: Hanson-Wright inequality and sub-Gaussian concentration. Electron. Commun. Probab. 18(82), 1–9 (2013)
- Sampford, M.R.: Some inequalities on Mill's ratio and related functions. Ann. Math. Stat. 24(1), 130–132 (1953)
- Spokoiny, V.: Penalized maximum likelihood estimation and effective dimension (2012). ArXiv preprint arXiv:1205.0498
- Su, W., Bogdan, M., Candes, E.: False discoveries occur early on the Lasso path. Ann. Stat. 45, 2133–2150 (2015)
- Sur, P., Candès, E.J.: Additional supplementary materials for: a modern maximum-likelihood theory for high-dimensional logistic regression. https://statweb.stanford.edu/~candes/papers/proofs\_ LogisticAMP.pdf (2018)

- Sur, P., Candès, E.J.: A modern maximum-likelihood theory for high-dimensional logistic regression (2018). ArXiv preprint arXiv:1803.06964
- Sur, P., Chen, Y., Candès, E.: Supplemental materials for "the likelihood ratio test in high-dimensional logistic regression is asymptotically a rescaled chi-square". http://statweb.stanford.edu/~candes/ papers/supplement\_LRT.pdf (2017)
- 59. Tang, C.Y., Leng, C.: Penalized high-dimensional empirical likelihood. Biometrika 97, 905–919 (2010)
- 60. Tao, T.: Topics in Random Matrix Theory, vol. 132. American Mathematical Society, Providence (2012)
- Thrampoulidis, C., Abbasi, E., Hassibi, B.: Precise error analysis of regularized m-estimators in highdimensions (2016). ArXiv preprint arXiv:1601.06233
- Van de Geer, S., Bühlmann, P., Ritov, Y., Dezeure, R., et al.: On asymptotically optimal confidence regions and tests for high-dimensional models. Ann. Stat. 42(3), 1166–1202 (2014)
- Van de Geer, S.A., et al.: High-dimensional generalized linear models and the lasso. Ann. Stat. 36(2), 614–645 (2008)
- 64. Van der Vaart, A.W.: Asymptotic Statistics, vol. 3. Cambridge University Press, Cambridge (2000)
- Vershynin, R.: Introduction to the non-asymptotic analysis of random matrices. In: Compressed Sensing: Theory and Applications, pp. 210–268 (2012)
- 66. Wilks, S.S.: The large-sample distribution of the likelihood ratio for testing composite hypotheses. Ann. Math. Stat. 9(1), 60–62 (1938)
- Yan, T., Li, Y., Xu, J., Yang, Y., Zhu, J.: High-dimensional Wilks phenomena in some exponential random graph models (2012). ArXiv preprint arXiv:1201.0058

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## Affiliations

# Pragya Sur<sup>1</sup> · Yuxin Chen<sup>2</sup> · Emmanuel J. Candès<sup>1,3</sup>

- <sup>1</sup> Department of Statistics, Stanford University, Stanford, CA 94305, USA
- <sup>2</sup> Department of Electrical Engineering, Princeton University, Princeton, NJ 08544, USA
- <sup>3</sup> Department of Mathematics, Stanford University, Stanford, CA 94305, USA