

Weierstraß-Institut
für Angewandte Analysis und Stochastik
Leibniz-Institut im Forschungsverbund Berlin e. V.

Preprint

ISSN 2198-5855

Bootstrap confidence sets under model misspecification

Vladimir Spokoiny ^{1,2} Mayya Zhilova ³

submitted: November 18, 2014

¹ HU Berlin and Weierstrass Institute,
Mohrenstr. 39
10117 Berlin
Germany
E-Mail: vladimir.spokoiny@wias-berlin.de

² Moscow Institute of Physics and Technology
9 Institutskiy per., Dolgoprudny
Moscow Region, 141700
Russian Federation

³ Weierstrass Institute
Mohrenstr. 39
10117 Berlin
Germany
E-Mail: mayya.zhilova@wias-berlin.de

No. 1992
Berlin 2014



2010 *Mathematics Subject Classification.* 62F25 (Primary) 62F40, 62E17 (Secondary).

Key words and phrases. likelihood-based bootstrap confidence set, misspecified model, finite sample size, multiplier bootstrap, weighted bootstrap, Gaussian approximation, Pinsker's inequality.

¹ The author is partially supported by Laboratory for Structural Methods of Data Analysis in Predictive Modeling, MIPT, RF government grant, ag. 11.G34.31.0073. Financial support by the German Research Foundation (DFG) through the Research Unit 1735 is gratefully acknowledged.

² Financial support by the German Research Foundation (DFG) through the Collaborative Research Center 649 "Economic Risk" is gratefully acknowledged.

Edited by
Weierstraß-Institut für Angewandte Analysis und Stochastik (WIAS)
Leibniz-Institut im Forschungsverbund Berlin e. V.
Mohrenstraße 39
10117 Berlin
Germany

Fax: +49 30 20372-303
E-Mail: preprint@wias-berlin.de
World Wide Web: <http://www.wias-berlin.de/>

Abstract

A multiplier bootstrap procedure for construction of likelihood-based confidence sets is considered for finite samples and a possible model misspecification. Theoretical results justify the bootstrap consistency for a small or moderate sample size and allow to control the impact of the parameter dimension p : the bootstrap approximation works if p^3/n is small. The main result about bootstrap consistency continues to apply even if the underlying parametric model is misspecified under the so called Small Modeling Bias condition. In the case when the true model deviates significantly from the considered parametric family, the bootstrap procedure is still applicable but it becomes a bit conservative: the size of the constructed confidence sets is increased by the modeling bias. We illustrate the results with numerical examples for misspecified constant and logistic regressions.

1 Introduction

Since introducing in 1979 by [Efron \(1979\)](#) the bootstrap procedure became one of the most powerful and common tools in statistical confidence estimation and hypothesis testing. Many versions and extensions of the original bootstrap method have been proposed in the literature; see e.g. [Wu \(1986\)](#), [Newton and Raftery \(1994\)](#); [Barbe and Bertail \(1995\)](#); [Horowitz \(2001\)](#); [Chatterjee and Bose \(2005\)](#); [Ma and Kosorok \(2005\)](#); [Chen and Pouzo \(2009\)](#); [Lavergne and Patilea \(2013\)](#); [Chen and Pouzo \(2014\)](#) among many others. This paper focuses on the multiplier bootstrap procedure which attracted a lot of attention last time due to its nice theoretical properties and numerical performance. We mention the papers [Chatterjee and Bose \(2005\)](#), [Arlot et al. \(2010\)](#) and [Chernozhukov et al. \(2013\)](#) for the most advanced recent results. [Chatterjee and Bose \(2005\)](#) showed some results on asymptotic bootstrap consistency in a very general framework: for estimators obtained by solving estimating equations. [Chernozhukov et al. \(2013\)](#) presented a number of non-asymptotic results on bootstrap validity with applications to special problems like testing many moment restrictions or parameter choice for a LASSO procedure. [Arlot et al. \(2010\)](#) constructed a non-asymptotical confidence bound in ℓ_s norm ($s \in [1, \infty]$) for the mean of a sample of high dimensional i.i.d. Gaussian vectors (or with a symmetric and bounded distribution), using the generalized weighted bootstrap for resampling of the quantiles.

This paper makes a further step in studying the multiplier bootstrap method in the problem of confidence estimation by a quasi maximum likelihood method. For a rather general parametric model, we consider likelihood-based confidence sets with the radius determined by a multiplier bootstrap. The aim of the study is to check the validity of the bootstrap procedure in situations with a large parameter dimension, a limited sample size, and a possible misspecification of the parametric assumption. The main result of the paper explicitly describes the error term of the bootstrap approximation. This particularly allows to track the impact of the parameter

dimension p and of the sample size n in the quality of the bootstrap procedure. As one of the corollaries, we show bootstrap validity under the constraint “ p^3/n -small”. Chatterjee and Bose (2005) stated results under the condition “ p/n -small” but their results only apply to low dimensional projections of the MLE vector. In the likelihood based approach, the construction involves the Euclidean norm of the MLE which leads to completely different tools and results. Chernozhukov et al. (2013) allowed a huge parameter dimension with “ $\log(p)/n$ small” but they essentially work with a family of univariate tests which again differs essentially from the maximum likelihood approach.

Another interesting and important issue is the impact of the model misspecification on the accuracy of bootstrap approximation. A surprising corollary of our error bounds is that the bootstrap confidence set can be used even if the underlying parametric model is slightly misspecified under the so called *small modeling bias (SmB)* condition. If the modeling bias becomes large, the bootstrap confidence sets are still applicable, but they become more and more conservative. (SmB) condition is given in Section 4 and it is consistent with classical bias-variance relation in nonparametric estimation.

Our theoretical study uses the square-root Wilks (sq-Wilks) expansion from Spokoiny (2012a), Spokoiny (2013) which approximates the square root likelihood ratio statistic by the norm of the standardized score vector. Further we extend the sq-Wilks expansion to the bootstrap log-likelihood and adopt the Gaussian approximation theory (GAR) to the special case when the distribution of the Euclidean norm of a non-Gaussian vector is approximated by the distribution of the norm of a Gaussian one with the same first and second moments. The Gaussian comparison technique based on the Pinsker inequality completes the study and allows to bridge the real unknown coverage probability and the conditional bootstrap coverage probability under (SmB) condition. In the case of a large modeling bias we state a one-sided bound: the bootstrap quantiles are uniformly larger than the real ones. This effect is nicely confirmed by our simulation study.

Now consider the problem and the approach in more detail. Let the data sample $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ consist of *independent* random observations and belong to the probability space $(\Omega, \mathcal{F}, \mathbb{P})$. We do not assume that the observations Y_i are identically distributed, moreover, no specific parametric structure of \mathbb{P} is being required. In order to explain the idea of the approach we start here with a parametric case, however the assumption (1.1) below is not required for the results. Let \mathbb{P} belong to some known regular parametric family $\{\mathbb{P}_\theta\} \stackrel{\text{def}}{=} \{\mathbb{P}_\theta \ll \mu_0, \theta \in \Theta \subset \mathbb{R}^p\}$. In this case the true parameter $\theta^* \in \Theta$ is such that

$$\mathbb{P} \equiv \mathbb{P}_{\theta^*} \in \{\mathbb{P}_\theta\}, \quad (1.1)$$

and the initial problem of finding the properties of unknown distribution \mathbb{P} is reduced to the equivalent problem for the finite-dimensional parameter θ^* . The parametric family $\{\mathbb{P}_\theta\}$ induces the log-likelihood process $L(\theta)$ of the sample \mathbf{Y} :

$$L(\theta) = L(\mathbf{Y}, \theta) \stackrel{\text{def}}{=} \log \left\{ \frac{d\mathbb{P}_\theta}{d\mu_0}(\mathbf{Y}) \right\}$$

and the maximum likelihood estimate (MLE) of $\boldsymbol{\theta}^*$:

$$\tilde{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}). \quad (1.2)$$

The asymptotic Wilks phenomenon [Wilks \(1938\)](#) states that for the case of i.i.d. observations with the sample size tending to the infinity the likelihood ratio statistic converges in distribution to $\chi_p^2/2$, where p is the parameter dimension:

$$2\{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\} \xrightarrow{w} \chi_p^2, \quad n \rightarrow \infty.$$

Define the likelihood-based confidence set as

$$\mathcal{E}(\mathfrak{z}) \stackrel{\text{def}}{=} \left\{ \boldsymbol{\theta} : L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}) \leq \mathfrak{z}^2/2 \right\}, \quad (1.3)$$

then the Wilks phenomenon implies

$$\mathbb{P} \left\{ \boldsymbol{\theta}^* \in \mathcal{E}(\mathfrak{z}_{\alpha, \chi_p^2}) \right\} \rightarrow \alpha, \quad n \rightarrow \infty,$$

where $\mathfrak{z}_{\alpha, \chi_p^2}^2$ is the $(1 - \alpha)$ -quantile for the χ_p^2 distribution. This result is very important and useful under the parametric assumption, i.e. when (1.1) holds. In this case the limit distribution of the likelihood ratio is independent of the model parameters or in other words it is *pivotal*. By this result a sufficiently large sample size allows to construct the confidence sets for $\boldsymbol{\theta}^*$ with a given coverage probability. However, a possibly low speed of convergence of the likelihood ratio statistic makes the asymptotic Wilks result hardly applicable to the case of small or moderate samples. Moreover, the asymptotical pivotality breaks down if the parametric assumption (1.1) does not hold (see [Huber \(1967\)](#)), and, therefore, the whole approach may be misleading if the model is considerably misspecified. If the assumption (1.1) does not hold, then the “true” parameter is defined by the projection of the true measure \mathbb{P} on the parametric family $\{\mathbb{P}_{\boldsymbol{\theta}}\}$:

$$\boldsymbol{\theta}^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \mathbb{E}L(\boldsymbol{\theta}). \quad (1.4)$$

The recent results by [Spokoiny \(2012a\)](#), [Spokoiny \(2013\)](#) provide a non-asymptotic version of square-root Wilks phenomenon for the case of misspecified model. It holds with an exponentially high probability

$$\left| \sqrt{2\{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\}} - \|\boldsymbol{\xi}\| \right| \leq \Delta_w \simeq \frac{p}{\sqrt{n}}, \quad (1.5)$$

where $\boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*)$, $D_0^2 \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E}L(\boldsymbol{\theta}^*)$. The bound is non-asymptotical, the approximation error term Δ_w has an explicit form (the precise statement is given in [Theorem A.2](#), [Section A.1](#), and it depends on the parameter dimension p , sample size n , and the probability of the random set on which the result holds.

Due to this bound, the original problem of finding a quantile of the LR test statistic $L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)$ is reduced to a similar question for the approximating quantity $\|\boldsymbol{\xi}\|$. The difficulty here is that in general $\|\boldsymbol{\xi}\|$ is non-pivotal, it depends on the unknown distribution \mathbb{P} and the target parameter

θ^* . Another result by [Spokoiny \(2012b\)](#) gives the following non-asymptotical deviation bound for $\|\xi\|^2$: for some explicit constant $C > 0$ it holds for $x \geq \sqrt{p}$

$$\mathbb{P}(\|\xi\|^2 \geq \mathbb{E}\|\xi\|^2 + Cx) \leq 2e^{-x}$$

(the precise statement is given in [Theorem A.3](#). This is a non-asymptotic deviation bound, sharp in leading approximating terms, however, the critical values yielded by it are too conservative for a valuable confidence set.

In the present work we study the *multiplier bootstrap* (or *weighted bootstrap*) procedure for estimation of the quantiles of the likelihood ratio statistic. The idea of the procedure is to mimic a distribution of the likelihood ratio statistic by reweighing its summands with random multipliers independent of the data:

$$L^\circ(\theta) \stackrel{\text{def}}{=} \sum_{i=1}^n \log \left\{ \frac{d\mathbb{P}_\theta}{d\mu_0}(Y_i) \right\} u_i.$$

Here the probability distribution is taken conditionally on the data \mathbf{Y} , which is denoted by the sign $^\circ$. The random weights u_1, \dots, u_n are i.i.d. with continuous c.d.f., independent of \mathbf{Y} and it holds for them: $\mathbb{E}(u_i) = 1$, $\text{Var}(u_i) = 1$, $\mathbb{E} \exp(u_i) < \infty$. Therefore, the multiplier bootstrap induces the probability space conditional on the data \mathbf{Y} . A simple but important observation is that $\mathbb{E}^\circ L^\circ(\theta) \equiv \mathbb{E}[L^\circ(\theta) | \mathbf{Y}] = L(\theta)$, and hence,

$$\text{argmax}_\theta \mathbb{E}^\circ L^\circ(\theta) = \text{argmax}_\theta L(\theta) = \tilde{\theta}.$$

This means that the target parameter in the bootstrap world is precisely known and it coincides with the maximum likelihood estimator $\tilde{\theta}$ conditioned on \mathbf{Y} , therefore, the bootstrap likelihood ratio statistic $L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta}) \stackrel{\text{def}}{=} \sup_{\theta \in \Theta} L^\circ(\theta) - L^\circ(\tilde{\theta})$ is fully computable and leads to a simple computational procedure for the approximation of the distribution of $L(\tilde{\theta}) - L(\theta^*)$.

The goal of the present study is to show in a non-asymptotic way the consistency of the described multiplier bootstrap procedure and to obtain an explicit bound on the error of coverage probability. In other words, we are interested in non-asymptotic approximation of the distribution of $\{L(\tilde{\theta}) - L(\theta^*)\}^{1/2}$ with the distribution of $\{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}^{1/2}$. So far there exist very few theoretical non-asymptotic results about bootstrap validity. Important contributions are given in the works by [Chernozhukov et al. \(2013\)](#) and [Arlot et al. \(2010\)](#). Finite sample methods for study of the bootstrap validity are essentially different from the asymptotic ones which are mainly based on weak convergence arguments. The main steps of our theoretical study are illustrated by the following scheme:

$$\begin{array}{ccccccc} & & & \text{sq-Wilks} & & \text{Gauss.} & \\ & & & \text{theorem} & & \text{approx.} & \\ \mathbf{Y}\text{-world:} & \sqrt{2L(\tilde{\theta}) - 2L(\theta^*)} & \approx & \|\xi\| & \stackrel{w}{\approx} & \|\bar{\xi}\| & \\ & & & & & \Downarrow w & \text{Gauss.} \\ & & & & & & \text{compar.} & (1.6) \\ \text{Bootstrap} & \sqrt{2L^\circ(\tilde{\theta}^\circ) - 2L^\circ(\tilde{\theta})} & \approx & \|\xi^\circ\| & \stackrel{w}{\approx} & \|\bar{\xi}^\circ\|, & \\ \text{world:} & & & & & & \end{array}$$

where $\xi^\circ \stackrel{\text{def}}{=} D_0^{-1} \nabla_\theta [L^\circ(\theta^*) - \mathbb{E} \{L^\circ(\theta^*) | \mathbf{Y}\}]$; compare with the definition (1.5) of the vector ξ in the \mathbf{Y} -world. The vectors $\bar{\xi}$ and $\bar{\xi}^\circ$ are zero mean Gaussian and they mimic the covariance structure of the vectors ξ and ξ° : $\bar{\xi} \sim \mathcal{N}(0, \text{Var} \xi)$, $\bar{\xi}^\circ \sim \mathcal{N}(0, \text{Var} \{\xi^\circ | \mathbf{Y}\})$.

The upper line of the scheme corresponds to the \mathbf{Y} -world, the lower line - to the bootstrap world. In both lines we apply two steps for approximating the corresponding likelihood ratio statistics. The first approximating step is the non-asymptotic square-root Wilks theorem: the bound (1.5) for the \mathbf{Y} case and a similar statement for the bootstrap case, which is obtained in Theorem A.4, Section A.2.

The next step is called *Gaussian approximation* (GAR) which means that the distribution of the Euclidean norm $\|\xi\|$ of a centered random vector ξ is close to the distribution of the similar norm of a Gaussian vector $\|\bar{\xi}\|$ with the same covariance matrix as ξ . A similar statement holds for the vector ξ° . Thus, the initial problem of comparing the distributions of the likelihood ratio statistics is reduced to the comparison of the distributions of the Euclidean norms of two centered normal vectors $\bar{\xi}$ and $\bar{\xi}^\circ$ (Gaussian comparison). This last step links their distributions and encloses the approximating scheme. The Gaussian comparison step is done by computing the Kullback-Leibler divergence between two multivariate Gaussian distributions (i.e. by comparison of the covariance matrices of $\nabla_\theta L(\theta^*)$ and $\nabla_\theta L^\circ(\theta^*)$) and applying Pinsker's inequality (Lemma 5.7). At this point we need to introduce the "small modeling bias" condition (**SmB**) from Section 4.2. It is formulated in terms of the following nonnegative-definite $p \times p$ symmetric matrices:

$$H_0^2 \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} [\nabla_\theta \ell_i(\theta^*) \nabla_\theta \ell_i(\theta^*)^\top], \quad (1.7)$$

$$B_0^2 \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} [\nabla_\theta \ell_i(\theta^*)] \mathbb{E} [\nabla_\theta \ell_i(\theta^*)]^\top, \quad (1.8)$$

so that $\text{Var} \{\nabla_\theta L(\theta^*)\} = H_0^2 - B_0^2$. If the parametric assumption (1.1) is true or if the data \mathbf{Y} are i.i.d., then it holds $\mathbb{E} [\nabla_\theta \ell_i(\theta^*)] \equiv 0$ and $B_0^2 = 0$. The (**SmB**) condition roughly means that the bias term B_0^2 is small relative to H_0^2 . Below we show that the Kullback-Leibler distance between the distributions of two Gaussian vectors $\bar{\xi}$ and $\bar{\xi}^\circ$ is bounded by $p \|H_0^{-1} B_0^2 H_0^{-1}\|^2 / 2$. The (**SmB**) condition precisely means that this quantity is small. We consider two situations: when the condition (**SmB**) is fulfilled and when it is not. Theorem 2.1 in Section 2 deals with the first case, it provides the cumulative error term for the coverage probability of the confidence set (1.3), taken at the $(1-\alpha)$ -quantile computed with the multiplier bootstrap procedure. The proof of this result (see Section A.3) summarizes the steps of scheme (1.6). The biggest term in the full error is induced by Gaussian approximation and requires the ratio p^3/n to be small. In the case of a "large modelling bias" i.e., when (**SmB**) does not hold, the multiplier bootstrap procedure continues to apply. It turns out that the bootstrap quantile increases with the growing modelling bias, hence, the confidence set based on it remains valid, however, it may become conservative. This result is given in Theorem 2.4 of Section 2. The problems of Gaussian approximation and comparison for the Euclidean norm are considered in Sections 5.2 and 5.4 in general terms independently of the statistical setting of the paper, and might be interesting by themselves. Section 5.4 presents also an anti-concentration inequality for the Euclidean norm of a Gaussian vector. This inequality shows how the deviation probability changes with a threshold. The general results on GAR are summarized in Theorem 5.1 and

restated in Proposition A.9 for the setting of scheme (1.6). These results are also non-asymptotic with explicit errors and apply under the condition that the ratio p^3/n to be small.

In Theorem 2.3 we consider the case of a scalar parameter $p = 1$ with an improved error term. Furthermore in Section 2.1 we propose a modified version of a quantile function based on a smoothed probability distribution. In this case the obtained error term is also better, than in the general result.

Notations: $\|\cdot\|$ denotes Euclidean norm for vectors and spectral norm for matrices; C is a generic constant. The value $x > 0$ describes our tolerance level: all the results will be valid on a random set of probability $(1 - Ce^{-x})$ for an explicit constant C . Everywhere we give explicit error bounds and show how they depend on p and n for the case of the i.i.d. observations Y_1, \dots, Y_n and $x \leq C \log n$. More details on it are given in Section 4.3.

The paper is organized as follows: the main results are stated in Section 2, their proofs are given in Sections A.3, A.4 and A.5; Section 3 contains numerical results for misspecified constant and logistic regressions. In Section 4 we give all the necessary conditions and provide an information about dependency of the involved terms on n and p . Section 5 collects some useful statements on Gaussian approximation and Gaussian comparison.

2 Multiplier bootstrap procedure

Let $\ell_i(\boldsymbol{\theta})$ denote the parametric log-density of the i -th observation:

$$\ell_i(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \log \left\{ \frac{d\mathbb{P}_{\boldsymbol{\theta}}}{d\mu_0}(Y_i) \right\},$$

then $L(\boldsymbol{\theta}) = \sum_{i=1}^n \ell_i(\boldsymbol{\theta})$. Consider i.i.d. scalar random variables u_i independent of \mathbf{Y} with continuous c.d.f., $\mathbb{E}u_i = 1$, $\text{Var} u_i = 1$, $\mathbb{E} \exp(u_i) < \infty$ for all $i = 1, \dots, n$. Multiply the summands of the likelihood function $L(\boldsymbol{\theta})$ with the new random variables:

$$L^\circ(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{i=1}^n \ell_i(\boldsymbol{\theta}) u_i,$$

then it holds $\mathbb{E}^\circ L^\circ(\boldsymbol{\theta}) = L(\boldsymbol{\theta})$, where \mathbb{E}° stands for the conditional expectation given \mathbf{Y} :

$$\mathbb{E}^\circ(\cdot) \stackrel{\text{def}}{=} \mathbb{E}(\cdot | \mathbf{Y}), \quad \mathbb{P}^\circ(\cdot) \stackrel{\text{def}}{=} \mathbb{P}(\cdot | \mathbf{Y}).$$

Therefore, the quasi MLE for the \mathbf{Y} -world is a target parameter for the bootstrap world:

$$\text{argmax}_{\boldsymbol{\theta} \in \Theta} \mathbb{E}^\circ L^\circ(\boldsymbol{\theta}) = \text{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}) = \tilde{\boldsymbol{\theta}}.$$

The corresponding quasi MLE under the conditional measure \mathbb{P}° is defined

$$\tilde{\boldsymbol{\theta}}^\circ \stackrel{\text{def}}{=} \text{argmax}_{\boldsymbol{\theta} \in \Theta} L^\circ(\boldsymbol{\theta}).$$

The likelihood ratio statistic in the bootstrap world is equal to $L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}})$, where all the elements: the function $L^\circ(\boldsymbol{\theta})$ and the arguments $\tilde{\boldsymbol{\theta}}^\circ$, $\tilde{\boldsymbol{\theta}}$ are known and available for computation.

Let $1 - \alpha \in (0, 1)$ be a fixed desirable confidence level of the set $\mathcal{E}(\mathfrak{z})$:

$$\mathbb{P}(\boldsymbol{\theta}^* \in \mathcal{E}(\mathfrak{z})) \geq 1 - \alpha. \quad (2.1)$$

Here the parameter $\mathfrak{z} \geq 0$ determines the size of the confidence set. Usually we are interested in finding a set of the smallest possible diameter satisfying this property. This leads to the problem of fixing the minimal possible value of \mathfrak{z} such that (2.1) is fulfilled. Let \mathfrak{z}_α denote the upper α -quantile of the square-root likelihood ratio statistic:

$$\mathfrak{z}_\alpha \stackrel{\text{def}}{=} \min \left\{ \mathfrak{z} \geq 0 : \mathbb{P} \left(L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) > \mathfrak{z}^2/2 \right) \leq \alpha \right\}. \quad (2.2)$$

This means, that \mathfrak{z}_α is exactly the value of our interest. Estimation of \mathfrak{z}_α leads to recovering of the distribution of $L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)$. The multiplier bootstrap procedure consists of generating a large number of independent samples $\{u_1, \dots, u_n\}$ and computing from them the empirical distribution function of $L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}})$. By this procedure we can estimate $\mathfrak{z}_\alpha^\circ$, the upper α -quantile of $\sqrt{2L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - 2L^\circ(\tilde{\boldsymbol{\theta}})}$:

$$\mathfrak{z}_\alpha^\circ \stackrel{\text{def}}{=} \min \left\{ \mathfrak{z} \geq 0 : \mathbb{P}^\circ \left(L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}}) > \mathfrak{z}^2/2 \right) = \alpha \right\}. \quad (2.3)$$

Theorem 2.1 (Validity of the bootstrap under a small modeling bias). *Let the conditions of Section 4 be fulfilled. It holds with probability $\geq 1 - 12e^{-x}$ for $\mathfrak{z}_\alpha^\circ \geq \max\{2, \sqrt{p}\} + \mathbb{C}(p + x)/\sqrt{n}$:*

$$\left| \mathbb{P} \left(L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) > (\mathfrak{z}_\alpha^\circ)^2/2 \right) - \alpha \right| \leq \Delta_{\text{full}}, \quad (2.4)$$

where $\Delta_{\text{full}} \leq \mathbb{C}\{(p + x)^3/n\}^{1/8}$ in the case 4.3. An explicit definition of the error term Δ_{full} is given in the proof (see (A.26), (A.27) in Section A.3).

The term Δ_{full} can be viewed as a sum of the error terms corresponding to each step in the scheme (1.6). The largest error term equal to $\mathbb{C}\{(p + x)^3/n\}^{1/8}$ is induced by GAR. This error rate is not always optimal for GAR, e.g. in the case of $p = 1$ or for the i.i.d. observations (see Remark 5.2). In Theorems 2.3 and 2.5 the rate is $\mathbb{C}\{(p + x)^3/n\}^{1/2}$.

In view of definition (1.3) of the likelihood-based confidence set Theorem 2.1 implies the following

Corollary 2.2 (Coverage probability error). *Under the conditions of Theorem 2.1 it holds:*

$$|\mathbb{P} \{ \boldsymbol{\theta}^* \in \mathcal{E}(\mathfrak{z}_\alpha^\circ) \} - (1 - \alpha)| \leq \Delta_{\text{full}}.$$

REMARK 2.1 (Critical dimension). The error term Δ_{full} depends on the ratio p^3/n . The bootstrap validity can be only stated if this ratio is small. The obtained error bound seems to be mainly of theoretical interest, because the condition “ $(p^3/n)^{1/8}$ is small” may require a huge sample. However, it provides some qualitative information about the bootstrap behavior as the parameter dimension grows. Our numerical results show that the accuracy of bootstrap approximation is very reasonable in a variety of examples.

In the following theorem we consider the case of a scalar parameter $p = 1$. The obtained error rate is $1/\sqrt{n}$, which is sharper, than $1/n^{1/8}$. Instead of the GAR for the Euclidean norm from Section 5 we use here Berry-Esseen theorem (see also Remark 5.2).

Theorem 2.3 (The case of $p = 1$, using Berry-Esseen theorem). *Let the conditions of Section 4 be fulfilled. It holds with probability $\geq 1 - 12e^{-x}$ for $\mathfrak{z}_\alpha^\circ \geq 1 + C(1 + \mathfrak{x})/\sqrt{n}$:*

$$\left| \mathbb{P} \left(L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) > (\mathfrak{z}_\alpha^\circ)^2/2 \right) - \alpha \right| \leq \Delta_{B.E., \text{full}}, \quad (2.5)$$

where $\Delta_{B.E., \text{full}} \leq C(1 + \mathfrak{x})/\sqrt{n}$ in the case 4.3. An explicit definition of $\Delta_{B.E., \text{full}}$ is given in (A.28) in Section A.3.

REMARK 2.2 (Bootstrap validity and weak convergence). The standard way of proving the bootstrap validity is based on weak convergence arguments; see e.g. Mammen (1992), van der Vaart and Wellner (1996), Janssen and Pauls (2003), Chatterjee and Bose (2005). If the statistic $L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)$ weakly converges to a χ^2 -type distribution, one can state an asymptotic version of the results (2.4), (2.5). Our way is based on a kind of non-asymptotic Gaussian approximation and Gaussian comparison for random vectors and allows to get explicit error terms.

REMARK 2.3 (Use of Edgeworth expansion). The classical results on confidence sets for the mean of population states the accuracy of order $1/n$ based on the second order Edgeworth expansion Hall (1992). Unfortunately, if the considered parametric model can be misspecified, even the leading term is affected by the modeling bias, and the use of Edgeworth expansion cannot help in improving the bootstrap accuracy.

REMARK 2.4 (Choice of the weights). In our construction, similarly to Chatterjee and Bose (2005), we apply a general distribution of the bootstrap weights u_i under some moment conditions. One particularly can use Gaussian multipliers as suggested by Chernozhukov et al. (2013). This leads to the exact Gaussian distribution of the vectors $\boldsymbol{\xi}^\circ$ and is helpful to avoid one step of Gaussian approximation for these vectors.

Now we discuss the impact of modeling bias, which comes from a possible misspecification of the parametric model. As explained by the approximating diagram (1.6), the distance between the distributions of the likelihood ratio statistics can be characterized via the distance between two multivariate normal distributions. To state the result let us recall the definition of the full Fisher information matrix $D_0^2 \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E} L(\boldsymbol{\theta}^*)$. For the matrices H_0^2 and B_0^2 , given in (1.7) and (1.8), it holds $H_0^2 > B_0^2 \geq 0$. If the parametric assumption (1.1) is true or in the case of an i.i.d. sample \mathbf{Y} , $B_0^2 = 0$. Under the condition (SmB) $\|H_0^{-1} B_0^2 H_0^{-1}\|$ enters linearly in the error term Δ_{full} in Theorem 2.1.

The first statement in Theorem 2.4 below says that the effective coverage probability of the confidence set based on the multiplier bootstrap is *larger* than the nominal coverage probability up to the error term $\Delta_{b, \text{full}} \leq C\{(p + \mathfrak{x})^3/n\}^{1/8}$. The inequalities in the second part of Theorem 2.4 prove the *conservativeness of the bootstrap quantiles*: the quantity $\sqrt{\text{tr}\{D_0^{-1} H_0^2 D_0^{-1}\}} - \sqrt{\text{tr}\{D_0^{-1} (H_0^2 - B_0^2) D_0^{-1}\}} \geq 0$ increases with the growing modeling bias.

Theorem 2.4 (Performance of the bootstrap for a large modeling bias). *Under the conditions of Section 4 except for (SmB) it holds with probability $\geq 1 - 14e^{-x}$ for $\mathfrak{z}, \mathfrak{z}_\alpha^\circ \geq \max\{2, \sqrt{p}\} + C(p + \mathfrak{x})/\sqrt{n}$*

$$\begin{aligned}
1. \quad & \mathbb{P} \left(L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) > \mathfrak{z}^2/2 \right) \leq \mathbb{P}^\circ \left(L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}}) > \mathfrak{z}^2/2 \right) + \Delta_{\text{b, full}}. \\
2. \quad & \mathfrak{z}_\alpha^\circ \geq \mathfrak{z}_{(\alpha + \Delta_{\text{b, full}})} \\
& \quad + \sqrt{\text{tr}\{D_0^{-1}H_0^2D_0^{-1}\}} - \sqrt{\text{tr}\{D_0^{-1}(H_0^2 - B_0^2)D_0^{-1}\}} - \Delta_{\text{qf},1}, \\
& \mathfrak{z}_\alpha^\circ \leq \mathfrak{z}_{(\alpha - \Delta_{\text{b, full}})} \\
& \quad + \sqrt{\text{tr}\{D_0^{-1}H_0^2D_0^{-1}\}} - \sqrt{\text{tr}\{D_0^{-1}(H_0^2 - B_0^2)D_0^{-1}\}} + \Delta_{\text{qf},2}.
\end{aligned}$$

The term $\Delta_{\text{b, full}} \leq C\{(p + \mathbf{x})^3/n\}^{1/8}$ is given in (A.30) in Section A.4. The positive values $\Delta_{\text{qf},1}, \Delta_{\text{qf},2}$ are given in (A.34), (A.33) in Section A.4, they are bounded from above with $(\alpha^2 + \alpha_B^2)(\sqrt{8\mathbf{x}p} + 6\mathbf{x})$ for the constants $\alpha^2, \alpha_B^2 > 0$ from conditions $(\mathcal{I}), (\mathcal{I}_B)$.

REMARK 2.5. There exists some literature on robust (and heteroscedasticity robust) bootstrap procedures; see e.g. Mammen (1993), Aerts and Claeskens (2001), Kline and Santos (2012). However, up to our knowledge there are no robust bootstrap procedures for the likelihood ratio statistic, most of the results compare the distribution of the estimator obtained from estimating equations, or Wald / score test statistics with their bootstrap counterparts in the i.i.d. setup. In our context this would correspond to the noise misspecification in the log-likelihood function and it is addressed automatically by the multiplier bootstrap. Our notion of modeling bias includes the situation when the target value $\boldsymbol{\theta}^*$ from (1.4) only defines a projection (the best parametric fit) of the data distribution. In particular, the quantities $\mathbb{E}\nabla_{\boldsymbol{\theta}}\ell_i(\boldsymbol{\theta}^*)$ for different i do not necessarily vanish yielding a significant modeling bias. Similar notion of misspecification is used in the literature on Generalized Method of Moments; see e.g. Hall (2005). Chapter 5 therein considers the hypothesis testing problem with two kinds of misspecification: local and non-local, which would correspond to our small and large modeling bias cases.

An interesting message of Theorem 2.4 is that the multiplier bootstrap procedure ensures a prescribed coverage level for this target value $\boldsymbol{\theta}^*$ even without small modeling bias restriction, however, in this case the method is somehow conservative because the modeling bias is transferred into the additional variance in the bootstrap world. The numerical experiments in Section 3 agree with this result.

2.1 Smoothed version of a quantile function

This section briefly discusses the use of a smoothed quantile function. The $(1 - \alpha)$ -quantile of $\sqrt{2L(\tilde{\boldsymbol{\theta}}) - 2L(\boldsymbol{\theta}^*)}$ is defined as

$$\begin{aligned}
\mathfrak{z}_\alpha &\stackrel{\text{def}}{=} \min \left\{ \mathfrak{z} \geq 0 : \mathbb{P} \left(L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) > \mathfrak{z}^2/2 \right) \leq \alpha \right\} \\
&= \min \left\{ \mathfrak{z} \geq 0 : \mathbb{E} \mathbb{1} \left\{ L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) > \mathfrak{z}^2/2 \right\} \leq \alpha \right\}.
\end{aligned}$$

Introduce for $x \geq 0$ and $z, \Delta > 0$ the following function

$$g_{\Delta}(x, z) \stackrel{\text{def}}{=} g\left(\frac{1}{2\Delta z}(x^2 - z^2)\right), \quad (2.6)$$

where $g(\cdot) \in C^2(\mathbb{R})$ is a non-negative function, which grows monotonously from 0 to 1, $g(x) = 0$ for $x \leq 0$ and $g(x) = 1$ for $x \geq 1$, therefore:

$$\mathbb{I}\{x \geq 1\} \leq g(x) \leq \mathbb{I}\{x \geq 0\} \leq g(x + 1).$$

An example of such function is given in (5.9). In (5.10) it is shown

$$\mathbb{I}\{x - z \geq \Delta\} \leq g_{\Delta}(x, z) \leq \mathbb{I}\{x - z \geq 0\} \leq g_{\Delta}(x, z + \Delta).$$

This approximation is used in the proofs of Theorems 2.1 and 2.4 in the part of Gaussian approximation of Euclidean norm of a sum of independent vectors (see Section 5.2) yielding the error rate $(p^3/n)^{1/8}$ in the final bound (Theorems 2.1, 5.1). The next result shows that the use of a smoothed quantile function helps to improve the accuracy of bootstrap approximation: it becomes $(p^3/n)^{1/2}$ instead of $(p^3/n)^{1/8}$. The reason is that we do not need to account for the error induced by a smooth approximation of the indicator function.

Theorem 2.5 (Validity of the bootstrap in the smoothed case under (SmB) condition). *Let the conditions of Section 4 be fulfilled. It holds with probability $\geq 1 - 12e^{-x}$ for $\mathfrak{z} \geq \max\{2, \sqrt{p}\} + C(p + x)/\sqrt{n}$ and $\Delta \in (0, 0.22]$:*

$$\left| \mathbb{E} g_{\Delta}\left(\sqrt{2L(\tilde{\theta}) - 2L(\theta^*)}, \mathfrak{z}\right) - \mathbb{E}^{\circ} g_{\Delta}\left(\sqrt{2L^{\circ}(\tilde{\theta}^{\circ}) - 2L^{\circ}(\tilde{\theta})}, \mathfrak{z}\right) \right| \leq \Delta_{\text{sm}},$$

where $\Delta_{\text{sm}} \leq C\{(p + x)^3/n\}^{1/2} \Delta^{-3}$ in the case 4.3. An explicit definition of Δ_{sm} is given in (A.38), (A.39) in Section A.5.

The modified bootstrap quantile function reads as

$$\mathfrak{z}_{\Delta, \alpha}^{\circ} \stackrel{\text{def}}{=} \min \left\{ \mathfrak{z} \geq 0 : \mathbb{E}^{\circ} g_{\Delta}\left(\sqrt{2L^{\circ}(\tilde{\theta}^{\circ}) - 2L^{\circ}(\tilde{\theta})}, \mathfrak{z}\right) = \alpha \right\}.$$

3 Numerical results

This section illustrates the performance of the multiplier bootstrap for some artificial examples. We especially aim to address the issues of noise misspecification and of increasing modeling bias. In all the experiments we took 10^4 data samples for estimation of empirical c.d.f. of $\sqrt{2L(\tilde{\theta}) - 2L(\theta^*)}$, 10^4 $\{u_1, \dots, u_n\}$ samples and 10^4 data samples for the estimation of the quantiles of $\sqrt{2L^{\circ}(\tilde{\theta}^{\circ}) - 2L^{\circ}(\tilde{\theta})}$. All sample sizes are $n = 50$. It should be mentioned that the obtained results are nicely consistent with the theoretical statements.

3.1 Computational error

Here we check numerically, how well the multiplier procedure works in the case of the correct model. Let the i.i.d. data follow the distribution $Y_i \sim \mathcal{N}(2, 1)$, $i = 1, \dots, n$. The true likelihood function is $L(\theta) = -\sum_{i=1}^n (Y_i - \theta)^2/2$.

Table 1 shows the effective coverage probabilities of the quantiles estimated using the multiplier bootstrap. The second line contains the range of the nominal confidence levels: 0.99, . . . , 0.75. The first left column describes the distribution of the bootstrap weights: $\mathcal{N}(1, 1)$ or $exp(1)$. The 3-d and the 4-th lines show the frequency of the event: “the real likelihood ratio \leq the quantile of the bootstrap likelihood ratio”.

Table 1: Coverage probabilities for the correct i.i.d. model

	Confidence levels					
$\mathcal{L}(u_i)$	0.99	0.95	0.90	0.85	0.80	0.75
$exp(1)$	0.99	0.94	0.89	0.83	0.78	0.73
$\mathcal{N}(1, 1)$	0.99	0.95	0.89	0.84	0.80	0.75

3.2 Constant regression with misspecified heteroscedastic errors

Here we show on a constant regression model that the quality of the confidence sets obtained by the multiplier bootstrap procedure is not significantly deteriorated by misspecified heteroscedastic errors. Let the data be defined as $Y_i = 2 + \sigma_i \varepsilon_i$, $i = 1, \dots, n$. The i.i.d. random variables $\varepsilon_i \sim Lap(0, 2^{-1/2})$ are s.t. $\mathbb{E}(\varepsilon_i) = 0$, $\text{Var}(\varepsilon_i) = 1$. The coefficients σ_i are deterministic: $\sigma_i \stackrel{\text{def}}{=} 0.5 \{4 - i \pmod{4}\}$. The quasi-likelihood function is the same as in the previous experiment: $L(\theta) = -\sum_{i=1}^n (Y_i - \theta)^2/2$, i.e. it is misspecified, since it corresponds to the i.i.d. standard normal distribution. Table 2 describes the 2-nd experiment’s results similarly to the Table 1.

Table 2: Coverage probabilities for the misspecified heteroscedastic noise

	Confidence levels					
$\mathcal{L}(u_i)$	0.99	0.95	0.90	0.85	0.80	0.75
$exp(1)$	0.98	0.93	0.87	0.82	0.77	0.72
$\mathcal{N}(1, 1)$	0.98	0.94	0.88	0.83	0.78	0.73

3.3 Biased constant regression with misspecified errors

In the third experiment we consider biased regression with misspecified i.i.d. errors:

$$Y_i = \beta \sin(X_i) + \varepsilon_i, \quad \varepsilon_i \sim Lap(0, 2^{-1/2}), \text{ i.i.d.},$$

$$X_i \text{ are equidistant in } [0, 2\pi].$$

Taking the likelihood function $L(\theta) = -\sum_{i=1}^n (Y_i - \theta)^2/2$ yields $\theta^* = 0$. Therefore, the larger is the deterministic amplitude $\beta > 0$, the bigger is bias of the mean constant regression. We consider two cases: $\beta = 0.25$ with fulfilled **(SmB)** condition and $\beta = 1.25$ when **(SmB)** does not hold. Table 3 shows that for the large bias quantiles yielded by the multiplier bootstrap are conservative. This conservative property of the multiplier bootstrap quantiles is

Table 3: Coverage probabilities for the misspecified biased regression

$\mathcal{L}(u_i)$	β	Confidence levels					
		0.99	0.95	0.90	0.85	0.80	0.75
$\mathcal{N}(1, 1)$	0.25	0.98	0.94	0.89	0.84	0.79	0.74
	1.25	1.0	0.99	0.97	0.94	0.91	0.87

also illustrated with the graphs in Figure 3.1. They show the empirical distribution functions of the likelihood ratio statistics $L(\tilde{\theta}) - L(\theta^*)$ and $L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})$ for $\beta = 0.25$ and $\beta = 1.25$. On the right graph for $\beta = 1.25$ the empirical distribution functions for the bootstrap case are smaller than the one for the \mathbf{Y} case. It means that for the large bias the bootstrap quantiles are bigger than the \mathbf{Y} quantiles, which increases the diameter of the confidence set based on the bootstrap quantiles. This confidence set remains valid, since it still contains the true parameter with a given confidence level.

Figure 3.2 shows the growth of the difference between the quantiles of $L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})$ and $L(\tilde{\theta}) - L(\theta^*)$ with increasing β for the range of the confidence levels: 0.75, 0.8, ..., 0.99.

3.4 Logistic regression with bias

In this example we consider logistic regression. Let the data come from the following distribution:

$$Y_i \sim Bernoulli(\beta X_i), \quad X_i \text{ are equidistant in } [0, 2], \quad \beta \in (0, 1/2].$$

Consider the likelihood function corresponding to the i.i.d. observations:

$$L(\theta) = \sum_{i=1}^n \{Y_i \theta - \log(1 + e^\theta)\}.$$

By definition (1.4) $\theta^* = \log\{\beta/(1 - \beta)\}$, bigger values of β induce larger modeling bias. The graphs below demonstrate the conservativeness of bootstrap quantiles. Here we consider two cases: $\beta = 0.1$ and $\beta = 0.5$. Similarly to the Example 3.3 in the case of the bigger β on

Figure 3.1: Empirical distribution functions of the likelihood ratios

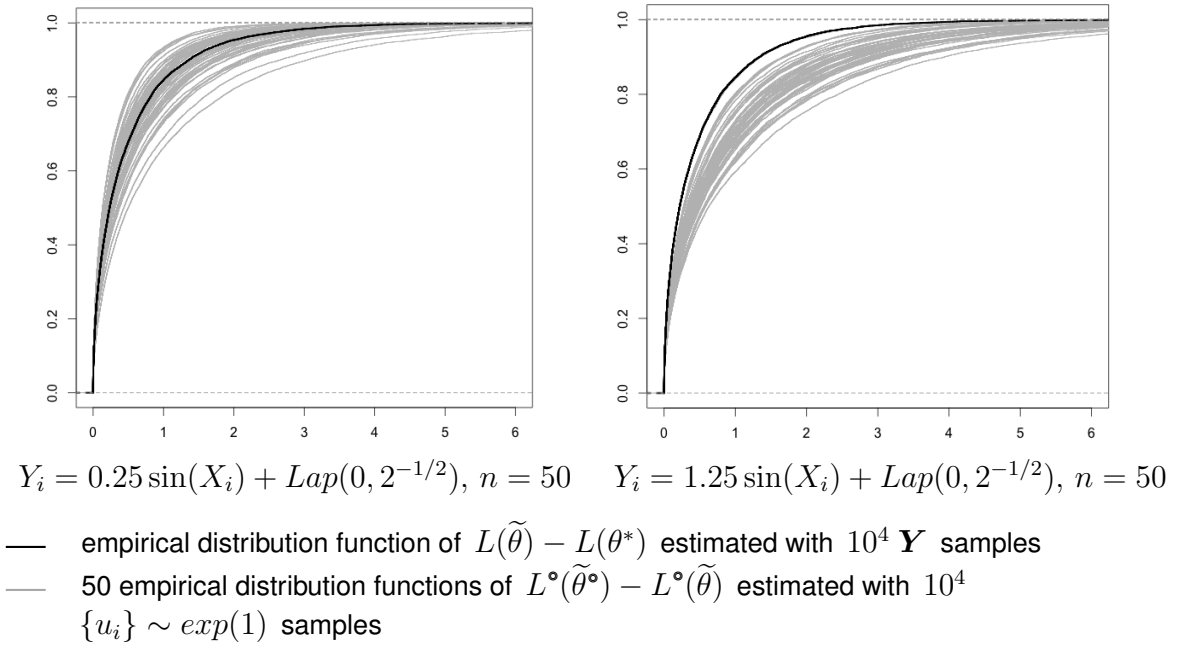
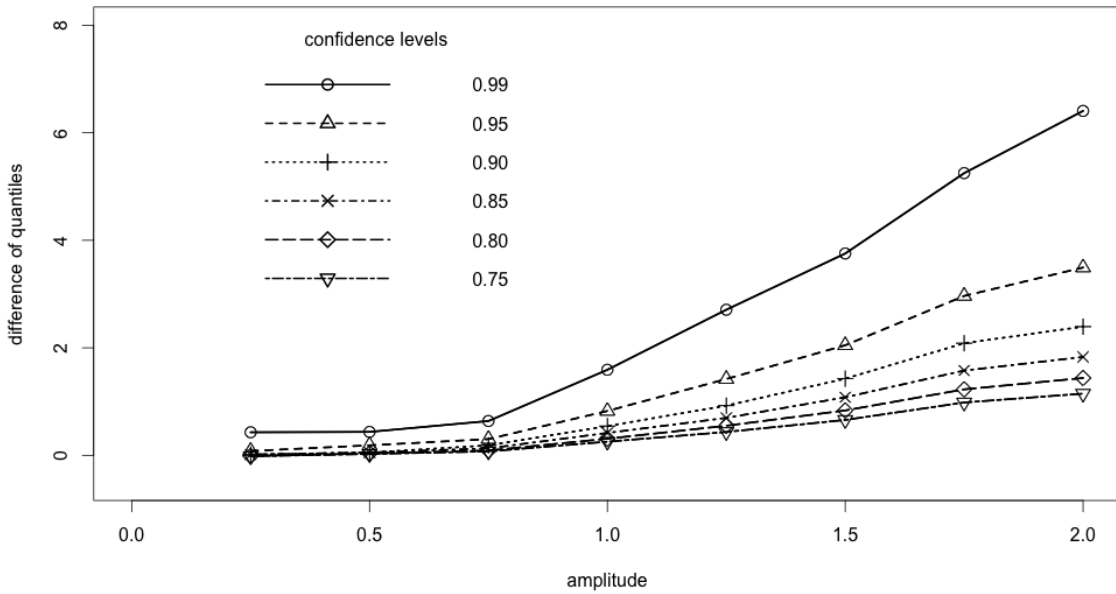
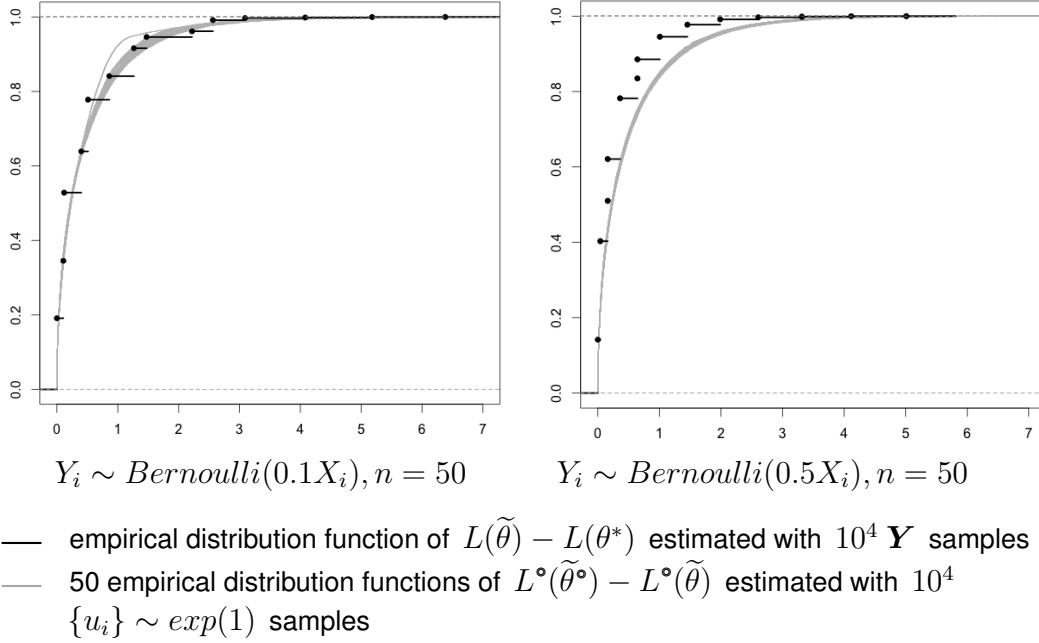


Figure 3.2: The difference (“Bootstrap quantile” – “ \mathbf{Y} -quantile”) growing with modeling bias



the right graph in Figure 3.3 the empirical distribution functions of $L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})$ are smaller than the one for $L(\tilde{\theta}) - L(\theta^*)$.

Figure 3.3:



4 Conditions

Here we state the conditions necessary for the main results. The conditions in Section 4.1 come from the general finite sample theory by Spokoiny (2012a), they are required for the results of Sections A.1 and A.2. Spokoiny (2012a) considers the examples of i.i.d. setup, generalized linear model and linear median regression providing a check of conditions from Section 4.1. The conditions in Section 4.2 are necessary to prove the results on multiplier bootstrap from Section 2.

4.1 Basic conditions

Introduce the stochastic part of the likelihood process: $\zeta(\boldsymbol{\theta}) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta})$, and its marginal summand: $\zeta_i(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \ell_i(\boldsymbol{\theta}) - \mathbb{E}\ell_i(\boldsymbol{\theta})$.

(ED₀) There exist a positive-definite symmetric matrix V_0^2 and constants $g > 0, \nu_0 \geq 1$ such that $\text{Var} \{ \nabla_{\boldsymbol{\theta}} \zeta(\boldsymbol{\theta}^*) \} \leq V_0^2$ and

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla_{\boldsymbol{\theta}} \zeta(\boldsymbol{\theta}^*)}{\|V_0 \boldsymbol{\gamma}\|} \right\} \leq \nu_0^2 \lambda^2 / 2, \quad |\lambda| \leq g.$$

(ED₂) There exists a constant $\omega \geq 0$ and for each $\mathfrak{r} > 0$ a constant $g_2(\mathfrak{r})$ such that it

holds for all $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ and for $j = 1, 2$

$$\sup_{\substack{\gamma_j \in \mathbb{R}^p \\ \|\gamma_j\| \leq 1}} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \gamma_1^\top D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \zeta(\boldsymbol{\theta}) D_0^{-1} \gamma_2 \right\} \leq \nu_0^2 \lambda^2 / 2, \quad |\lambda| \leq \mathfrak{g}_2(\mathbf{r}).$$

(\mathcal{L}_0) For each $\mathbf{r} > 0$ there exists a constant $\delta(\mathbf{r}) \geq 0$ such that for $\mathbf{r} \leq \mathbf{r}_0$ (\mathbf{r}_0 comes from condition (A.1) of Theorem A.1 in Section A.1) $\delta(\mathbf{r}) \leq 1/2$, and for all $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ it holds

$$\|D_0^{-1} D^2(\boldsymbol{\theta}) D_0^{-1} - \mathbf{I}_p\| \leq \delta(\mathbf{r}),$$

where $D^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E} L(\boldsymbol{\theta})$, $\Theta_0(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} : \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}\}$.

(\mathcal{I}) There exists a constant $\mathfrak{a} > 0$ s.t. $\mathfrak{a}^2 D_0^2 \geq V_0^2$.

($\mathcal{L}_\mathbf{r}$) For each $\mathbf{r} \geq \mathbf{r}_0$ there exists a value $\mathfrak{b}(\mathbf{r}) > 0$ s.t. $\mathbf{r}\mathfrak{b}(\mathbf{r}) \rightarrow \infty$ for $\mathbf{r} \rightarrow \infty$ and $\forall \boldsymbol{\theta} : \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| = \mathbf{r}$ it holds

$$-2 \{\mathbb{E} L(\boldsymbol{\theta}) - \mathbb{E} L(\boldsymbol{\theta}^*)\} \geq \mathbf{r}^2 \mathfrak{b}(\mathbf{r}).$$

4.2 Conditions required for the bootstrap validity

(SmB) There exists a constant $\delta_{\text{smb}}^2 \in [0, 1/8]$ such that it holds for all $i = 1, \dots, n$ and the matrices H_0^2, B_0^2 defined in (1.7) and (1.8).

$$\|H_0^{-1} B_0^2 H_0^{-1}\| \leq \delta_{\text{smb}}^2 \leq \mathfrak{C} p n^{-1/2},$$

(\mathcal{ED}_{2m}) For each $\mathbf{r} > 0$, $i = 1, \dots, n$, $j = 1, 2$ and for all $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ it holds for the values $\omega \geq 0$ and $\mathfrak{g}_2(\mathbf{r})$ from the condition (\mathcal{ED}_2):

$$\sup_{\substack{\gamma_j \in \mathbb{R}^p \\ \|\gamma_j\| \leq 1}} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \gamma_1^\top D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \zeta_i(\boldsymbol{\theta}) D_0^{-1} \gamma_2 \right\} \leq \frac{\nu_0^2 \lambda^2}{2n}, \quad |\lambda| \leq \mathfrak{g}_2(\mathbf{r}),$$

(\mathcal{L}_{0m}) For each $\mathbf{r} > 0$, $i = 1, \dots, n$ and for all $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ there exists a constant $\mathfrak{C}_m(\mathbf{r}) \geq 0$ such that

$$\|D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \mathbb{E} \ell_i(\boldsymbol{\theta}) D_0^{-1}\| \leq \mathfrak{C}_m(\mathbf{r}) n^{-1}.$$

(\mathcal{L}_{3m}) For all $\boldsymbol{\theta} \in \Theta$ and $i = 1, \dots, n$ it holds $\|D_0^{-1} \nabla_{\boldsymbol{\theta}}^3 \mathbb{E} \ell_i(\boldsymbol{\theta}) D_0^{-1}\| \leq \mathfrak{C}$.

(\mathcal{I}_B) There exists a constant $\mathfrak{a}_B^2 > 0$ s.t. $\mathfrak{a}_B^2 D_0^2 \geq B_0^2$.

(SD₁) There exists a constant $0 \leq \delta_v \leq Cp/n$. such that it holds for all $i = 1, \dots, n$ with exponentially high probability

$$\|H_0^{-1} \{ \nabla_{\theta} \ell_i(\theta^*) \nabla_{\theta} \ell_i(\theta^*)^\top - \mathbb{E} [\nabla_{\theta} \ell_i(\theta^*) \nabla_{\theta} \ell_i(\theta^*)^\top] \} H_0^{-1}\| \leq \delta_v^2.$$

(Eb) The i.i.d. bootstrap weights u_i have continuous c.d.f., and it holds for all $i = 1, \dots, n$: $\mathbb{E}^\circ u_i = 1$, $\text{Var}^\circ u_i = 1$,

$$\log \mathbb{E}^\circ \exp \{ \lambda(u_i - 1) \} \leq \nu_0^2 \lambda^2 / 2, \quad |\lambda| \leq g.$$

4.3 Dependence of the involved terms on the sample size and parameter dimension

Here we consider the case of the i.i.d. observations Y_1, \dots, Y_n and $\mathbf{x} = C \log n$ in order to specify the dependence of the non-asymptotic bounds on n and p . Example 5.1 in Spokoiny (2012a) demonstrates that in this situation $g = C\sqrt{n}$ and $\omega = C/\sqrt{n}$. then $\mathfrak{Z}(\mathbf{x}) = C\sqrt{p + \mathbf{x}}$ for some constant $C \geq 1.85$, for the function $\mathfrak{Z}(\mathbf{x})$ given in (A.3) in Section A.1. Similarly it can be checked that $g_2(\mathbf{r})$ from condition (ED₂) is proportional to \sqrt{n} : due to independency of the observations

$$\begin{aligned} & \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \gamma_1^\top D_0^{-1} \nabla_{\theta}^2 \zeta(\theta) D_0^{-1} \gamma_2 \right\} \\ &= \sum_{i=1}^n \log \mathbb{E} \exp \left\{ \frac{\lambda}{\sqrt{n}} \frac{1}{\omega/\sqrt{n}} \gamma_1^\top d_0^{-1} \nabla_{\theta}^2 \zeta_i(\theta) d_0^{-1} \gamma_2 \right\} \\ &\leq n \frac{\lambda^2}{n} C \quad \text{for } |\lambda| \leq \bar{g}_2(\mathbf{r}) \sqrt{n}, \end{aligned}$$

where $\zeta_i(\theta) \stackrel{\text{def}}{=} \ell_i(\theta) - \mathbb{E} \ell_i(\theta)$, $d_0^2 \stackrel{\text{def}}{=} -\nabla_{\theta}^2 \mathbb{E} \ell_i(\theta^*)$ and $D_0^2 = n d_0^2$ in the i.i.d. case. Function $\bar{g}_2(\mathbf{r})$ denotes the marginal analog of $g_2(\mathbf{r})$.

Let us show, that for the value $\delta(\mathbf{r})$ from the condition (L₀) it holds $\delta(\mathbf{r}) = C\mathbf{r}/\sqrt{n}$. For some $\bar{\theta}$

$$\begin{aligned} \|D_0^{-1} D^2(\theta) D_0^{-1} - \mathbf{I}_p\| &= \|D_0^{-1}(\theta^* - \theta)^\top \nabla_{\theta}^3 \mathbb{E} L(\bar{\theta}) D_0^{-1}\| \\ &= \|D_0^{-1}(\theta^* - \theta)^\top D_0 D_0^{-1} \nabla_{\theta}^3 \mathbb{E} L(\bar{\theta}) D_0^{-1}\| \\ &\leq \mathbf{r} \|D_0^{-1}\| \|D_0^{-1} \nabla_{\theta}^3 \mathbb{E} L(\bar{\theta}) D_0^{-1}\| \leq C\mathbf{r}/\sqrt{n} \quad (\text{by condition } (\mathcal{L}_{3m})). \end{aligned}$$

Similarly $C_m(\mathbf{r}) \leq C\mathbf{r}$ in condition (L_{0m}).

If $\delta(\mathbf{r}) = C\mathbf{r}/\sqrt{n}$ is sufficiently small, then the value $b(\mathbf{r})$ from condition (L_r) can be taken as $C\{1 - \delta(\mathbf{r})\}^2$. Indeed, by (L₀) and (L_r) for $\theta : \|D_0(\theta - \theta^*)\| = \mathbf{r}$

$$-2 \{ \mathbb{E} L(\theta) - \mathbb{E} L(\theta^*) \} \geq \mathbf{r}^2 \{ 1 - \delta^2(\mathbf{r}) \}.$$

Therefore, if $\delta(\mathbf{r})$ is small, then $\mathfrak{b}(\mathbf{r}) \stackrel{\text{def}}{=} \mathbb{C}\{1 - \delta^2(\mathbf{r})\} \approx \text{const}$. Due to the obtained orders the conditions (A.1) and (A.17) of Theorems A.1 and A.6 on concentration of the MLEs $\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\theta}}^\circ$ require $r_0 \geq \mathbb{C}\sqrt{p+x}$.

5 Approximation of distributions of ℓ_2 norms of sums of independent vectors

Consider two samples ϕ_1, \dots, ϕ_n and ψ_1, \dots, ψ_n , each consists of centered independent random vectors in \mathbb{R}^p with nearly the same second moments. This section explains how one can quantify the closeness in distribution between the norms of $\boldsymbol{\phi} = \sum_i \phi_i$ and of $\boldsymbol{\psi} = \sum_i \psi_i$. Suppose that

$$\mathbb{E}\phi_i = \mathbb{E}\psi_i = 0, \quad \text{Var } \phi_i = \Sigma_i, \quad \text{Var } \psi_i = \check{\Sigma}_i, \quad i = 1, \dots, n.$$

Let also

$$\boldsymbol{\phi} \stackrel{\text{def}}{=} \sum_{i=1}^n \phi_i, \quad \boldsymbol{\psi} \stackrel{\text{def}}{=} \sum_{i=1}^n \psi_i, \quad (5.1)$$

$$\Sigma \stackrel{\text{def}}{=} \text{Var } \boldsymbol{\phi} = \sum_{i=1}^n \Sigma_i, \quad \check{\Sigma} \stackrel{\text{def}}{=} \text{Var } \boldsymbol{\psi} = \sum_{i=1}^n \check{\Sigma}_i. \quad (5.2)$$

Also introduce multivariate Gaussian vectors $\bar{\phi}_i, \bar{\psi}_i$ which are mutually independent for $i = 1, \dots, n$ and

$$\begin{aligned} \bar{\phi}_i &\sim \mathcal{N}(0, \Sigma_i), \quad \bar{\psi}_i \sim \mathcal{N}(0, \check{\Sigma}_i), \\ \bar{\boldsymbol{\phi}} \stackrel{\text{def}}{=} \sum_{i=1}^n \bar{\phi}_i &\sim \mathcal{N}(0, \Sigma), \quad \bar{\boldsymbol{\psi}} \stackrel{\text{def}}{=} \sum_{i=1}^n \bar{\psi}_i \sim \mathcal{N}(0, \check{\Sigma}). \end{aligned} \quad (5.3)$$

The bar sign for a vector stands here for a normal distribution. The following theorem gives the conditions on Σ and $\check{\Sigma}$ which ensure that $\|\boldsymbol{\phi}\|$ and $\|\boldsymbol{\psi}\|$ are close to each other in distribution. It also presents a general result on Gaussian approximation of $\|\boldsymbol{\phi}\|$ with $\|\bar{\boldsymbol{\phi}}\|$.

Introduce the following deterministic values, which are supposed to be finite:

$$\delta_n \stackrel{\text{def}}{=} \frac{1}{2} \sum_{i=1}^n \mathbb{E} (\|\phi_i\|^3 + \|\bar{\phi}_i\|^3), \quad \check{\delta}_n \stackrel{\text{def}}{=} \frac{1}{2} \sum_{i=1}^n \mathbb{E} (\|\psi_i\|^3 + \|\bar{\psi}_i\|^3). \quad (5.4)$$

Theorem 5.1. Assume for the covariance matrices defined in (5.2) that

$$\|\check{\Sigma}^{-1/2} \Sigma \check{\Sigma}^{-1/2} - \mathbf{I}_p\| \leq 1/2, \quad \text{and} \quad \text{tr}\{(\check{\Sigma}^{-1/2} \Sigma \check{\Sigma}^{-1/2} - \mathbf{I}_p)^2\} \leq \delta_\Sigma^2 \quad (5.5)$$

for some $\delta_\Sigma^2 \geq 0$. The sign $\|\cdot\|$ for matrices denotes the spectral norm. Let also for $z, \bar{z} \geq 2$ and some $\delta_z \geq 0$ $|z - \bar{z}| \leq \delta_z$, then it holds for all $0 < \Delta \leq 0.22$

$$\begin{aligned} 1.1. \quad |IP(\|\boldsymbol{\phi}\| \geq z) - IP(\|\bar{\boldsymbol{\psi}}\| \geq \bar{z})| &\leq 16\delta_n \Delta^{-3} + \frac{\Delta + \delta_z}{\bar{z}} \sqrt{p/2} + \delta_\Sigma/2 \\ &\leq 16\delta_n \Delta^{-3} + (\Delta + \delta_z)/\sqrt{2} + \delta_\Sigma/2 \\ &\quad \text{for } \bar{z} \geq \sqrt{p}, \end{aligned}$$

$$\begin{aligned}
1.2. \quad |\mathbb{P}(\|\phi\| \geq z) - \mathbb{P}(\|\psi\| \geq \bar{z})| &\leq 16\Delta^{-3}(\delta_n + \check{\delta}_n) + \frac{2\Delta + \delta_z}{\bar{z}}\sqrt{p/2} + \delta_\Sigma/2 \\
&\leq 16\Delta^{-3}(\delta_n + \check{\delta}_n) + (2\Delta + \delta_z)/\sqrt{2} + \delta_\Sigma/2 \\
&\quad \text{for } \bar{z} \geq \sqrt{p}.
\end{aligned}$$

Moreover, if $z, \bar{z} \geq \max\{2, \sqrt{p}\}$ and $\max\{\delta_n^{1/4}, \check{\delta}_n^{1/4}\} \leq 0.11$, then

$$2.1. \quad |\mathbb{P}(\|\phi\| \geq z) - \mathbb{P}(\|\bar{\psi}\| \geq \bar{z})| \leq 1.55\delta_n^{1/4} + \delta_z/\sqrt{2} + \delta_\Sigma/2,$$

$$2.2. \quad |\mathbb{P}(\|\phi\| \geq z) - \mathbb{P}(\|\psi\| \geq \bar{z})| \leq 1.55(\delta_n^{1/4} + \check{\delta}_n^{1/4}) + \delta_z/\sqrt{2} + \delta_\Sigma/2.$$

Proof of Theorem 5.1. The inequality 1.1 is based on the results of Lemmas 5.3, 5.6 and 5.7:

$$\begin{aligned}
\mathbb{P}(\|\phi\| \geq z) &\stackrel{\text{by L. 5.3}}{\leq} \mathbb{P}(\|\bar{\phi}\| \geq z - \Delta) + 16\Delta^{-3}\delta_n \\
&\stackrel{\text{by L. 5.7}}{\leq} \mathbb{P}(\|\bar{\psi}\| \geq z - \Delta) + 16\Delta^{-3}\delta_n + \delta_\Sigma/2 \\
&\stackrel{\text{by L. 5.6}}{\leq} \mathbb{P}(\|\bar{\psi}\| \geq \bar{z}) + 16\Delta^{-3}\delta_n + \delta_\Sigma/2 + (\delta_z + \Delta)\bar{z}^{-1}\sqrt{p/2}.
\end{aligned}$$

The inequality 1.2 is implied by the triangle inequality and the sum of two bounds: the bound 1.1 for $|\mathbb{P}(\|\phi\| \geq z) - \mathbb{P}(\|\bar{\psi}\| \geq \bar{z})|$ and the bound

$$|\mathbb{P}(\|\psi\| \geq \bar{z}) - \mathbb{P}(\|\bar{\psi}\| \geq \bar{z})| \leq 16\check{\delta}_n\Delta^{-3} + \Delta\bar{z}^{-1}\sqrt{p/2},$$

which also follows from 1.1 by taking $\phi := \psi$, $z := \bar{z}$. In this case $\Sigma = \check{\Sigma}$ and $\delta_\Sigma = \delta_z = 0$.

The second part of the statement follows from the the first part by balancing the error term $16\delta_n\Delta^{-3} + \Delta/\sqrt{2}$ w.r.t. Δ . \square

REMARK 5.1. The approximation error in the statements of Theorem 5.1 includes three terms, each of them is responsible for a step of derivation: Gaussian approximation, Gaussian comparison and anti-concentration. The value δ_Σ bounds the relation between covariance matrices, δ_z corresponds to the difference between quantiles. $\delta_n^{1/4}$ comes from the Gaussian approximation, under certain conditions this is the biggest term in the expressions 2.1, 2.2 (cf. the proof of Theorem 2.1).

REMARK 5.2. Here we briefly comment how our results can be compared with what is available in the literature. In the case of i.i.d. vectors ϕ_i and $\text{Var } \phi_i \equiv \mathbf{I}_p$ Bentkus (2003) obtained the rate $\mathbb{E}\|\phi_i\|^3/\sqrt{n}$ for the error of approximation $\sup_{A \in \mathcal{A}} |\mathbb{P}(\phi \in A) - \mathbb{P}(\bar{\phi} \in A)|$, where \mathcal{A} is a class of all Euclidean balls in \mathbb{R}^p . Götze (1991) showed for independent vectors ϕ_i and their standardized sum ϕ :

$$\delta_{GAR} \leq \begin{cases} C_1\sqrt{p} \sum_{i=1}^n \mathbb{E}\|\phi_i\|^3/\sqrt{n}, & p \in [2, 5], \\ C_2p \sum_{i=1}^n \mathbb{E}\|\phi_i\|^3/\sqrt{n}, & p \geq 6, \end{cases}$$

where $\delta_{GAR} \stackrel{\text{def}}{=} \sup_{B \in \mathcal{B}} |\mathbb{P}(\phi \in B) - \mathbb{P}(\bar{\phi} \in B)|$ and \mathcal{B} is a class of all measurable convex sets in \mathbb{R}^p , the constants $C_1, C_2 > 150$. [Bhattacharya and Holmes \(2010\)](#) argued that the results by [Götze \(1991\)](#) might require more thorough derivation, they obtained the rate $p^{5/2} \sum_{i=1}^n \mathbb{E} \|\phi_i\|^3$ for the previous bound (and $p^{5/2} \mathbb{E} \|\phi_1\|^3 / n^{1/2}$ in the i.i.d. case). [Chen and Fang \(2011\)](#) prove that $\delta_{GAR} \leq 115\sqrt{p} \sum_{i=1}^n \mathbb{E} \|\phi_i\|^3$ for independent vectors ϕ_i with a standardized sum. [Götze and Zaitsev \(2014\)](#) obtained the rate $\mathbb{E} \|\phi_i\|^4 / n$ for i.i.d. vectors ϕ_i with a standardized sum but only for $p \geq 5$. See also [Prokhorov and Ulyanov \(2013\)](#) for the review of the results about normal approximation of quadratic forms.

Our results ensure the error of the Gaussian approximation of order $1.55\delta_n^{1/4} \leq 1.31 \left\{ \sum_{i=1}^n \mathbb{E} (\|\phi_i\|^3 + \|\bar{\phi}_i\|^3) \right\}^{1/4}$. The technique used here is much simpler than in the previous works, and the obtained bounding terms are explicit and only use independence of the ϕ_i and ψ_i . However, for some special cases, the use of more advanced results on Gaussian approximation may lead to sharper bounds. For instance, for an i.i.d. sample, the GAR error rate $\delta_{GAR} = \sqrt{p^3/n}$ by [Bentkus \(2003\)](#) is better than ours $(p^3/n)^{1/8}$, and in the one-dimensional case Berry-Esseen's theorem would also work better (see [Section 5.1](#)). In those cases one can improve the overall error bound of the bootstrap approximation by putting δ_{GAR} in place of the sum $16\delta_n \Delta^{-3} + \Delta/\sqrt{2}$. [Section 5.3](#) comments how our results can be used to obtain the error rate $\sqrt{p^3/n}$ by using a smoothed quantile function.

5.1 The case of $p = 1$ using Berry-Esseen theorem

Let us consider how the results of [Theorem 5.1](#) can be refined in the case $p = 1$ using Berry-Esseen theorem. Introduce similarly to δ_n and $\check{\delta}_n$ from [\(5.4\)](#) the bounded values

$$\delta_{n,\text{B.E.}} \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} |\phi_i|^3, \quad \check{\delta}_{n,\text{B.E.}} \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} |\psi_i|^3. \quad (5.6)$$

Due to Berry-Esseen theorem by [Berry \(1941\)](#) and [Esseen \(1942\)](#) it holds

$$\begin{aligned} \sup_{z \in \mathbb{R}} |\mathbb{P}(|\phi| \geq z) - \mathbb{P}(|\bar{\phi}| \geq z)| &\leq 2C_0 \frac{\delta_{n,\text{B.E.}}}{(\text{Var } \phi)^{3/2}}, \\ \sup_{z \in \mathbb{R}} |\mathbb{P}(|\psi| \geq z) - \mathbb{P}(|\bar{\psi}| \geq z)| &\leq 2C_0 \frac{\check{\delta}_{n,\text{B.E.}}}{(\text{Var } \psi)^{3/2}}, \end{aligned} \quad (5.7)$$

for the constant $C_0 \in [0.4097, 0.560]$ by [Esseen \(1956\)](#) and [Shevtsova \(2010\)](#).

Lemma 5.2. *Under the conditions of [Theorem 5.1](#) it holds*

$$\begin{aligned} 1. \quad |\mathbb{P}(|\phi| \geq z) - \mathbb{P}(|\bar{\psi}| \geq \bar{z})| &\leq 2C_0 \frac{\delta_{n,\text{B.E.}}}{(\text{Var } \phi)^{3/2}} + \frac{\delta_\Sigma}{2} + \frac{\delta_z}{\sqrt{2}} \frac{1}{\bar{z}} \\ &\leq 2C_0 \frac{\delta_{n,\text{B.E.}}}{(\text{Var } \phi)^{3/2}} + \frac{\delta_\Sigma}{2} + \frac{\delta_z}{\sqrt{2}} \\ &\quad \text{for } \bar{z} \geq 1, \end{aligned}$$

$$\begin{aligned}
2. \quad & |\mathbb{P}(|\phi| \geq z) - \mathbb{P}(|\psi| \geq \bar{z})| \\
& \leq 2C_0 \left\{ \frac{\delta_{n,B.E.}}{(\text{Var } \phi)^{3/2}} + \frac{\check{\delta}_{n,B.E.}}{(\text{Var } \psi)^{3/2}} \right\} + \frac{\delta_\Sigma}{2} + \frac{\delta_z}{\sqrt{2}} \frac{1}{\bar{z}} \\
& \leq 2C_0 \left\{ \frac{\delta_{n,B.E.}}{(\text{Var } \phi)^{3/2}} + \frac{\check{\delta}_{n,B.E.}}{(\text{Var } \psi)^{3/2}} \right\} + \frac{\delta_\Sigma}{2} + \frac{\delta_z}{\sqrt{2}} \quad \text{for } \bar{z} \geq 1. \quad (5.8)
\end{aligned}$$

Proof of Lemma 5.2. Similarly to the proof of Theorem 5.1:

$$\begin{aligned}
\mathbb{P}(|\phi| \geq z) & \stackrel{\text{by (5.7)}}{\leq} \mathbb{P}(|\bar{\phi}| \geq z) + 2C_0(\text{Var } \phi)^{-3/2}\delta_{n,B.E.} \\
& \stackrel{\text{by L. 5.7}}{\leq} \mathbb{P}(|\bar{\psi}| \geq z) + 2C_0(\text{Var } \phi)^{-3/2}\delta_{n,B.E.} + \delta_\Sigma/2 \\
& \stackrel{\text{by L. 5.6}}{\leq} \mathbb{P}(|\bar{\psi}| \geq \bar{z}) + 2C_0(\text{Var } \phi)^{-3/2}\delta_{n,B.E.} + \delta_\Sigma/2 + \delta_z \bar{z}^{-1} 2^{-1/2}.
\end{aligned}$$

The analogous chain in the inverse direction finishes the proof of the first part of the statement. The second part is implied by the triangle inequality applied to the first part and again to it with $\phi := \psi$ and $z := \bar{z}$. \square

5.2 Gaussian approximation of ℓ_2 norm of a sum of independent vectors

Lemma 5.3 (GAR with equal covariance matrices). *For the random vectors ϕ and $\bar{\phi}$ defined in (5.1), (5.3), s.t. $\text{Var } \phi = \text{Var } \bar{\phi}$, and for δ_n given in (5.4), it holds for all $z \geq 2$ and $\Delta \in (0, 0.22]$:*

$$\begin{aligned}
\mathbb{P}(\|\phi\| \geq z) & \leq \mathbb{P}(\|\bar{\phi}\| \geq z - \Delta) + 16\Delta^{-3}\delta_n, \\
\mathbb{P}(\|\phi\| \geq z) & \geq \mathbb{P}(\|\bar{\phi}\| \geq z + \Delta) - 16\Delta^{-3}\delta_n.
\end{aligned}$$

Proof of Lemma 5.3. It holds for $z \in \mathbb{R}$ $\mathbb{P}(\|\phi\| \geq z) = \mathbb{E} \mathbb{I}\{\|\phi\| \geq z\}$. The main idea of the proof is to approximate the discontinuous function $\mathbb{I}\{\|\phi\| \geq z\}$ by a smooth function $f_\Delta(\phi, z)$ and then to apply the Lindeberg's telescopic sum device. Let us introduce a non-negative function $g(\cdot) \in C^2(\mathbb{R})$, which grows monotonously from 0 to 1:

$$g(x) \stackrel{\text{def}}{=} \begin{cases} 0, & x \leq 0, \\ 16x^3/3, & x \in [0, 1/4], \\ 0.5 + 2(x - 0.5) - 16(x - 0.5)^3/3, & x \in [1/4, 3/4], \\ 1 + 16(x - 1)^3/3, & x \in [3/4, 1], \\ 1, & x \geq 1. \end{cases} \quad (5.9)$$

It holds for all $x \in \mathbb{R}$ $\mathbb{I}\{x \geq 1\} \leq g(x) \leq \mathbb{I}\{x \geq 0\}$. Hence, for the function $f_\Delta(\phi, z) \stackrel{\text{def}}{=} g((\|\phi\|^2 - z^2)/(2z\Delta))$ with $z, \Delta > 0$, it holds due to $\mathbb{I}\{\|\phi\| \geq z\} = \mathbb{I}\{(\|\phi\|^2 - z^2)/2 \geq 0\}$:

$$\mathbb{I}\{\|\phi\| \geq z + \Delta\} \leq \mathbb{I}\{\|\phi\|^2 \geq z^2 + 2\Delta z\} \leq f_\Delta(\phi, z) \leq \mathbb{I}\{\|\phi\| \geq z\}. \quad (5.10)$$

Due to Lemma 5.4 one can apply the Lindeberg's telescopic sum device (see Lindeberg (1922)) in order to approximate $\mathbb{E}f_\Delta(\phi, z)$ with $\mathbb{E}f_\Delta(\bar{\phi}, z)$. Define for $k = 2, \dots, n-1$ the following random sums

$$S_k \stackrel{\text{def}}{=} \sum_{i=1}^{k-1} \bar{\phi}_i + \sum_{i=k+1}^n \phi_i, \quad S_1 \stackrel{\text{def}}{=} \sum_{i=2}^n \phi_i, \quad S_n \stackrel{\text{def}}{=} \sum_{i=1}^{n-1} \bar{\phi}_i.$$

The difference $f_\Delta(\phi, z) - f_\Delta(\bar{\phi}, z)$ can be represented as the telescopic sum:

$$f_\Delta(\phi, z) - f_\Delta(\bar{\phi}, z) = \sum_{k=1}^n \{f_\Delta(S_k + \phi_k, z) - f_\Delta(S_k + \bar{\phi}_k, z)\}.$$

Due to Lemma 5.4 and the third order Taylor expansions of $f_\Delta(S_k + \phi_k, z)$ and $f_\Delta(S_k + \bar{\phi}_k, z)$ w.r.t. the first argument at S_k , it holds for each $k = 1, \dots, n$:

$$\begin{aligned} & \left| f_\Delta(S_k + \phi_k, z) - f_\Delta(S_k + \bar{\phi}_k, z) - \nabla_{\phi} f_\Delta(S_k, z)^\top (\phi_k - \bar{\phi}_k) \right. \\ & \left. - \frac{1}{2} (\phi_k - \bar{\phi}_k)^\top \nabla_{\phi}^2 f_\Delta(S_k, z) (\phi_k - \bar{\phi}_k) \right| \leq \mathbf{C}(\Delta, z) (\|\phi_k\|^3 + \|\bar{\phi}_k\|^3) / 6, \end{aligned}$$

where the value $\mathbf{C}(\Delta, z)$ is defined in (5.14). As S_k and $\phi_k - \bar{\phi}_k$ are independent, $\mathbb{E}\phi_k = \mathbb{E}\bar{\phi}_k = 0$ and $\text{Var} \phi_k = \text{Var} \bar{\phi}_k$, we derive

$$\begin{aligned} \left| \mathbb{E}f_\Delta(\phi, z) - \mathbb{E}f_\Delta(\bar{\phi}, z) \right| &= \left| \sum_{k=1}^n \{ \mathbb{E}f_\Delta(S_k + \phi_k, z) - \mathbb{E}f_\Delta(S_k + \bar{\phi}_k, z) \} \right| \\ &\leq \mathbf{C}(\Delta, z) \sum_{k=1}^n \mathbb{E} (\|\phi_k\|^3 + \|\bar{\phi}_k\|^3) / 6 \\ &\text{(by Def. (5.4))} = \mathbf{C}(\Delta, z) \delta_n / 3. \end{aligned} \tag{5.11}$$

Combining the derived bounds, we obtain:

$$\begin{aligned} \mathbb{P}(\|\phi\| \geq z + \Delta) &\stackrel{\text{by (5.10)}}{\leq} \mathbb{E}f_\Delta(\phi, z) \stackrel{\text{by (5.11)}}{\leq} \mathbb{E}f_\Delta(\bar{\phi}, z) + \frac{\mathbf{C}(\Delta, z)}{3} \delta_n \\ &\stackrel{\text{by (5.10)}}{\leq} \mathbb{P}(\|\bar{\phi}\| \geq z) + \frac{\mathbf{C}(\Delta, z)}{3} \delta_n, \end{aligned}$$

or $\mathbb{P}(\|\phi\| \geq z) \leq \mathbb{P}(\|\bar{\phi}\| \geq z - \Delta) + \mathbf{C}(\Delta, z - \Delta) \delta_n / 3$. Interchanging the arguments ϕ and $\bar{\phi}$ implies the inequality in the inverse direction:

$$\mathbb{P}(\|\phi\| \geq z) \geq \mathbb{P}(\|\bar{\phi}\| \geq z + \Delta) - \mathbf{C}(\Delta, z) \delta_n / 3.$$

Let us bound the constants $\mathbf{C}(\Delta, z)$ and $\mathbf{C}(\Delta, z - \Delta)$ for the function $g(x)$ given above in (5.9). $|g''(x)| \leq 8$ and $|g'''(x)| \leq 32$ for all $x \in \mathbb{R}$. By definition (5.14) it holds for $0 < \Delta \leq 0.22$ and $z \geq 2$:

$$\mathbf{C}(\Delta, z) \leq \mathbf{C}(\Delta, z - \Delta) \leq \Delta^{-3} 48. \tag{5.12}$$

□

Lemma 5.4 (A property of the smooth approximant of the indicator). *Let a function $g(\cdot) \in C^2(\mathbb{R})$ be non-negative, monotonously increasing from 0 to 1 s.t. $g(x) = 0$ for $x < 0$, $g(x) = 1$ for $x \geq 1$. It holds for all $\phi, \phi_0 \in \mathbb{R}^p$, $z, \Delta \geq 0$, for the Euclidean norm $\|\cdot\|$ and for the function*

$$f_\Delta(\phi, z) \stackrel{\text{def}}{=} g\left(\frac{1}{2z\Delta}(\|\phi\|^2 - z^2)\right) \quad (5.13)$$

$$\begin{aligned} |f_\Delta(\phi_0 + \phi, z) - f_\Delta(\phi_0, z) - \phi^\top \nabla_\phi f_\Delta(\phi_0, z) - \phi^\top \nabla_\phi^2 f_\Delta(\phi_0, z) \phi / 2| \\ \leq c(\Delta, z) \|\phi\|^3 / 3!, \end{aligned}$$

where

$$c(\Delta, z) \stackrel{\text{def}}{=} \frac{1}{\Delta^3} \left(1 + 2\frac{\Delta}{z}\right)^{1/2} \left\{ \left(1 + 2\frac{\Delta}{z}\right) \|g'''\|_\infty + 3\frac{\Delta}{z} \|g''\|_\infty \right\}. \quad (5.14)$$

Proof of Lemma 5.4. By the Taylor formula:

$$f_\Delta(\phi_0 + \phi, z) = f_\Delta(\phi_0, z) + \phi^\top \nabla_\phi f_\Delta(\phi_0, z) + \phi^\top \nabla_\phi^2 f_\Delta(\phi_0, z) \phi / 2 + R_3,$$

where R_3 is the 3-d order remainder term. Consider for $\gamma \in \mathbb{R}^p : \|\gamma\| = 1$ and $t \in \mathbb{R}$ the function $f_\Delta(\phi_0 + t\gamma, z) = g\left(\frac{1}{2z\Delta}(\|\phi_0 + t\gamma\|^2 - z^2)\right)$. It holds

$$|R_3| \leq \frac{\|\phi\|^3}{3!} \sup_{\gamma \in \mathbb{R}^p, \|\gamma\|=1} \sup_{t \in \mathbb{R}} \left| \frac{d^3 f_\Delta(\phi_0 + t\gamma, z)}{dt^3} \right|.$$

Now let us bound the third derivative $\frac{d^3}{dt^3} f_\Delta(\phi_0 + t\gamma, z)$:

$$\begin{aligned} \frac{df_\Delta(\phi_0 + t\gamma, z)}{dt} &= \frac{\gamma^\top(\phi_0 + t\gamma)}{z\Delta} g' \left(\frac{1}{2z\Delta}(\|\phi_0 + t\gamma\|^2 - z^2) \right), \\ \frac{d^2 f_\Delta(\phi_0 + t\gamma, z)}{dt^2} &= \frac{\{\gamma^\top(\phi_0 + t\gamma)\}^2}{(z\Delta)^2} g'' \left(\frac{1}{2z\Delta}(\|\phi_0 + t\gamma\|^2 - z^2) \right) \\ &\quad + \frac{1}{z\Delta} g' \left(\frac{1}{2z\Delta}(\|\phi_0 + t\gamma\|^2 - z^2) \right), \\ \frac{d^3 f_\Delta(\phi_0 + t\gamma, z)}{dt^3} &= \frac{\{\gamma^\top(\phi_0 + t\gamma)\}^3}{(z\Delta)^3} g''' \left(\frac{1}{2z\Delta}(\|\phi_0 + t\gamma\|^2 - z^2) \right) \\ &\quad + 3 \frac{\gamma^\top(\phi_0 + t\gamma)}{(z\Delta)^2} g'' \left(\frac{1}{2z\Delta}(\|\phi_0 + t\gamma\|^2 - z^2) \right). \end{aligned}$$

Now we use that $g''(x)$ and $g'''(x)$ vanish if $x < 0$ or $x \geq 1$. The inequality $\frac{1}{2z\Delta}(\|\phi_0 + t\gamma\|^2 - z^2) \leq 1$ implies in view of $\|\gamma\| = 1$ that

$$\gamma^\top(\phi_0 + t\gamma) \leq \|\phi_0 + t\gamma\| \leq (2z\Delta + z^2)^{1/2}.$$

Therefore

$$\left| \frac{d^3 f_\Delta(\phi_0 + t\gamma, z)}{dt^3} \right| \leq \frac{1}{\Delta^3} \left(1 + 2\frac{\Delta}{z} \right)^{1/2} \left\{ \left(1 + 2\frac{\Delta}{z} \right) \|g'''\|_\infty + 3\frac{\Delta}{z} \|g''\|_\infty \right\}.$$

□

5.3 Results for the smoothed indicator function

Theorem 5.5 (Theorem 5.1 for a smoothed indicator function). *Under the conditions of Theorem 5.1 it holds for all $\delta_z \in [0, 1]$ and the function $f_\Delta(\phi, z)$ defined in (5.13):*

1.
$$\begin{aligned} |\mathbb{E} f_\Delta(\phi, z) - \mathbb{E} f_\Delta(\bar{\psi}, \bar{z})| &\leq \frac{16}{\Delta^3} \delta_n + 2\sqrt{p} \frac{\delta_z}{z} + \sqrt{p} \frac{\delta_z^2}{z^2} + \delta_\Sigma \\ &\leq \frac{16}{\Delta^3} \delta_n + \sqrt{5} \delta_z + \delta_\Sigma \quad \text{for } z \geq \sqrt{p}. \end{aligned}$$
2.
$$\begin{aligned} |\mathbb{E} f_\Delta(\phi, z) - \mathbb{E} f_\Delta(\psi, \bar{z})| &\leq \frac{16}{\Delta^3} (\delta_n + \check{\delta}_n) + 2\sqrt{p} \frac{\delta_z}{z} + \sqrt{p} \frac{\delta_z^2}{z^2} + \delta_\Sigma \\ &\leq \frac{16}{\Delta^3} (\delta_n + \check{\delta}_n) + \sqrt{5} \delta_z + \delta_\Sigma \quad \text{for } z \geq \sqrt{p}. \end{aligned}$$

REMARK 5.3. The approximating bounds above do not contain the term proportional to Δ unlike the bound in Theorem 5.1. This yields the smaller error terms for the case of the smoothed indicator.

Proof of Theorem 5.5. The following inequality is proved in Lemma 5.3 (see the expression (5.11)): $|\mathbb{E} f_\Delta(\phi, z) - \mathbb{E} f_\Delta(\bar{\phi}, z)| \leq \mathbf{C}(\Delta, z) \delta_n / 3$.

The function $f_\Delta(\phi, z)$ is non-increasing in z :

$$\frac{df_\Delta(\phi, z)}{dz} = -\frac{1}{2\Delta} \left(1 + \frac{\|\phi\|^2}{z^2} \right) g' \left(\frac{1}{2\Delta z} (\|\phi\|^2 - z^2) \right) \leq 0.$$

The definition of $f_\Delta(\phi, z)$ yields for $\bar{z} \geq z$, $a \stackrel{\text{def}}{=} \bar{z}/z \geq 1$ and any ϕ

$$\begin{aligned} f_\Delta(\phi, \bar{z}) &\leq f_\Delta(\phi, z) \leq f_\Delta(a\phi, \bar{z}), \\ 0 &\leq f_\Delta(\phi, z) - f_\Delta(\phi, \bar{z}) \leq f_\Delta(a\phi, \bar{z}) - f_\Delta(\phi, \bar{z}). \end{aligned} \tag{5.15}$$

Lemma 5.8 yields for $\delta_z \leq z(\sqrt{3/2} - 1)$:

$$\begin{aligned} |\mathbb{E} f_\Delta(a\bar{\phi}, \bar{z}) - \mathbb{E} f_\Delta(\bar{\phi}, \bar{z})| &\leq \sqrt{p} \left(\frac{\bar{z}^2}{z^2} - 1 \right) \leq 2\sqrt{p} \frac{\delta_z}{z} + \sqrt{p} \frac{\delta_z^2}{z^2} \\ &\leq (1 + \sqrt{3/2}) \delta_z \leq \sqrt{5} \delta_z \quad \text{for } z \geq \sqrt{p}. \end{aligned}$$

Inequalities similar to (5.15) hold for $\bar{z} \leq z$ and $a \stackrel{\text{def}}{=} z/\bar{z}$, therefore, by triangle inequality, bound (5.12) on $\mathbb{C}(\Delta, z)$ and Lemma 5.8:

$$\begin{aligned} |\mathbb{E}f_{\Delta}(\phi, z) - \mathbb{E}f_{\Delta}(\bar{\psi}, \bar{z})| &\leq \frac{16}{\Delta^3}\delta_n + 2\sqrt{p}\frac{\delta_z}{z} + \sqrt{p}\frac{\delta_z^2}{z^2} + \delta_{\Sigma} \\ &\leq \frac{16}{\Delta^3}\delta_n + \sqrt{5}\delta_z + \delta_{\Sigma} \quad \text{for } z \geq \sqrt{p}. \end{aligned}$$

The second part of the statement follows from triangle inequality applied to the first inequality and again to the same one with $\phi := \psi$ and $z := \bar{z}$. \square

5.4 Gaussian anti-concentration and comparison by Pinsker's inequality

Lemma 5.6 (Anti-concentration bound for ℓ_2 norm of a Gaussian vector). *Let $\bar{\phi} \sim \mathcal{N}(0, \Sigma)$, $\bar{\phi} \in \mathbb{R}^p$, then it holds for all $z > 0$ and $0 \leq \Delta \leq z$:*

$$\begin{aligned} |\mathbb{P}(\|\bar{\phi}\| \geq z + \Delta) - \mathbb{P}(\|\bar{\phi}\| \geq z)| &\leq \Delta\sqrt{p}/(z\sqrt{2}) \\ &\leq \Delta/\sqrt{2} \quad \text{for } z \geq \sqrt{p}. \end{aligned}$$

Proof of Lemma 5.6. It holds $\mathbb{P}(\|\bar{\phi}\| \geq z + \Delta) = \mathbb{P}(\|\bar{\phi}_{\Delta}\| \geq z)$, where $\bar{\phi}_{\Delta} \stackrel{\text{def}}{=} \bar{\phi}_{z+\Delta}$. The Kullback-Leibler divergence between $\mathbb{P}_1 \stackrel{\text{def}}{=} \mathcal{N}(0, \Sigma)$ and $\mathbb{P}_2 \stackrel{\text{def}}{=} \mathcal{N}(0, \Sigma \frac{z^2}{(z+\Delta)^2})$ is equal to

$$\begin{aligned} \text{KL}(\mathbb{P}_1, \mathbb{P}_2) &= p \{ (\Delta/z)^2 + 2(\Delta/z) - 2\log(1 + \Delta/z) \} / 2 \\ &\leq p(\Delta/z)^2 \quad \text{for } 0 \leq \Delta \leq z. \end{aligned}$$

We use Pinsker's inequality in the following form (see the book by [Tsybakov \(2009\)](#), pp. 88, 132): for a measurable space (Ω, \mathcal{F}) and two measures on it $\mathbb{P}_1, \mathbb{P}_2$:

$$\sup_{A \in \mathcal{F}} |\mathbb{P}_1(A) - \mathbb{P}_2(A)| \leq \sqrt{\text{KL}(\mathbb{P}_1, \mathbb{P}_2)}/2. \quad (5.16)$$

Therefore, it holds:

$$|\mathbb{P}(\|\bar{\phi}\| \geq z + \Delta) - \mathbb{P}(\|\bar{\phi}\| \geq z)| \leq \sqrt{\text{KL}(\mathbb{P}_1, \mathbb{P}_2)}/2 \leq \Delta\sqrt{p}/(z\sqrt{2}).$$

\square

Lemma 5.7 (Comparison of the Euclidian norms of Gaussian vectors). *Let $\bar{\psi}_1 \sim \mathcal{N}(0, \Sigma_1)$ and $\bar{\psi}_2 \sim \mathcal{N}(0, \Sigma_2)$ belong to \mathbb{R}^p , and*

$$\|\Sigma_2^{-1/2}\Sigma_1\Sigma_2^{-1/2} - \mathbf{I}_p\| \leq 1/2, \quad \text{and} \quad \text{tr}\{(\Sigma_2^{-1/2}\Sigma_1\Sigma_2^{-1/2} - \mathbf{I}_p)^2\} \leq \delta_{\Sigma}^2,$$

for some $\delta_{\Sigma}^2 \geq 0$. Then it holds

$$\sup_{z \in \mathbb{R}} |\mathbb{P}(\|\bar{\psi}_1\| \geq z) - \mathbb{P}(\|\bar{\psi}_2\| \geq z)| \leq \delta_{\Sigma}/2.$$

Proof of Lemma 5.7. Let $\mathbb{P}_1 = \mathcal{N}(0, \Sigma_1)$ and $\mathbb{P}_2 = \mathcal{N}(0, \Sigma_2)$. Denote $G \stackrel{\text{def}}{=} \Sigma_2^{-1/2} \Sigma_1 \Sigma_2^{-1/2}$, then the Kullback-Leibler divergence between \mathbb{P}_1 and \mathbb{P}_2 is equal to

$$\begin{aligned} \text{KL}(\mathbb{P}_1, \mathbb{P}_2) &= -0.5 \log\{\det(G)\} + 0.5 \text{tr}\{G - \mathbf{I}_p\} \\ &= 0.5 \sum_{j=1}^p \{\lambda_j - \log(\lambda_j + 1)\}, \end{aligned}$$

where $\lambda_p \leq \dots \leq \lambda_1$ are the eigenvalues the matrix $G - \mathbf{I}_p$. By conditions of the lemma $|\lambda_1| \leq 1/2$, and it holds:

$$\text{KL}(\mathbb{P}_1, \mathbb{P}_2) \leq 0.5 \sum_{j=1}^p \lambda_j^2 = 0.5 \text{tr}\{(G - \mathbf{I}_p)^2\} \leq \delta_\Sigma^2/2, \quad (5.17)$$

which finishes the proof due to the Pinsker's inequality (5.16). \square

Lemma 5.8 (Gaussian comparison, smoothed version). Let $\bar{\boldsymbol{\psi}}_1 \sim \mathcal{N}(0, \Sigma_1)$ and $\bar{\boldsymbol{\psi}}_2 \sim \mathcal{N}(0, \Sigma_2)$ belong to \mathbb{R}^p , and for some $\delta_\Sigma^2 \geq 0$:

$$\|\Sigma_2^{-1/2} \Sigma_1 \Sigma_2^{-1/2} - \mathbf{I}_p\| \leq 1/2, \quad \text{and} \quad \text{tr}\{(\Sigma_2^{-1/2} \Sigma_1 \Sigma_2^{-1/2} - \mathbf{I}_p)^2\} \leq \delta_\Sigma^2.$$

Then it holds for any function $f(\mathbf{x}) : \mathbb{R}^p \mapsto \mathbb{R}$ s.t. $|f(\mathbf{x})| \leq 1$:

$$|\mathbb{E}f(\bar{\boldsymbol{\psi}}_1) - \mathbb{E}f(\bar{\boldsymbol{\psi}}_2)| \leq \delta_\Sigma.$$

Proof of Lemma 5.8. Let $\mathbb{P}_1 = \mathcal{N}(0, \Sigma_1)$ and $\mathbb{P}_2 = \mathcal{N}(0, \Sigma_2)$. Due to $|f(\mathbf{x})| \leq 1$ and Pinsker's inequality (5.16) it holds:

$$\begin{aligned} |\mathbb{E}f(\bar{\boldsymbol{\psi}}_1) - \mathbb{E}f(\bar{\boldsymbol{\psi}}_2)| &\leq \int_{\mathbb{R}^p} |f(\mathbf{x})| \cdot |d\mathbb{P}_1(\mathbf{x}) - d\mathbb{P}_2(\mathbf{x})| \\ &\leq \int_{\mathbb{R}^p} |d\mathbb{P}_1(\mathbf{x}) - d\mathbb{P}_2(\mathbf{x})| \leq 2 \sqrt{\text{KL}(\mathbb{P}_1, \mathbb{P}_2)/2}. \end{aligned}$$

Finally, as in (5.17), $2\sqrt{\text{KL}(\mathbb{P}_1, \mathbb{P}_2)/2} \leq \delta_\Sigma$. \square

A Appendix

This section contains proofs of the main results from Section 2. Due to the scheme (1.6) the key ingredients are:

- the square-root Wilks approximation for the \mathbf{Y} -world (Theorem A.2),
- the square-root Wilks approximation for the bootstrap world (Theorem A.4),
- the statement about closeness in distribution of the approximating terms $\|\boldsymbol{\xi}\|$ and $\|\boldsymbol{\xi}^\circ\|$ (Proposition A.9).

In Section A.1 we recall some results from the general finite sample theory by Spokoiny (2012a), Spokoiny (2012b) and Spokoiny (2013), including the square-root Wilks approximation in \mathbf{Y} case. In Section A.2 we derive the necessary results from the finite sample theory for the bootstrap world (including the square-root Wilks approximation). In Section A.3 we adapt Theorem 5.1 (GAR for ℓ_2 norm of a sum of independent vectors) to the setting of maximum likelihood estimation (Proposition A.9). The proofs of the main results are given in Sections A.3, A.4 and A.5.

A.1 Finite sample theory

Let us use the notations given in the introduction: $L(\boldsymbol{\theta})$ is the log-likelihood process, which depends on the data \mathbf{Y} and corresponds to the regular parametric family of probability distributions $\{P_{\boldsymbol{\theta}}\}$. The general finite sample approach by Spokoiny (2012a) does not require the true measure P to belong to $\{P_{\boldsymbol{\theta}}\}$. The target parameter $\boldsymbol{\theta}^*$ is defined as in (1.4) by projection of the true measure P on $\{P_{\boldsymbol{\theta}}\}$. D_0^2 denotes the full Fisher information $p \times p$ matrix, which is deterministic, symmetric and positive-definite:

$$D_0^2 \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E}L(\boldsymbol{\theta}^*).$$

A centered p -dimensional random vector $\boldsymbol{\xi}$ denotes the normalised score:

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*).$$

Introduce the following elliptic vicinity around the true point $\boldsymbol{\theta}^*$:

$$\Theta_0(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} : \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}\}.$$

The non-asymptotic Wilks approximating bound by Spokoiny (2012a), Spokoiny (2013) requires that the maximum likelihood estimate $\tilde{\boldsymbol{\theta}}$ gets into the local vicinity $\Theta_0(\mathbf{r}_0)$ of some radius $\mathbf{r}_0 > 0$ with probability $\geq 1 - 3e^{-\mathbf{x}}$, $\mathbf{x} > 0$. This is guaranteed by the following concentration result:

Theorem A.1 (Concentration of MLE, Spokoiny (2013)). *Let the conditions (ED_0) , (ED_2) , (\mathcal{L}_0) , (\mathcal{I}) and $(\mathcal{L}_\mathbf{r})$ be fulfilled. If for the constant $\mathbf{r}_0 > 0$ and for the function $b(\mathbf{r})$ from $(\mathcal{L}_\mathbf{r})$:*

$$b(\mathbf{r})\mathbf{r} \geq 2 \{ \mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathbb{B}) + 6\omega\nu_0\mathfrak{Z}(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0)) \}, \quad \mathbf{r} > \mathbf{r}_0 \quad (\text{A.1})$$

where the functions $\mathfrak{Z}(\mathbf{x})$ and $\mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathbb{B})$ are defined respectively in (A.3) and (A.4), then it holds

$$P \left(\tilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r}_0) \right) \leq 3e^{-\mathbf{x}}.$$

The constants ω, ν_0 and α come from the imposed conditions $(ED_0) - (\mathcal{I})$ (from Section 4). In the case 4.3 $\mathbf{r}_0 \geq C\sqrt{p + \mathbf{x}}$.

The following result is one of the central in the general finite sample theory and is crucial for the present study due to the scheme (1.6):

Theorem A.2 (Wilks approximation, Spokoiny (2013)). *Under the conditions of Theorem A.1 for some $r_0 > 0$ s.t. (A.1) is fulfilled, it holds with probability $\geq 1 - 5e^{-x}$*

$$\begin{aligned} \left| 2\{L(\tilde{\theta}) - L(\theta^*)\} - \|\xi\|^2 \right| &\leq \Delta_{W^2}(r_0, x), \\ \left| \sqrt{2\{L(\tilde{\theta}) - L(\theta^*)\}} - \|\xi\| \right| &\leq \Delta_W(r_0, x) \end{aligned}$$

for

$$\Delta_W(r, x) \stackrel{\text{def}}{=} 3r \{ \delta(r) + 6\nu_0 \mathfrak{Z}(x) \omega \}, \quad (\text{A.2})$$

$$\Delta_{W^2}(r, x) \stackrel{\text{def}}{=} \frac{2}{3} \{ 2r + \mathfrak{Z}_{\text{qf}}(x, \mathbb{B}) \} \Delta_W(r, x),$$

$$\mathfrak{Z}(x) \stackrel{\text{def}}{=} 2\sqrt{p} + \sqrt{2x} + 4p(xg^{-2} + 1)g^{-1}. \quad (\text{A.3})$$

In the case 4.3 it holds for $r \leq r_0$:

$$\Delta_W(r, x) = C \frac{p+x}{\sqrt{n}}, \quad \Delta_{W^2}(r, x) = C \sqrt{\frac{(p+x)^3}{n}}.$$

The constants g and $\delta(r)$ come from the imposed conditions (\mathbf{ED}_0) , (\mathcal{L}_0) (from Section 4), and the function $\mathfrak{Z}_{\text{qf}}(x, \mathbb{B})$, defined in (A.4), corresponds to the quantile function of deviations of the random value $\|\xi\|$ (see Theorem A.3 below).

The following theorem characterizes the tail behaviour of the approximating term $\|\xi\|^2$. It means that with a bounded exponential moment of the vector ξ (condition (\mathbf{ED}_0)) its squared Euclidean norm $\|\xi\|^2$ has three regimes of deviations: sub-Gaussian, Poissonian and large-deviations' zone.

Theorem A.3 (Deviation bound for a random quadratic form, Spokoiny (2012b)). *Let condition (\mathbf{ED}_0) be fulfilled, then for $g \geq \sqrt{2 \text{tr}(\mathbb{B}^2)}$ it holds:*

$$\mathbb{P}(\|\xi\|^2 \geq \mathfrak{Z}_{\text{qf}}^2(x, \mathbb{B})) \leq 2e^{-x} + 8.4e^{-x_c},$$

where $\mathbb{B}^2 \stackrel{\text{def}}{=} D_0^{-1} V_0^2 D_0^{-1}$, $\lambda(\mathbb{B})$ is a maximum eigenvalue of \mathbb{B}^2 ,

$$\mathfrak{Z}_{\text{qf}}^2(x, \mathbb{B}) \stackrel{\text{def}}{=} \begin{cases} \text{tr}(\mathbb{B}^2) + \sqrt{8 \text{tr}(\mathbb{B}^4)x}, & x \leq \sqrt{2 \text{tr}(\mathbb{B}^4)} / \{18\lambda(\mathbb{B})\}, \\ \text{tr}(\mathbb{B}^2) + 6x\lambda(\mathbb{B}), & \sqrt{2 \text{tr}(\mathbb{B}^4)} / \{18\lambda(\mathbb{B})\} < x \leq x_c, \\ |z_c + 2(x - x_c)/g_c|^2 \lambda(\mathbb{B}), & x > x_c, \end{cases} \quad (\text{A.4})$$

$$\begin{aligned}
2\mathbf{x}_c &\stackrel{\text{def}}{=} 2\mathbf{x}_c(\mathcal{B}) \stackrel{\text{def}}{=} \mu_c \mathbf{z}_c^2 + \log \det (\mathbf{I}_p - \mu_c \mathcal{B}^2 / \lambda(\mathcal{B})), \\
\mathbf{z}_c^2 &\stackrel{\text{def}}{=} \{ \mathbf{g}^2 / \mu_c^2 - \text{tr}(\mathcal{B}^2) / \mu_c \} / \lambda(\mathcal{B}), \\
\mathbf{g}_c &\stackrel{\text{def}}{=} \sqrt{\mathbf{g}^2 - \mu_c \text{tr}(\mathcal{B}^2)} / \sqrt{\lambda(\mathcal{B})}, \\
\mu_c &\stackrel{\text{def}}{=} 2/3.
\end{aligned} \tag{A.5}$$

The matrix V_0^2 comes from condition (\mathbf{ED}_0) and can be defined as

$$V_0^2 \stackrel{\text{def}}{=} \text{Var} \{ \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*) \}.$$

By condition (\mathcal{I}) $\text{tr}(\mathcal{B}^2) \leq \mathfrak{a}^2 p$, $\text{tr}(\mathcal{B}^4) \leq \mathfrak{a}^4 p$ and $\lambda(\mathcal{B}) \leq \mathfrak{a}^2$. In the case 4.3 $\mathbf{g} = C\sqrt{n}$, hence $\mathbf{x}_c = Cn$, and for $\mathbf{x} \leq \mathbf{x}_c$ it holds:

$$\mathfrak{Z}_{\text{qf}}^2(\mathbf{x}, \mathcal{B}) \leq \mathfrak{a}^2(p + 6\mathbf{x}). \tag{A.6}$$

A.2 Finite sample theory for the bootstrap world

Let us introduce the bootstrap score vector at a point $\boldsymbol{\theta} \in \Theta$:

$$\begin{aligned}
\xi^\circ(\boldsymbol{\theta}) &\stackrel{\text{def}}{=} D_0^{-1} \nabla_{\boldsymbol{\theta}} \zeta^\circ(\boldsymbol{\theta}) \\
&= \sum_{i=1}^n D_0^{-1} \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta})(u_i - 1).
\end{aligned}$$

Theorem A.4 (Bootstrap Wilks approximation). *Under the conditions of Theorems A.1 and A.6 for some $\mathbf{r}_0^2 \geq 0$ s.t. (A.1) and (A.17) are fulfilled, it holds with \mathbb{P} -probability $\geq 1 - 5e^{-\mathbf{x}}$*

$$\begin{aligned}
\mathbb{P}^\circ \left(\left| \sup_{\boldsymbol{\theta} \in \Theta} 2 \{ L^\circ(\boldsymbol{\theta}) - L^\circ(\tilde{\boldsymbol{\theta}}) \} - \|\xi^\circ(\tilde{\boldsymbol{\theta}})\|^2 \right| \leq \Delta_{W^2}^\circ(\mathbf{r}_0, \mathbf{x}) \right) &\geq 1 - 4e^{-\mathbf{x}}, \\
\mathbb{P}^\circ \left(\left| \sqrt{\sup_{\boldsymbol{\theta} \in \Theta} 2 \{ L^\circ(\boldsymbol{\theta}) - L^\circ(\tilde{\boldsymbol{\theta}}) \}} - \|\xi^\circ(\tilde{\boldsymbol{\theta}})\| \right| \leq \Delta_W^\circ(\mathbf{r}_0, \mathbf{x}) \right) &\geq 1 - 4e^{-\mathbf{x}}.
\end{aligned}$$

where the error terms $\Delta_W^\circ(\mathbf{r}, \mathbf{x})$, $\Delta_{W^2}^\circ(\mathbf{r}, \mathbf{x})$ are deterministic and

$$\begin{aligned}
\Delta_W^\circ(\mathbf{r}, \mathbf{x}) &\stackrel{\text{def}}{=} 2\Delta_W(\mathbf{r}, \mathbf{x}) + 36\nu_0 \mathbf{r} \omega_1(\mathbf{r}) \mathfrak{Z}(\mathbf{x}), \\
\Delta_{W^2}^\circ(\mathbf{r}, \mathbf{x}) &\stackrel{\text{def}}{=} \frac{1}{18} \{ 12\mathbf{r} \Delta_W^\circ(\mathbf{r}, \mathbf{x}) + \Delta_W^\circ(\mathbf{r}, \mathbf{x})^2 \}.
\end{aligned}$$

$\Delta_W(\mathbf{r}, \mathbf{x})$ and $\mathfrak{Z}(\mathbf{x})$ are defined respectively in (A.2) and (A.3), and $\omega_1(\mathbf{r})$ is given in (A.12). For the case 4.3 and $\mathbf{r} \leq \mathbf{r}_0$ it holds:

$$\Delta_W^\circ(\mathbf{r}, \mathbf{x}) = C \frac{p + \mathbf{x}}{\sqrt{n}} \sqrt{\mathbf{x}}, \quad \Delta_{W^2}^\circ(\mathbf{r}, \mathbf{x}) = C \sqrt{\frac{(p + \mathbf{x})^3}{n}} \sqrt{\mathbf{x}}.$$

Proof of Theorem A.4. Let us consider the following random process in the bootstrap world for $\boldsymbol{\theta}, \boldsymbol{\theta}_1 \in \Theta_0(\mathbf{r})$:

$$\mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1) \stackrel{\text{def}}{=} L^\circ(\boldsymbol{\theta}) - L^\circ(\boldsymbol{\theta}_1) - (\boldsymbol{\theta} - \boldsymbol{\theta}_1)^\top \nabla_{\boldsymbol{\theta}} L^\circ(\boldsymbol{\theta}_1) + \frac{1}{2} \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}_1)\|^2.$$

It holds $\mathcal{A}^\circ(\boldsymbol{\theta}_1, \boldsymbol{\theta}_1) = 0$. Taylor expansion w.r.t. $\boldsymbol{\theta}$ around $\boldsymbol{\theta}_1$ implies:

$$\mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1) = (\boldsymbol{\theta} - \boldsymbol{\theta}_1)^\top \nabla_{\boldsymbol{\theta}} \mathcal{A}^\circ(\bar{\boldsymbol{\theta}}_1, \boldsymbol{\theta}_1),$$

where $\bar{\boldsymbol{\theta}}_1$ is some convex combination of the vectors $\boldsymbol{\theta}$ and $\boldsymbol{\theta}_1$. Therefore,

$$|\mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)| \leq \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}_1)\| \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \|D_0^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)\| \quad (\text{A.7})$$

$$\leq 2\mathbf{r} \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \|D_0^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)\|. \quad (\text{A.8})$$

Now let us consider the normalized gradient process:

$$D_0^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1) = D_0^{-1} \{ \nabla_{\boldsymbol{\theta}} L^\circ(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} L^\circ(\boldsymbol{\theta}_1) \} + D_0(\boldsymbol{\theta} - \boldsymbol{\theta}_1).$$

The deterministic part of it reads as:

$$D_0^{-1} \nabla_{\boldsymbol{\theta}} \mathbb{E}^\circ \mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1) = D_0^{-1} \{ \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_1) \} + D_0(\boldsymbol{\theta} - \boldsymbol{\theta}_1).$$

Proposition 3.1 in Spokoiny (2013) implies due to the conditions (\mathcal{L}_0) , (ED_2) , that the following random event holds with \mathbb{P} -probability at least $1 - e^{-\mathbf{x}}$ for all $\boldsymbol{\theta}, \boldsymbol{\theta}_1 \in \Theta_0(\mathbf{r})$ and $\mathbf{r} \leq \mathbf{r}_0$:

$$\begin{aligned} \|D_0^{-1} \nabla_{\boldsymbol{\theta}} \mathbb{E}^\circ \mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)\| &= \|D_0^{-1} \{ \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_1) \} + D_0(\boldsymbol{\theta} - \boldsymbol{\theta}_1)\| \\ &\leq \frac{2}{3} \Delta_{\mathbf{w}}(\mathbf{r}, \mathbf{x}), \end{aligned} \quad (\text{A.9})$$

where the deterministic error term $\Delta_{\mathbf{w}}(\mathbf{r}, \mathbf{x})$ is given in (A.2).

Denote the stochastic part of $D_0^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)$ as follows:

$$\begin{aligned} \mathcal{Y}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1) &\stackrel{\text{def}}{=} D_0^{-1} \{ \nabla_{\boldsymbol{\theta}} \mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1) - \nabla_{\boldsymbol{\theta}} \mathbb{E}^\circ \mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1) \} \\ &= \sum_{i=1}^n D_0^{-1} \{ \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}_1) \} (u_i - 1). \end{aligned}$$

In order to bound its norm's supremum w.r.t. $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ for $\mathbf{r} \leq \mathbf{r}_0$ we use the idea from the proof of Proposition 3.1 in Spokoiny (2013). Let us introduce the new parameters $\mathbf{v} \stackrel{\text{def}}{=} D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$ and $\mathbf{v}_1 \stackrel{\text{def}}{=} D_0(\boldsymbol{\theta}_1 - \boldsymbol{\theta}^*)$, then

$$\nabla_{\mathbf{v}} \mathcal{Y}^\circ(\mathbf{v}, \mathbf{v}_1) = \sum_{i=1}^n D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}) D_0^{-1} (u_i - 1).$$

Thus, we obtain a proper normalisation for $\nabla_{\mathbf{v}} \mathcal{Y}^\circ(\mathbf{v}, \mathbf{v}_1)$. Independency of u_1, \dots, u_n and Lemma A.5 imply with probability $\geq 1 - e^{-x}$ for $j = 1, 2$ and $\omega_1(\mathbf{r})$ given in (A.12):

$$\sup_{\substack{\gamma_j \in \mathbb{R}^p \\ \|\gamma_j\|=1}} \log \mathbb{E}^\circ \exp \left\{ \frac{\lambda}{\omega_1(\mathbf{r})} \gamma_1^\top \nabla_{\mathbf{v}} \mathcal{Y}^\circ(\mathbf{v}, \mathbf{v}_1) \gamma_2 \right\} \leq \frac{\lambda^2 \nu_0^2}{2}, \quad |\lambda| \leq \mathfrak{g}_2(\mathbf{r}).$$

This allows to apply Theorem A.3 from Spokoiny (2013) on a uniform bound for the norm of stochastic process to $\omega_1^{-1}(\mathbf{r}) \mathcal{Y}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)$. By the triangle inequality it holds for $\mathbf{r} \leq \mathbf{r}_0$:

$$\mathbb{P}^\circ \left(\sup_{\boldsymbol{\theta}, \boldsymbol{\theta}_1 \in \Theta_0(\mathbf{r})} \|\mathcal{Y}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)\| \leq 12\nu_0 \mathbf{r} \omega_1(\mathbf{r}) \mathfrak{Z}(\mathbf{x}) \right) \geq 1 - e^{-x}, \quad (\text{A.10})$$

where $\mathfrak{Z}(\mathbf{x})$ is defined in (A.3). Collecting together the bounds (A.8), (A.9) and (A.10) we obtain that the following bound holds with \mathbb{P} -probability at least $1 - e^{-x}$:

$$\mathbb{P}^\circ \left(\sup_{\boldsymbol{\theta}, \boldsymbol{\theta}_1 \in \Theta_0(\mathbf{r})} |\mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)| \leq 4\mathbf{r} \{ \Delta_{\mathbf{w}}(\mathbf{r}, \mathbf{x})/3 + 6\nu_0 \mathbf{r} \omega_1(\mathbf{r}) \mathfrak{Z}(\mathbf{x}) \} \right) \geq 1 - e^{-x}$$

for $\mathbf{r} \leq \mathbf{r}_0$.

Theorems A.6 and A.1 say that the maximum likelihood estimators $\tilde{\boldsymbol{\theta}}^\circ$ and $\tilde{\boldsymbol{\theta}}$ get into the local vicinity $\Theta_0(\mathbf{r}_0)$ with exponentially high \mathbb{P}° - and \mathbb{P} -probabilities correspondingly. Therefore, taking $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}^\circ$ and $\boldsymbol{\theta}_1 = \tilde{\boldsymbol{\theta}}$ in the last bound, we obtain with dominating probability:

$$\begin{aligned} & \left| L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}}) - (\tilde{\boldsymbol{\theta}}^\circ - \tilde{\boldsymbol{\theta}})^\top \nabla_{\boldsymbol{\theta}} L^\circ(\tilde{\boldsymbol{\theta}}) + \frac{1}{2} \|D_0(\tilde{\boldsymbol{\theta}}^\circ - \tilde{\boldsymbol{\theta}})\|^2 \right| \\ & \leq 4\mathbf{r} \{ \Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x})/3 + 6\nu_0 \mathbf{r}_0 \omega_1(\mathbf{r}) \mathfrak{Z}(\mathbf{x}) \}. \end{aligned}$$

Similarly bounds (A.9) and (A.10) imply:

$$\begin{aligned} & \frac{1}{2} \left| \|\boldsymbol{\xi}^\circ(\tilde{\boldsymbol{\theta}})\|^2 - 2(\tilde{\boldsymbol{\theta}}^\circ - \tilde{\boldsymbol{\theta}})^\top \nabla_{\boldsymbol{\theta}} L^\circ(\tilde{\boldsymbol{\theta}}) + \|D_0(\tilde{\boldsymbol{\theta}}^\circ - \tilde{\boldsymbol{\theta}})\|^2 \right| \\ & = \frac{1}{2} \|D_0^{-1} \nabla_{\boldsymbol{\theta}} L^\circ(\tilde{\boldsymbol{\theta}}) - D_0(\tilde{\boldsymbol{\theta}}^\circ - \tilde{\boldsymbol{\theta}})\|^2 \\ & \leq 2 \{ \Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x})/3 + 6\nu_0 \mathbf{r}_0 \omega_1(\mathbf{r}) \mathfrak{Z}(\mathbf{x}) \}^2. \end{aligned} \quad (\text{A.11})$$

Therefore it holds with \mathbb{P} -probability at least $1 - 4e^{-x}$:

$$\begin{aligned} & \mathbb{P}^\circ \left(\left| L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}}) - \frac{1}{2} \|\boldsymbol{\xi}^\circ(\tilde{\boldsymbol{\theta}})\|^2 \right| \leq \Delta_{\mathbf{w}^2}^\circ(\mathbf{r}_0, \mathbf{x}) \right) \geq 1 - 4e^{-x}, \\ & \Delta_{\mathbf{w}^2}^\circ(\mathbf{r}_0, \mathbf{x}) \stackrel{\text{def}}{=} 4\mathbf{r} \{ \Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x})/3 + 6\nu_0 \mathbf{r}_0 \omega_1(\mathbf{r}) \mathfrak{Z}(\mathbf{x}) \} \\ & \quad + 2 \{ \Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x})/3 + 6\nu_0 \mathbf{r}_0 \omega_1(\mathbf{r}) \mathfrak{Z}(\mathbf{x}) \}^2. \end{aligned}$$

For the second bound of the statement we use the similar approach as in Theorem 2.3 in Spokoiny (2013).

$$\begin{aligned}
& \left| \sqrt{2 \left\{ L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}}) \right\}} - \|D_0(\tilde{\boldsymbol{\theta}}^\circ - \tilde{\boldsymbol{\theta}})\| \right| \\
& \leq \frac{\left| 2 \left\{ L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}}) \right\} - \|D_0(\tilde{\boldsymbol{\theta}}^\circ - \tilde{\boldsymbol{\theta}})\|^2 \right|}{\|D_0(\tilde{\boldsymbol{\theta}}^\circ - \tilde{\boldsymbol{\theta}})\|} \\
& = \frac{|2\mathcal{A}^\circ(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\theta}}^\circ)|}{\|D_0(\tilde{\boldsymbol{\theta}}^\circ - \tilde{\boldsymbol{\theta}})\|} \leq \sup_{\boldsymbol{\theta}, \boldsymbol{\theta}_1 \in \Theta_0(\mathbf{r}_0)} \frac{|2\mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)|}{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}_1)\|} \\
& \stackrel{\text{by (A.7)}}{\leq} \sup_{\boldsymbol{\theta}, \boldsymbol{\theta}_1 \in \Theta_0(\mathbf{r}_0)} 2 \|D_0^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{A}^\circ(\boldsymbol{\theta}, \boldsymbol{\theta}_1)\| \\
& \stackrel{\text{by (A.9), (A.10)}}{\leq} 4\Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x})/3 + 24\nu_0\mathbf{r}_0\omega_1(\mathbf{r})\mathfrak{J}(\mathbf{x}).
\end{aligned}$$

This together with (A.11) imply the final statement. \square

Lemma A.5 (Check of the bootstrap equivalent of (ED₂)). *Conditions (Eb), (L_{0m}) and (ED_{2m}) imply for each $\mathbf{r} > 0$, $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$, $\|\boldsymbol{\gamma}_j\| = 1$, $j = 1, 2$ and all $|\lambda| \leq \mathfrak{g}_2(\mathbf{r})$ with probability $\geq 1 - e^{-\mathbf{x}}$:*

$$\sup_{\substack{\boldsymbol{\gamma}_j \in \mathbb{R}^p \\ \|\boldsymbol{\gamma}_j\|=1}} \sum_{i=1}^n \log \mathbb{E}^\circ \exp \left\{ \frac{\lambda}{\omega_1(\mathbf{r})} \boldsymbol{\gamma}_1^\top D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}) D_0^{-1} \boldsymbol{\gamma}_2 (u_i - 1) \right\} \leq \frac{\lambda^2 \nu_0^2}{2}.$$

where

$$\omega_1(\mathbf{r}) = \omega_1 \stackrel{\text{def}}{=} \frac{\mathbf{C}_m(\mathbf{r})}{\sqrt{n}} + 2\omega\nu_0\sqrt{2\mathbf{x}} \tag{A.12}$$

In the case 4.3 it holds for $\mathbf{r} \leq \mathbf{r}_0$ $\omega_1(\mathbf{r}) = \mathbf{C}\mathbf{r}/n + \mathbf{C}\sqrt{\mathbf{x}/n}$.

Proof of Lemma A.5. Introduce the independent random scalar values for $i = 1, \dots, n$ and $j = 1, 2$:

$$\mu_i(\boldsymbol{\theta}, \boldsymbol{\gamma}_j) \stackrel{\text{def}}{=} \boldsymbol{\gamma}_1^\top D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}) D_0^{-1} \boldsymbol{\gamma}_2.$$

It holds

$$\begin{aligned}
& \sum_{i=1}^n \log \mathbb{E}^\circ \exp \left\{ \frac{\lambda}{\omega_1} \boldsymbol{\gamma}_1^\top D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}) D_0^{-1} \boldsymbol{\gamma}_2 (u_i - 1) \right\} \\
&= \sum_{i=1}^n \log \mathbb{E}^\circ \exp \left\{ \frac{\lambda}{\omega_1} \mu_i(\boldsymbol{\theta}, \boldsymbol{\gamma}_j) (u_i - 1) \right\} \\
&\leq \frac{\lambda^2 \nu_0^2}{2\omega_1^2} \sum_{i=1}^n \mu_i^2(\boldsymbol{\theta}, \boldsymbol{\gamma}_j), \tag{A.13}
\end{aligned}$$

here the inequality (A.13) follows from condition (Eb) if $|\lambda \mu_i(\boldsymbol{\theta}, \boldsymbol{\gamma}_j)| \leq g\omega_1$ for all $i = 1, \dots, n$, which is true due to the arguments below. Let us consider $\mu_i(\boldsymbol{\theta}, \boldsymbol{\gamma}_j)$, for each $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$, $i = 1, \dots, n$ it holds:

$$\begin{aligned}
|\mu_i(\boldsymbol{\theta}, \boldsymbol{\gamma}_j)| &\leq \|D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \mathbb{E} \ell_i(\boldsymbol{\theta}) D_0^{-1}\| \\
&\quad + \|D_0^{-1} \{ \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}}^2 \mathbb{E} \ell_i(\boldsymbol{\theta}) \} D_0^{-1}\|. \tag{A.14}
\end{aligned}$$

Condition (ED_{2m}), which is a stronger version of (ED₂), implies that for all $i = 1, \dots, n$, $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ and each $\mathbf{x} > 0$ it holds with \mathbb{P} -probability $\geq 1 - e^{-\mathbf{x}}$

$$\|D_0^{-1} \{ \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}}^2 \mathbb{E} \ell_i(\boldsymbol{\theta}) \} D_0^{-1}\| \leq 2\omega\nu_0 \left(\frac{2\mathbf{x}}{n} \right)^{1/2}. \tag{A.15}$$

Indeed, by the exponential Chebyshev inequality for $\lambda > 0$

$$\begin{aligned}
& \mathbb{P} \left(\omega^{-1} \|D_0^{-1} \{ \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}}^2 \mathbb{E} \ell_i(\boldsymbol{\theta}) \} D_0^{-1}\| \geq t \right) \\
&\leq \mathbb{E} \exp \left[-\lambda t + \omega^{-1} \lambda \|D_0^{-1} \{ \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}}^2 \mathbb{E} \ell_i(\boldsymbol{\theta}) \} D_0^{-1}\| \right] \\
&\stackrel{\text{by (ED}_{2m}\text{)}}{\leq} \exp \left\{ -\lambda t + \lambda^2 \nu_0^2 / (2n) \right\}, \quad 0 < \lambda < g_2(\mathbf{r}) \\
&\leq \exp \{-\mathbf{x}\},
\end{aligned}$$

here the last inequality holds under the assumption, that $g_2(\mathbf{r})$ is large enough. In the case 4.3 it holds $g_2(\mathbf{r}) = Cn^{1/2}$, $\omega = Cn^{-1/2}$ and $\mathbf{x} = C \log(n)$; $t^2 := 8\nu_0^2 \mathbf{x} / n$ implies $\lambda t - \lambda^2 \nu_0^2 / (2n) - \mathbf{x} \geq 0$ for $0 < \lambda < g_2(\mathbf{r})$. For the deterministic term in (A.14) condition (L_{0m}) reads as:

$$\|D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \mathbb{E} \ell_i(\boldsymbol{\theta}) D_0^{-1}\| \leq \frac{\mathbf{C}_m(\mathbf{r})}{n}. \tag{A.16}$$

Collecting the inequalities (A.13), (A.14), (A.15) and (A.16), we obtain:

$$\begin{aligned}
& \sum_{i=1}^n \log \mathbb{E}^\circ \exp \left\{ \frac{\lambda}{\omega_1} \boldsymbol{\gamma}_1^\top D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}) D_0^{-1} \boldsymbol{\gamma}_2 (u_i - 1) \right\} \\
&\leq \frac{\lambda^2 \nu_0^2}{2} \frac{1}{\omega_1^2} \left\{ \frac{\mathbf{C}_m(\mathbf{r})}{\sqrt{n}} + 2\omega\nu_0 \sqrt{2\mathbf{x}} \right\}^2
\end{aligned}$$

Taking $\omega_1 = \omega_1(\mathbf{r})$ as in (A.12) implies the necessary statement. \square

Theorem A.6 (Concentration of bootstrap MLE). *Let the conditions of Theorems A.1 and A.8, (\mathcal{L}_{0m}) and (ED_{2m}) be fulfilled. If the following holds for $\omega_1(\mathbf{r})$ defined in (A.12) and the \mathbb{P} -random matrix $\mathcal{B}^2 \stackrel{\text{def}}{=} D_0^{-1} \text{Var}^\circ \{ \nabla_\theta L^\circ(\theta^*) \} D_0^{-1}$*

$$\begin{aligned} \mathfrak{b}(\mathbf{r})\mathbf{r} &\geq 2 \{ \mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathbb{B}) + \mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathcal{B}) + 6\nu_0 \mathfrak{Z}(\mathbf{x}) \omega_1(\mathbf{r}_0) \mathbf{r}_0 \} \\ &\quad + 12\nu_0(\omega + \omega_1(\mathbf{r})) \mathfrak{Z}(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0)) \quad \text{for } \mathbf{r} > \mathbf{r}_0, \end{aligned} \quad (\text{A.17})$$

then it holds with \mathbb{P} -probability $\geq 1 - 3e^{-x}$

$$\mathbb{P}^\circ \left(\tilde{\theta}^\circ \notin \Theta_0(\mathbf{r}_0) \right) \leq 3e^{-x}.$$

Proof of Theorem A.6. We use the idea by Spokoiny (2013): if $\sup_{\theta \in \Theta \setminus \Theta_0(\mathbf{r}_0)} \{ L(\theta) - L(\theta^*) \} < 0$, then $\tilde{\theta} \in \Theta_0(\mathbf{r}_0)$. We apply it here for the the bootstrap objects: $L^\circ(\theta) - L^\circ(\tilde{\theta})$ and $\tilde{\theta}^\circ$. Denote the stochastic part of the bootstrap likelihood process as $\zeta^\circ(\theta) \stackrel{\text{def}}{=} L^\circ(\theta) - \mathbb{E}^\circ L^\circ(\theta)$. It holds

$$\begin{aligned} L^\circ(\theta) - L^\circ(\tilde{\theta}) &= \zeta^\circ(\theta) - \zeta^\circ(\tilde{\theta}) + \mathbb{E}^\circ L^\circ(\theta) - \mathbb{E}^\circ L^\circ(\tilde{\theta}) \\ &= \zeta^\circ(\theta) - \zeta^\circ(\tilde{\theta}) + L(\theta) - L(\tilde{\theta}) \\ &= \{ \zeta^\circ(\theta) - \zeta^\circ(\tilde{\theta}) \} + \{ L(\theta) - L(\theta^*) \} + \{ L(\theta^*) - L(\tilde{\theta}) \}. \end{aligned}$$

Here the last summand $\{ L(\theta^*) - L(\tilde{\theta}) \}$ is non-positive by definition (1.2) of $\tilde{\theta}$. The following bound follows from the proof of Theorem 2.1 in Spokoiny (2013):

$$\begin{aligned} \mathbb{P} \left(\sup_{\theta \in \Theta \setminus \Theta_0(\mathbf{r}_0)} \{ L(\theta) - L(\theta^*) \} < \varrho(\mathbf{r}, \mathbf{x})\mathbf{r} + \mathbf{r} \mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathbb{B}) - \mathbf{r}^2 \mathfrak{b}(\mathbf{r})/2 \right) &\geq 1 - 3e^{-x}, \\ \varrho(\mathbf{r}, \mathbf{x}) &\stackrel{\text{def}}{=} 6\nu_0 \mathfrak{Z}(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0))\omega. \end{aligned}$$

Due to Lemma A.5 the process $\zeta^\circ(\theta) - \zeta^\circ(\tilde{\theta})$ satisfies the necessary conditions of Theorem A.1 in Spokoiny (2013), and it holds for $\mathbf{r} \geq \mathbf{r}_0$

$$\begin{aligned} \mathbb{P}^\circ \left(\sup_{\theta \in \Theta_0(\mathbf{r})} \left| \zeta^\circ(\theta) - \zeta^\circ(\tilde{\theta}) - (\theta - \tilde{\theta})^\top \nabla_\theta \zeta^\circ(\tilde{\theta}) \right| \leq \varrho_1(\mathbf{r}, \mathbf{x})\mathbf{r} \right) &\geq 1 - e^{-x}, \\ \varrho_1(\mathbf{r}, \mathbf{x}) &\stackrel{\text{def}}{=} 6\nu_0 \mathfrak{Z}(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0))\omega_1(\mathbf{r}). \end{aligned}$$

By Lemma A.7 and Theorem A.8 it holds with dominating probability

$$\begin{aligned} \sup_{\theta \in \Theta_0(\mathbf{r})} \left| (\theta - \tilde{\theta})^\top \nabla_\theta \zeta^\circ(\tilde{\theta}) \right| &\leq \mathbf{r} \|\xi^\circ(\tilde{\theta})\| \\ &\leq \mathbf{r} \left\{ \|\xi^\circ(\theta^*)\| + \|\xi^\circ(\tilde{\theta}) - \xi^\circ(\theta^*)\| \right\} \\ &\leq \mathbf{r} \{ \mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathcal{B}) + 6\nu_0 \mathfrak{Z}(\mathbf{x}) \omega_1(\mathbf{r}_0) \mathbf{r}_0 \}. \end{aligned}$$

Finally we have:

$$\begin{aligned}
\sup_{\boldsymbol{\theta} \in \Theta \setminus \Theta_0(r_0)} \left\{ L^\circ(\boldsymbol{\theta}) - L^\circ(\tilde{\boldsymbol{\theta}}) \right\} &\leq \sup_{\boldsymbol{\theta} \in \Theta \setminus \Theta_0(r_0)} \left\{ L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*) \right\} \\
&\quad + \sup_{\substack{\boldsymbol{\theta} \in \Theta_0(r), \\ r \geq r_0}} \left\{ \zeta^\circ(\boldsymbol{\theta}) - \zeta^\circ(\tilde{\boldsymbol{\theta}}) \right\} \\
&\leq r \mathfrak{J}_{\text{qf}}(\mathbf{x}, \mathcal{B}) + r \mathfrak{J}_{\text{qf}}(\mathbf{x}, \mathcal{B}) + \varrho_1(r, \mathbf{x})r \\
&\quad + \varrho(r, \mathbf{x})r - r^2 \mathfrak{b}(r)/2 + 6\nu_0 \mathfrak{J}(\mathbf{x}) \omega_1(r_0) r r_0,
\end{aligned}$$

which implies the condition (A.17) in the statement. \square

REMARK A.1. Condition (A.17) imposed for the bootstrap MLE concentration result is stronger, than condition (A.1) for the concentration of \mathbf{Y} - MLE, and (A.17) implies the latter one.

The following lemma had already been derived in the proof of Theorem A.4: see the bound (A.10). We formulate it separately, since it is used again in another statements.

Lemma A.7. *Let the conditions of Lemma A.5 be fulfilled, then it holds with \mathbb{P} -probability $\geq 1 - e^{-x}$*

$$\mathbb{P}^\circ \left(\sup_{\boldsymbol{\theta} \in \Theta_0(r)} \|\boldsymbol{\xi}^\circ(\boldsymbol{\theta}) - \boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\| \leq \Delta_\xi^\circ(r, \mathbf{x}) \right) \geq 1 - e^{-x},$$

where

$$\Delta_\xi^\circ(r, \mathbf{x}) \stackrel{\text{def}}{=} 6\nu_0 \mathfrak{J}(\mathbf{x}) \omega_1(r)r$$

In the case 4.3 it holds for the bounding term.

$$\Delta_\xi^\circ(r_0, \mathbf{x}) \leq C \frac{p + \mathbf{x}}{\sqrt{n}} \sqrt{\mathbf{x}}.$$

Theorem A.8 (Deviation bound for the bootstrap quadratic form). *Let conditions (Eb), (I), (SD₁), (I_B) be fulfilled, then for $g \geq \sqrt{2 \text{tr}(\mathcal{B}^2)}$ it holds:*

$$\mathbb{P}^\circ \left(\|\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\|^2 \leq \mathfrak{J}_{\text{qf}}^2(\mathbf{x}, \mathcal{B}) \right) \geq 1 - 2e^{-x} - 8.4e^{-x_c(\mathcal{B})},$$

where

$$\mathcal{B}^2 \stackrel{\text{def}}{=} D_0^{-1} \mathcal{V}^2(\boldsymbol{\theta}^*) D_0^{-1}, \quad \mathcal{V}^2(\boldsymbol{\theta}^*) \stackrel{\text{def}}{=} \text{Var}^\circ \nabla_{\boldsymbol{\theta}} L^\circ(\boldsymbol{\theta}^*), \tag{A.18}$$

$\mathfrak{J}_{\text{qf}}(\mathbf{x}, \cdot)$ and $x_c(\cdot)$ are defined respectively in (A.4) and (A.5). Similarly to (A.6) it holds for $\mathbf{x} \leq x_c(\mathcal{B})$:

$$\mathfrak{J}_{\text{qf}}^2(\mathbf{x}, \mathcal{B}) \leq \mathbf{a}^{\circ 2} (p + 6\mathbf{x}) \tag{A.19}$$

$$\text{for } \mathbf{a}^{\circ 2} \stackrel{\text{def}}{=} (1 + \delta_{\mathcal{V}}^2) (\mathbf{a}^2 + \mathbf{a}_B^2).$$

Proof of Theorem A.8. This result is the bootstrap equivalent of Theorem A.3. For the \mathbf{Y} -world it demands condition (ED_0) to be fulfilled. Let us check whether the bootstrap equivalent of (ED_0) holds. It reads as follows: *there exist constants $g^\circ > 0$, $\nu_0^\circ \geq 1$ such that for the positive-definite symmetric matrix $\mathcal{V}^2(\boldsymbol{\theta}^*)$ it holds for all $|\lambda| \leq g^\circ$*

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \log \mathbb{E}^\circ \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \{ \nabla_{\boldsymbol{\theta}} L^\circ(\boldsymbol{\theta}^*) - \nabla_{\boldsymbol{\theta}} \mathbb{E}^\circ L^\circ(\boldsymbol{\theta}^*) \}}{\| \mathcal{V}(\boldsymbol{\theta}^*) \boldsymbol{\gamma} \|} \right\} \leq \nu_0^{\circ 2} \lambda^2 / 2.$$

By definition $\mathcal{V}^2(\boldsymbol{\theta}^*) = \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)^\top$. Let us introduce the independent \mathbb{P} -random variables $s_i(\boldsymbol{\gamma}) \stackrel{\text{def}}{=} \boldsymbol{\gamma}^\top \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*) / \| \mathcal{V}(\boldsymbol{\theta}^*) \boldsymbol{\gamma} \|$ for $i = 1, \dots, n$. It holds $\sum_{i=1}^n s_i^2(\boldsymbol{\gamma}) = 1$, hence $\max_{1 \leq i \leq n} |s_i| \leq 1$. Condition (Eb) implies:

$$\begin{aligned} & \log \mathbb{E}^\circ \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \{ \nabla_{\boldsymbol{\theta}} L^\circ(\boldsymbol{\theta}^*) - \nabla_{\boldsymbol{\theta}} \mathbb{E}^\circ L^\circ(\boldsymbol{\theta}^*) \}}{\| \mathcal{V}(\boldsymbol{\theta}^*) \boldsymbol{\gamma} \|} \right\} \\ &= \sum_{i=1}^n \log \mathbb{E}^\circ \exp \{ \lambda s_i(\boldsymbol{\gamma})(u_i - 1) \} \\ &\leq \frac{\nu_0^{\circ 2} \lambda^2}{2} \sum_{i=1}^n s_i^2(\boldsymbol{\gamma}) = \nu_0^{\circ 2} \lambda^2 / 2, \quad |\lambda| \leq g. \end{aligned}$$

Thus the bootstrap equivalent for the condition (ED_0) is fulfilled with the same constants ν_0, g , and the theorem's statements holds as well as for Theorem A.3.

The inequality (A.19) follows from conditions (\mathcal{I}) , (\mathcal{I}_B) , (SD_1) and Bernstein matrix inequality by Tropp (2012) (see Section A.6):

$$\| D_0^{-1} \mathcal{V}_0^2(\boldsymbol{\theta}^*) D_0^{-1} \| \leq \| D_0^{-1} H_0 \|^2 (1 + \delta_V^2) \leq (1 + \delta_V^2) (\boldsymbol{\alpha}^2 + \boldsymbol{\alpha}_B^2).$$

□

A.3 Proofs of Theorems 2.1 and 2.3

In order to justify theoretically the multiplier bootstrap procedure it has to be shown that the approximating terms $\| \boldsymbol{\xi} \|$ and $\| \boldsymbol{\xi}^\circ(\tilde{\boldsymbol{\theta}}) \|$ from the Wilks Theorems A.2 and A.4 have nearly the same distributions. By Lemma A.7 the random values $\| \boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*) \|$ and $\| \boldsymbol{\xi}^\circ(\tilde{\boldsymbol{\theta}}) \|$ are close to each other within the error term $\leq C(p+x) \sqrt{x/n}$ with exponentially high probability, therefore, it is sufficient to compare the distributions of $\| \boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*) \|$ and $\| \boldsymbol{\xi} \|$. This is done in Proposition A.9 using the results on Gaussian approximation for Euclidean norms from Section 5.

Let us introduce the multivariate normal vectors similarly to (5.3):

$$\bar{\boldsymbol{\xi}} \sim \mathcal{N}(0, \text{Var } \boldsymbol{\xi}), \quad \bar{\boldsymbol{\xi}}^\circ(\boldsymbol{\theta}^*) \sim \mathcal{N}(0, \text{Var}^\circ \{ \boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*) \}). \quad (\text{A.20})$$

Let us also represent the vectors $\boldsymbol{\xi}$ and $\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)$ as sums of the marginal score vectors $\boldsymbol{\xi}_i$ and

$\xi_i^\circ(\theta^*)$ s.t. $\mathbb{E}\xi_i = \mathbb{E}^\circ\xi_i^\circ = 0$:

$$\begin{aligned}\xi_i &\stackrel{\text{def}}{=} D_0^{-1} \{ \nabla_{\theta} \ell_i(\theta^*) - \nabla_{\theta} \mathbb{E} \ell_i(\theta^*) \}, \\ \xi_i^\circ(\theta^*) &\stackrel{\text{def}}{=} D_0^{-1} \nabla_{\theta} \ell_i(\theta^*) \{u_i - 1\}.\end{aligned}$$

Their Gaussian analogs are

$$\bar{\xi}_i \sim \mathcal{N}(0, \text{Var } \xi_i) \quad \text{and} \quad \bar{\xi}_i^\circ \sim \mathcal{N}(0, \text{Var}^\circ \{ \xi_i^\circ(\theta^*) \}).$$

Similarly to (5.4) denote

$$\begin{aligned}\delta_n &\stackrel{\text{def}}{=} \frac{1}{2} \sum_{i=1}^n \mathbb{E} (\|\xi_i\|^3 + \|\bar{\xi}_i\|^3), \\ \check{\delta}_n &\stackrel{\text{def}}{=} \frac{1}{2} \sum_{i=1}^n \mathbb{E}^\circ (\|\xi_i^\circ(\theta^*)\|^3 + \|\bar{\xi}_i^\circ(\theta^*)\|^3).\end{aligned}\tag{A.21}$$

Proposition A.9 (Closeness of the c.d.f. of $\|\xi\|$ and $\|\xi^\circ(\theta^*)\|$). *If conditions **(SmB)** and **(SD₁)** are fulfilled, then it holds with probability $\geq 1 - e^{-x}$ for all $0 < \Delta \leq 0.22$ and for all $z, \bar{z} > 2$ s.t. $|z - \bar{z}| \leq \delta_z$ for some $\delta_z \geq 0$:*

$$\begin{aligned}& \left| \mathbb{P} (\|\xi\| \geq z) - \mathbb{P}^\circ (\|\xi^\circ(\theta^*)\| \geq \bar{z}) \right| \\ & \leq 16\Delta^{-3} (\delta_n + \check{\delta}_n) + \frac{2\Delta + \delta_z}{\bar{z}} \sqrt{\frac{p}{2}} + \frac{\sqrt{p}}{2} \frac{\delta_{\mathcal{V}}^2(\mathbf{x}) + \delta_{\text{smb}}^2}{1 - \delta_{\mathcal{V}}^2(\mathbf{x})} \\ & \leq 16\Delta^{-3} (\delta_n + \check{\delta}_n) + \frac{2\Delta + \delta_z}{\sqrt{2}} + \frac{2\sqrt{p}}{3} (\delta_{\mathcal{V}}^2(\mathbf{x}) + \delta_{\text{smb}}^2) \\ & \quad \text{for } \bar{z} \geq \sqrt{p}, \delta_{\mathcal{V}}^2(\mathbf{x}) \leq 1/4.\end{aligned}$$

Moreover, if $z, \bar{z} \geq \max\{2, \sqrt{p}\}$ and $\max\{\delta_n^{1/4}, \check{\delta}_n^{1/4}\} \leq 0.11$, then

$$\begin{aligned}& \left| \mathbb{P} (\|\xi\| \geq z) - \mathbb{P}^\circ (\|\xi^\circ(\theta^*)\| \geq \bar{z}) \right| \\ & \leq 1.55 (\delta_n^{1/4} + \check{\delta}_n^{1/4}) + \frac{\delta_z}{\sqrt{2}} + \frac{2\sqrt{p}}{3} (\delta_{\mathcal{V}}^2(\mathbf{x}) + \delta_{\text{smb}}^2).\end{aligned}\tag{A.22}$$

Proof of Proposition A.9. We use Theorem 5.1 taking $\phi := \xi$ and $\psi := \xi^\circ(\theta^*)$. Let us check that the conditions (5.5) on the covariance matrices are fulfilled. By definitions (1.7), (1.8) and (A.18)

$$\begin{aligned}\text{Var } \xi &= D_0^{-1} H_0^2 D_0^{-1} - D_0^{-1} B_0^2 D_0^{-1}, \\ \text{Var}^\circ \{ \xi^\circ(\theta^*) \} &= D_0^{-1} \mathcal{V}^2(\theta^*) D_0^{-1}.\end{aligned}$$

Due to Theorem A.13 by Tropp (2012) (see Section A.6) it holds with probability $\geq 1 - e^{-x}$

$$\|H_0^{-1} \mathcal{V}^2(\theta^*) H_0^{-1} - \mathbf{I}_p\| \leq \delta_{\mathcal{V}}^2(\mathbf{x}),\tag{A.23}$$

therefore, by Cauchy-Schwarz inequality

$$\|\mathcal{V}^{-1}(\boldsymbol{\theta}^*)H_0^2\mathcal{V}^{-1}(\boldsymbol{\theta}^*) - \mathbf{I}_p\| \leq \delta_{\mathcal{V}}^2(\mathbf{x})(1 - \delta_{\mathcal{V}}^2(\mathbf{x}))^{-1}.$$

Condition **(SmB)** says that $\|H_0^{-1}B_0^2H_0^{-1}\| \leq \delta_{\text{smb}}^2$, therefore, by the triangle inequality it holds:

$$\begin{aligned} \left\| [\text{Var}^\circ\{\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\}]^{-1/2} \text{Var}\{\boldsymbol{\xi}\} [\text{Var}^\circ\{\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\}]^{-1/2} - \mathbf{I}_p \right\| &\leq \frac{\delta_{\mathcal{V}}^2(\mathbf{x}) + \delta_{\text{smb}}^2}{1 - \delta_{\mathcal{V}}^2(\mathbf{x})} \\ &\leq 1/2 \\ &\text{for } \delta_{\text{smb}}^2 \leq 1/8, \delta_{\mathcal{V}}^2(\mathbf{x}) \leq 1/4. \end{aligned}$$

□

Now we are ready to collect all the obtained bounds together for the following

Proof of Theorem 2.1. On a random set of probability $\geq 1 - 12e^{-x}$ it holds:

$$\begin{aligned} \alpha &\stackrel{\text{(Def. (2.3))}}{=} \mathbb{P}^\circ \left(\sqrt{2\{L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}})\}} > \mathfrak{z}_\alpha^\circ \right) \\ &\stackrel{\text{(Th. A.4)}}{\geq} \mathbb{P}^\circ \left(\|\boldsymbol{\xi}^\circ(\tilde{\boldsymbol{\theta}})\| > \mathfrak{z}_\alpha^\circ + \Delta_{\mathbf{W}}^\circ(\mathbf{r}_0, \mathbf{x}) \right) \\ &\stackrel{\text{(L. A.7)}}{\geq} \mathbb{P}^\circ \left(\|\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\| > \mathfrak{z}_\alpha^\circ + \Delta_{\mathbf{W}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_\xi^\circ(\mathbf{r}_0, \mathbf{x}) \right) \end{aligned} \quad (\text{A.24})$$

$$\begin{aligned} &\stackrel{\text{(Prop. A.9)}}{\geq} \mathbb{P} \left(\|\boldsymbol{\xi}\| > \mathfrak{z}_\alpha^\circ - \Delta_{\mathbf{W}}(\mathbf{r}_0, \mathbf{x}) \right) - \Delta_{\text{full}} \\ &\stackrel{\text{(Th. A.2)}}{\geq} \mathbb{P} \left(\sqrt{2\{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\}} > \mathfrak{z}_\alpha^\circ \right) - \Delta_{\text{full}}, \end{aligned} \quad (\text{A.25})$$

where the value Δ_{full} comes from the bound (A.22) with $\delta_z := \Delta_{\mathbf{W}}(\mathbf{r}_0, \mathbf{x}) + \Delta_{\mathbf{W}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_\xi^\circ(\mathbf{r}_0, \mathbf{x})$:

$$\begin{aligned} \Delta_{\text{full}} &\stackrel{\text{def}}{=} 1.55(\delta_n^{1/4} + \check{\delta}_n^{1/4}) + \frac{2\sqrt{p}}{3} (\delta_{\mathcal{V}}^2(\mathbf{x}) + \delta_{\text{smb}}^2) \\ &\quad + \{\Delta_{\mathbf{W}}(\mathbf{r}_0, \mathbf{x}) + \Delta_{\mathbf{W}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_\xi^\circ(\mathbf{r}_0, \mathbf{x})\} / \sqrt{2} \end{aligned} \quad (\text{A.26})$$

By the similar arguments in the inverse direction we obtain the following inequality:

$$\left| \mathbb{P} \left(\sqrt{2\{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\}} > \mathfrak{z}_\alpha^\circ \right) - \alpha \right| \leq \Delta_{\text{full}}.$$

Notice that inequality (A.22) from Proposition A.9, that we use here, requires $\max\{\delta_n^{1/4}, \check{\delta}_n^{1/4}\} \leq 0.11$.

Let us quantify, how the error term Δ_{full} depends on p and n . In the case 4.3 random vectors ξ_i and $\xi_i^\circ(\theta^*)$ satisfy the conditions of Theorems A.3 and A.8 correspondingly. Hence $\|\xi_i\|, \|\xi_i^\circ(\theta^*)\| \leq C\sqrt{(p+x)/n}$ and $\delta_n, \check{\delta}_n \leq C\sqrt{(p+x)^3/n}$. Finally we have in the case 4.3

$$\Delta_{full} = C \left\{ \frac{(p+x)^3}{n} \right\}^{1/8} + C \frac{p+x}{\sqrt{n}} \sqrt{x} + C \frac{p+x}{\sqrt{n}}. \quad (\text{A.27})$$

□

REMARK A.2. It is clear from expression (A.27), that the impact of the error term, induced by the Gaussian approximation, is the biggest. The requirement for the ratio $(p+x)^3/n$ to be small is imposed by our Gaussian approximation results (see also Remark 5.2 about the multivariate GAR).

Let us introduce for $p = 1$ similarly to (5.6) and (A.21)

$$\delta_{n,\text{B.E.}} \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} |\xi_i|^3, \quad \check{\delta}_{n,\text{B.E.}} \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E}^\circ |\xi_i^\circ(\theta^*)|^3.$$

Proof of Theorem 2.3. On a random set of probability $\geq 1 - 12e^{-x}$ it holds:

$$\begin{aligned} \alpha &\stackrel{(\text{Def. (2.3)})}{=} \mathbb{P}^\circ \left(\sqrt{2 \{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}} > \mathfrak{z}_\alpha^\circ \right) \\ &\stackrel{(\text{Th. A.4})}{\geq} \mathbb{P}^\circ \left(\|\xi^\circ(\tilde{\theta})\| > \mathfrak{z}_\alpha^\circ + \Delta_{\mathbb{W}}^\circ(\mathbf{r}_0, \mathbf{x}) \right) \\ &\stackrel{(\text{L. A.7})}{\geq} \mathbb{P}^\circ \left(\|\xi^\circ(\theta^*)\| > \mathfrak{z}_\alpha^\circ + \Delta_{\mathbb{W}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_{\xi}^\circ(\mathbf{r}_0, \mathbf{x}) \right) \\ &\stackrel{(\text{L. 5.2, Prop. A.9})}{\geq} \mathbb{P} \left(\|\xi\| > \mathfrak{z}_\alpha^\circ - \Delta_{\mathbb{W}}(\mathbf{r}_0, \mathbf{x}) - \Delta_{\text{B.E., full}} \right) \\ &\stackrel{(\text{Th. A.2})}{\geq} \mathbb{P} \left(\sqrt{2 \{L(\tilde{\theta}) - L(\theta^*)\}} > \mathfrak{z}_\alpha^\circ \right) - \Delta_{\text{B.E., full}}, \end{aligned}$$

where the value $\Delta_{\text{B.E., full}}$ comes from the bound (5.8) with $\delta_z := \Delta_{\mathbb{W}}(\mathbf{r}_0, \mathbf{x}) + \Delta_{\mathbb{W}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_{\xi}^\circ(\mathbf{r}_0, \mathbf{x})$, $C_0 \in [0.4097, 0.560]$ and

$$\begin{aligned} \text{Var}^\circ \{ \xi^\circ(\theta^*) \} &\geq \{1 - \delta_V^2(\mathbf{x})\} \mathbb{E} \text{Var}^\circ \{ \xi^\circ(\theta^*) \} \\ &\geq \frac{3}{4} D_0^{-1} H_0^2 D_0^{-1} \quad \text{for } \delta_V^2(\mathbf{x}) \leq 1/4 \end{aligned}$$

with probability $\geq 1 - e^{-x}$ (due to the bound (A.23)):

$$\begin{aligned} \Delta_{\text{B.E., full}} &\stackrel{\text{def}}{=} 2C_0 \left\{ \frac{\delta_{n,\text{B.E.}}}{(\text{Var } \boldsymbol{\xi})^{3/2}} + \frac{\check{\delta}_{n,\text{B.E.}}}{(\mathbb{E} \text{Var}^\circ\{\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\})^{3/2}} \left(\frac{2}{\sqrt{3}}\right)^3 \right\} \\ &\quad + \frac{1}{\sqrt{2}} \left\{ \Delta_{\text{W}}(\mathbf{r}_0, \mathbf{x}) + \Delta_{\text{W}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_{\xi}^\circ(\mathbf{r}_0, \mathbf{x}) \right\} + \frac{2}{3} \left\{ \delta_{\mathcal{V}}^2(\mathbf{x}) + \delta_{\xi}^2 \right\} \\ &\leq C \frac{1+x}{\sqrt{n}} \quad \text{in the case 4.3.} \end{aligned} \tag{A.28}$$

The similar inequalities in the inverse direction finish the proof with the error term \square

A.4 Proof of Theorem 2.4 (large modeling bias)

Lemma A.10 (Lower bound for deviations of a Gaussian quadratic form). *Let $\boldsymbol{\phi} \sim \mathcal{N}(0, \mathbf{I}_p)$ and Σ is any symmetric non-negative definite matrix, then it holds for any $x > 0$*

$$\mathbb{P} \left(\text{tr } \Sigma - \|\Sigma^{1/2} \boldsymbol{\phi}\|^2 \geq 2\sqrt{x \text{tr}(\Sigma^2)} \right) \leq \exp(-x).$$

Proof of Lemma A.10. It is sufficient to consider w.l.o.g. only the case of diagonal matrix Σ , since it can be represented as $\Sigma = U^\top \text{diag}\{a_1, \dots, a_p\}U$ for an orthogonal matrix U and the eigenvalues $a_1 \geq \dots \geq a_p$; $U\boldsymbol{\phi} \sim \mathcal{N}(0, \mathbf{I}_p)$.

By the exponential Chebyshev inequality it holds for $\mu > 0$, $\Delta > 0$

$$\begin{aligned} \mathbb{P} \left(\text{tr } \Sigma - \|\Sigma^{1/2} \boldsymbol{\phi}\|^2 \geq \Delta \right) &\leq \exp(-\mu\Delta/2) \mathbb{E} \exp \left(\mu \left\{ \text{tr } \Sigma - \|\Sigma^{1/2} \boldsymbol{\phi}\|^2 \right\} / 2 \right). \\ \log \mathbb{E} \exp \left(\mu \left\{ \text{tr } \Sigma - \|\Sigma^{1/2} \boldsymbol{\phi}\|^2 \right\} / 2 \right) &\leq \frac{1}{2} \sum_{j=1}^p \left\{ \mu a_j - \log(1 + a_j \mu) \right\}, \end{aligned}$$

therefore

$$\begin{aligned} \mathbb{P} \left(\text{tr } \Sigma - \|\Sigma^{1/2} \boldsymbol{\phi}\|^2 \geq \Delta \right) &\leq \exp \left(-\frac{1}{2} \left[\mu\Delta + \sum_{j=1}^p \left\{ \log(1 + a_j \mu) - \mu a_j \right\} \right] \right) \\ &\leq \exp \left(-\frac{1}{2} \left[\mu\Delta - \mu^2 \sum_{j=1}^p a_j^2 / 2 \right] \right) \\ &\leq \exp \left(-\Delta^2 / \left\{ 4 \sum_{j=1}^p a_j^2 \right\} \right). \end{aligned}$$

If $x := \Delta^2 / \left\{ 4 \sum_{j=1}^p a_j^2 \right\}$, then $\Delta = 2\sqrt{x \sum_{j=1}^p a_j^2}$. \square

Proof of Theorem 2.4. Due to the bound (A.24) it holds for $\mathfrak{z} \geq \max\{2, \sqrt{p}\} + \mathsf{C}(p + \mathbf{x})/\sqrt{n}$ with probability $\geq 1 - 5e^{-x}$

$$\begin{aligned} & \mathbb{P}^\circ \left(\sqrt{2 \{L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}})\}} > \mathfrak{z} \right) \\ & \geq \mathbb{P}^\circ \left(\|\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\| > \mathfrak{z} + \Delta_{\mathbf{w}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_{\boldsymbol{\xi}}^\circ(\mathbf{r}_0, \mathbf{x}) \right). \end{aligned}$$

Let us introduce the random vector $\boldsymbol{\xi}_0 \stackrel{\text{def}}{=} (D_0^{-1} H_0^2 D_0^{-1})^{1/2} (\text{Var } \boldsymbol{\xi})^{-1/2} \boldsymbol{\xi}$. The bound (A.23) implies with probability $\geq 1 - e^{-x}$

$$\text{tr} \left\{ \left((\text{Var } \boldsymbol{\xi}_0)^{-1/2} \text{Var}^\circ \{ \boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*) \} (\text{Var } \boldsymbol{\xi}_0)^{-1/2} - \mathbf{I}_p \right)^2 \right\} \leq p \delta_{\mathcal{V}}^4(\mathbf{x}). \quad (\text{A.29})$$

Applying statement 2.2 of Theorem 5.1 to the vectors $\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)$ and $\boldsymbol{\xi}_0$, we have with probability $\geq 1 - e^{-x}$

$$\begin{aligned} & \mathbb{P}^\circ \left(\|\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\| > \mathfrak{z} + \Delta_{\mathbf{w}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_{\boldsymbol{\xi}}^\circ(\mathbf{r}_0, \mathbf{x}) \right) \\ & \geq \mathbb{P} \left(\|\boldsymbol{\xi}_0\| > \mathfrak{z} - \Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x}) \right) - \Delta_{\mathbf{b}, \text{full}} \end{aligned}$$

where

$$\begin{aligned} \Delta_{\mathbf{b}, \text{full}} & \stackrel{\text{def}}{=} 1.55 \left(\delta_n^{1/4} + \check{\delta}_n^{1/4} \right) + \frac{\sqrt{p}}{2} \delta_{\mathcal{V}}^2(\mathbf{x}) \\ & + \frac{\Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x}) + \Delta_{\mathbf{w}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_{\boldsymbol{\xi}}^\circ(\mathbf{r}_0, \mathbf{x})}{\sqrt{2}}. \end{aligned} \quad (\text{A.30})$$

By the definition of $\boldsymbol{\xi}_0$ it holds $\|\boldsymbol{\xi}_0\| \geq \|\boldsymbol{\xi}\| \| (D_0^{-1} H_0^2 D_0^{-1})^{-1/2} \|^{-1}$. Consider the following matrix

$$\begin{aligned} \tilde{V}^2 & \stackrel{\text{def}}{=} (D_0^{-1} H_0^2 D_0^{-1})^{-1/2} (\text{Var } \boldsymbol{\xi}) (D_0^{-1} H_0^2 D_0^{-1})^{-1/2} \\ & = (D_0^{-1} H_0^2 D_0^{-1})^{1/2} (D_0 H_0^{-2} V_0^2 H_0^{-2} D_0) (D_0^{-1} H_0^2 D_0^{-1})^{1/2} \\ & \leq (D_0^{-1} H_0^2 D_0^{-1})^{1/2} (D_0 H_0^{-2} D_0) (D_0^{-1} H_0^2 D_0^{-1})^{1/2} \\ & = \mathbf{I}_p, \end{aligned} \quad (\text{A.31})$$

here $V_0^2 \stackrel{\text{def}}{=} \text{Var} \{ \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*) \}$; the inequality (A.31) holds due to the definitions (1.7), (1.8) and $V_0^2 = H_0^2 - B_0^2 > 0$. Therefore $\|\tilde{V}^2\| \leq 1$ and $\|\boldsymbol{\xi}_0\| \geq \|\boldsymbol{\xi}\|$. By (A.25)

$$\begin{aligned} & \mathbb{P}^\circ \left(\sqrt{2 \{L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}})\}} > \mathfrak{z} \right) \\ & \geq \mathbb{P} \left(\|\boldsymbol{\xi}\| > \mathfrak{z} - \Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x}) \right) - \Delta_{\mathbf{b}, \text{full}} \\ & \geq \mathbb{P} \left(\sqrt{2 \{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\}} > \mathfrak{z} \right) - \Delta_{\mathbf{b}, \text{full}} \end{aligned}$$

with probability $\geq 1 - 12e^{-x}$, which finishes the proof of the first part. For the second part let us introduce $\bar{\xi}_0 \sim \mathcal{N}(0, D_0^{-1}H_0^2D_0^{-1})$ s.t. $\text{Var } \bar{\xi}_0 = \text{Var } \xi_0$. Applying statement 2.1 of Theorem 5.1 to the vectors $\xi^\circ(\theta^*)$ and $\bar{\xi}_0$, using the bound (A.29), we have with probability $\geq 1 - e^{-x}$

$$\begin{aligned} P^\circ(\|\xi^\circ(\theta^*)\| > \mathfrak{z} + \Delta_w^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_\xi^\circ(\mathbf{r}_0, \mathbf{x})) \\ \geq P(\|\bar{\xi}_0\| > \mathfrak{z}) - \Delta_{G,1}, \end{aligned}$$

where

$$\Delta_{G,1} \stackrel{\text{def}}{=} 1.55\delta_n^{1/4} + \frac{\Delta_w^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_\xi^\circ(\mathbf{r}_0, \mathbf{x})}{\sqrt{2}} + \frac{\sqrt{p}}{2}\delta_V^2(\mathbf{x}).$$

By definition (A.20) $\bar{\xi} \sim \mathcal{N}(0, \text{Var } \xi)$. Lemma A.10 and Theorem 1.2 by Spokoiny (2012b) imply

$$\begin{aligned} P\left(\|\bar{\xi}\| - \|\bar{\xi}_0\| \geq \sqrt{\text{tr}(\text{Var } \xi)} - \sqrt{\text{tr}(\text{Var } \bar{\xi}_0)} + \Delta_{\text{qf},1}\right) &\leq 2e^{-x}, \\ P\left(\|\bar{\xi}\| - \|\bar{\xi}_0\| \leq \sqrt{\text{tr}(\text{Var } \xi)} - \sqrt{\text{tr}(\text{Var } \bar{\xi}_0)} - \Delta_{\text{qf},2}\right) &\leq 2e^{-x}, \end{aligned} \quad (\text{A.32})$$

where

$$\begin{aligned} \Delta_{\text{qf},1} &\stackrel{\text{def}}{=} \left[4x \text{tr}\{(\text{Var } \bar{\xi}_0)^2\}\right]^{1/4} \\ &\quad + \max\left[2\sqrt{2x \text{tr}\{(\text{Var } \xi)^2\}}, 6x\|\text{Var } \xi\|\right]^{1/2}, \\ \Delta_{\text{qf},2} &\stackrel{\text{def}}{=} \left[4x \text{tr}\{(\text{Var } \xi)^2\}\right]^{1/4} \\ &\quad + \max\left[2\sqrt{2x \text{tr}\{(\text{Var } \bar{\xi}_0)^2\}}, 6x\|\text{Var } \bar{\xi}_0\|\right]^{1/2}. \end{aligned} \quad (\text{A.33})$$

By conditions (\mathcal{I}) , (\mathcal{I}_B)

$$\begin{aligned} \Delta_{\text{qf},1} &\leq \left\{\sqrt{4xp}(\mathbf{a}^2 + \mathbf{a}_B^2)\right\}^{1/2} + \mathbf{a} \max\left\{\sqrt{8xp}, 6x\right\}^{1/2}, \\ \Delta_{\text{qf},2} &\leq \left\{4xp\mathbf{a}^4\right\}^{1/4} + \sqrt{\mathbf{a}^2 + \mathbf{a}_B^2} \max\left\{\sqrt{8xp}, 6x\right\}^{1/2}. \end{aligned} \quad (\text{A.34})$$

Further, it holds on a random set with probability $\geq 1 - 2e^{-x}$

$$\begin{aligned}
& \mathbb{P}(\|\bar{\boldsymbol{\xi}}_0\| > \mathfrak{z}) - \Delta_{G,1} \\
&= \mathbb{P}(\|\bar{\boldsymbol{\xi}}\| > \mathfrak{z} + \|\bar{\boldsymbol{\xi}}\| - \|\bar{\boldsymbol{\xi}}_0\|) - \Delta_{G,1} \\
&\stackrel{\text{(by (A.32))}}{\geq} \mathbb{P}\left(\|\bar{\boldsymbol{\xi}}\| > \mathfrak{z} + \sqrt{\text{tr}(\text{Var } \boldsymbol{\xi})} - \sqrt{\text{tr}(\text{Var } \bar{\boldsymbol{\xi}}_0)} + \Delta_{\text{qf},1}\right) - \Delta_{G,1} \\
&\stackrel{\text{(Th. 5.1)}}{\geq} \mathbb{P}\left(\|\boldsymbol{\xi}\| > \mathfrak{z} - \Delta_{\text{w}}(\mathbf{r}_0, \mathbf{x}) + \sqrt{\text{tr}(\text{Var } \boldsymbol{\xi})} - \sqrt{\text{tr}(\text{Var } \bar{\boldsymbol{\xi}}_0)} + \Delta_{\text{qf},1}\right) \\
&\quad - \Delta_{G,1} - \Delta_{G,2} \\
&\stackrel{\text{(Th. A.2)}}{\geq} \mathbb{P}\left(\sqrt{2\{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\}} > \mathfrak{z} + \sqrt{\text{tr}(\text{Var } \boldsymbol{\xi})} - \sqrt{\text{tr}(\text{Var } \bar{\boldsymbol{\xi}}_0)} + \Delta_{\text{qf},1}\right) \\
&\quad - \Delta_{\text{b, full}},
\end{aligned}$$

where

$$\begin{aligned}
\Delta_{G,2} &\stackrel{\text{def}}{=} 1.55\delta_n^{1/4} + \frac{\Delta_{\text{w}}(\mathbf{r}_0, \mathbf{x})}{\sqrt{2}}, \\
\Delta_{\text{b, full}} &= \Delta_{G,1} + \Delta_{G,2}.
\end{aligned}$$

Hence, we obtain

$$\begin{aligned}
& \mathbb{P}^\circ\left(\sqrt{2\{L^\circ(\tilde{\boldsymbol{\theta}}^\circ) - L^\circ(\tilde{\boldsymbol{\theta}})\}} > \mathfrak{z}\right) \\
&\geq \mathbb{P}\left(\sqrt{2\{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\}} > \mathfrak{z} + \sqrt{\text{tr}(\text{Var } \boldsymbol{\xi})} - \sqrt{\text{tr}(\text{Var } \bar{\boldsymbol{\xi}}_0)} + \Delta_{\text{qf},1}\right) \\
&\quad - \Delta_{\text{b, full}}.
\end{aligned}$$

By definition (2.2) of $(1 - \alpha)$ -quantile \mathfrak{z}_α it holds:

$$\mathfrak{z}_{(\alpha + \Delta_{\text{b, full}})} \leq \mathfrak{z}_\alpha + \sqrt{\text{tr}(\text{Var } \boldsymbol{\xi})} - \sqrt{\text{tr}(\text{Var } \bar{\boldsymbol{\xi}}_0)} + \Delta_{\text{qf},1},$$

and in addition

$$\sqrt{\text{tr}(\text{Var } \boldsymbol{\xi})} - \sqrt{\text{tr}(\text{Var } \bar{\boldsymbol{\xi}}_0)} \leq -\frac{\text{tr}(D_0^{-1}B_0^2D_0^{-1})}{2\sqrt{\text{tr}(D_0^{-1}H_0^2D_0^{-1})}} \leq 0.$$

The inverse inequalities are implied with the similar arguments:

$$\begin{aligned} & \mathbb{P}^\circ \left(\sqrt{2 \{L^\circ(\tilde{\theta}^\circ) - L^\circ(\tilde{\theta})\}} > \mathfrak{z} \right) \\ & \leq \mathbb{P} \left(\sqrt{2 \{L(\tilde{\theta}) - L(\theta^*)\}} > \mathfrak{z} + \sqrt{\text{tr}(\text{Var } \boldsymbol{\xi})} - \sqrt{\text{tr}(\text{Var } \bar{\boldsymbol{\xi}}_0)} - \Delta_{\text{qf},2} \right) \\ & \quad + \Delta_{\text{b, full}}. \end{aligned}$$

And

$$\mathfrak{z}_{(\alpha - \Delta_{\text{b, full}})} \geq \mathfrak{z}_\alpha^\circ + \sqrt{\text{tr}(\text{Var } \boldsymbol{\xi})} - \sqrt{\text{tr}(\text{Var } \bar{\boldsymbol{\xi}}_0)} - \Delta_{\text{qf},2}.$$

□

A.5 Proof of Theorem 2.5 (the smoothed version)

Lemma A.11. For the function $g_\Delta(x, z)$ defined in (2.6), all $\Delta_1 \in [0, x]$ and all $C \geq 1$ it holds

$$g_\Delta(x - \Delta_1, z) \geq g_\Delta(x, z + \Delta_1 C)$$

Proof of Lemma A.11. By definition (5.9) of $g(x)$

$$\begin{aligned} \max_{x \geq 0} \{g_\Delta(x - \Delta_1, z) = 0\} &= z + \Delta_1, \\ \max_{x \geq 0} \{g_\Delta(x, z + \Delta_1 C) = 0\} &= z + \Delta_1 C. \end{aligned}$$

For $x \geq z + \Delta_1 C$ it holds

$$\begin{aligned} g_\Delta(x - \Delta_1, z) &= g \left(\frac{1}{2\Delta z} \{(x - \Delta_1)^2 - z^2\} \right) \\ &\geq g \left(\frac{1}{2\Delta(z + \Delta_1 C)} \{x^2 - (z + \Delta_1 C)^2\} \right) \\ &= g_\Delta(x, z + \Delta_1 C). \end{aligned} \tag{A.35}$$

Indeed, the comparison in (A.35) reads as

$$\begin{aligned} & (z + \Delta_1 C)(x - \Delta_1 + z)(x - \Delta_1 - z) \\ & \vee z(x + z + \Delta_1 C)(x - z - \Delta_1 C). \end{aligned} \tag{A.36}$$

Since $C \geq 1$, $(x - \Delta_1 - z) \geq (x - \Delta_1 C - z)$ and it holds for the left side of (A.36):

$$\begin{aligned} (z + \Delta_1 C)(x - \Delta_1 + z) &= (zx + z^2 + 2\Delta_1 C) + \Delta_1(xC - \Delta_1 C - z) \\ &\geq (zx + z^2 + 2\Delta_1 C), \end{aligned}$$

which is equal to the multiplier $z(x + \Delta_1 C + z)$ in right side. \square

Proposition A.12 (Smooth analog of Proposition A.9). *If conditions (SmB) and (SD₁) are fulfilled, then it holds for all $0 < \Delta \leq 0.22$ and for all $z, \bar{z} > 2$ s.t. $|z - \bar{z}| \leq \delta_z$ for some $\delta_z \in [0, 1]$ with probability $\geq 1 - e^{-x}$:*

$$\begin{aligned}
& \left| \mathbb{E} g_\Delta (\|\boldsymbol{\xi}\|, z) - \mathbb{E}^\circ g_\Delta (\|\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\|, \bar{z}) \right| \\
& \leq \frac{16}{\Delta^3} (\delta_n + \check{\delta}_n) + 2\sqrt{p} \frac{\delta_z}{z} + \sqrt{p} \frac{\delta_z^2}{z^2} + \sqrt{p} \frac{\delta_V^2(\mathbf{x}) + \delta_{\text{smb}}^2}{1 - \delta_V^2(\mathbf{x})} \\
& \leq \frac{16}{\Delta^3} (\delta_n + \check{\delta}_n) + \sqrt{5} \delta_z + \frac{4\sqrt{p}}{3} \{\delta_V^2(\mathbf{x}) + \delta_{\text{smb}}^2\} \tag{A.37} \\
& \text{for } \bar{z} \geq \sqrt{p}, \delta_V^2(\mathbf{x}) \leq 1/4.
\end{aligned}$$

Proof of Proposition A.12. The conditions of Theorem 5.5 are fulfilled with the value $\delta_\Sigma = \sqrt{p} \{\delta_V^2(\mathbf{x}) + \delta_{\text{smb}}^2\} / \{1 - \delta_V^2(\mathbf{x})\}$ due to the proof of Proposition A.9. \square

Proof of Theorem 2.5. The following holds on a random set of probability $\geq 1 - 12e^{-x}$:

$$\begin{aligned}
& \mathbb{E}^\circ g_\Delta \left(\sqrt{2 \{L^\circ(\tilde{\boldsymbol{\theta}}) - L^\circ(\boldsymbol{\theta})\}}, \mathfrak{z} \right) \\
& \stackrel{(\text{Th. A.4})}{\geq} \mathbb{E}^\circ g_\Delta \left(\|\boldsymbol{\xi}^\circ(\tilde{\boldsymbol{\theta}})\| - \Delta_{\mathbf{w}}^\circ(\mathbf{r}_0, \mathbf{x}), \mathfrak{z} \right) \\
& \stackrel{(\text{L. A.7})}{\geq} \mathbb{E}^\circ g_\Delta \left(\|\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\| - \Delta_{\mathbf{w}}^\circ(\mathbf{r}_0, \mathbf{x}) - \Delta_{\xi}^\circ(\mathbf{r}_0, \mathbf{x}), \mathfrak{z} \right) \\
& \stackrel{(\text{L. A.11})}{\geq} \mathbb{E}^\circ g_\Delta \left(\|\boldsymbol{\xi}^\circ(\boldsymbol{\theta}^*)\|, \mathfrak{z} + \Delta_{\mathbf{w}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_{\xi}^\circ(\mathbf{r}_0, \mathbf{x}) \right) \\
& \stackrel{(\text{Prop. A.12})}{\geq} \mathbb{E} g_\Delta (\|\boldsymbol{\xi}\|, \mathfrak{z} - \Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x})) - \Delta_{\text{sm}} \\
& \stackrel{(\text{Th. A.2, L. A.11})}{\geq} \mathbb{E} g_\Delta \left(\sqrt{2 \{L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)\}}, \mathfrak{z} \right) - \Delta_{\text{sm}},
\end{aligned}$$

where the term Δ_{sm} comes from (A.37) with $\delta_z := \Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x}) + \Delta_{\mathbf{w}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_{\xi}^\circ(\mathbf{r}_0, \mathbf{x})$:

$$\begin{aligned}
\Delta_{\text{sm}} & \stackrel{\text{def}}{=} \frac{16}{\Delta^3} (\delta_n + \check{\delta}_n) + \frac{4\sqrt{p}}{3} \{\delta_V^2(\mathbf{x}) + \delta_{\text{smb}}^2\} \\
& \quad + \sqrt{5} \{\Delta_{\mathbf{w}}(\mathbf{r}_0, \mathbf{x}) + \Delta_{\mathbf{w}}^\circ(\mathbf{r}_0, \mathbf{x}) + \Delta_{\xi}^\circ(\mathbf{r}_0, \mathbf{x})\}. \tag{A.38}
\end{aligned}$$

By the similar inequalities in the inverse direction we get the statement proved. Due to the arguments in the end of the proof of Theorem 2.1 it holds in the case 4.3

$$\Delta_{\text{sm}} = C \frac{1}{\Delta^3} \left\{ \frac{(p+x)^3}{n} \right\}^{1/2} + C \frac{p+x}{\sqrt{n}} \sqrt{x} + C \frac{p+x}{\sqrt{n}}. \tag{A.39}$$

\square

A.6 Bernstein matrix inequality

Consider the following symmetric $p \times p$ \mathbb{P} -random matrix and its expected value:

$$\begin{aligned}\mathcal{V}^2(\boldsymbol{\theta}^*) &\stackrel{\text{def}}{=} \text{Var}^\circ(\nabla_{\boldsymbol{\theta}} L^\circ(\boldsymbol{\theta}^*)) = \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)^\top, \\ H_0^2 &\stackrel{\text{def}}{=} \mathbb{E} \mathcal{V}^2(\boldsymbol{\theta}^*) = \sum_{i=1}^n \mathbb{E} [\nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)^\top].\end{aligned}$$

Matrix $\mathcal{V}^2(\boldsymbol{\theta}^*)$ equals to a sum of the independent random matrices $\nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)^\top$. Assuming the condition **(SD₁)** to be fulfilled, we can refer to the result by [Tropp \(2012\)](#) in order to get the concentration bound below. Let us previously introduce some notations.

$$v_i^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} H_0^{-1} \{ \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta})^\top - \mathbb{E} [\nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta})^\top] \} H_0^{-1},$$

then

$$H_0^{-1} \mathcal{V}^2(\boldsymbol{\theta}^*) H_0^{-1} = \sum_{i=1}^n v_i^2(\boldsymbol{\theta}^*).$$

Define also

$$\chi_v^2 \stackrel{\text{def}}{=} \left\| \sum_{i=1}^n \mathbb{E} v_i^4(\boldsymbol{\theta}^*) \right\|.$$

Theorem A.13 (Bernstein inequality for $\mathcal{V}^2(\boldsymbol{\theta}^*)$). *Let the condition **(SD₁)** be fulfilled, then it holds with probability $\geq 1 - e^{-x}$:*

$$\|H_0^{-1} \mathcal{V}^2(\boldsymbol{\theta}^*) H_0^{-1} - \mathbf{I}_p\| \leq \delta_{\mathcal{V}}^2(\mathbf{x}),$$

where the error term is defined as

$$\delta_{\mathcal{V}}^2(\mathbf{x}) \stackrel{\text{def}}{=} \sqrt{2\chi_v^2 \{\log(p) + \mathbf{x}\}} + \frac{2}{3} \delta_v^2 \{\log(p) + \mathbf{x}\}$$

and is proportional to $\sqrt{\{\log(p) + \mathbf{x}\}/n}$ in the case [4.3](#).

Proof. Due to Theorem 1.4 by [Tropp \(2012\)](#):

$$\mathbb{P}(\|H_0^{-1} \mathcal{V}^2(\boldsymbol{\theta}^*) H_0^{-1} - \mathbf{I}_p\| \geq t) \leq p \exp\left(\frac{-t^2}{2\chi_v^2 + 2\delta_v^2 t/3}\right).$$

For

$$\mathbf{x} = \frac{t^2}{2\chi_v^2 + 2\delta_v^2 t/3} - \log(p)$$

it holds:

$$\mathbb{P} (\|H_0^{-1}\mathcal{V}^2(\boldsymbol{\theta}^*)H_0^{-1} - \mathbf{I}_p\| \geq \delta_{\mathcal{V}}^2(\mathbf{x})) \leq e^{-x}.$$

□

References

- Aerts, M. and Claeskens, G. (2001). Bootstrap tests for misspecified models, with application to clustered binary data. *Computational statistics & data analysis*, 36(3):383–401.
- Arlot, S., Blanchard, G., and Roquain, E. (2010). Some nonasymptotic results on resampling in high dimension. I. Confidence regions. *Ann. Statist.*, 38(1):51–82.
- Barbe, P. and Bertail, P. (1995). *The weighted bootstrap*, volume 98. Springer.
- Bentkus, V. (2003). On the dependence of the Berry–Esseen bound on dimension. *Journal of Statistical Planning and Inference*, 113(2):385–402.
- Berry, A. C. (1941). The accuracy of the Gaussian approximation to the sum of independent variates. *Transactions of the american mathematical society*, 49(1):122–136.
- Bhattacharya, R. and Holmes, S. (2010). An exposition of Götze’s estimation of the rate of convergence in the multivariate central limit theorem. *arXiv preprint arXiv:1003.4254*.
- Chatterjee, S. and Bose, A. (2005). Generalized bootstrap for estimating equations. *The Annals of Statistics*, 33(1):414–436.
- Chen, L. H. and Fang, X. (2011). Multivariate normal approximation by Stein’s method: The concentration inequality approach. *arXiv preprint arXiv:1111.4073*.
- Chen, X. and Pouzo, D. (2009). Efficient estimation of semiparametric conditional moment models with possibly nonsmooth residuals. *Journal of Econometrics*, 152(1):46–60.
- Chen, X. and Pouzo, D. (2014). Sieve Wald and QLR Inferences on semi/nonparametric conditional moment models.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2013). Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *The Annals of Statistics*, 41(6):2786–2819.
- Efron, B. (1979). Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, 7(1):1–26.
- Esseen, C.-G. (1942). *On the Liapounoff limit of error in the theory of probability*. Almqvist & Wiksell.
- Esseen, C. G. (1956). A moment inequality with an application to the central limit theorem. *Scandinavian Actuarial Journal*, 1956(2):160–170.

- Götze, F. (1991). On the rate of convergence in the multivariate CLT. *The Annals of Probability*, pages 724–739.
- Götze, F. and Zaitsev, A. Y. (2014). Explicit rates of approximation in the CLT for quadratic forms. *The Annals of Probability*, 42(1):354–397.
- Hall, A. R. (2005). *Generalized method of moments*. Oxford University Press Oxford.
- Hall, P. (1992). *The bootstrap and Edgeworth expansion*. Springer.
- Horowitz, J. L. (2001). The bootstrap. *Handbook of econometrics*, 5:3159–3228.
- Huber, P. (1967). The behavior of maximum likelihood estimates under nonstandard conditions. Proc. 5th Berkeley Symp. Math. Stat. Probab., Univ. Calif. 1965/66, 1, 221–233 (1967).
- Janssen, A. and Pauls, T. (2003). How do bootstrap and permutation tests work? *Annals of statistics*, pages 768–806.
- Kline, P. and Santos, A. (2012). Higher order properties of the wild bootstrap under misspecification. *Journal of Econometrics*, 171(1):54–70.
- Lavergne, P. and Patilea, V. (2013). Smooth minimum distance estimation and testing with conditional estimating equations: uniform in bandwidth theory. *Journal of Econometrics*, 177(1):47–59.
- Lindeberg, J. W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeitsrechnung. *Mathematische Zeitschrift*, 15(1):211–225.
- Ma, S. and Kosorok, M. R. (2005). Robust semiparametric M-estimation and the weighted bootstrap. *Journal of Multivariate Analysis*, 96(1):190–217.
- Mammen, E. (1992). *When does bootstrap work?*, volume 77. Springer.
- Mammen, E. (1993). Bootstrap and wild bootstrap for high dimensional linear models. *The Annals of Statistics*, pages 255–285.
- Newton, M. A. and Raftery, A. E. (1994). Approximate bayesian inference with the weighted likelihood bootstrap. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 3–48.
- Prokhorov, Y. V. and Ulyanov, V. V. (2013). Some approximation problems in statistics and probability. In *Limit Theorems in Probability, Statistics and Number Theory*, pages 235–249. Springer.
- Shevtsova, I. (2010). An improvement of convergence rate estimates in the Lyapunov theorem. In *Doklady Mathematics*, volume 82, pages 862–864. Springer.
- Spokoiny, V. (2012a). Parametric estimation. Finite sample theory. *The Annals of Statistics*, 40(6):2877–2909.
- Spokoiny, V. (2012b). Supplement to “Parametric estimation. Finite sample theory”.

- Spokoiny, V. (2013). Bernstein-von Mises Theorem for growing parameter dimension. *arXiv preprint arXiv:1302.3430*.
- Tropp, J. A. (2012). User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12(4):389–434.
- Tsybakov, A. B. (2009). *Introduction to Nonparametric Estimation*. Springer, New York.
- van der Vaart, A. W. and Wellner, J. A. (1996). *Weak Convergence and Empirical processes*. Springer, New York.
- Wilks, S. S. (1938). The large-sample distribution of the likelihood ratio for testing composite hypotheses. *The Annals of Mathematical Statistics*, 9(1):60–62.
- Wu, C. F. J. (1986). Jackknife, bootstrap and other resampling methods in regression analysis. *The Annals of Statistics*, 14(4):1261–1295+.