

SFB 649 Discussion Paper 2015-031

# Simultaneous likelihood-based bootstrap confidence sets for a large number of models

Mayya Zhilova\*



\* Humboldt-Universität zu Berlin, Germany

This research was supported by the Deutsche  
Forschungsgemeinschaft through the SFB 649 "Economic Risk".

<http://sfb649.wiwi.hu-berlin.de>  
ISSN 1860-5664

SFB 649, Humboldt-Universität zu Berlin  
Spandauer Straße 1, D-10178 Berlin



SFB 649 ECONOMIC RISK BERLIN

# Simultaneous likelihood-based bootstrap confidence sets for a large number of models

Mayya Zhilova<sup>\*,†</sup>

Weierstrass-Institute,  
Mohrenstr. 39,  
10117 Berlin, Germany,  
zhilova@wias-berlin.de

June 18, 2015

## Abstract

The paper studies a problem of constructing simultaneous likelihood-based confidence sets. We consider a simultaneous multiplier bootstrap procedure for estimating the quantiles of the joint distribution of the likelihood ratio statistics, and for adjusting the confidence level for multiplicity. Theoretical results state the bootstrap validity in the following setting: the sample size  $n$  is fixed, the maximal parameter dimension  $p_{\max}$  and the number of considered parametric models  $K$  are s.t.  $(\log K)^{12} p_{\max}^3/n$  is small. We also consider the situation when the parametric models are misspecified. If the models' misspecification is significant, then the bootstrap critical values exceed the true ones and the simultaneous bootstrap confidence set becomes conservative. Numerical experiments for local constant and local quadratic regressions illustrate the theoretical results.

*JEL classification codes:* C13, C15

*Keywords:* simultaneous inference, correction for multiplicity, family-wise error, misspecified model, multiplier/weighted bootstrap

---

<sup>\*</sup>I am very grateful to Vladimir Spokoiny for many helpful discussions and comments.

<sup>†</sup>Financial support by the German Research Foundation (DFG) through the Collaborative Research Center 649 "Economic Risk" is gratefully acknowledged.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>The multiplier bootstrap procedure</b>	<b>9</b>
<b>3</b>	<b>Theoretical justification of the bootstrap procedure</b>	<b>11</b>
3.1	Overview of the theoretical approach . . . . .	11
3.2	Main results . . . . .	14
<b>4</b>	<b>Numerical experiments</b>	<b>16</b>
4.1	Local constant regression . . . . .	16
4.2	Local quadratic regression . . . . .	16
4.3	Simulated data . . . . .	17
4.4	Effect of the modeling bias on a width of a bootstrap confidence band . . . . .	17
4.5	Effective coverage probability (local constant estimate) . . . . .	18
4.6	Correction for multiplicity . . . . .	21
<b>5</b>	<b>Conditions</b>	<b>22</b>
5.1	Basic conditions . . . . .	22
5.2	Conditions required for the bootstrap validity . . . . .	23
5.3	Dependence of the involved terms on the sample size and cardinality of the parameters' set . . . . .	24
<b>A</b>	<b>Approximation of the joint distributions of <math>\ell_2</math>-norms</b>	<b>25</b>
A.1	Joint Gaussian approximation of $\ell_2$ -norm of sums of independent vectors by Lindeberg's method . . . . .	27
A.2	Gaussian comparison . . . . .	35
A.3	Simultaneous anti-concentration for $\ell_2$ -norms of Gaussian vectors . . . . .	37
A.4	Proof of Proposition A.1 . . . . .	39
<b>B</b>	<b>Square-root Wilks approximations</b>	<b>42</b>
B.1	Finite sample theory . . . . .	42
B.2	Finite sample theory for the bootstrap world . . . . .	44
B.3	Simultaneous square-root Wilks approximations . . . . .	46
<b>C</b>	<b>Proofs of the main results</b>	<b>47</b>
C.1	Bernstein matrix inequality . . . . .	47
C.2	Bootstrap validity for the case of one parametric model . . . . .	48
C.3	Proof of Theorem 3.1 . . . . .	49
C.4	Proof of Theorem 3.2 . . . . .	51
C.5	Proof of Theorem 3.3 . . . . .	55
	<b>References</b>	<b>56</b>

## 1 Introduction

The problem of simultaneous confidence estimation appears in numerous practical applications when a confidence statement has to be made simultaneously for a collection of objects, e.g. in safety analysis in clinical trials, gene expression analysis, population biology, functional magnetic resonance imaging and many others. See e.g. [Miller \(1981\)](#); [Westfall \(1993\)](#); [Manly \(2006\)](#); [Benjamini \(2010\)](#); [Dickhaus \(2014\)](#), and references therein. This problem is also closely related to construction of simultaneous confidence bands in curve estimation, which goes back to [Working and Hotelling \(1929\)](#). For an extensive literature review about constructing the simultaneous confidence bands we refer to [Hall and Horowitz \(2013\)](#), [Liu \(2010\)](#), and [Wasserman \(2006\)](#).

A simultaneous confidence set requires a probability bound to be constructed jointly for several possibly dependent statistics. Therefore, the critical values of the corresponding statistics should be chosen in such a way that the joint probability distribution achieves a required family-wise confidence level. This choice can be made by multiplicity correction of the marginal confidence levels. The Bonferroni correction method ([Bonferroni \(1936\)](#)) uses a probability union bound, the corrected marginal significance levels are taken equal to the total level divided by the number of models. This procedure can be very conservative if the considered statistics are positively correlated and if their number is large. The Šidák correction method ([Šidák \(1967\)](#)) is more powerful than Bonferroni correction, however, it also becomes conservative in the case of large number of dependent statistics.

Most of the existing results about simultaneous bootstrap confidence sets and resampling-based multiple testing are asymptotic (with sample size tending to infinity), see e.g. [Beran \(1988, 1990\)](#); [Hall and Pittelkow \(1990\)](#); [Härdle and Marron \(1991\)](#); [Shao and Tu \(1995\)](#); [Hall and Horowitz \(2013\)](#), and [Westfall \(1993\)](#); [Dickhaus \(2014\)](#). The results based on asymptotic distribution of maximum of an approximating Gaussian process (see [Bickel and Rosenblatt \(1973\)](#); [Johnston \(1982\)](#); [Härdle \(1989\)](#)) require a huge sample size  $n$ , since they yield a coverage probability error of order  $(\log(n))^{-1}$  (see [Hall \(1991\)](#)). Some papers considered an alternative approach in context of confidence band estimation based on the approximation of the underlying empirical processes by its bootstrap counterpart. In particular, [Hall \(1993\)](#) showed that such an approach leads to a significant improvement of the error rate (see also [Neumann and Polzehl \(1998\)](#); [Claeskens and Van Keilegom \(2003\)](#)). [Chernozhukov et al. \(2014a\)](#) constructed honest confidence bands for nonparametric density estimators without requiring the existence of limit distribution of the supremum of the studentized empirical process: instead, they used an approximation between sup-norms of an empirical and Gaussian processes, and anti-concentration

property of suprema of Gaussian processes.

In many modern applications the sample size cannot be large, and/or can be smaller than a parameter dimension, for example, in genomics, brain imaging, spatial epidemiology and microarray data analysis, see [Leek and Storey \(2008\)](#); [Kim and van de Wiel \(2008\)](#); [Arlot et al. \(2010\)](#); [Cao and Kosorok \(2011\)](#), and references therein.

For the recent results on resampling-based simultaneous confidence sets in high-dimensional finite sample set-up we refer to the papers by [Arlot et al. \(2010\)](#) and [Chernozhukov et al. \(2013a, 2014a,b\)](#). [Arlot et al. \(2010\)](#) considered i.i.d. observations of a Gaussian vector with a dimension possibly much larger than the sample size, and with unknown covariance matrix. They examined multiple testing problems for the mean values of its coordinates and provided non-asymptotic control for the family-wise error rate using resampling-type procedures. [Chernozhukov et al. \(2013a\)](#) presented a number of non-asymptotic results on Gaussian approximation and multiplier bootstrap for maxima of sums of high-dimensional vectors (with a dimension possibly much larger than a sample size) in a very general set-up. As an application the authors considered the problem of multiple hypothesis testing in the framework of approximate means. They derived non-asymptotic results for the general stepdown procedure by [Romano and Wolf \(2005\)](#) with improved error rates and in high-dimensional setting. [Chernozhukov et al. \(2014a\)](#) showed how this technique applies to the problem of constructing an honest confidence set in nonparametric density estimation. [Chernozhukov et al. \(2014b\)](#) extended the results from maxima to the class of sparsely convex sets.

The present paper studies simultaneous likelihood-based bootstrap confidence sets in the following setting:

1. the sample size  $n$  is fixed;
2. the parametric models can be misspecified;
3. the number  $K$  of the parametric models can be exponentially large w.r.t.  $n$ ;
4. the maximal dimension  $p_{\max}$  of the considered parametric models can be dependent on the sample size  $n$ .

This set-up, in contrast with the paper by [Chernozhukov et al. \(2014b\)](#), does not require the sparsity condition, in particular the dimension  $p_1, \dots, p_K$  of each parametric family may grow with the sample size. Moreover, the simultaneous likelihood-based confidence sets are not necessarily convex, and the parametric assumption can be violated.

The considered simultaneous multiplier bootstrap procedure involves two main steps: estimation of the quantile functions of the likelihood ratio statistics, and multiplicity correction of the marginal confidence level. Theoretical results of the paper state the

bootstrap validity in the setting 1-4 taking in account the multiplicity correction. The resulting approximation bound requires the quantity  $(\log K)^{12} p_{\max}^3/n$  to be small. The log-factor here is suboptimal and can probably be improved. The paper particularly focuses on the impact of the model misspecification. We distinguish between slight and strong misspecifications. Under the so called small modeling bias condition (**SmB**) given in Section 5.2 the bootstrap approximation is accurate. This condition roughly means that all the parametric models are close to the true distribution. If the (**SmB**) condition is not fulfilled, then the simultaneous bootstrap confidence set is still applicable, however, it becomes conservative. This property is nicely confirmed by the numerical experiments in Section 4.

Let the random data

$$\mathbf{Y} \stackrel{\text{def}}{=} (Y_1, \dots, Y_n)^\top \quad (1.1)$$

consist of *independent* observations  $Y_i$ , and belong to the probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . The sample size  $n$  is *fixed*.  $\mathbb{P}$  is an *unknown probability distribution* of the sample  $\mathbf{Y}$ . Consider  $K$  regular parametric families of probability distributions:

$$\{\mathbb{P}_k(\boldsymbol{\theta})\} \stackrel{\text{def}}{=} \{\mathbb{P}_k(\boldsymbol{\theta}) \ll \mu_0, \boldsymbol{\theta} \in \Theta_k \subset \mathbb{R}^{p_k}\}, \quad k = 1, \dots, K.$$

Each parametric family induces the quasi log-likelihood function for  $\boldsymbol{\theta} \in \Theta_k \subset \mathbb{R}^{p_k}$

$$\begin{aligned} L_k(\mathbf{Y}, \boldsymbol{\theta}) &\stackrel{\text{def}}{=} \log \left( \frac{d\mathbb{P}_k(\boldsymbol{\theta})}{d\mu_0}(\mathbf{Y}) \right) \\ &= \sum_{i=1}^n \log \left( \frac{d\mathbb{P}_k(\boldsymbol{\theta})}{d\mu_0}(Y_i) \right). \end{aligned} \quad (1.2)$$

It is important that we *do not require* that  $\mathbb{P}$  belongs to any of the known parametric families  $\{\mathbb{P}_k(\boldsymbol{\theta})\}$ , that is why the term *quasi* log-likelihood is used here. Below in this section we consider two popular examples of simultaneous confidence sets in terms of the quasi log-likelihood functions (1.2). Namely, the simultaneous confidence band for local constant regression, and multiple quantiles regression.

The target of estimation for the misspecified log-likelihood  $L_k(\boldsymbol{\theta})$  is such a parameter  $\boldsymbol{\theta}_k^*$ , that minimises the Kullback-Leibler distance between the unknown true measure  $\mathbb{P}$  and the parametric family  $\{\mathbb{P}_k(\boldsymbol{\theta})\}$ :

$$\boldsymbol{\theta}_k^* \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta} \in \Theta_k}{\operatorname{argmax}} \mathbb{E} L_k(\boldsymbol{\theta}). \quad (1.3)$$

The maximum likelihood estimator is defined as:

$$\tilde{\boldsymbol{\theta}}_k \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta} \in \Theta_k}{\operatorname{argmax}} L_k(\boldsymbol{\theta}).$$

The parametric sets  $\Theta_k$  have dimensions  $p_k$ , therefore,  $\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_k^* \in \mathbb{R}^{p_k}$ . For  $1 \leq k, j \leq K$  and  $k \neq j$  the numbers  $p_k$  and  $p_j$  can be unequal.

The likelihood-based confidence set for the target parameter  $\boldsymbol{\theta}_k^*$  is

$$\mathcal{E}_k(\mathfrak{z}) \stackrel{\text{def}}{=} \left\{ \boldsymbol{\theta} \in \Theta_k : L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}) \leq \mathfrak{z}^2/2 \right\} \subset \mathbb{R}^{p_k}. \quad (1.4)$$

Let  $\mathfrak{z}_k(\alpha)$  denote the  $(1 - \alpha)$ -quantile of the corresponding square-root likelihood ratio statistic:

$$\mathfrak{z}_k(\alpha) \stackrel{\text{def}}{=} \inf \left\{ \mathfrak{z} \geq 0 : \mathbb{P} \left( L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) > \mathfrak{z}^2/2 \right) \leq \alpha \right\}. \quad (1.5)$$

Together with (1.4) this implies for each  $k = 1, \dots, K$ :

$$\mathbb{P} \left( \boldsymbol{\theta}_k^* \in \mathcal{E}_k(\mathfrak{z}_k(\alpha)) \right) \geq 1 - \alpha. \quad (1.6)$$

Thus  $\mathcal{E}_k(\mathfrak{z})$  and the quantile function  $\mathfrak{z}_k(\alpha)$  fully determine the marginal  $(1 - \alpha)$ -confidence set. The simultaneous confidence set requires a correction for multiplicity. Let  $\mathfrak{c}(\alpha)$  denote a maximal number  $c \in (0, \alpha]$  s.t.

$$\mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} > \mathfrak{z}_k(c) \right\} \right) \leq \alpha. \quad (1.7)$$

This is equivalent to

$$\mathfrak{c}(\alpha) \stackrel{\text{def}}{=} \sup \left\{ c \in (0, \alpha] : \mathbb{P} \left( \max_{1 \leq k \leq K} \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} - \mathfrak{z}_k(c) \right\} > 0 \right) \leq \alpha \right\}. \quad (1.8)$$

Therefore, taking the marginal confidence sets with the same confidence levels  $1 - \mathfrak{c}(\alpha)$  yields the simultaneous confidence bound of the total level  $1 - \alpha$ . The value  $\mathfrak{c}(\alpha) \in (0, \alpha]$  is the correction for multiplicity. In order to construct the simultaneous confidence set using this correction, one has to estimate the values  $\mathfrak{z}_k(\mathfrak{c}(\alpha))$  for all  $k = 1, \dots, K$ . By its definition this problem splits into two subproblems:

1. **Marginal step.** Estimation of the marginal quantile functions  $\mathfrak{z}_1(\alpha), \dots, \mathfrak{z}_K(\alpha)$  given in (1.5).
2. **Correction for multiplicity.** Estimation of the correction for multiplicity  $\mathfrak{c}(\alpha)$  given in (1.8).

If the 1-st problem is solved for any  $\alpha \in (0, 1)$ , the 2-nd problem can be treated by calibrating the value  $\alpha$  s.t. (1.8) holds. It is important to take into account the correlation between the likelihood ratio statistics  $L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*)$ ,  $k = 1, \dots, K$ , otherwise the estimate of the correction  $\mathfrak{c}(\alpha)$  can be too conservative. For instance, the Bonferroni

correction would lead to the marginal confidence level  $1 - \alpha/K$ , which may be very conservative if  $K$  is large and the statistics  $L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*)$  are highly correlated.

In Section 2 we suggest a multiplier bootstrap procedure, which performs the steps 1 and 2 described above. Theoretical justification of the procedure is given in Section 3. The proofs are based on several approximation bounds: non-asymptotic square-root Wilks theorem, simultaneous Gaussian approximation for  $\ell_2$ -norms, Gaussian comparison, and simultaneous Gaussian anti-concentration inequality.

Spokoiny and Zhilova (2014) considered the 1-st subproblem for the case of a single parametric model ( $K = 1$ ): a multiplier bootstrap procedure was applied for construction of a likelihood-based confidence set, and justified theoretically for a fixed sample size and for possibly misspecified parametric model. In the present paper we extend that approach for the case of simultaneously many parametric models.

Below we illustrate the definitions (1.2)-(1.8) of the simultaneous likelihood-based confidence sets with two popular examples.

**Example 1 (Simultaneous confidence band for local constant regression):**

Let  $Y_1, \dots, Y_n$  be independent random scalar observations and  $X_1, \dots, X_n$  some deterministic design points. Consider the following quadratic likelihood function reweighted with the kernel functions  $K(\cdot)$ :

$$L(\boldsymbol{\theta}, x, h) \stackrel{\text{def}}{=} -\frac{1}{2} \sum_{i=1}^n (Y_i - \boldsymbol{\theta})^2 w_i(x, h),$$

$$w_i(x, h) \stackrel{\text{def}}{=} K(\{x - X_i\}/h),$$

$$K(x) \in [0, 1], \int_{\mathbb{R}} K(x) dx = 1, K(x) = K(-x).$$

Here  $h > 0$  denotes bandwidth, the local smoothing parameter. The target point and the local MLE read as:

$$\boldsymbol{\theta}^*(x, h) \stackrel{\text{def}}{=} \frac{\sum_{i=1}^n w_i(x, h) \mathbb{E}Y_i}{\sum_{i=1}^n w_i(x, h)}, \quad \tilde{\boldsymbol{\theta}}(x, h) \stackrel{\text{def}}{=} \frac{\sum_{i=1}^n w_i(x, h) Y_i}{\sum_{i=1}^n w_i(x, h)}.$$

$\tilde{\boldsymbol{\theta}}(x, h)$  is also known as Nadaraya-Watson estimate. Fix a bandwidth  $h$  and consider the range of points  $x_1, \dots, x_K$ . They yield  $K$  local constant models with the target parameters  $\boldsymbol{\theta}_k^* \stackrel{\text{def}}{=} \boldsymbol{\theta}^*(x_k, h)$  and the likelihood functions  $L_k(\boldsymbol{\theta}) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}, x_k, h)$  for  $k = 1, \dots, K$ . The confidence intervals for each model are defined as

$$\mathcal{E}_k(\mathfrak{z}, h) \stackrel{\text{def}}{=} \left\{ \boldsymbol{\theta} \in \Theta : L(\tilde{\boldsymbol{\theta}}(x_k, h), x_k, h) - L(\boldsymbol{\theta}, x_k, h) \leq \mathfrak{z}^2/2 \right\},$$

for the quintile functions  $\mathfrak{z}_k(\alpha)$  and for the multiplicity correction  $\mathfrak{c}(\alpha)$  from (1.5) and (1.8) they form the following simultaneous confidence band:

$$\mathbb{P} \left( \bigcap_{k=1}^K \left\{ \boldsymbol{\theta}_k^* \in \mathcal{E}_k(\mathfrak{z}_k(\mathfrak{c}(\alpha))) \right\} \right) \geq 1 - \alpha.$$



In Section 4 we provide results of numerical experiments for this model.

**Example 2 (Multiple quantiles regression):** Quantile regression is an important method of statistical analysis, widely used in various applications. It aims at estimating conditional quantile functions of a response variable, see [Koenker \(2005\)](#). Multiple quantiles regression model considers simultaneously several quantile regression functions based on a range of quantile indices, see e.g. [Liu and Wu \(2011\)](#); [Qu \(2008\)](#); [He \(1997\)](#). Let  $Y_1, \dots, Y_n$  be independent random scalar observations and  $X_1, \dots, X_n \in \mathbb{R}^d$  some deterministic design points, as in Example 1. Consider the following quantile regression models for  $k = 1, \dots, K$ :

$$Y_i = g_k(X_i) + \varepsilon_{k,i}, \quad i = 1, \dots, n,$$

where  $g_k(\mathbf{x}) : \mathbb{R}^d \mapsto \mathbb{R}$  are unknown functions, the random values  $\varepsilon_{k,1}, \dots, \varepsilon_{k,n}$  are independent for each fixed  $k$ , and

$$\mathbb{P}(\varepsilon_{k,i} < 0) = \tau_k \quad \text{for all } i = 1, \dots, n.$$

The range of quantile indices  $\tau_1, \dots, \tau_K \in (0, 1)$  is known and fixed. We are interested in simultaneous parametric confidence sets for the functions  $g_1(\cdot), \dots, g_K(\cdot)$ . Let  $f_k(\mathbf{x}, \boldsymbol{\theta}) : \mathbb{R}^d \times \mathbb{R}^{p_k} \mapsto \mathbb{R}$  be known regression functions. Using the quantile regression approach by [Koenker and Bassett Jr \(1978\)](#), this problem can be treated with the quasi maximum likelihood method and the following log-likelihood functions:

$$L_k(\boldsymbol{\theta}) = - \sum_{i=1}^n \rho_{\tau_k}(Y_i - f_k(X_i, \boldsymbol{\theta})),$$

$$\rho_{\tau_k}(x) \stackrel{\text{def}}{=} x(\tau_k - \mathbb{I}\{x < 0\}).$$

for  $k = 1, \dots, K$ . This quasi log-likelihood function corresponds to the Asymmetric Laplace distribution with the density  $\tau_k(1 - \tau_k)e^{-\rho_{\tau_k}(x-a)}$ . If  $\tau = 1/2$ , then  $\rho_{1/2}(x) = |x|/2$  and  $L(\boldsymbol{\theta}) = - \sum_{i=1}^n |Y_i - f_k(X_i, \boldsymbol{\theta})|/2$ , which corresponds to the median regression.

The paper is organised as follows: Section 2 describes the multiplier bootstrap procedure, Section 3 explains the ideas of the theoretical approach and provides main results in Sections 3.1 and 3.2 correspondingly. All the necessary conditions are given in Section 5. In Section 5.3 and in statements of the main theoretical results we provide information about dependence of the involved terms on the sample size and parametric dimensions in the case of i.i.d. observations. Proofs of the main results are given in Section C. Statements from Sections A and B are used for the proofs in Section C. Numerical experiments are described in Section 4: we construct simultaneous confidence corridors for local constant and local quadratic regressions using both bootstrap and Monte Carlo

procedures. The quality of the bootstrap procedure is checked by computing the effective simultaneous coverage probabilities of the bootstrap confidence sets. We also compare the widths of the confidence bands and the values of multiplicity correction obtained with bootstrap and with Monte Carlo procedures. The experiments confirm that the multiplier bootstrap and the bootstrap multiplicity correction become conservative if the local parametric model is considerably misspecified.

The results given here are valid on a random set of probability  $1 - Ce^{-x}$  for some explicit constant  $C > 0$ . The number  $x > 0$  determines this dominating probability level. For the case of the i.i.d. observations (see Secion 5.3) we take  $x = C \log n$ . Throughout the text  $\|\cdot\|$  denotes the Euclidean norm for a vector and spectral norm for a matrix.  $\|\cdot\|_{\max}$  is the maximal absolute value of elements of a vector (or a matrix),  $p_{\text{sum}} \stackrel{\text{def}}{=} p_1 + \dots + p_K$ ,  $p_{\max} \stackrel{\text{def}}{=} \max_{1 \leq k \leq K} p_k$ .

## 2 The multiplier bootstrap procedure

Let  $\ell_{i,k}(\boldsymbol{\theta})$  denote the log-density from the  $k$ -th parametric distribution family evaluated at the  $i$ -th observation:

$$\ell_{i,k}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \log \left( \frac{d\mathbb{P}_k(\boldsymbol{\theta})}{d\mu_0}(Y_i) \right), \quad (2.1)$$

then due to independence of  $Y_1, \dots, Y_n$

$$L_k(\boldsymbol{\theta}) = \sum_{i=1}^n \ell_{i,k}(\boldsymbol{\theta}) \quad \forall k = 1, \dots, K.$$

Consider i.i.d. scalar random variables  $u_i$  independent of the data  $\mathbf{Y}$ , s.t.  $\mathbb{E}u_i = 1$ ,  $\text{Var } u_i = 1$ ,  $\mathbb{E} \exp(u_i) < \infty$  (e.g.  $u_i \sim \mathcal{N}(1, 1)$  or  $u_i \sim \exp(1)$  or  $u_i \sim 2\text{Bernoulli}(0.5)$ ). Multiply the summands of the likelihood function  $L_k(\boldsymbol{\theta})$  with the new random variables:

$$L_k^\circ(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{i=1}^n \ell_{i,k}(\boldsymbol{\theta}) u_i, \quad (2.2)$$

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

Therefore, the quasi MLE for the  $\mathbf{Y}$ -world is a target parameter for the bootstrap world for each  $k = 1, \dots, K$ :

$$\text{argmax}_{\boldsymbol{\theta} \in \Theta_k} \mathbb{E}^\circ L_k^\circ(\boldsymbol{\theta}) = \text{argmax}_{\boldsymbol{\theta} \in \Theta_k} L_k(\boldsymbol{\theta}) = \tilde{\boldsymbol{\theta}}_k.$$

The corresponding bootstrap MLE is:

$$\tilde{\boldsymbol{\theta}}_k \stackrel{\text{def}}{=} \text{argmax}_{\boldsymbol{\theta} \in \Theta_k} L_k^\circ(\boldsymbol{\theta}).$$

The  $k$ -th likelihood ratio statistic in the bootstrap world equals to  $L_k^\circ(\tilde{\theta}_k^\circ) - L_k^\circ(\tilde{\theta}_k)$ , where all the elements: the function  $L_k^\circ(\theta)$  and the arguments  $\tilde{\theta}_k^\circ, \tilde{\theta}_k$  are known and available for computation. This means, that given the data  $\mathbf{Y}$ , one can estimate the distribution or quantiles of the statistic  $L_k^\circ(\tilde{\theta}_k^\circ) - L_k^\circ(\tilde{\theta}_k)$  by generating many independent samples of the bootstrap weights  $u_1, \dots, u_n$  and computing with them the bootstrap likelihood ratio.

Let us introduce similarly to (1.5) the  $(1 - \alpha)$ -quantile for the bootstrap square-root likelihood ratio statistic:

$$\mathfrak{z}_k^\circ(\alpha) \stackrel{\text{def}}{=} \inf \left\{ \mathfrak{z} \geq 0 : \mathbb{P}^\circ \left( L_k^\circ(\tilde{\theta}_k^\circ) - L_k^\circ(\tilde{\theta}_k) > \mathfrak{z}^2/2 \right) \leq \alpha \right\}, \quad (2.3)$$

here  $\mathbb{P}^\circ$  denotes probability measure conditional on the data  $\mathbf{Y}$ , therefore,  $\mathfrak{z}_k^\circ(\alpha)$  is a random value dependent on  $\mathbf{Y}$ .

Spokoiny and Zhilova (2014) considered the case of a single parametric model ( $K = 1$ ), and showed that the bootstrap quantile  $\mathfrak{z}_k^\circ(\alpha)$  is close to the true one  $\mathfrak{z}_k(\alpha)$  under a so called ‘‘Small Modeling Bias’’ (SmB) condition, which is fulfilled when the true distribution is close to the parametric family or when the observations are i.i.d. When the SmB condition does not hold, the bootstrap quantile is still valid, however, it becomes conservative. Therefore, for each fixed  $k = 1, \dots, K$  the bootstrap quantiles  $\mathfrak{z}_k^\circ(\alpha)$  are rather good estimates for the true unknown ones  $\mathfrak{z}_k(\alpha)$ , however, they are still ‘‘pointwise’’ in  $k$ , i.e. the confidence bounds (1.6) hold for each  $k$  separately. Our goal here is to estimate  $\mathfrak{z}_1(\alpha), \dots, \mathfrak{z}_K(\alpha)$  and  $\mathfrak{c}(\alpha)$  according to (1.7) and (1.8). Let us introduce the bootstrap correction for multiplicity:

$$\mathfrak{c}^\circ(\alpha) \stackrel{\text{def}}{=} \sup \left\{ c \in (0, \alpha] : \mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^\circ(\tilde{\theta}_k^\circ) - 2L_k^\circ(\tilde{\theta}_k)} > \mathfrak{z}_k^\circ(c) \right\} \right) \leq \alpha \right\}. \quad (2.4)$$

By its definition  $\mathfrak{c}^\circ(\alpha)$  depends on the random sample  $\mathbf{Y}$ .

The multiplier bootstrap procedure below explains how to estimate the bootstrap quantile functions  $\mathfrak{z}_k^\circ(\mathfrak{c}^\circ(\alpha))$  corrected for multiplicity.

---

### The simultaneous bootstrap procedure:

---

**Input:** The data  $\mathbf{Y}$  (as in (1.1)) and a fixed confidence level  $(1 - \alpha) \in (0, 1)$ .

**Step 1:** Generate  $B$  independent samples of i.i.d. bootstrap weights  $\{u_1^{(b)}, \dots, u_n^{(b)}\}$ ,  $b = 1, \dots, B$ . For the bootstrap likelihood processes

$$L_k^{\circ(b)}(\theta) \stackrel{\text{def}}{=} \sum_{i=1}^n \ell_{i,k}(\theta) u_i^{(b)}. \quad (2.5)$$

compute the bootstrap likelihood ratios  $L_k^{\circ(b)}(\theta_k^{\circ(b)}) - L_k^{\circ(b)}(\tilde{\theta}_k)$ . For each fixed  $b$  the bootstrap likelihoods  $L_1^{\circ(b)}(\theta), \dots, L_K^{\circ(b)}(\theta)$  are computed using

the same bootstrap sample  $\{u_i^{(b)}\}$ , s.t. the  $i$ -th summand  $\ell_{i,k}(\boldsymbol{\theta})$  is always multiplied with the  $i$ -th weight  $u_i^{(b)}$  as in (2.5).

**Step 2:** Estimate the marginal quantile functions  $\mathfrak{z}_k^\circ(\alpha)$  defined in (2.3) separately for each  $k = 1, \dots, K$ , using  $B$  bootstrap realisations of  $L_k^\circ(\tilde{\boldsymbol{\theta}}_k^\circ) - L_k^\circ(\tilde{\boldsymbol{\theta}}_k)$  from **Step 1**.

**Step 3:** Find by an iterative procedure the maximum value  $c \in (0, \alpha]$  s.t.

$$\mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^\circ(\tilde{\boldsymbol{\theta}}_k^\circ) - 2L_k^\circ(\tilde{\boldsymbol{\theta}}_k)} \geq \mathfrak{z}_k^\circ(c) \right\} \right) \leq \alpha.$$

**Output:** The resulting critical values are  $\mathfrak{z}_k^\circ(c)$ ,  $k = 1, \dots, K$ .

**Remark 2.1.** The requirement in Step 1 to use the same bootstrap sample  $\{u_i^{(b)}\}$  for generation of the bootstrap likelihood ratios  $L_k^{\circ(b)}(\boldsymbol{\theta}_k^{\circ(b)}) - L_k^{\circ(b)}(\tilde{\boldsymbol{\theta}}_k)$ ,  $k = 1, \dots, K$  allows to preserve the correlation structure between the ratios and, therefore, to make a sharper simultaneous adjustment in Step 3.

This procedure is justified theoretically in the next section.

### 3 Theoretical justification of the bootstrap procedure

Before stating the main results in Section 3.2 we introduce in Section 3.1 the basic ingredients of the proofs. The general scheme of the theoretical approach here is taken from Spokoiny and Zhilova (2014). In the present work we extend that approach for the case of simultaneously many parametric models.

#### 3.1 Overview of the theoretical approach

For justification of the described multiplier bootstrap procedure for simultaneous inference it has to be checked that the joint distributions of the sets of likelihood ratio statistics  $\{L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) : k = 1, \dots, K\}$  and  $\{L_k^\circ(\tilde{\boldsymbol{\theta}}_k^\circ) - L_k^\circ(\tilde{\boldsymbol{\theta}}_k) : k = 1, \dots, K\}$  are close to each other. These joint distributions are approximated using several non-asymptotic

steps given in the following scheme:

$$\begin{array}{l}
\mathbf{Y}\text{-world: } \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \\
\text{Bootstrap world: } \sqrt{2L_k^\circ(\tilde{\boldsymbol{\theta}}_k^\circ) - 2L_k^\circ(\tilde{\boldsymbol{\theta}}_k)}
\end{array}
\begin{array}{l}
\approx \\
\approx
\end{array}
\begin{array}{l}
\frac{p_k + \log K}{\sqrt{n}} \\
\frac{p_k + \log K}{\sqrt{n}}
\end{array}
\begin{array}{l}
\|\boldsymbol{\xi}_k\| \\
\|\boldsymbol{\xi}_k^\circ\|
\end{array}
\begin{array}{l}
\approx \\
\approx
\end{array}
\begin{array}{l}
\|\bar{\boldsymbol{\xi}}_k\| \\
\|\bar{\boldsymbol{\xi}}_k^\circ\|,
\end{array}
\begin{array}{l}
\text{uniform} \\
\text{sq-Wilks} \\
\text{theorem}
\end{array}
\begin{array}{l}
\text{joint Gauss.} \\
\text{approx. \&} \\
\text{anti-concentr.*}
\end{array}
\begin{array}{l}
\bigcap_{1 \leq k \leq K} \\
\bigcap_{1 \leq k \leq K}
\end{array}
\begin{array}{l}
\rightsquigarrow w \\
\rightsquigarrow w
\end{array}
\begin{array}{l}
\text{simultaneous} \\
\text{Gauss. compar.**}
\end{array}
\quad (3.1)$$

\* the accuracy of these approximating steps is  $\mathcal{C} \left\{ \frac{p_{\max}^3}{n} \log^9(K) \log^3(np_{\text{sum}}) \right\}^{1/8}$ ;

\*\* Gaussian comparison step yields an approximation error proportional to  $\hat{\delta}_{\text{smb}}^2 \left( \frac{p_{\max}^3}{n} \right)^{1/4} p_{\max} \log^2(K) \log^{3/4}(np_{\text{sum}})$ , where  $\hat{\delta}_{\text{smb}}^2$  comes from condition **(SmB)**, see also (3.4) below.

Here  $\boldsymbol{\xi}_k$  and  $\boldsymbol{\xi}_k^\circ$  denote normalized score vectors for the  $\mathbf{Y}$  and bootstrap likelihood processes:

$$\boldsymbol{\xi}_k \stackrel{\text{def}}{=} D_k^{-1} \nabla_{\boldsymbol{\theta}} L_k(\boldsymbol{\theta}_k^*), \quad \boldsymbol{\xi}_k^\circ \stackrel{\text{def}}{=} \boldsymbol{\xi}_k^\circ(\boldsymbol{\theta}_k^*) \stackrel{\text{def}}{=} D_k^{-1} \nabla_{\boldsymbol{\theta}} L_k(\boldsymbol{\theta}_k^*), \quad (3.2)$$

$D_k^2$  is the full Fisher information matrix for the corresponding  $k$ -th likelihood:

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

$\bar{\boldsymbol{\xi}}_k \sim \mathcal{N}(0, \text{Var } \boldsymbol{\xi}_k)$  and  $\bar{\boldsymbol{\xi}}_k^\circ \sim \mathcal{N}(0, \text{Var}^\circ \boldsymbol{\xi}_k^\circ)$  denote approximating Gaussian vectors, which have the same covariance matrices as  $\boldsymbol{\xi}$  and  $\boldsymbol{\xi}^\circ$ . Moreover the vectors  $(\bar{\boldsymbol{\xi}}_1^\top, \dots, \bar{\boldsymbol{\xi}}_K^\top)^\top$  and  $(\bar{\boldsymbol{\xi}}_1^{\circ\top}, \dots, \bar{\boldsymbol{\xi}}_K^{\circ\top})^\top$  are normally distributed and have the same covariance matrices as the vectors  $(\boldsymbol{\xi}_1^\top, \dots, \boldsymbol{\xi}_K^\top)^\top$  and  $(\boldsymbol{\xi}_1^{\circ\top}, \dots, \boldsymbol{\xi}_K^{\circ\top})^\top$  correspondingly.  $\text{Var}^\circ$  and  $\text{Cov}^\circ$  denote variance and covariance operators w.r.t. the probability measure  $\mathbb{P}^\circ$  conditional on  $\mathbf{Y}$ .

The first two approximating steps: square root Wilks and Gaussian approximations are performed in parallel for both  $\mathbf{Y}$  and bootstrap worlds, which is shown in the corresponding lines of the scheme (3.1). The two worlds are connected in the last step: Gaussian comparison for  $\ell_2$ -norms of Gaussian vectors. All the approximations are performed simultaneously for  $K$  parametric models.

Let us consider each step in more details. Non-asymptotic square-root Wilks approximation result had been obtained recently by Spokoiny (2012a, 2013). It says that for

a fixed sample size and misspecified parametric assumption:  $\mathbb{P} \notin \{\mathbb{P}_k\}$ , it holds with exponentially high probability:

$$\left| \sqrt{2\{L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*)\}} - \|\boldsymbol{\xi}_k\| \right| \leq \Delta_{k,W} \simeq \frac{p_k}{\sqrt{n}},$$

here the index  $k$  is fixed, i.e. this statement is for one parametric model. The precise statement of this result is given in Section B.1, and its simultaneous version – in Section B.3. The approximating value  $\|\boldsymbol{\xi}_k\|$  is  $\ell_2$ -norm of the score vector  $\boldsymbol{\xi}_k$  given in (3.2). The next approximating step is between the joint distributions of  $\|\boldsymbol{\xi}_1\|, \dots, \|\boldsymbol{\xi}_K\|$  and  $\|\bar{\boldsymbol{\xi}}_1\|, \dots, \|\bar{\boldsymbol{\xi}}_K\|$ . This is done in Section A.1 for general centered random vectors under bounded exponential moments assumptions. The main tools for the simultaneous Gaussian approximation are: Lindeberg's telescopic sum, smooth maximum function and three times differentiable approximation of the indicator function  $\mathbb{I}\{x \in \mathbb{R} : x > 0\}$ . The simultaneous anti-concentration inequality for the  $\ell_2$ -norms of Gaussian vectors is obtained in Section A.3. The result is based on approximation of the  $\ell_2$ -norm with a maximum over a finite grid on a hypersphere, and on the anti-concentration inequality for maxima of a Gaussian random vector by Chernozhukov et al. (2014c). The same approximating steps are performed for the bootstrap world, the square-root bootstrap Wilks approximation is given in Sections B.2, B.3. The last step in the scheme (3.1) is comparison of the joint distributions of the sets of  $\ell_2$ -norms of Gaussian vectors:  $\|\bar{\boldsymbol{\xi}}_1\|, \dots, \|\bar{\boldsymbol{\xi}}_K\|$  and  $\|\bar{\boldsymbol{\xi}}_1^\circ\|, \dots, \|\bar{\boldsymbol{\xi}}_K^\circ\|$  by Slepian interpolation (see Section A.2 for the result in a general setting). The error of approximation is proportional to

$$\max_{1 \leq k_1, k_2 \leq K} \left\| \text{Cov}(\boldsymbol{\xi}_{k_1}, \boldsymbol{\xi}_{k_2}) - \text{Cov}^\circ(\boldsymbol{\xi}_{k_1}^\circ, \boldsymbol{\xi}_{k_2}^\circ) \right\|_{\max}. \quad (3.3)$$

It is shown, using Bernstein matrix inequality (Sections C.1 and C.3), that the value (3.3) is bounded from above (up to a constant) on a random set of dominating probability with

$$\max_{1 \leq k \leq K} \|H_k^{-1} B_k^2 H_k^{-1}\| \leq \widehat{\delta}_{\text{smb}}^2 \quad (3.4)$$

for

$$\begin{aligned} B_k^2 &\stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} \{ \nabla_{\boldsymbol{\theta}} \ell_{i,k}(\boldsymbol{\theta}_k^*) \} \mathbb{E} \{ \nabla_{\boldsymbol{\theta}} \ell_{i,k}(\boldsymbol{\theta}_k^*) \}^\top, \\ H_k^2 &\stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} \left\{ \nabla_{\boldsymbol{\theta}} \ell_{i,k}(\boldsymbol{\theta}_k^*) \nabla_{\boldsymbol{\theta}} \ell_{i,k}(\boldsymbol{\theta}_k^*)^\top \right\}. \end{aligned} \quad (3.5)$$

The value  $\|H_k^{-1} B_k^2 H_k^{-1}\|$  is responsible for the modelling bias of the  $k$ -th model. If the parametric family  $\{\mathbb{P}_k(\boldsymbol{\theta})\}$  contains the true distribution  $\mathbb{P}$  or if the observations  $Y_i$  are i.i.d., then  $B_k^2$  equals to zero. Condition  $(\widehat{\text{SmB}})$  assumes that all the values  $\|H_k^{-1} B_k^2 H_k^{-1}\|$  are rather small.

### 3.2 Main results

The following theorem shows the closeness of the joint cumulative distribution functions (c.d.f.s.) of  $\left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)}, k = 1, \dots, K \right\}$  and  $\left\{ \sqrt{2L_k^\circ(\tilde{\boldsymbol{\theta}}_k^\circ) - 2L_k^\circ(\tilde{\boldsymbol{\theta}}_k)}, k = 1, \dots, K \right\}$ . The approximating error term  $\Delta_{\text{total}}$  equals to a sum of the errors from all the steps in the scheme (3.1).

**Theorem 3.1.** *Under the conditions of Section 5 it holds with probability  $\geq 1 - 12e^{-x}$  for  $z_k \geq C\sqrt{p_k}$ ,  $1 \leq C < 2$*

$$\left| \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} > z_k \right\} \right) - \mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^\circ(\tilde{\boldsymbol{\theta}}_k^\circ) - 2L_k^\circ(\tilde{\boldsymbol{\theta}}_k)} > z_k \right\} \right) \right| \leq \Delta_{\text{total}}.$$

The approximating total error  $\Delta_{\text{total}} \geq 0$  is deterministic and in the case of i.i.d. observations (see Section 5.3) it holds:

$$\Delta_{\text{total}} \leq \mathfrak{c} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \left\{ (\hat{\mathbf{a}}^2 + \hat{\mathbf{a}}_B^2) (1 + \delta_{\mathbb{V}}^2(\mathbf{x})) \right\}^{3/8}, \quad (3.6)$$

where the deterministic terms  $\hat{\mathbf{a}}^2$ ,  $\hat{\mathbf{a}}_B^2$  and  $\delta_{\mathbb{V}}^2(\mathbf{x})$  come from the conditions  $(\mathcal{I})$ ,  $(\mathcal{I}_B)$  and  $(\widehat{\mathcal{SD}}_1)$ .  $\Delta_{\text{total}}$  is defined in (C.5).

**Remark 3.1.** The obtained approximation bound is mainly of theoretical interest, although it shows the impact of  $p_{\max}$ ,  $K$  and  $n$  on the quality of the bootstrap procedure. For more details on the error term see Remark A.1.

The next theorem justifies the bootstrap procedure under the  $(\widehat{\mathcal{SmB}})$  condition. The theorem says that the bootstrap quantile functions  $\mathfrak{z}_k^\circ(\cdot)$  with the bootstrap-corrected for multiplicity confidence levels  $1 - \mathfrak{c}^\circ(\alpha)$  can be used for construction of the simultaneous confidence set in the  $\mathbf{Y}$ -world.

**Theorem 3.2** (Bootstrap validity for a small modeling bias). *Assume the conditions of Theorem 3.1, and  $\mathfrak{c}(\alpha), 0.5\mathfrak{c}^\circ(\alpha) \geq \Delta_{\text{full, max}}$ , then for  $\alpha \leq 1 - 8e^{-x}$  it holds with probability  $1 - 12e^{-x}$*

$$\begin{aligned} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq \mathfrak{z}_k^\circ(\mathfrak{c}^\circ(\alpha) - 2\Delta_{\text{full, max}}) \right\} \right) - \alpha &\leq \Delta_{\mathfrak{z}, \text{total}}, \\ \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq \mathfrak{z}_k^\circ(\mathfrak{c}^\circ(\alpha) + 2\Delta_{\text{full, max}}) \right\} \right) - \alpha &\geq -\Delta_{\mathfrak{z}, \text{total}}, \end{aligned}$$

where  $\Delta_{\text{full, max}} \leq \mathfrak{C}\{(p_{\max} + \mathbf{x})^3/n\}^{1/8}$  in the case of i.i.d. observations (see Section 5.3), and  $\Delta_{\mathfrak{z}, \text{total}} \leq 3\Delta_{\text{total}}$ ; their explicit definitions are given in (C.11) and (C.14).

Moreover

$$\begin{aligned} \mathbf{c}^\circ(\alpha) &\leq \mathbf{c}(\alpha + \Delta_{\mathbf{c}}) + \Delta_{\text{full, max}}, \\ \mathbf{c}^\circ(\alpha) &\geq \mathbf{c}(\alpha - \Delta_{\mathbf{c}}) - \Delta_{\text{full, max}}, \end{aligned}$$

for  $0 \leq \Delta_{\mathbf{c}} \leq 2\Delta_{\text{total}}$ , defined in (C.15).

The following theorem does not assume the  $(\widehat{\mathbf{SmB}})$  condition to be fulfilled. It turns out that in this case the bootstrap procedure becomes conservative, and the bootstrap critical values corrected for the multiplicity  $\mathfrak{z}_k^\circ(\mathbf{c}^\circ(\alpha))$  are increased with the modelling bias  $\sqrt{\text{tr}\{D_k^{-1}H_k^2D_k^{-1}\}} - \sqrt{\text{tr}\{D_k^{-1}(H_k^2 - B_k^2)D_k^{-1}\}}$ , therefore, the confidence set based on the bootstrap estimates can be conservative.

**Theorem 3.3** (Bootstrap conservativeness for a large modeling bias). *Under the conditions of Section 5 except for  $(\widehat{\mathbf{SmB}})$  it holds with probability  $\geq 1 - 14e^{-x}$  for  $z_k \geq C\sqrt{p_k}$ ,  $1 \leq C < 2$*

$$\begin{aligned} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k)} - 2L_k(\boldsymbol{\theta}_k^*) > z_k \right\} \right) \\ \leq \mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^\circ(\tilde{\boldsymbol{\theta}}_k^\circ)} - 2L_k^\circ(\tilde{\boldsymbol{\theta}}_k) > z_k \right\} \right) + \Delta_{\text{b, total}}. \end{aligned}$$

The deterministic value  $\Delta_{\text{b, total}} \in [0, \Delta_{\text{total}}]$  (see (3.6) in the case 5.3). Moreover, the bootstrap-corrected for multiplicity confidence level  $1 - \mathbf{c}^\circ(\alpha)$  is conservative in comparison with the true corrected confidence level:

$$1 - \mathbf{c}^\circ(\alpha) \geq 1 - \mathbf{c}(\alpha + \Delta_{\text{b,c}}) - \Delta_{\text{full, max}},$$

and it holds for all  $k = 1, \dots, K$  and  $\alpha \leq 1 - 8e^{-x}$

$$\begin{aligned} \mathfrak{z}_k^\circ(\mathbf{c}^\circ(\alpha)) &\geq \mathfrak{z}_k(\mathbf{c}(\alpha + \Delta_{\text{b,c}}) + \Delta_{\text{full, max}}) \\ &\quad + \sqrt{\text{tr}\{D_k^{-1}H_k^2D_k^{-1}\}} - \sqrt{\text{tr}\{D_k^{-1}(H_k^2 - B_k^2)D_k^{-1}\}} - \Delta_{\text{qf},1,k}, \end{aligned}$$

for  $0 \leq \Delta_{\text{b,c}} \leq 2\Delta_{\text{total}}$ , defined in (C.18), and the positive value  $\Delta_{\text{qf},1,k}$  is bounded from above with  $(\mathbf{a}_k^2 + \mathbf{a}_{B,k}^2)(\sqrt{8xp_k} + 6x)$  for the constants  $\mathbf{a}_k^2 > 0$ ,  $\mathbf{a}_{B,k}^2 \geq 0$  from conditions  $(\mathcal{I})$ ,  $(\mathcal{I}_B)$ .

The  $(\widehat{\mathbf{SmB}})$  condition is automatically fulfilled if all the parametric models are correct or in the case of i.i.d. observations. This condition is checked for generalised linear model and linear quantile regression in Spokoiny and Zhilova (2014) (the version of 2015).



## 4 Numerical experiments

Here we check the performance of the bootstrap procedure by constructing simultaneous confidence sets based on the local constant and local quadratic estimates, the former one is also known as Nadaraya-Watson estimate [Nadaraya \(1964\)](#); [Watson \(1964\)](#). Let  $Y_1, \dots, Y_n$  be independent random scalar observations and  $X_1, \dots, X_n$  some deterministic design points. In [Sections 4.1-4.3](#) below we introduce the models and the data, [Sections 4.4-4.6](#) present the results of the experiments.

### 4.1 Local constant regression

Consider the following quadratic likelihood function reweighted with the kernel functions  $K(\cdot)$ :

$$\begin{aligned} L(\boldsymbol{\theta}, x, h) &\stackrel{\text{def}}{=} -\frac{1}{2} \sum_{i=1}^n (Y_i - \boldsymbol{\theta})^2 w_i(x, h), \\ w_i(x, h) &\stackrel{\text{def}}{=} K(\{x - X_i\}/h), \\ K(x) &\in [0, 1], \quad \int_{\mathbb{R}} K(x) dx = 1, \quad K(x) = K(-x). \end{aligned}$$

Here  $h > 0$  denotes bandwidth, the local smoothing parameter. The target point and the local MLE read as:

$$\boldsymbol{\theta}^*(x, h) \stackrel{\text{def}}{=} \frac{\sum_{i=1}^n w_i(x, h) \mathbb{E}Y_i}{\sum_{i=1}^n w_i(x, h)}, \quad \tilde{\boldsymbol{\theta}}(x, h) \stackrel{\text{def}}{=} \frac{\sum_{i=1}^n w_i(x, h) Y_i}{\sum_{i=1}^n w_i(x, h)}.$$

Let us fix a bandwidth  $h$  and consider the range of points  $x_1, \dots, x_K$ . They yield  $K$  local constant models with the target parameters  $\boldsymbol{\theta}_k^* \stackrel{\text{def}}{=} \boldsymbol{\theta}^*(x_k, h)$  and the likelihood functions  $L_k(\boldsymbol{\theta}) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}, x_k, h)$  for  $k = 1, \dots, K$ .

The bootstrap local likelihood function is defined similarly to the global one [\(2.2\)](#), by reweighting  $L(\boldsymbol{\theta}, x, h)$  with the bootstrap multipliers  $u_1, \dots, u_n$ :

$$\begin{aligned} L_k^\circ(\boldsymbol{\theta}) &\stackrel{\text{def}}{=} L^\circ(\boldsymbol{\theta}, x_k, h) \stackrel{\text{def}}{=} -\frac{1}{2} \sum_{i=1}^n (Y_i - \boldsymbol{\theta})^2 w_i(x_k, h) u_i, \\ \tilde{\boldsymbol{\theta}}_k^\circ &\stackrel{\text{def}}{=} \tilde{\boldsymbol{\theta}}^\circ(x_k, h) \stackrel{\text{def}}{=} \frac{\sum_{i=1}^n w_i(x_k, h) u_i Y_i}{\sum_{i=1}^n w_i(x_k, h) u_i}. \end{aligned}$$

### 4.2 Local quadratic regression

Here the local likelihood function reads as

$$\begin{aligned} L(\boldsymbol{\theta}, x, h) &\stackrel{\text{def}}{=} -\frac{1}{2} \sum_{i=1}^n (Y_i - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta})^2 w_i(x, h), \\ \boldsymbol{\theta}, \boldsymbol{\Psi}_i &\in \mathbb{R}^3, \quad \boldsymbol{\Psi}_i \stackrel{\text{def}}{=} (1, X_i, X_i^2)^\top, \end{aligned}$$

and

$$\begin{aligned}\boldsymbol{\theta}^*(x, h) &\stackrel{\text{def}}{=} \left( \Psi W(x, h) \Psi^\top \right)^{-1} \Psi W(x, h) \mathbb{E} \mathbf{Y}, \\ \tilde{\boldsymbol{\theta}}(x, h) &\stackrel{\text{def}}{=} \left( \Psi W(x, h) \Psi^\top \right)^{-1} \Psi W(x, h) \mathbf{Y},\end{aligned}$$

where

$$\begin{aligned}\mathbf{Y} &\stackrel{\text{def}}{=} (Y_1, \dots, Y_n)^\top, \quad \Psi \stackrel{\text{def}}{=} (\Psi_1, \dots, \Psi_n) \in \mathbb{R}^{3 \times n}, \\ W(x, h) &\stackrel{\text{def}}{=} \text{diag} \{w_1(x, h), \dots, w_n(x, h)\}.\end{aligned}$$

And similarly for the bootstrap objects

$$\begin{aligned}L^\circ(\boldsymbol{\theta}, x, h) &\stackrel{\text{def}}{=} -\frac{1}{2} \sum_{i=1}^n (Y_i - \Psi_i^\top \boldsymbol{\theta})^2 w_i(x, h) u_i, \\ \tilde{\boldsymbol{\theta}}^\circ(x, h) &\stackrel{\text{def}}{=} \left( \Psi U W(x, h) \Psi^\top \right)^{-1} \Psi U W(x, h) \mathbf{Y},\end{aligned}$$

for  $U \stackrel{\text{def}}{=} \text{diag} \{u_1, \dots, u_n\}$ .

### 4.3 Simulated data

In the numerical experiments we constructed two 90% simultaneous confidence bands: using Monte Carlo (MC) samples and bootstrap procedure with Gaussian weights ( $u_i \sim \mathcal{N}(1, 1)$ ), in each case we used  $10^4$   $\{Y_i\}$  and  $10^4$   $\{u_i\}$  independent samples. The sample size  $n = 400$ .  $K(x)$  is Epanechnikov's kernel function. The independent random observations  $Y_i$  are generated as follows:

$$Y_i = f(X_i) + \mathcal{N}(0, 1), \quad X_i \text{ are equidistant on } [0, 1], \quad (4.1)$$

$$f(x) = \begin{cases} 5, & x \in [0, 0.25] \cup [0.65, 1]; \\ 5 + 3.8\{1 - 100(x - 0.35)^2\}, & x \in [0.25, 0.45]; \\ 5 - 3.8\{1 - 100(x - 0.55)^2\}, & x \in [0.45, 0.65]. \end{cases} \quad (4.2)$$

The number of local models  $K = 71$ , the points  $x_1, \dots, x_{71}$  are equidistant on  $[0, 1]$ . For the bandwidth we considered two cases:  $h = 0.12$  and  $h = 0.3$ .

### 4.4 Effect of the modeling bias on a width of a bootstrap confidence band

The function  $f(x)$  defined in (4.2) should yield a considerable modeling bias for both mean constant and mean quadratic estimators. Figures 4.1, 4.2 demonstrate that the bootstrap confidence bands become conservative (i.e. wider than the MC confidence

band) when the local model is misspecified. The top graphs on Figures 4.1, 4.2 show the 90% confidence bands, the middle graphs show their width, and the bottom graphs show the value of the modelling bias for  $K = 71$  local models (see formulas (4.3) and (4.4) below). For the local constant estimate (Figure 4.1) the width of the bootstrap confidence sets is considerably increased by the modeling bias when  $x \in [0.25, 0.65]$ . In this case the expression for the modeling bias term for the  $k$ -th model (see also **(SmB)** condition) reads as:

$$\begin{aligned} |H_k^{-1}B_k^2H_k^{-1}| &= \frac{\sum_{i=1}^n \{\mathbb{E}Y_i - \boldsymbol{\theta}^*(x_k)\}^2 w_i^2(x_k, h)}{\sum_{i=1}^n \mathbb{E}\{Y_i - \boldsymbol{\theta}^*(x_k)\}^2 w_i^2(x_k, h)} \\ &= 1 - \left(1 + \frac{\sum_{i=1}^n w_i^2(x_k, h) \{f(X_i) - \boldsymbol{\theta}^*(x_k)\}^2}{\sum_{i=1}^n w_i^2(x_k, h)}\right)^{-1}. \end{aligned} \quad (4.3)$$

And for the local quadratic estimate it holds:

$$\|H_k^{-1}B_k^2H_k^{-1}\| = \left\| \mathbf{I}_p - H_k^{-1} \left\{ \sum_{i=1}^n \Psi_i \Psi_i^\top w_i^2(x_k, h) \right\} H_k^{-1} \right\|, \quad (4.4)$$

where  $\mathbf{I}_p$  is the identity matrix of dimension  $p \times p$  (here  $p = 3$ ), and

$$\begin{aligned} H_k^2 &= \sum_{i=1}^n \Psi_i \Psi_i^\top w_i^2(x_k, h) \mathbb{E}\{Y_i - \boldsymbol{\theta}^*(x_k)\}^2 \\ &= \sum_{i=1}^n \Psi_i \Psi_i^\top w_i^2(x_k, h) \{f(X_i) - \boldsymbol{\theta}^*(x_k)\}^2 + \sum_{i=1}^n \Psi_i \Psi_i^\top w_i^2(x_k, h). \end{aligned} \quad (4.5)$$

Therefore, if  $\max_{1 \leq k \leq K} \{f(X_i) - \boldsymbol{\theta}^*(x_k)\}^2 = 0$ , then  $\|H_k^{-1}B_k^2H_k^{-1}\| = 0$ . On the Figure 4.1 both the modelling bias and the difference between the widths of the bootstrap and MC confidence bands are close to zero in the regions where the true function  $f(x)$  is constant. On Figure 4.2 the modelling bias for  $h = 0.12$  is overall smaller than the corresponding value on Figure 4.1. For the bigger bandwidth  $h = 0.3$  the modelling biases on Figures 4.1 and 4.2 are comparable with each other.

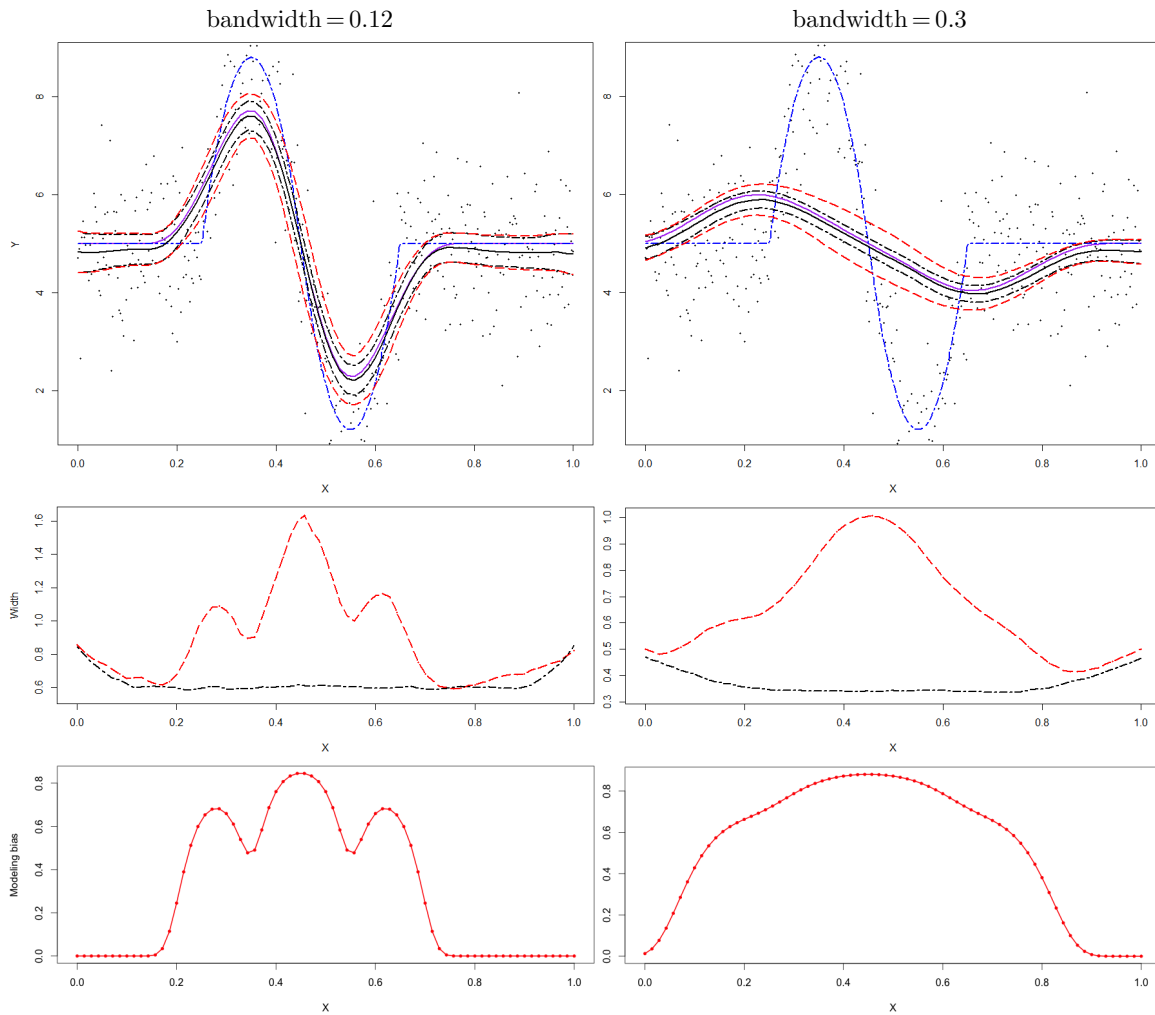
Thus the numerical experiment is consistent with the theoretical results from Section 3.2, and confirm that in the case when a (local) parametric model is close to the true distribution the simultaneous bootstrap confidence set is valid. Otherwise the bootstrap procedure is conservative: the modelling bias widens the simultaneous bootstrap confidence set.

#### 4.5 Effective coverage probability (local constant estimate)

In this part of the experiment we check the bootstrap validity by computing the effective coverage probability values. This requires to perform many independent experiments: for each of independent 5000  $\{Y_i\} \sim (4.1)$  samples we took  $10^4$  independent bootstrap samples  $\{u_i\} \sim \mathcal{N}(1, 1)$ , and constructed simultaneous bootstrap confidence sets for a range of confidence levels. The second row of Table 4.1 contains this range  $(1 - \alpha) =$

Figure 4.1: **Local constant regression:**

Confidence bands, their widths, and the modeling bias



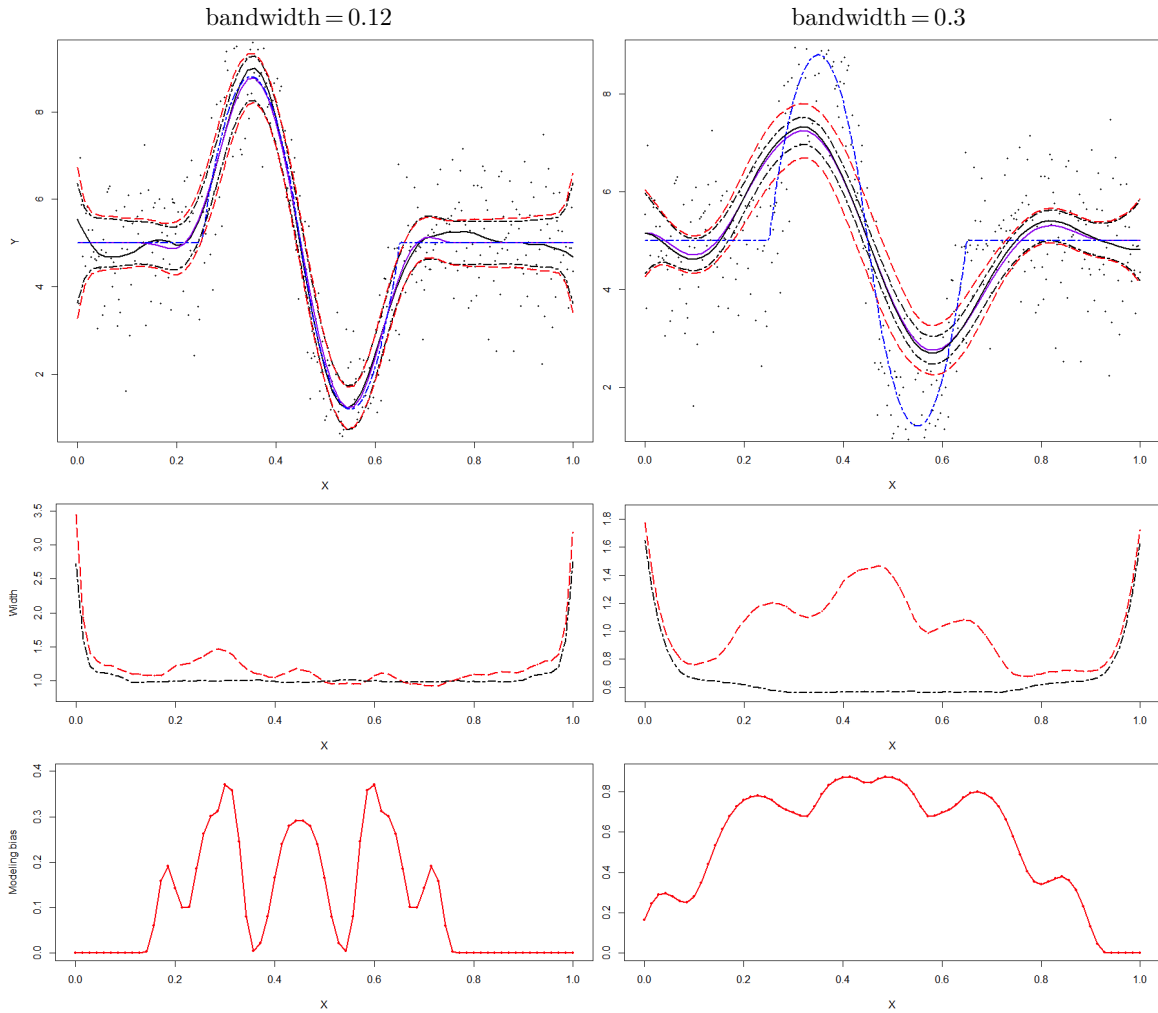
Legend for the top graphs:

- 90% bootstrap simultaneous confidence band
- 90% MC simultaneous confidence band
- smoothed target function
- - - - the true function  $f(x)$
- local constant MLE

Legend for the middle and the bottom graphs:

- width of the 90% bootstrap confidence bands from the upper graphs
- width of the 90% MC confidence bands from the upper graphs
- modeling bias from the expression (4.3)

Figure 4.2: **Local quadratic regression:**  
Confidence bands, their widths, and the modeling bias



Legend for the top graphs:

- 90% bootstrap simultaneous confidence band
- 90% MC simultaneous confidence band
- smoothed target function
- - - - the true function  $f(x)$
- local constant MLE

Legend for the middle and the bottom graphs:

- width of the 90% bootstrap confidence bands from the upper graphs
- width of the 90% MC confidence bands from the upper graphs
- modeling bias from the expression (4.4)

0.95, 0.9, ..., 0.5. The third and the fourth rows of Table 4.1 show the frequencies of the event

$$\max_{1 \leq k \leq K} \left\{ L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) - \mathfrak{z}_k^\circ(\mathbf{c}^\circ(\alpha)) \right\} \leq 0$$

among 5000 data samples, for the bandwidths  $h = 0.12, 0.3$ , and for the range of  $(1 - \alpha)$ . The results show that the bootstrap procedure is rather conservative for both  $h = 0.12$  and  $h = 0.3$ , however, the larger bandwidth yields bigger coverage probabilities.

Table 1: Effective coverage probabilities for the local constant regression

$h$	Confidence levels									
	<b>0.95</b>	<b>0.90</b>	<b>0.85</b>	<b>0.80</b>	<b>0.75</b>	<b>0.70</b>	<b>0.65</b>	<b>0.60</b>	<b>0.55</b>	<b>0.50</b>
0.12	0.971	0.947	0.917	0.888	0.863	0.830	0.800	0.769	0.738	0.702
0.3	0.982	0.963	0.942	0.918	0.895	0.868	0.842	0.815	0.784	0.750

#### 4.6 Correction for multiplicity

Here we compare the  $\mathbf{Y}$  and the bootstrap corrections for multiplicity, i.e. the values  $\mathbf{c}(\alpha)$  and  $\mathbf{c}^\circ(\alpha)$  defined in (1.8) and (2.4). The numerical results in Tables 2, 3 are based on  $10^4 \{Y_i\} \sim (4.1)$  independent samples and  $10^4$  independent bootstrap samples  $\{u_i\} \sim \mathcal{N}(1, 1)$ . The second line in Tables 2, 3 contains the range of the nominal confidence levels  $(1 - \alpha) = 0.95, 0.9, \dots, 0.5$  (similarly to the Table 1). The first column contains the values of the bandwidth  $h = 0.12, 0.3$ , and the second column – the resampling scheme: Monte Carlo (MC) or bootstrap (B). The Monte Carlo experiment yields the corrected confidence levels  $1 - \mathbf{c}(\alpha)$ , and the bootstrap yields  $1 - \mathbf{c}^\circ(\alpha)$ . The lines 3–6 contain the average values of  $1 - \mathbf{c}(\alpha)$  and  $1 - \mathbf{c}^\circ(\alpha)$  over all the experiments. The results show that for the smaller bandwidth both the MC and bootstrap corrections are bigger than the ones for the larger bandwidth. In the case of a smaller bandwidth the local models have less intersections with each other, and hence, the corrections for multiplicity are closer to the Bonferroni's bound.

**Remark 4.1.** The theoretical results of this paper can be extended to the case when a set of considered local models has cardinality of the continuum, and the confidence bands are uniform w.r.t. the local parameter. This extension would require some uniform statements such as locally uniform square-root Wilks approximation (see e.g. Spokoiny and Zhilova (2013)).

**Remark 4.2.** The use of the bootstrap procedure in the problem of choosing an optimal bandwidth is considered in Spokoiny and Willrich (2015).

Table 2: **Local constant regression:**

MC vs Bootstrap confidence levels corrected for multiplicity

		Confidence levels									
$h$	r.m.	<b>0.95</b>	<b>0.90</b>	<b>0.85</b>	<b>0.80</b>	<b>0.75</b>	<b>0.70</b>	<b>0.65</b>	<b>0.60</b>	<b>0.55</b>	<b>0.50</b>
0.12	MC	0.997	0.994	0.989	0.985	0.980	0.975	0.969	0.963	0.956	0.949
	B	0.998	0.995	0.991	0.988	0.984	0.979	0.975	0.969	0.963	0.957
0.3	MC	0.993	0.983	0.973	0.962	0.949	0.936	0.922	0.906	0.891	0.873
	B	0.994	0.986	0.977	0.968	0.958	0.947	0.935	0.922	0.908	0.893

Table 3: **Local quadratic regression:**

MC vs Bootstrap confidence levels corrected for multiplicity

		Confidence levels									
$h$	r.m.	<b>0.95</b>	<b>0.90</b>	<b>0.85</b>	<b>0.80</b>	<b>0.75</b>	<b>0.70</b>	<b>0.65</b>	<b>0.60</b>	<b>0.55</b>	<b>0.50</b>
0.12	MC	0.997	0.993	0.989	0.985	0.979	0.974	0.968	0.961	0.954	0.946
	B	0.998	0.995	0.991	0.988	0.984	0.979	0.974	0.969	0.963	0.956
0.3	MC	0.993	0.983	0.973	0.961	0.949	0.936	0.921	0.904	0.887	0.868
	B	0.996	0.991	0.985	0.978	0.971	0.963	0.954	0.944	0.934	0.923

## 5 Conditions

Here we show necessary conditions for the main results. The conditions in Section 5.1 come from the general finite sample theory by Spokoiny (2012a), they are required for the results of Sections B.1 and B.2. The conditions in Section 5.2 are necessary to prove the statements on multiplier bootstrap validity.

### 5.1 Basic conditions

Introduce the stochastic part of the  $k$ -th likelihood process:  $\zeta_k(\boldsymbol{\theta}) \stackrel{\text{def}}{=} L_k(\boldsymbol{\theta}) - \mathbb{E}L_k(\boldsymbol{\theta})$ , and its marginal summand:  $\zeta_{i,k}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \ell_{i,k}(\boldsymbol{\theta}) - \mathbb{E}\ell_{i,k}(\boldsymbol{\theta})$  for  $\ell_{i,k}(\boldsymbol{\theta})$  defined in (2.1).

**(ED<sub>0</sub>)** For each  $k = 1, \dots, K$  there exist a positive-definite  $p_k \times p_k$  symmetric matrix  $V_k^2$  and constants  $\mathbf{g}_k > 0, \nu_k \geq 1$  such that  $\text{Var} \{ \nabla_{\boldsymbol{\theta}} \zeta_k(\boldsymbol{\theta}_k^*) \} \leq V_k^2$  and

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^{p_k}} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla_{\boldsymbol{\theta}} \zeta_k(\boldsymbol{\theta}_k^*)}{\|V_k \boldsymbol{\gamma}\|} \right\} \leq \nu_k^2 \lambda^2 / 2, \quad |\lambda| \leq \mathbf{g}_k.$$

**(ED<sub>2</sub>)** For each  $k = 1, \dots, K$  there exist a constant  $\omega_k > 0$  and for each  $\mathbf{r} > 0$  a

constant  $\mathfrak{g}_{2,k}(\mathbf{r})$  such that it holds for all  $\boldsymbol{\theta} \in \Theta_{0,k}(\mathbf{r})$  and for  $j = 1, 2$

$$\sup_{\substack{\boldsymbol{\gamma}_j \in \mathbb{R}^{p_k} \\ \|\boldsymbol{\gamma}_j\| \leq 1}} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega_k} \boldsymbol{\gamma}_1^\top D_k^{-1} \nabla_{\boldsymbol{\theta}}^2 \zeta_k(\boldsymbol{\theta}) D_k^{-1} \boldsymbol{\gamma}_2 \right\} \leq \nu_k^2 \lambda^2 / 2, \quad |\lambda| \leq \mathfrak{g}_{2,k}(\mathbf{r}).$$

**( $\mathcal{L}_0$ )** For each  $k = 1, \dots, K$  and for each  $\mathbf{r} > 0$  there exists a constant  $\delta_k(\mathbf{r}) \geq 0$  such that for  $\mathbf{r} \leq \mathbf{r}_{0,k}$  ( $\mathbf{r}_{0,k}$  come from condition **(B.1)** of Theorem **B.1** in Section **B.1**)  $\delta(\mathbf{r}) \leq 1/2$ , and for all  $\boldsymbol{\theta} \in \Theta_{0,k}(\mathbf{r})$  it holds

$$\|D_k^{-1} \check{D}_k^2(\boldsymbol{\theta}) D_k^{-1} - \mathbf{I}_{p_k}\| \leq \delta_k(\mathbf{r}),$$

where  $\check{D}_k^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E} L_k(\boldsymbol{\theta})$  and  $\Theta_{0,k}(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} \in \Theta_k : \|D_k(\boldsymbol{\theta} - \boldsymbol{\theta}_k^*)\| \leq \mathbf{r}\}$ .

**( $\mathcal{I}$ )** There exist constants  $\mathbf{a}_k > 0$  for all  $k = 1, \dots, K$  s.t.

$$\mathbf{a}_k^2 D_k^2 \geq V_k^2.$$

Denote  $\widehat{\mathbf{a}}^2 \stackrel{\text{def}}{=} \max_{1 \leq k \leq K} \mathbf{a}_k^2$ .

**( $\mathcal{L}_r$ )** For each  $k = 1, \dots, K$  and  $\mathbf{r} \geq \mathbf{r}_{0,k}$  there exists a value  $\mathbf{b}_k(\mathbf{r}) > 0$  s.t.

$\mathbf{r} \mathbf{b}_k(\mathbf{r}) \rightarrow \infty$  for  $\mathbf{r} \rightarrow \infty$  and  $\forall \boldsymbol{\theta} \in \Theta_k : \|D_k(\boldsymbol{\theta} - \boldsymbol{\theta}_k^*)\| = \mathbf{r}$  it holds

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

## 5.2 Conditions required for the bootstrap validity

**( $\widehat{\text{SmB}}$ )** There exists a constant  $\widehat{\delta}_{\text{smb}} \geq 0$  such that it holds for the matrices  $B_k^2$  and  $H_k^2$  defined in (3.5):

$$\begin{aligned} \max_{1 \leq k \leq K} \|H_k^{-1} B_k^2 H_k^{-1}\| &\leq \widehat{\delta}_{\text{smb}}^2, \\ \widehat{\delta}_{\text{smb}}^2 &\leq \mathfrak{c} \left( \frac{n}{p_{\max}^{13}} \right)^{1/8} \log^{-7/8}(K) \log^{-3/8}(np_{\text{sum}}). \end{aligned}$$

**( $\text{ED}_{2m}$ )** For each  $k = 1, \dots, K$ ,  $\mathbf{r} > 0$ ,  $i = 1, \dots, n$ ,  $j = 1, 2$  and for all  $\boldsymbol{\theta} \in \Theta_{0,k}(\mathbf{r})$  it holds for the values  $\omega_k \geq 0$  and  $\mathfrak{g}_{2,k}(\mathbf{r})$  from the condition **(ED<sub>2</sub>)**:

$$\sup_{\substack{\boldsymbol{\gamma}_j \in \mathbb{R}^{p_k} \\ \|\boldsymbol{\gamma}_j\| \leq 1}} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega_k} \boldsymbol{\gamma}_1^\top D_k^{-1} \nabla_{\boldsymbol{\theta}}^2 \zeta_{i,k}(\boldsymbol{\theta}) D_k^{-1} \boldsymbol{\gamma}_2 \right\} \leq \frac{\nu_0^2 \lambda^2}{2n}, \quad |\lambda| \leq \mathfrak{g}_{2,k}(\mathbf{r}),$$

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

$$\|D_k^{-1} \nabla_{\boldsymbol{\theta}}^2 \mathbb{E} \ell_{i,k}(\boldsymbol{\theta}) D_k^{-1}\| \leq \mathfrak{C}_{m,k}(\mathbf{r}) n^{-1}.$$



( $\mathcal{I}_B$ ) For each  $k = 1, \dots, K$  there exists a constant  $\mathfrak{a}_{B,k}^2 > 0$  s.t.

$$\mathfrak{a}_{B,k}^2 D_k^2 \geq B_k^2.$$

Denote  $\widehat{\mathfrak{a}}_B^2 \stackrel{\text{def}}{=} \max_{1 \leq k \leq K} \mathfrak{a}_{B,k}^2$ .

( $\widehat{SD}_1$ ) There exists a constant  $0 \leq \delta_{v^*}^2 \leq \mathfrak{C} p_{\text{sum}}/n$  such that it holds for all  $i = 1, \dots, n$  with exponentially high probability

$$\left\| \widehat{H}^{-1} \left\{ \mathbf{g}_i \mathbf{g}_i^\top - \mathbb{E} \left[ \mathbf{g}_i \mathbf{g}_i^\top \right] \right\} \widehat{H}^{-1} \right\| \leq \delta_{v^*}^2,$$

where

$$\mathbf{g}_i \stackrel{\text{def}}{=} \left( \nabla_{\boldsymbol{\theta}} \ell_{i,1}(\boldsymbol{\theta}_1^*)^\top, \dots, \nabla_{\boldsymbol{\theta}} \ell_{i,K}(\boldsymbol{\theta}_K^*)^\top \right)^\top \in \mathbb{R}^{p_{\text{sum}}},$$

$$\widehat{H}^2 \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} \left\{ \mathbf{g}_i \mathbf{g}_i^\top \right\},$$

$$p_{\text{sum}} \stackrel{\text{def}}{=} p_1 + \dots + p_K.$$

( $\mathbf{Eb}$ ) The i.i.d. bootstrap weights  $u_i$  are independent of  $\mathbf{Y}$ , and for all  $i = 1, \dots, n$  it holds for some constants  $\mathfrak{g}_k > 0, \nu_k \geq 1$

$$\mathbb{E} u_i = 1, \quad \text{Var } u_i = 1,$$

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

### 5.3 Dependence of the involved terms on the sample size and cardinality of the parameters' set

Here we consider the case of the i.i.d. observations  $Y_1, \dots, Y_n$  and  $\mathbf{x} = \mathfrak{C} \log n$  in order to specify the dependence of the non-asymptotic bounds on  $n$  and  $p$ . In the paper by [Spokoiny and Zhilova \(2014\)](#) (the version of 2015) this is done in detail for the i.i.d. case, generalized linear model and quantile regression.

Example 5.1 in [Spokoiny \(2012a\)](#) demonstrates that in this situation  $\mathfrak{g}_k = \mathfrak{C} \sqrt{n}$  and  $\omega_k = \mathfrak{C} / \sqrt{n}$ . then  $\mathfrak{z}_k(\mathbf{x}) = \mathfrak{C} \sqrt{p_k + \mathbf{x}}$  for some constant  $\mathfrak{C} \geq 1.85$ , for the function  $\mathfrak{z}_k(\mathbf{x})$  given in (B.3) in Section B.1. Similarly it can be checked that  $\mathfrak{g}_{2,k}(\mathbf{r})$  from condition ( $\mathbf{ED}_2$ ) is proportional to  $\sqrt{n}$ : due to independence of the observations

$$\begin{aligned} & \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega_k} \boldsymbol{\gamma}_1^\top D_k^{-1} \nabla_{\boldsymbol{\theta}}^2 \zeta_k(\boldsymbol{\theta}) D_k^{-1} \boldsymbol{\gamma}_2 \right\} \\ &= \sum_{i=1}^n \log \mathbb{E} \exp \left\{ \frac{\lambda}{\sqrt{n}} \frac{1}{\omega_k \sqrt{n}} \boldsymbol{\gamma}_1^\top d_k^{-1} \nabla_{\boldsymbol{\theta}}^2 \zeta_{i,k}(\boldsymbol{\theta}) d_k^{-1} \boldsymbol{\gamma}_2 \right\} \\ &\leq n \frac{\lambda^2}{n} \mathfrak{C} \quad \text{for } |\lambda| \leq \bar{\mathfrak{g}}_{2,k}(\mathbf{r}) \sqrt{n}, \end{aligned}$$

where  $\zeta_{i,k}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \ell_{i,k}(\boldsymbol{\theta}) - \mathbb{E}\ell_{i,k}(\boldsymbol{\theta})$ ,  $d_k^2 \stackrel{\text{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E}\ell_{i,k}(\boldsymbol{\theta}_k^*)$  and  $D_k^2 = nd_k^2$  in the i.i.d. case. Function  $\bar{\mathbf{g}}_{2,k}(\mathbf{r})$  denotes the marginal analog of  $\mathbf{g}_{2,k}(\mathbf{r})$ .

Let us show, that for the value  $\delta_k(\mathbf{r})$  from the condition  $(\mathcal{L}_0)$  it holds  $\delta_k(\mathbf{r}) = \mathbf{C}\mathbf{r}/\sqrt{n}$ . Suppose for all  $\boldsymbol{\theta} \in \Theta_{0,k}(\mathbf{r})$  and  $\boldsymbol{\gamma} \in \mathbb{R}^{p_k} : \|\boldsymbol{\gamma}\| = 1$   $\|D_k^{-1}\boldsymbol{\gamma}^\top \nabla_{\boldsymbol{\theta}}^3 \mathbb{E}L_k(\boldsymbol{\theta})D_k^{-1}\| \leq \mathbf{C}$ , then it holds for some  $\bar{\boldsymbol{\theta}} \in \Theta_{0,k}(\mathbf{r})$ :

$$\begin{aligned} \|D_k^{-1}D^2(\boldsymbol{\theta})D_k^{-1} - \mathbf{I}_{p_k}\| &= \|D_k^{-1}(\boldsymbol{\theta}_k^* - \boldsymbol{\theta})^\top \nabla_{\boldsymbol{\theta}}^3 \mathbb{E}L_k(\bar{\boldsymbol{\theta}})D_k^{-1}\| \\ &= \|D_k^{-1}(\boldsymbol{\theta}_k^* - \boldsymbol{\theta})^\top D_k D_k^{-1} \nabla_{\boldsymbol{\theta}}^3 \mathbb{E}L_k(\bar{\boldsymbol{\theta}})D_k^{-1}\| \\ &\leq \mathbf{r} \|D_k^{-1}\| \|D_k^{-1}\boldsymbol{\gamma}^\top \nabla_{\boldsymbol{\theta}}^3 \mathbb{E}L_k(\bar{\boldsymbol{\theta}})D_k^{-1}\| \leq \mathbf{C}\mathbf{r}/\sqrt{n}. \end{aligned}$$

Similarly  $\mathbf{C}_{m,k}(\mathbf{r}) \leq \mathbf{C}\mathbf{r}/\sqrt{n} + \mathbf{C}$  in condition  $(\mathcal{L}_{0m})$ .

The next remark helps to check the global identifiability condition  $(\mathcal{L}\mathbf{r})$  in many situations. Suppose that the parameter domain  $\Theta_k$  is compact and  $n$  is sufficiently large, then the value  $\mathbf{b}_k(\mathbf{r})$  from condition  $(\mathcal{L}\mathbf{r})$  can be taken as  $\mathbf{C}\{1 - \mathbf{r}/\sqrt{n}\} \approx \mathbf{C}$ . Indeed, for  $\boldsymbol{\theta} : \|D_k(\boldsymbol{\theta} - \boldsymbol{\theta}_k^*)\| = \mathbf{r}$

$$\begin{aligned} -2\{\mathbb{E}L_k(\boldsymbol{\theta}) - \mathbb{E}L_k(\boldsymbol{\theta}_k^*)\} &\geq \mathbf{r}^2 \left\{ 1 - \mathbf{r} \|D_k^{-1}\| \|D_k^{-1}\boldsymbol{\gamma}^\top \nabla_{\boldsymbol{\theta}}^3 \mathbb{E}L_k(\bar{\boldsymbol{\theta}})D_k^{-1}\| \right\} \\ &\geq \mathbf{r}^2(1 - \mathbf{C}\mathbf{r}/\sqrt{n}). \end{aligned}$$

Due to the obtained orders, the conditions (B.1) and (B.9) of Theorems B.1 and B.5 on concentration of the MLEs  $\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_k^\circ$  require  $\mathbf{r}_{0,k} \geq \mathbf{C}\sqrt{p_k + \mathbf{x}}$ .

## A Approximation of the joint distributions of $\ell_2$ -norms

Let us previously introduce some notations:

$$\mathbf{1}_K \stackrel{\text{def}}{=} (1, \dots, 1)^\top \in \mathbb{R}^K;$$

$\|\cdot\|$  is the Euclidean norm for a vector and spectral norm for a matrix;

$\|\cdot\|_{\max}$  is the maximum of absolute values of elements of a vector or of a matrix;

$\|\cdot\|_1$  is the sum of absolute values of elements of a vector or of a matrix.

Consider  $K$  random centered vectors  $\boldsymbol{\phi}_k \in \mathbb{R}^{p_k}$  for  $k = 1, \dots, K$ . Each vector equals to a sum of  $n$  centered independent vectors:

$$\begin{aligned} \boldsymbol{\phi}_k &= \boldsymbol{\phi}_{k,1} + \dots + \boldsymbol{\phi}_{k,n}, \\ \mathbb{E}\boldsymbol{\phi}_k &= \mathbb{E}\boldsymbol{\phi}_{k,i} = \mathbf{0} \quad \forall 1 \leq i \leq n. \end{aligned} \tag{A.1}$$

Introduce similarly the vectors  $\boldsymbol{\psi}_k \in \mathbb{R}^{p_k}$  for  $k = 1, \dots, K$ :

$$\begin{aligned} \boldsymbol{\psi}_k &= \boldsymbol{\psi}_{k,1} + \dots + \boldsymbol{\psi}_{k,n}, \\ \mathbb{E}\boldsymbol{\psi}_k &= \mathbb{E}\boldsymbol{\psi}_{k,i} = \mathbf{0} \quad \forall 1 \leq i \leq n, \end{aligned} \tag{A.2}$$

with the same independence properties as  $\phi_{k,i}$ , and also independent of all  $\phi_{k,i}$ .

The goal of this section is to compare the joint distributions of the  $\ell_2$ -norms of the sets of vectors  $\phi_k$  and  $\psi_k$ ,  $k = 1, \dots, K$  (i.e. the probability laws  $\mathcal{L}(\|\phi_1\|, \dots, \|\phi_K\|)$  and  $\mathcal{L}(\|\psi_1\|, \dots, \|\psi_K\|)$ ), assuming that their correlation structures are close to each other.

Denote

$$\begin{aligned} p_{\max} &\stackrel{\text{def}}{=} \max_{1 \leq k \leq K} p_k, & p_{\text{sum}} &\stackrel{\text{def}}{=} p_1 + \dots + p_K, \\ \lambda_{\phi, \max}^2 &\stackrel{\text{def}}{=} \max_{1 \leq k \leq K} \|\text{Var}(\phi_j)\|, & \lambda_{\psi, \max}^2 &\stackrel{\text{def}}{=} \max_{1 \leq k \leq K} \|\text{Var}(\psi_j)\|, \\ z_{\max} &\stackrel{\text{def}}{=} \max_{1 \leq k \leq K} z_k, & z_{\min} &\stackrel{\text{def}}{=} \min_{1 \leq k \leq K} z_k, \\ \delta_{z, \max} &\stackrel{\text{def}}{=} \max_{1 \leq k \leq K} \delta_{z_k}, & \delta_{z, \min} &\stackrel{\text{def}}{=} \min_{1 \leq k \leq K} \delta_{z_k}, \end{aligned}$$

let also

$$\begin{aligned} \Delta_\varepsilon &\stackrel{\text{def}}{=} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/16}(K) \log^{3/8}(np_{\text{sum}}) z_{\min}^{1/8} \\ &\quad \times \max \{ \lambda_{\phi, \max}, \lambda_{\psi, \max} \}^{3/4} \log^{-1/8}(5n^{1/2}). \end{aligned} \tag{A.3}$$

The following conditions are necessary for the Proposition A.1

(C1) For some  $\mathbf{g}_k, \nu_k, \mathbf{c}_\phi, \mathbf{c}_\psi > 0$  and for all  $i = 1, \dots, n$ ,  $k = 1, \dots, K$

$$\begin{aligned} \sup_{\substack{\gamma_k \in \mathbb{R}^{p_k}, \\ \|\gamma_k\|=1}} \log \mathbb{E} \exp \left\{ \lambda \sqrt{n} \gamma_k^\top \phi_{k,i} / \mathbf{c}_\phi \right\} &\leq \lambda^2 \nu_k^2 / 2, \quad |\lambda| < \mathbf{g}_k, \\ \sup_{\substack{\gamma_k \in \mathbb{R}^{p_k}, \\ \|\gamma_k\|=1}} \log \mathbb{E} \exp \left\{ \lambda \sqrt{n} \gamma_k^\top \psi_{k,i} / \mathbf{c}_\psi \right\} &\leq \lambda^2 \nu_k^2 / 2, \quad |\lambda| < \mathbf{g}_k, \end{aligned}$$

where  $\mathbf{c}_\phi \geq \mathbf{C} \lambda_{\phi, \max}$  and  $\mathbf{c}_\psi \geq \mathbf{C} \lambda_{\psi, \max}$ .

(C2) For some  $\delta_\Sigma^2 \geq 0$

$$\max_{1 \leq k_1, k_2 \leq K} \|\text{Cov}(\phi_{k_1}, \phi_{k_2}) - \text{Cov}(\psi_{k_1}, \psi_{k_2})\|_{\max} \leq \delta_\Sigma^2. \tag{A.4}$$

**Proposition A.1** (Approximation of the joint distributions of  $\ell_2$ -norms). *Consider the centered random vectors  $\phi_1, \dots, \phi_K$  and  $\psi_1, \dots, \psi_K$  given in (A.1), (A.2). Let the conditions (C1) and (C2) be fulfilled, and the values  $z_k \geq \sqrt{p_k + \Delta_\varepsilon}$  and  $\delta_{z_k} \geq 0$  be s.t.*

$\mathfrak{C} \max\{n^{-1/2}, \delta_{z, \max}\} \leq \Delta_\varepsilon \leq \mathfrak{C} z_{\max}^{-1}$ , then it holds with dominating probability

$$\begin{aligned} \mathbb{P} \left( \bigcup_{k=1}^K \{\|\phi_k\| > z_k\} \right) - \mathbb{P} \left( \bigcup_{k=1}^K \{\|\psi_k\| > z_k - \delta_{z_k}\} \right) &\geq -\Delta_{\ell_2}, \\ \mathbb{P} \left( \bigcup_{k=1}^K \{\|\phi_k\| > z_k\} \right) - \mathbb{P} \left( \bigcup_{k=1}^K \{\|\psi_k\| > z_k + \delta_{z_k}\} \right) &\leq \Delta_{\ell_2} \end{aligned}$$

for the deterministic non-negative value

$$\begin{aligned} \Delta_{\ell_2} &\leq 12.5\mathfrak{C} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \max\{\lambda_{\phi, \max}, \lambda_{\psi, \max}\}^{3/4} \\ &\quad + 3.2\mathfrak{C} \delta_\Sigma^2 \left( \frac{p_{\max}^3}{n} \right)^{1/4} p_{\max} z_{\min}^{1/2} \log^2(K) \log^{3/4}(np_{\text{sum}}) \max\{\lambda_{\phi, \max}, \lambda_{\psi, \max}\}^{7/2} \\ &\leq 25\mathfrak{C} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \max\{\lambda_{\phi, \max}, \lambda_{\psi, \max}\}^{3/4}, \end{aligned}$$

where the last inequality holds for

$$\delta_\Sigma^2 \leq 4\mathfrak{C} \left( \frac{n}{p_{\max}^{13}} \right)^{1/8} \log^{-7/8}(K) \log^{-3/8}(np_{\text{sum}}) (\max\{\lambda_{\phi, \max}, \lambda_{\psi, \max}\})^{-11/4}.$$

**Remark A.1.** The approximating error term  $\Delta_{\ell_2}$  consists of three errors, which correspond to: the Gaussian approximation result (Lemma A.2), Gaussian comparison (Lemma A.7), and anti-concentration inequality (Lemma A.8). The bound on  $\Delta_{\ell_2}$  above implies that the number  $K$  of the random vectors  $\phi_1, \dots, \phi_K$  should satisfy  $\log K \ll (n/p_{\max}^3)^{1/12}$  in order to keep the approximating error term  $\Delta_{\ell_2}$  small. This condition can be relaxed by using a sharper Gaussian approximation result. For instance, using in Lemma A.2 the Slepian-Stein technique plus induction argument from the recent paper by Chernozhukov et al. (2014b) instead of the Lindeberg's approach, would lead to the improved bound:  $\mathfrak{C} \left( \frac{p_{\max}^3}{n} \right)^{1/6}$  multiplied by a logarithmic term.

## A.1 Joint Gaussian approximation of $\ell_2$ -norm of sums of independent vectors by Lindeberg's method

Introduce the following random vectors from  $\mathbb{R}^{p_{\text{sum}}}$ :

$$\begin{aligned} \Phi &\stackrel{\text{def}}{=} \left( \phi_1^\top, \dots, \phi_K^\top \right)^\top, & \Phi_i &\stackrel{\text{def}}{=} \left( \phi_{1,i}^\top, \dots, \phi_{K,i}^\top \right)^\top, \quad i = 1, \dots, n, \\ \Phi &= \sum_{i=1}^n \Phi_i, & \mathbb{E}\Phi &= \mathbb{E}\Phi_i = 0. \end{aligned} \tag{A.5}$$

Define their Gaussian analogs as follows:

$$\bar{\Phi}_i \stackrel{\text{def}}{=} \left( \bar{\phi}_{1,i}^\top, \dots, \bar{\phi}_{K,i}^\top \right)^\top, \quad \bar{\Phi} \stackrel{\text{def}}{=} \left( \bar{\phi}_1^\top, \dots, \bar{\phi}_K^\top \right)^\top = \sum_{i=1}^n \bar{\Phi}_i, \quad (\text{A.6})$$

$$\bar{\Phi}_i \sim \mathcal{N}(0, \text{Var } \Phi_i), \quad \bar{\Phi} \sim \mathcal{N}(0, \text{Var } \Phi), \quad (\text{A.7})$$

$$\bar{\phi}_{k,i} \sim \mathcal{N}(0, \text{Var } \phi_{k,i}), \quad \bar{\phi}_k \stackrel{\text{def}}{=} \sum_{i=1}^n \bar{\phi}_{k,i} \sim \mathcal{N}(0, \text{Var } \phi_k). \quad (\text{A.8})$$

**Lemma A.2** (Joint GAR with equal covariance matrices). *Consider the sets of random vectors  $\phi_j$  and  $\bar{\phi}_j$ ,  $j = 1, \dots, K$  defined in (A.1), and (A.5)–(A.8). If the conditions of Lemmas A.4 are A.5 are fulfilled, then it holds for all  $\Delta, \beta > 0$ ,  $z_j \geq \max \{ \Delta + \sqrt{p_j}, 2.25 \log(K)/\beta \}$  with dominating probability*

$$\begin{aligned} \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\phi_j\| > z_j \} \right) &\leq \mathbb{P} \left( \bigcup_{j=1}^K \left\{ \|\bar{\phi}_j\| > z_j - \Delta - \frac{3 \log(K)}{2\beta} \right\} \right) + \delta_{3,\phi}(\Delta, \beta), \\ \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\phi_j\| > z_j \} \right) &\geq \mathbb{P} \left( \bigcup_{j=1}^K \left\{ \|\bar{\phi}_j\| > z_j + \Delta + \frac{3 \log(K)}{2\beta} \right\} \right) - \delta_{3,\phi}(\Delta, \beta) \end{aligned}$$

for  $\delta_{3,\phi}(\Delta, \beta) \leq \mathbf{c} \left( \frac{1}{\Delta^3} + \frac{\beta}{\Delta^2} + \frac{\beta^2}{\Delta} \right) \left\{ \frac{p_{\max}^3}{n} \log(K) \log^3(np_{\text{sum}}) \right\}^{1/2}$  given in (A.15).

*Proof of Lemma A.2.*

$$\mathbb{P} \left( \bigcup_{j=1}^K \{ \|\phi_j\| > z_j \} \right) = \mathbb{E} \mathbb{I}(\max_{1 \leq j \leq K} \{ \|\phi_j\|^2 - z_j^2 \} > 0).$$

Let us approximate the  $\max_{1 \leq j \leq K}$  function using the smooth maximum:

$$\begin{aligned} h_\beta(\{x_j\}) &\stackrel{\text{def}}{=} \beta^{-1} \log \left( \sum_{j=1}^K e^{\beta x_j} \right) \text{ for } \beta > 0, x_j \in \mathbb{R}, \\ h_\beta(\{x_j\}) - \beta^{-1} \log(K) &\leq \max_{1 \leq j \leq K} \{x_j\} \leq h_\beta(\{x_j\}). \end{aligned} \quad (\text{A.9})$$

The indicator function  $\mathbb{I}\{x > 0\}$  is approximated with the three times differentiable function  $g(x)$  growing monotonously from 0 to 1:

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

It holds for all  $x \in \mathbb{R}$  and  $\Delta > 0$

$$\mathbb{I}\{x > \Delta\} \leq g(x/\Delta) \leq \mathbb{I}\{x/\Delta > 0\}.$$

Therefore

$$\begin{aligned}
& \mathbb{P} \left( \max_{1 \leq j \leq K} \{ \|\phi_j\| - z_j \} > \Delta \right) \\
& \leq \mathbb{E} \mathbb{I} \left( \max_{1 \leq j \leq K} \left\{ \frac{\|\phi_j\|^2 - z_j^2}{2z_j} \right\} > \Delta \right) \\
& \leq \mathbb{E} g \left( \max_{1 \leq j \leq K} \left\{ \frac{\|\phi_j\|^2 - z_j^2}{2z_j \Delta} \right\} \right) \\
& \leq \mathbb{E} g \left( \frac{1}{\Delta \beta} \log \left\{ \sum_{j=1}^K \exp \left[ \beta \frac{\|\phi_j\|^2 - z_j^2}{2z_j} \right] \right\} \right) \tag{A.10}
\end{aligned}$$

$$\begin{aligned}
& \leq \mathbb{E} g \left( \max_{1 \leq j \leq K} \left\{ \frac{\|\phi_j\|^2 - z_j^2}{2z_j \Delta} \right\} + \frac{\log(K)}{\beta \Delta} \right) \\
& \leq \mathbb{E} \mathbb{I} \left( \max_{1 \leq j \leq K} \left\{ \frac{\|\phi_j\|^2 - z_j^2}{2z_j} \right\} > -\frac{\log(K)}{\beta} \right) \\
& \leq \mathbb{P} \left( \max_{1 \leq j \leq K} \{ \|\phi_j\| - z_j \} > -1.5 \frac{\log(K)}{\beta} \right), \tag{A.11}
\end{aligned}$$

where the last inequality holds for  $z_j \geq 2.25 \log(K)/\beta$ . Denote

$$\mathbf{z} \stackrel{\text{def}}{=} (z_1, \dots, z_K)^\top \in \mathbb{R}^K, \quad z_j > 0.$$

Introduce the function  $F_{\Delta, \beta}(\Phi, \mathbf{z}) : \mathbb{R}^{p_{\text{sum}}} \times \mathbb{R}^K \mapsto \mathbb{R}$ :

$$F_{\Delta, \beta}(\Phi, \mathbf{z}) \stackrel{\text{def}}{=} g \left( \frac{1}{\Delta \beta} \log \left\{ \sum_{j=1}^K \exp \left[ \beta \frac{\|\phi_j\|^2 - z_j^2}{2z_j} \right] \right\} \right) \tag{A.12}$$

Then by (A.10) and (A.11)

$$\begin{aligned}
& \mathbb{P} \left( \max_{1 \leq j \leq K} \{ \|\phi_j\| - z_j \} > \Delta \right) \\
& \leq \mathbb{E} F_{\Delta, \beta}(\Phi, \mathbf{z}) \tag{A.13}
\end{aligned}$$

$$\leq \mathbb{P} \left( \max_{1 \leq j \leq K} \{ \|\phi_j\| - z_j \} > -\frac{3 \log(K)}{2\beta} \right). \tag{A.14}$$

Lemma A.6 checks that  $F_{\Delta, \beta}(\cdot, \mathbf{z})$  admits applying the Lindeberg's telescopic sum device (see Lindeberg (1922)) in order to approximate  $\mathbb{E} F_{\Delta, \beta}(\Phi, \mathbf{z})$  with  $\mathbb{E} F_{\Delta, \beta}(\bar{\Phi}, \mathbf{z})$ . Define for  $q = 2, \dots, n-1$  the following  $\mathbb{R}^{p_{\text{sum}}}$ -valued random sums:

$$S_q \stackrel{\text{def}}{=} \sum_{i=1}^{q-1} \bar{\Phi}_i + \sum_{i=q+1}^n \Phi_i, \quad S_1 \stackrel{\text{def}}{=} \sum_{i=2}^n \Phi_i, \quad S_n \stackrel{\text{def}}{=} \sum_{i=1}^{n-1} \bar{\Phi}_i.$$

The difference  $F_{\Delta, \beta}(\Phi, \mathbf{z}) - F_{\Delta, \beta}(\bar{\Phi}, \mathbf{z})$  can be represented as the telescopic sum:

$$F_{\Delta, \beta}(\Phi, \mathbf{z}) - F_{\Delta, \beta}(\bar{\Phi}, \mathbf{z}) = \sum_{i=1}^n \{F_{\Delta, \beta}(S_i + \Phi_i, \mathbf{z}) - F_{\Delta, \beta}(S_i + \bar{\Phi}_i, \mathbf{z})\}.$$

The third order Taylor expansions of  $F_{\Delta, \beta}(S_i + \Phi_i, \mathbf{z})$  and  $F_{\Delta, \beta}(S_i + \bar{\Phi}_i, \mathbf{z})$  w.r.t. the first argument at  $S_i$ , and Lemma A.6 imply for each  $i = 1, \dots, n$ :

$$\begin{aligned} & \left| F_{\Delta, \beta}(S_i + \Phi_i, \mathbf{z}) - F_{\Delta, \beta}(S_i + \bar{\Phi}_i, \mathbf{z}) - \nabla_{\Phi} F_{\Delta, \beta}(S_i, \mathbf{z})^{\top} (\Phi_i - \bar{\Phi}_i) \right. \\ & \quad \left. - \frac{1}{2} (\Phi_i - \bar{\Phi}_i)^{\top} \nabla_{\Phi}^2 F_{\Delta, \beta}(S_i, \mathbf{z}) (\Phi_i + \bar{\Phi}_i) \right| \\ & \leq \frac{\mathbf{C}_3(\Delta, \beta)}{6} \left( \max_{1 \leq j \leq K} \{\|S_{j,i} + \phi_{j,i}\|^3\} \|\Phi_i\|_{\max}^3 + \max_{1 \leq j \leq K} \{\|S_{j,i} + \bar{\phi}_{j,i}\|^3\} \|\bar{\Phi}_i\|_{\max}^3 \right), \end{aligned}$$

where the value  $\mathbf{C}_3(\Delta, \beta)$  is defined in Lemma A.6, and the random vectors  $S_{j,i} \in \mathbb{R}^{p_j}$  for  $j = 1, \dots, K$  are s.t. for all  $i = 1, \dots, n$

$$S_i = \left( S_{1,i}^{\top}, S_{2,i}^{\top}, \dots, S_{K,i}^{\top} \right)^{\top}.$$

By their construction  $S_i$  and  $\Phi_i - \bar{\Phi}_i$  are independent,  $\mathbb{E}\Phi_i = \mathbb{E}\bar{\Phi}_i = 0$  and  $\text{Var}\Phi_i = \text{Var}\bar{\Phi}_i$ , therefore

$$\begin{aligned} & \left| \mathbb{E}F_{\Delta, \beta}(\Phi, \mathbf{z}) - \mathbb{E}F_{\Delta, \beta}(\bar{\Phi}, \mathbf{z}) \right| \\ & = \left| \sum_{i=1}^n \{ \mathbb{E}H_{\Delta}(S_i + \Phi_i, \mathbf{z}) - \mathbb{E}H_{\Delta}(S_i + \bar{\Phi}_i, \mathbf{z}) \} \right| \\ & \leq \frac{\mathbf{C}_3(\Delta, \beta)}{6} \sum_{i=1}^n \mathbb{E} \left( \max_{1 \leq j \leq K} \{\|S_{j,i} + \phi_{j,i}\|^3\} \|\Phi_i\|_{\max}^3 + \max_{1 \leq j \leq K} \{\|S_{j,i} + \bar{\phi}_{j,i}\|^3\} \|\bar{\Phi}_i\|_{\max}^3 \right). \end{aligned}$$

Lemma A.5 implies for all  $i = 1, \dots, n$  with probability  $\geq 1 - 2e^{-\mathbf{x}}$

$$\left( \mathbb{E} \max_{1 \leq j \leq K} \{\|S_{j,i} + \phi_{j,i}\|^6\} \right)^{1/2} \leq \mathbf{C}\nu_0 \max_{1 \leq j \leq K} \|\text{Var}^{1/2}(\phi_j)\|^3 \sqrt{p_{\max} \log(K)} (p_{\max} + 6\mathbf{x}),$$

and the same bound holds for  $(\mathbb{E} \max_{1 \leq j \leq K} \{\|S_{j,i} + \bar{\phi}_{j,i}\|^6\})^{1/2}$ . Denote

$$\delta_{\max, \phi} \stackrel{\text{def}}{=} \frac{1}{2} \sum_{i=1}^n \{ \mathbb{E} (\|\Phi_i\|_{\max}^6) \}^{1/2} + \{ \mathbb{E} (\|\bar{\Phi}_i\|_{\max}^6) \}^{1/2}.$$

By Lemma A.4 it holds for  $t = (\mathbf{x} + \log(p_{\text{sum}}))^3 (\sqrt{2}\mathbf{c}_{\phi}\nu_0)^6 n^{-3}$  with probability  $\geq 1 - e^{-\mathbf{x}}$

$$\|\Phi_i\|_{\max}^6 \leq t, \quad \|\bar{\Phi}_i\|_{\max}^6 \leq t.$$

If  $\mathbf{x} = \mathbf{C} \log n$ , then the last bound on  $|\mathbb{E}F_{\Delta, \beta}(\Phi, \mathbf{z}) - \mathbb{E}F_{\Delta, \beta}(\bar{\Phi}, \mathbf{z})|$  continues with

probability  $\geq 1 - 6 \exp(-x)$  as follows

$$\begin{aligned}
& |\mathbb{E}F_{\Delta, \beta}(\Phi, \mathbf{z}) - \mathbb{E}F_{\Delta, \beta}(\bar{\Phi}, \mathbf{z})| \\
& \leq \mathbf{C} \frac{\mathbf{C}_3(\Delta, \beta)}{3} \sqrt{p_{\max}^3 \log(K)} \delta_{\max, \phi} \max_{1 \leq j \leq K} \|\text{Var}^{1/2}(\phi_j)\|^3 \\
& \leq \frac{\mathbf{C}}{3} \left( \frac{1}{\Delta^3} + \frac{\beta}{\Delta^2} + \frac{\beta^2}{\Delta} \right) \frac{p_{\max}^{3/2}}{n^{1/2}} \log^{1/2}(K) \log^{3/2}(np_{\text{sum}}) \max_{1 \leq j \leq K} \|\text{Var}^{1/2}(\phi_j)\|^3 (2\nu_0^2 \mathbf{c}_\phi^2)^{3/2} \\
& \stackrel{\text{def}}{=} \delta_{3, \phi}(\Delta, \beta). \tag{A.15}
\end{aligned}$$

The derived bounds imply:

$$\begin{aligned}
& \mathbb{P} \left( \bigcup_{j=1}^K \{\|\phi_j\| > z_j\} \right) \\
& \stackrel{\text{by (A.13)}}{\leq} \mathbb{E}F_{\Delta, \beta}(\Phi, \mathbf{z} - \Delta \mathbf{1}_K) \\
& \stackrel{\text{by (A.15)}}{\leq} \mathbb{E}F_{\Delta, \beta}(\bar{\Phi}, \mathbf{z} - \Delta \mathbf{1}_K) + \delta_{3, \phi}(\Delta, \beta) \tag{A.16} \\
& \stackrel{\text{by (A.14)}}{\leq} \mathbb{P} \left( \bigcup_{j=1}^K \left\{ \|\bar{\phi}_j\| > z_j - \Delta - \frac{3 \log(K)}{2\beta} \right\} \right) + \delta_{3, \phi}(\Delta, \beta),
\end{aligned}$$

and similarly

$$\begin{aligned}
& \mathbb{P} \left( \bigcup_{j=1}^K \{\|\phi_j\| > z_j\} \right) \\
& \geq \mathbb{P} \left( \bigcup_{j=1}^K \left\{ \|\bar{\phi}_j\| > z_j + \frac{3 \log(K)}{2\beta} + \Delta \right\} \right) - \delta_{3, \phi}(\Delta, \beta).
\end{aligned}$$

□

The next lemma is formulated separately, since it is used for a proof of another result.

**Lemma A.3** (Smooth uniform GAR). *Under the conditions of Lemma A.2 it holds with dominating probability for the function  $F_{\Delta, \beta}(\cdot, \mathbf{z})$  given in (A.12):*

- 1.1.  $\mathbb{P} \left( \bigcup_{j=1}^K \{\|\phi_j\| > z_j\} \right) \leq \mathbb{E}F_{\Delta, \beta}(\bar{\Phi}, \mathbf{z} - \Delta \mathbf{1}_K) + \delta_{3, \phi}(\Delta, \beta),$
- 1.2.  $\mathbb{P} \left( \bigcup_{j=1}^K \{\|\phi_j\| > z_j\} \right) \geq \mathbb{E}H_{\Delta, \beta} \left( \bar{\Phi}, \mathbf{z} + \frac{3 \log(K)}{2\beta} \mathbf{1}_K \right) - \delta_{3, \phi}(\Delta, \beta);$
- 2.1.  $\mathbb{E}F_{\Delta, \beta}(\Phi, \mathbf{z}) \leq \mathbb{P} \left( \bigcup_{j=1}^K \left\{ \|\phi_j\| > z_j - \frac{3 \log(K)}{2\beta} \right\} \right),$
- 2.2.  $\mathbb{E}F_{\Delta, \beta}(\Phi, \mathbf{z}) \geq \mathbb{P} \left( \bigcup_{j=1}^K \{\|\phi_j\| > z_j + \Delta\} \right).$



*Proof of Lemma A.3.* The first inequality 1.1 is obtained in (A.16), the second inequality 1.2 follows similarly from (A.14) and (A.15). The inequalities 2.1 and 2.2 are given in (A.13) and (A.14).  $\square$

**Lemma A.4.** *Let for some  $\mathbf{c}_\phi, \mathbf{g}_1, \nu_0 > 0$  and for all  $i = 1, \dots, n$ ,  $j = 1, \dots, p_{\text{sum}}$*

$$\log \mathbb{E} \exp \left\{ \lambda \sqrt{n} |\phi_i^j| / \mathbf{c}_\phi \right\} \leq \lambda^2 \nu_0^2 / 2, \quad |\lambda| < \mathbf{g}_1,$$

here  $\phi_i^j$  denotes the  $j$ -th coordinate of vector  $\phi_i$ . Then it holds for all  $i = 1, \dots, n$  and  $m, t > 0$

$$\mathbb{P} \left( \max_{1 \leq j \leq p_{\text{sum}}} |\phi_i^j|^m > t \right) \leq \exp \left\{ -\frac{nt^{2/m}}{2\mathbf{c}_\phi^2 \nu_0^2} + \log(p_{\text{sum}}) \right\}.$$

*Proof of Lemma A.4.* Let us bound the  $\max_j |\phi_i^j|$  using the following bound for the maximum:

$$\max_{1 \leq j \leq p_{\text{sum}}} |\phi_i^j| \leq \log \left\{ \sum_{j=1}^{p_{\text{sum}}} \exp(|\phi_i^j|) \right\}.$$

By the Lemma's condition

$$\mathbb{E} \exp \left\{ \max_{1 \leq j \leq p} \frac{\lambda \sqrt{n}}{\mathbf{c}_\phi} |\phi_i^j| \right\} \leq \exp(\lambda^2 \nu_0^2 / 2 + \log p_{\text{sum}}).$$

Thus, the statement follows from the exponential Chebyshev's inequality.  $\square$

**Lemma A.5.** *If for the centered random vectors  $\phi_j \in \mathbb{R}^{p_j}$   $j = 1, \dots, K$*

$$\sup_{\substack{\gamma \in \mathbb{R}^{p_j}, \\ \|\gamma\| \neq 0}} \log \mathbb{E} \exp \left\{ \lambda \frac{\gamma^\top \phi_j}{\|\text{Var}^{1/2}(\phi_j) \gamma\|} \right\} \leq \nu_0^2 \lambda^2 / 2, \quad |\lambda| \leq \mathbf{g}$$

for some constants  $\nu_0 > 0$  and  $\mathbf{g} \geq \nu_0^{-1} \max_{1 \leq j \leq K} \sqrt{2p_j \log(K)}$ , then

$$\begin{aligned} \mathbb{E} \max_{1 \leq j \leq K} \{\|\phi_j\|\} &\leq \mathbf{C} \nu_0 \max_{1 \leq j \leq K} \|\text{Var}^{1/2}(\phi_j)\| \sqrt{2p_{\max} \log(K)}, \\ \left( \mathbb{E} \max_{1 \leq j \leq K} \{\|\phi_j\|^6\} \right)^{1/2} &\leq \mathbf{C} \nu_0 \max_{1 \leq j \leq K} \|\text{Var}^{1/2}(\phi_j)\|^3 \sqrt{2p_{\max} \log(K)} (p_{\max} + 6\mathbf{x}), \end{aligned}$$

The second bound holds with probability  $\geq 1 - 2e^{-\mathbf{x}}$ .

*Proof of Lemma A.5.* Let us take for each  $j = 1, \dots, K$  finite  $\varepsilon_j$ -grids  $\mathbf{G}_j(\varepsilon) \subset \mathbb{R}^{p_j}$  on the  $(p_j - 1)$ -spheres of radius 1 s.t

$$\forall \gamma \in \mathbb{R}^{p_j} \text{ s.t. } \|\gamma\| = 1 \quad \exists \gamma_0 \in \mathbf{G}_j(\varepsilon) : \|\gamma - \gamma_0\| \leq \varepsilon, \quad \|\gamma_0\| = 1.$$

Then

$$\|\phi_j\| \leq (1 - \varepsilon_j)^{-1} \max_{\gamma \in G_j(\varepsilon_j)} \{\gamma^\top \phi_j\}.$$

Hence, by inequality (A.9) and the imposed condition it holds for all  $0 < \mu < \mathbf{g} / \max_{1 \leq j \leq K} \|\text{Var}^{1/2}(\phi_j)\|$ :

$$\begin{aligned} \mathbb{E} \max_{1 \leq j \leq K} \{\|\phi_j\|\} &\leq \max_{1 \leq j \leq K} \frac{1}{1 - \varepsilon_j} \mathbb{E} \max_{1 \leq j \leq K} \max_{\gamma \in G_j(\varepsilon_j)} \{\gamma^\top \phi_j\} \\ &\leq \mathbf{c} \frac{1}{\mu} \mathbb{E} \log \left\{ \sum_{1 \leq j \leq K} \sum_{\gamma \in G_j(\varepsilon_j)} \exp(\mu \gamma^\top \phi_j) \right\} \\ &\leq \mathbf{c} \frac{1}{\mu} \log \left\{ \sum_{1 \leq j \leq K} \sum_{\gamma \in G_j(\varepsilon_j)} \mathbb{E} \exp(\mu \gamma^\top \phi_j) \right\} \\ &\leq \mathbf{c} \max_{1 \leq j \leq K} \frac{\log(K \text{card}\{G_j(\varepsilon_j)\})}{\mu} + \mathbf{c} \frac{\mu \nu_0^2}{2} \max_{1 \leq j \leq K} \|\text{Var}(\phi_j)\| \\ &\leq \mathbf{c} \max_{1 \leq j \leq K} \{p_j\} \frac{\log(K)}{\mu} + \mathbf{c} \frac{\mu \nu_0^2}{2} \max_{1 \leq j \leq K} \|\text{Var}(\phi_j)\| \\ &= \mathbf{c} \nu_0 \max_{1 \leq j \leq K} \{\sqrt{p_j}\} \max_{1 \leq j \leq K} \|\text{Var}^{1/2}(\phi_j)\| \sqrt{2 \log(K)} \\ &\quad \text{for } \mu = \mathbf{c} \nu_0^{-1} \max_{1 \leq j \leq K} \{\sqrt{p_j}\} \sqrt{2 \log(K)} / \max_{1 \leq j \leq K} \|\text{Var}^{1/2}(\phi_j)\|. \end{aligned}$$

For the second part of the statement we combine the first part with the result of Theorem B.3 on deviation of a random quadratic form: it holds with dominating probability for  $V_{\phi_j}^2 \stackrel{\text{def}}{=} \text{Var} \phi_j$

$$\begin{aligned} \|\phi_j\|^2 &\leq 3_{\text{qf}}^2(\mathbf{x}, V_{\phi_j}) \\ &\leq \text{tr}(V_{\phi_j}^2) + 6\mathbf{x} \|V_{\phi_j}^2\| \leq \|V_{\phi_j}^2\| (p_j + 6\mathbf{x}). \end{aligned}$$

□

**Lemma A.6.** *Let  $\Gamma \in \mathbb{R}^{p_{\text{sum}}}$ ,  $\gamma_j \in \mathbb{R}^{p_j}$  for  $j = 1, \dots, K$  are s.t.  $\Gamma = (\gamma_1^\top, \dots, \gamma_K^\top)^\top$ , and  $\mathbf{z} \stackrel{\text{def}}{=} (z_1, \dots, z_K)^\top$  s.t.  $z_j \geq \sqrt{p_j}$ , then it holds for the function  $F_{\Delta, \beta}(\cdot, \mathbf{z})$  defined in (A.12):*

$$\begin{aligned} \|\nabla_{\Gamma}^2 F_{\Delta, \beta}(\Gamma, \mathbf{z})\|_1 &\leq \mathbf{c}_2(\Delta, \beta) \max_{1 \leq j \leq K} \{\|\gamma_j\|^2\}, \quad \mathbf{c}_2(\Delta, \beta) \stackrel{\text{def}}{=} \mathbf{c} \left( \frac{1}{\Delta^2} + \frac{\beta}{\Delta} \right), \\ \|\nabla_{\Gamma}^3 F_{\Delta, \beta}(\Gamma, \mathbf{z})\|_1 &\leq \mathbf{c}_3(\Delta, \beta) \max_{1 \leq j \leq K} \{\|\gamma_j\|^3\}, \quad \mathbf{c}_3(\Delta, \beta) \stackrel{\text{def}}{=} \mathbf{c} \left( \frac{1}{\Delta^3} + \frac{\beta}{\Delta^2} + \frac{\beta^2}{\Delta} \right). \end{aligned}$$

*Proof of Lemma A.6.* Denote

$$s(\Gamma) \stackrel{\text{def}}{=} \sum_{j=1}^K \exp\left(\beta \frac{\|\gamma_j\|^2 - z_j^2}{2z_j}\right), \quad h_\beta(s(\Gamma)) \stackrel{\text{def}}{=} \beta^{-1} \log\{s(\Gamma)\}, \quad (\text{A.17})$$

then  $F_{\beta,\Delta}(\Gamma, \mathbf{z}) = g(\Delta^{-1}h_\beta(s(\Gamma)))$ . Let  $\gamma^q$  denote the  $q$ -th coordinate of the vector  $\Gamma \in \mathbb{R}^{p_{\text{sum}}}$ . It holds for  $q, l, b, r = 1, \dots, p_{\text{sum}}$ :

$$\begin{aligned} \frac{d}{d\gamma^q} F_{\beta,\Delta}(\Gamma, \mathbf{z}) &= \frac{1}{\Delta} g' \{ \Delta^{-1} h_\beta(s(\Gamma)) \} \frac{d}{d\gamma^q} h_\beta(s(\Gamma)), \\ \frac{d^2}{d\gamma^q d\gamma^l} F_{\beta,\Delta}(\Gamma, \mathbf{z}) &= \frac{1}{\Delta^2} g'' \{ \Delta^{-1} h_\beta(s(\Gamma)) \} \frac{d}{d\gamma^q} h_\beta(s(\Gamma)) \frac{d}{d\gamma^l} h_\beta(s(\Gamma)) \\ &\quad + \frac{1}{\Delta} g' \{ \Delta^{-1} h_\beta(s(\Gamma)) \} \frac{d^2}{d\gamma^q d\gamma^l} h_\beta(s(\Gamma)), \\ \frac{d^3}{d\gamma^q d\gamma^l d\gamma^b} F_{\beta,\Delta}(\Gamma, \mathbf{z}) &= \frac{1}{\Delta^3} g''' \{ \Delta^{-1} h_\beta(s(\Gamma)) \} \frac{d}{d\gamma^q} h_\beta(s(\Gamma)) \frac{d}{d\gamma^l} h_\beta(s(\Gamma)) \frac{d}{d\gamma^b} h_\beta(s(\Gamma)) \\ &\quad + \frac{1}{\Delta^2} g'' \{ \Delta^{-1} h_\beta(s(\Gamma)) \} \left\{ \frac{d^2}{d\gamma^q d\gamma^b} h_\beta(s(\Gamma)) \frac{d}{d\gamma^l} h_\beta(s(\Gamma)) \right. \\ &\quad \left. + \frac{d}{d\gamma^q} h_\beta(s(\Gamma)) \frac{d^2}{d\gamma^l d\gamma^b} h_\beta(s(\Gamma)) + \frac{d}{d\gamma^b} h_\beta(s(\Gamma)) \frac{d^2}{d\gamma^q d\gamma^l} h_\beta(s(\Gamma)) \right\} \\ &\quad + \frac{1}{\Delta} g' \{ \Delta^{-1} h_\beta(s(\Gamma)) \} \frac{d^3}{d\gamma^q d\gamma^l d\gamma^b} h_\beta(s(\Gamma)). \end{aligned}$$

Let for  $1 \leq q \leq p_{\text{sum}}$   $j(q)$  denote an index from 1 to  $K$  s.t. the coordinate  $\gamma^q$  of the vector  $\Gamma = (\gamma_1^\top, \dots, \gamma_K^\top)^\top$  belongs to its sub-vector  $\gamma_{j(q)}$ .

$$\frac{d}{d\gamma^q} h_\beta(s(\Gamma)) = \frac{1}{\beta} \frac{1}{s(\Gamma)} \frac{d}{d\gamma^q} s(\Gamma) = \frac{1}{s(\Gamma)} \frac{\gamma^q}{z_{j(q)}} \exp\left(\beta \frac{\|\gamma_{j(q)}\|^2 - z_{j(q)}^2}{2z_{j(q)}}\right),$$

$$\begin{aligned}
\frac{d^2}{d\gamma^q d\gamma^l} h_\beta(s(\Gamma)) &= \frac{1}{\beta} \frac{1}{s(\Gamma)} \frac{d^2}{d\gamma^q d\gamma^l} s(\Gamma) - \frac{1}{\beta} \frac{1}{s^2(\Gamma)} \frac{d}{d\gamma^q} s(\Gamma) \frac{d}{d\gamma^l} s(\Gamma) \\
&= \begin{cases} \left\{ \frac{1}{z_{j(q)}} + \beta \left( \frac{\gamma^q}{z_{j(q)}} \right)^2 \right\} \frac{1}{s(\Gamma)} \exp \left( \beta \frac{\|\gamma_{j(q)}\|^2 - z_{j(q)}^2}{2z_{j(q)}} \right) \\ \quad - \frac{\beta}{s^2(\Gamma)} \left\{ \frac{\gamma^q}{z_{j(q)}} \right\}^2 \exp \left( 2\beta \frac{\|\gamma_{j(q)}\|^2 - z_{j(q)}^2}{2z_{j(q)}} \right), & q = l; \\ \frac{\beta}{s(\Gamma)} \frac{\gamma^q \gamma^l}{z_{j(q)}^2} \exp \left( \beta \frac{\|\gamma_{j(q)}\|^2 - z_{j(q)}^2}{2z_{j(q)}} \right) \\ \quad - \frac{\beta}{s^2(\Gamma)} \frac{\gamma^q \gamma^l}{z_{j(q)}^2} \exp \left( 2\beta \frac{\|\gamma_{j(q)}\|^2 - z_{j(q)}^2}{2z_{j(q)}} \right), & j(q) = j(l), q \neq l; \\ -\frac{\beta}{s^2(\Gamma)} \frac{\gamma^q \gamma^l}{z_{j(q)} z_{j(l)}} \exp \left( \beta \frac{\|\gamma_{j(q)}\|^2 - z_{j(q)}^2}{2z_{j(q)}} + \beta \frac{\|\gamma_{j(l)}\|^2 - z_{j(l)}^2}{2z_{j(l)}} \right), & j(q) \neq j(l). \end{cases}
\end{aligned}$$

By definition (A.17) of  $s(\Gamma)$  it holds for all  $\Gamma \in \mathbb{R}^{p_{\text{sum}}}$  :

$$\frac{1}{s(\Gamma)} \exp \left( \beta \frac{\|\gamma_j\|^2 - z_j^2}{2z_j} \right) \leq 1, \quad \sum_{j=1}^K \frac{1}{s(\Gamma)} \exp \left( \beta \frac{\|\gamma_j\|^2 - z_j^2}{2z_j} \right) = 1.$$

Therefore,

$$\begin{aligned}
\sum_{q,l=1}^{p_{\text{sum}}} \left| \frac{d}{d\gamma^q} h_\beta(s(\Gamma)) \frac{d}{d\gamma^l} h_\beta(s(\Gamma)) \right| &\leq \left\{ \sum_{j=1}^K \frac{1}{s(\Gamma) z_j} \exp \left( \beta \frac{\|\gamma_j\|^2 - z_j^2}{2z_j} \right) \sum_{q=1}^{p_j} \gamma^q \right\}^2 \\
&\leq \left| \max_{1 \leq j \leq K} \|\gamma_j\| \frac{\sqrt{p_j}}{z_j} \right|^2 \\
&\leq \max_{1 \leq j \leq K} \|\gamma_j\|^2 \quad \text{for } z_j \geq \sqrt{p_j}.
\end{aligned}$$

Similarly

$$\begin{aligned}
\sum_{q,l=1}^{p_{\text{sum}}} \left| \frac{d^2}{d\gamma^q d\gamma^l} h_\beta(s(\Gamma)) \right| &\leq \mathbf{c} \beta \max_{1 \leq j \leq K} \|\gamma_j\|^2, \\
\sum_{q,l,b=1}^{p_{\text{sum}}} \left| \frac{d^2}{d\gamma^q d\gamma^l} h_\beta(s(\Gamma)) \frac{d}{d\gamma^b} h_\beta(s(\Gamma)) + \frac{d^3}{d\gamma^q d\gamma^l d\gamma^b} h_\beta(s(\Gamma)) \right| &\leq \mathbf{c} (\beta + \beta^2) \max_{1 \leq j \leq K} \|\gamma_j\|^3.
\end{aligned}$$

□

## A.2 Gaussian comparison

The following Lemma shows how to compare the expected values of a twice differentiable function evaluated at the independent centered Gaussian vectors. This statement is used

for the Gaussian comparison step in the scheme (3.1). The proof of the result is based on the Gaussian interpolation method introduced by Stein (1981) and Slepian (1962) (see also Röllin (2013) and Chernozhukov et al. (2013b) and references therein). The proof is given here in order to keep the text self-contained.

**Lemma A.7** (Gaussian comparison using Slepian interpolation). *Let the  $\mathbb{R}^{p_{\text{sum}}}$ -dimensional random centered vectors  $\bar{\Phi}$  and  $\bar{\Psi}$  be independent and normally distributed,  $f(Z) : \mathbb{R}^{p_{\text{sum}}} \mapsto \mathbb{R}$  is any twice differentiable function s.t. the expected values in the expression below are bounded. Then it holds*

$$|\mathbb{E}f(\bar{\Phi}) - \mathbb{E}f(\bar{\Psi})| \leq \frac{1}{2} \|\text{Var} \bar{\Phi} - \text{Var} \bar{\Psi}\|_{\max} \sup_{t \in [0,1]} \left\| \mathbb{E} \nabla^2 f \left( \bar{\Phi} \sqrt{t} + \bar{\Psi} \sqrt{1-t} \right) \right\|_1.$$

*Proof of Lemma A.7.* Introduce for  $t \in [0, 1]$  the Gaussian vector process  $Z_t$  and the deterministic scalar-valued function  $\varkappa(t)$ :

$$\begin{aligned} Z_t &\stackrel{\text{def}}{=} \bar{\Phi} \sqrt{t} + \bar{\Psi} \sqrt{1-t} \in \mathbb{R}^{p_{\text{sum}}}, \\ \varkappa(t) &\stackrel{\text{def}}{=} \mathbb{E}f(Z(t)), \end{aligned}$$

then  $\mathbb{E}f(\bar{\Phi}) = \varkappa(1)$ ,  $\mathbb{E}f(\bar{\Psi}) = \varkappa(0)$  and

$$|\mathbb{E}f(\bar{\Phi}) - \mathbb{E}f(\bar{\Psi})| = |\varkappa(1) - \varkappa(0)| \leq \int_0^1 |\varkappa'(t)| dt.$$

Let us consider  $\varkappa'(t)$ :

$$\begin{aligned} \varkappa'(t) &= \frac{d}{dt} \mathbb{E}f(Z_t) = \mathbb{E} \left[ \{\nabla f(Z_t)\}^\top \frac{d}{dt} Z_t \right] \\ &= \frac{1}{2\sqrt{t}} \mathbb{E} \left\{ \bar{\Phi}^\top \nabla f(Z_t) \right\} - \frac{1}{2\sqrt{1-t}} \mathbb{E} \left\{ \bar{\Psi}^\top \nabla f(Z_t) \right\}. \end{aligned} \quad (\text{A.18})$$

Further we use the Gaussian integration by parts formula (see e.g Section A.6 in Talagrand (2003)): if  $(x_1, \dots, x_{p_{\text{sum}}})^\top$  is a centered Gaussian vector and  $f(x_1, \dots, x_{p_{\text{sum}}})$  is s.t. the integrals below exist, then it holds for all  $j = 1, \dots, p_{\text{sum}}$ :

$$\mathbb{E} \{x_j f(x_1, \dots, x_{p_{\text{sum}}})\} = \sum_{k=1}^{p_{\text{sum}}} \mathbb{E}(x_j x_k) \mathbb{E} \left\{ \frac{d}{dx_k} f(x_1, \dots, x_{p_{\text{sum}}}) \right\}. \quad (\text{A.19})$$

Let  $\bar{\Phi}^j, \bar{\Psi}^j$  denote the  $j$ -th coordinates of  $\bar{\Phi}$  and  $\bar{\Psi}$ . Let also  $\frac{d}{dx_j} f(Z_t)$  denote the partial derivative of the vectors  $f(Z_t)$  w.r.t. the  $j$ -th coordinate of  $Z_t$ . Then it holds

due to (A.19):

$$\begin{aligned} \mathbb{E} \left\{ \bar{\Phi}^\top \nabla f(Z_t) \right\} &= \sum_{j=1}^{p_{\text{sum}}} \mathbb{E} \left\{ \bar{\Phi}^j \frac{d}{d_j} f(Z_t) \right\} = \sum_{j,q=1}^{p_{\text{sum}}} \mathbb{E} \left( \bar{\Phi}^j \bar{\Phi}^q \right) \mathbb{E} \left\{ \frac{d}{d \bar{\Phi}^q} \frac{d}{d_j} f(Z_t) \right\} \\ &= \sqrt{t} \sum_{j,q=1}^{p_{\text{sum}}} \mathbb{E} \left( \bar{\Phi}^j \bar{\Phi}^q \right) \mathbb{E} \left\{ \frac{d^2}{d_q d_j} f(Z_t) \right\}. \end{aligned}$$

Similarly for the second term in (A.18):

$$\mathbb{E} \left\{ \bar{\Psi}^\top \nabla f(Z_t) \right\} = \sqrt{1-t} \sum_{j,q=1}^{p_{\text{sum}}} \mathbb{E} \left( \bar{\Psi}^j \bar{\Psi}^q \right) \mathbb{E} \left\{ \frac{d^2}{d_q d_j} f(Z_t) \right\},$$

therefore

$$\begin{aligned} \varkappa(t) &= \frac{1}{2} \sum_{j=1}^{p_{\text{sum}}} \sum_{q=1}^{p_{\text{sum}}} \left\{ \mathbb{E} \left( \bar{\Phi}^j \bar{\Phi}^q \right) - \mathbb{E} \left( \bar{\Psi}^j \bar{\Psi}^q \right) \right\} \mathbb{E} \left\{ \frac{d^2}{d_q d_j} f(Z_t) \right\} \\ &\leq \frac{1}{2} \left\| \text{Var} \bar{\Phi} - \text{Var} \bar{\Psi} \right\|_{\max} \sup_{t \in [0,1]} \left\| \mathbb{E} \nabla^2 f(Z_t) \right\|_1. \end{aligned}$$

□

### A.3 Simultaneous anti-concentration for $\ell_2$ -norms of Gaussian vectors

**Lemma A.8** (Simultaneous Gaussian anti-concentration). *Let  $(\bar{\Phi}_1^\top, \dots, \bar{\Phi}_K^\top)^\top \in \mathbb{R}^{p_{\text{sum}}}$  be centered normally distributed random vector, and  $\bar{\Phi}_j \in \mathbb{R}^{p_j}$ ,  $j = 1, \dots, K$ . It holds for all  $z_j \geq \sqrt{p_j}$  and  $0 < \Delta_j \leq z_j$ ,  $j = 1, \dots, K$ :*

$$\mathbb{P} \left( \bigcup_{j=1}^K \{ \|\bar{\Phi}_j\| > z_j \} \right) - \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\bar{\Phi}_j\| > z_j + \Delta_j \} \right) \leq \Delta_{\text{ac}}(\{\Delta_j\}),$$

where

$$\Delta_{\text{ac}}(\{\Delta_j\}) \leq \mathbf{C} \left\{ \varkappa \sqrt{1 \vee \log(K/2)} + \mathbf{C} \max_{1 \leq j \leq K} \{ \Delta_j \} \sqrt{\max_{1 \leq j \leq K} \log(2z_j/\Delta_j)} \right\},$$

and  $\varkappa \stackrel{\text{def}}{=} \max_{1 \leq j \leq K} \{ \Delta_j / z_j \} \leq 1$  is a deterministic positive constant. An explicit definition of  $\Delta_{\text{ac}}(\{\Delta_j\})$  is given in (A.22).

*Proof of Lemma A.8.*

$$\begin{aligned}
& \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\bar{\phi}_j\| > z_j \} \right) - \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\bar{\phi}_j\| > z_j + \Delta_j \} \right) \\
& \leq \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\bar{\phi}_j\| z_j^{-1} - 1 > 0 \} \right) - \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\bar{\phi}_j\| z_j^{-1} - 1 > \varkappa \} \right) \\
& = \mathbb{P} \left( \max_{1 \leq j \leq K} \{ \|\bar{\phi}_j\| z_j^{-1} - 1 \} > 0 \right) - \mathbb{P} \left( \max_{1 \leq j \leq K} \{ \|\bar{\phi}_j\| z_j^{-1} - 1 \} > \varkappa \right) \\
& \leq \mathbb{P} \left( 0 \leq \max_{1 \leq j \leq K} \{ \|\bar{\phi}_j\| z_j^{-1} - 1 \} \leq \varkappa \right). \tag{A.20}
\end{aligned}$$

It holds

$$\|\bar{\phi}_j\| = \sup_{\substack{\gamma \in \mathbb{R}^{p_j}, \\ \|\gamma\|=1}} \{ \gamma^\top \bar{\phi}_j \}.$$

Let  $G_j(\varepsilon_j) \subset \mathbb{R}^{p_j}$  (for  $1 \leq j \leq K$ ) denote a finite  $\varepsilon_j$ -net on  $(p_j - 1)$ -sphere of radius 1:

$$\forall \gamma \in \mathbb{R}^{p_j} \text{ s.t. } \|\gamma\| = 1 \quad \exists \gamma_0 \in G_j(\varepsilon_j) : \|\gamma - \gamma_0\| \leq \varepsilon_j, \|\gamma_0\| = 1.$$

This implies for all  $j = 1, \dots, K$

$$(1 - \varepsilon_j) \|\bar{\phi}_j\| \leq \max_{\gamma \in G_j(\varepsilon_j)} \{ \gamma^\top \bar{\phi}_j \} \leq \|\bar{\phi}_j\|.$$

Let us take  $\varepsilon_1, \dots, \varepsilon_K > 0$  s.t.  $\forall j = 1, \dots, K$

$$\varepsilon_j \|\bar{\phi}_j\| z_j^{-1} \leq \varkappa, \tag{A.21}$$

then

$$0 \leq \max_{1 \leq j \leq K} \left\{ \frac{\|\bar{\phi}_j\|}{z_j} \right\} - \max_{1 \leq j \leq K} \max_{\gamma \in G_j(\varepsilon_j)} \left\{ \frac{\gamma^\top \bar{\phi}_j}{z_j} \right\} \leq \varkappa,$$

and the inequality (A.20) continues as

$$\begin{aligned}
& \mathbb{P} \left( 0 \leq \max_{1 \leq j \leq K} \{ \|\bar{\phi}_j\| z_j^{-1} - 1 \} \leq \varkappa \right) \\
& \leq \mathbb{P} \left( \left| \max_{1 \leq j \leq K} \sup_{\gamma \in G_j(\varepsilon_j)} \left\{ \frac{\gamma^\top \bar{\phi}_j}{z_j} \right\} - 1 \right| \leq \varkappa \right).
\end{aligned}$$

The random values  $\gamma^\top \bar{\phi}_j z_j^{-1} \sim \mathcal{N}(0, z_j^{-2} \text{Var}\{\gamma^\top \bar{\phi}_j\})$ . The anti-concentration inequality by Chernozhukov et al. (2014c) for the maximum of a centered high-dimensional

Gaussian vector (see Theorem A.9 below), applied to  $\max_{1 \leq j \leq K} \sup_{\gamma \in G_j(\varepsilon_j)} \left\{ \gamma^\top \bar{\phi}_j z_j^{-1} \right\}$ , implies

$$\begin{aligned} \mathbb{P} \left( \left| \max_{1 \leq j \leq K} \sup_{\gamma \in G_j(\varepsilon_j)} \left\{ \frac{\gamma^\top \bar{\phi}_j}{z_j} \right\} - 1 \right| \leq \varkappa \right) \\ \leq \Delta_{\text{ac}} \stackrel{\text{def}}{=} \mathbf{C}_{\text{ac}} \varkappa \sqrt{1 \vee \log \left( \varkappa^{-1} \sum_{j=1}^K \{2/\varepsilon_j\}^{p_j} \right)}, \end{aligned} \quad (\text{A.22})$$

where the constant  $\mathbf{C}_{\text{ac}}$  depends on min and max of  $\text{Var}\{\gamma^\top \bar{\phi}_j z_j^{-1}\} \leq \mathbb{E}\|\bar{\phi}_j\|^2 z_j^{-2} \leq 1$ ; the sum  $\sum_{j=1}^K \{2/\varepsilon_j\}^{p_j}$  is proportional to cardinality of the set  $\{\gamma^\top \bar{\phi}_j z_j^{-1}, \gamma \in G_j(\varepsilon_j), j = 1, \dots, K\}$ . If one takes  $\varepsilon_j = 2\mathbf{C}\{\Delta_j/(2z_j)\}^{\frac{p_{\min}+1}{p_j+1}}$ , then (A.21) holds with exponentially high probability due to Gaussianity of the vectors  $\bar{\phi}_j$  and Theorem 1.2 in Spokoiny (2012b), hence

$$\begin{aligned} \Delta_{\text{ac}} &\leq \mathbf{C}_{\text{ac}} \varkappa \sqrt{1 \vee \mathbf{C} \log \left( \frac{1}{2} \sum_{j=1}^K \{2/\varepsilon_j\}^{p_j+1} \right)} \\ &\leq \mathbf{C}_{\text{ac}} \left\{ \varkappa \sqrt{1 \vee \log(K/2)} + \mathbf{C} \max_{1 \leq j \leq K} \{\Delta_j\} \sqrt{\max_{1 \leq j \leq K} \log(2z_j/\Delta_j)} \right\}. \end{aligned} \quad (\text{A.23})$$

□

**Theorem A.9** (Anti-concentration inequality for maxima of a Gaussian random vector, Chernozhukov et al. (2014c)). *Let  $(X_1, \dots, X_p)^\top$  be a centered Gaussian random vector with  $\sigma_j^2 \stackrel{\text{def}}{=} \mathbb{E}X_j^2 > 0$  for all  $1 \leq j \leq p$ . Let  $\underline{\sigma} \stackrel{\text{def}}{=} \min_{1 \leq j \leq p} \sigma_j$ ,  $\bar{\sigma} \stackrel{\text{def}}{=} \max_{1 \leq j \leq p} \sigma_j$ . Then for every  $\epsilon > 0$*

$$\sup_{x \in \mathbb{R}} \mathbb{P} \left( \left| \max_{1 \leq j \leq p} X_j - x \right| \leq \epsilon \right) \leq \mathbf{C}_{\text{ac}} \epsilon \sqrt{1 \vee \log(p/\epsilon)},$$

where  $\mathbf{C}_{\text{ac}}$  depends only on  $\underline{\sigma}$  and  $\bar{\sigma}$ . When the variances are all equal, namely  $\underline{\sigma} = \bar{\sigma} = \sigma$ ,  $\log(p/\epsilon)$  on the right side can be replaced by  $\log p$ .

#### A.4 Proof of Proposition A.1

*Proof of Proposition A.1.* Let  $\Phi \stackrel{\text{def}}{=} (\phi_1^\top, \dots, \phi_K^\top)^\top \in \mathbb{R}^{p_{\text{sum}}}$  for  $p_{\text{sum}} \stackrel{\text{def}}{=} p_1 + \dots + p_K$  (as in (A.5)), and similarly  $\Psi \stackrel{\text{def}}{=} (\psi_1^\top, \dots, \psi_K^\top)^\top \in \mathbb{R}^{p_{\text{sum}}}$ . Let also  $\bar{\Phi} \sim \mathcal{N}(0, \text{Var} \Phi)$  and  $\bar{\Psi} \sim \mathcal{N}(0, \text{Var} \Psi)$ . Introduce the following value, which comes from Lemma A.7 on Gaussian comparison:

$$\begin{aligned} \delta_2(\Delta, \beta) &\stackrel{\text{def}}{=} \mathbf{C}_2(\Delta, \beta) \max_{1 \leq j \leq K} \sup_{t \in [0,1]} \left\{ \mathbb{E} \|\bar{\phi}_j \sqrt{t} + \bar{\psi}_j \sqrt{1-t}\|^2 \right\} \\ &\leq \mathbf{C}_2(\Delta, \beta) \max_{1 \leq j \leq K} \max \left\{ \text{tr} \text{Var}(\bar{\phi}_j), \text{tr} \text{Var}(\bar{\psi}_j) \right\}. \end{aligned} \quad (\text{A.24})$$



It holds

$$\begin{aligned}
& \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\phi_j\| > z_j \} \right) \\
& \stackrel{\text{by L. A.3}}{\geq} \mathbb{E} H_{\Delta, \beta} \left( \bar{\Phi}, \mathbf{z} + \frac{3 \log(K)}{2\beta} \mathbf{1}_K \right) - \delta_{3, \phi}(\Delta, \beta) \\
& \stackrel{\text{by L. A.7, A.6}}{\geq} \mathbb{E} H_{\Delta, \beta} \left( \bar{\Psi}, \mathbf{z} + \frac{3 \log(K)}{2\beta} \mathbf{1}_K \right) - \frac{1}{2} \delta_{\Sigma}^2 \delta_2(\Delta, \beta) - \delta_{3, \phi}(\Delta, \beta) \\
& \stackrel{\text{by L. A.3}}{\geq} \mathbb{P} \left( \bigcup_{j=1}^K \left\{ \|\bar{\psi}_j\| > z_j + \Delta + \frac{3 \log(K)}{2\beta} \right\} \right) - \frac{1}{2} \delta_{\Sigma}^2 \delta_2(\Delta, \beta) - \delta_{3, \phi}(\Delta, \beta) \\
& \stackrel{\text{by L. A.8}}{\geq} \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\bar{\psi}_j\| > z_j - \delta_{z_j} - \Delta \} \right) - \frac{1}{2} \delta_{\Sigma}^2 \delta_2(\Delta, \beta) - \delta_{3, \phi}(\Delta, \beta) \\
& \quad - 2\Delta_{\text{ac}} \left( \{ \delta_{z_j} \} + 2\Delta + \frac{3 \log(K)}{\beta} \right) \tag{A.25}
\end{aligned}$$

$$\begin{aligned}
& \stackrel{\text{by L. A.2}}{\geq} \mathbb{P} \left( \bigcup_{j=1}^K \{ \|\psi_j\| > z_j - \delta_{z_j} \} \right) - \frac{1}{2} \delta_{\Sigma}^2 \delta_2(\Delta, \beta) \tag{A.26} \\
& \quad - \delta_{3, \phi}(\Delta, \beta) - \delta_{3, \psi}(\Delta, \beta) - 2\Delta_{\text{ac}} \left( \{ \delta_{z_j} \} + 2\Delta + \frac{3 \log(K)}{\beta} \right),
\end{aligned}$$

where  $\delta_{3, \psi}(\Delta, \beta)$  is defined similarly to  $\delta_{3, \phi}(\Delta, \beta)$  in (A.15):

$$\delta_{3, \psi}(\Delta, \beta) \stackrel{\text{def}}{=} \frac{\mathbf{C}_3(\Delta, \beta)}{3} \frac{p_{\max}^{3/2}}{n^{1/2}} \log^{1/2}(K) \log^{3/2}(np_{\text{sum}}) (2\nu_0^2 c_{\psi}^2 \lambda_{\psi, \max}^2)^{3/2}. \tag{A.27}$$

By Lemma A.8 inequality (A.25) requires the following:  $\delta_{z_j} + 2\Delta + \frac{3 \log(K)}{\beta} \leq z_j$ . The bound in the inverse direction is derived similarly. Denote the approximating error term obtained in (A.26) as

$$\Delta_{\ell_2} \stackrel{\text{def}}{=} \frac{1}{2} \delta_{\Sigma}^2 \delta_2(\Delta, \beta) + \delta_{3, \phi}(\Delta, \beta) + \delta_{3, \psi}(\Delta, \beta) + 2\Delta_{\text{ac}} \left( \{ \delta_{z_j} \} + 2\Delta + \frac{3 \log(K)}{\beta} \right).$$

Consider this term in more details, by inequality (A.23)

$$\begin{aligned}
\Delta_{\text{ac}} \left( \{ \delta_{z_j} \} + 2\Delta + \frac{3 \log(K)}{\beta} \right) & \leq \max_{1 \leq j \leq K} \left( \delta_{z_j} + 2\Delta + \frac{3 \log(K)}{\beta} \right) \\
& \quad \times \left\{ \mathbf{C} \frac{\log^{1/2}(K)}{z_j} + \log^{1/2}(2z_{\max}) - \log^{1/2} \left( \delta_{z_j} + 2\Delta + \frac{3 \log(K)}{\beta} \right) \right\}.
\end{aligned}$$

Let us take  $\beta = \frac{\log(K)}{\Delta}$ , then

$$\begin{aligned}
\Delta_{\text{ac}} &\leq 5\mathbf{C}\Delta \frac{\log^{1/2}(K)}{z_{\min}} + \mathbf{C} \max_{1 \leq j \leq K} \frac{\delta_{z_j}}{z_j} \log^{1/2}(K) \\
&\quad + \mathbf{C}(5\Delta + \delta_{z,\max}) \left( \log^{1/2}(2z_{\max}) + \sqrt{-\log(\delta_{z,\min} + 5\Delta)} \right), \\
&\leq 5\mathbf{C}\Delta \frac{\log^{1/2}(K)}{z_{\min}} + \mathbf{C} \max_{1 \leq j \leq K} \frac{\delta_{z_j}}{z_j} \log^{1/2}(K) \\
&\quad + 2\mathbf{C}(5\Delta + \delta_{z,\max}) \sqrt{-\log(\delta_{z,\min} + 5\Delta)} \\
&\leq 5\mathbf{C}\Delta \frac{\log^{1/2}(K)}{z_{\min}} + \mathbf{C} \max_{1 \leq j \leq K} \frac{\delta_{z_j}}{z_j} \log^{1/2}(K) + 2\mathbf{C}(5\Delta + \delta_{z,\max}) \sqrt{-\log(5\Delta)} \\
&\leq 5\mathbf{C}\Delta \left\{ \frac{\log^{1/2}(K)}{z_{\min}} + 2.4 \log^{1/2}(5n^{1/2}) \right\} + \mathbf{C} \max_{1 \leq j \leq K} \frac{\delta_{z_j}}{z_j} \log^{1/2}(K) \\
&\leq 6\mathbf{C}\Delta \left\{ \frac{\log^{1/2}(K)}{z_{\min}} + 0.4 \log^{1/2}(5n^{1/2}) \right\}, \tag{A.28}
\end{aligned}$$

where the second inequality holds for  $\delta_{z,\min} + 5\Delta \leq 1/(2z_{\max})$ , and the last one holds for  $\delta_{z,\max} \leq \Delta$  and  $\Delta \geq n^{-1/2}$ .

$$\begin{aligned}
\delta_{3,\phi}(\Delta, \beta) + \delta_{3,\psi}(\Delta, \beta) &\stackrel{\text{by (A.27)}}{\leq} \mathbf{C} \frac{\log^{5/2}(K) p_{\max}^{3/2}}{\Delta^3 n^{1/2}} \log^{3/2}(np_{\text{sum}}) (\lambda_{\phi,\max}^3 + \lambda_{\psi,\max}^3), \tag{A.29} \\
\delta_{\Sigma} \delta_2(\Delta, \beta) &\stackrel{\text{by (A.24)}}{\leq} \mathbf{C} \delta_{\Sigma}^2 \frac{\log(K)}{\Delta^2} \max_{1 \leq j \leq K} \max \{ \text{tr Var}(\bar{\phi}_j), \text{tr Var}(\bar{\psi}_j) \} \\
&\leq \mathbf{C} \delta_{\Sigma}^2 \frac{\log(K)}{\Delta^2} p_{\max} \max \{ \lambda_{\phi,\max}^2, \lambda_{\psi,\max}^2 \}.
\end{aligned}$$

After minimizing the sum of the expressions (A.28) and (A.29) w.r.t  $\Delta$ , we have

$$\begin{aligned}
\Delta_{\ell_2} &\leq 12.5\mathbf{C} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \max \{ \lambda_{\phi,\max}, \lambda_{\psi,\max} \}^{3/4} \\
&\quad + 3.2\mathbf{C} \delta_{\Sigma}^2 p_{\max} z_{\min}^{1/2} \left( \frac{p_{\max}^3}{n} \right)^{1/4} \log^2(K) \log^{3/4}(np_{\text{sum}}) \max \{ \lambda_{\phi,\max}, \lambda_{\psi,\max} \}^{7/2} \\
&\leq 25\mathbf{C} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \max \{ \lambda_{\phi,\max}, \lambda_{\psi,\max} \}^{3/4},
\end{aligned}$$

where the last inequality holds for

$$\delta_{\Sigma}^2 \leq 4\mathbf{C} p_{\max}^{-1} z_{\min}^{-1/2} \left( \frac{p_{\max}^3}{n} \right)^{-1/8} \log^{-7/8}(K) \log^{-3/8}(np_{\text{sum}}) (\max \{ \lambda_{\phi,\max}, \lambda_{\psi,\max} \})^{-11/4}.$$

□

## B Square-root Wilks approximations

This section's goal is to derive square root Wilks approximations simultaneously for  $K$  parametric models, for the  $\mathbf{Y}$  and bootstrap worlds. This is done in Section B.3 below. Both of the results are used in the approximating scheme (3.1) for the bootstrap justification. In order to make the text self-contained we recall in Section B.1 some results from the general finite sample theory by Spokoiny (2012a,b, 2013). In Section B.2 we recall similar finite sample results for the bootstrap world for a single parametric model, obtained in Spokoiny and Zhilova (2014).

### B.1 Finite sample theory

Let us use the notations given in the introduction:  $L_k(\boldsymbol{\theta})$ ,  $k = 1, \dots, K$  are the log-likelihood processes, which depend on the data  $\mathbf{Y}$  and correspond to the regular parametric families of probability distributions  $\{\mathbb{P}_k(\boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta_k \subset \mathbb{R}^{p_k}\}$ . The general finite sample approach by Spokoiny (2012a) does not require that the true distribution  $\mathbb{P}$  of the data  $\mathbf{Y}$  belongs to any of the parametric families  $\{\mathbb{P}_k(\boldsymbol{\theta})\}$ . The target parameters  $\boldsymbol{\theta}_k^*$  are defined as in (1.3) by projection of the true measure  $\mathbb{P}$  on  $\{\mathbb{P}_k(\boldsymbol{\theta})\}$ . Let  $D_k^2$  denote the full Fisher information  $p_k \times p_k$  matrices, which are deterministic, symmetric and positive-definite:

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

Centered  $p_k$ -dimensional random vectors  $\boldsymbol{\xi}_k$  denote the normalised scores:

$$\boldsymbol{\xi}_k \stackrel{\text{def}}{=} D_k^{-1} \nabla_{\boldsymbol{\theta}} L_k(\boldsymbol{\theta}_k^*).$$

Introduce the following elliptic vicinities around the true points  $\boldsymbol{\theta}_k^*$ :

$$\Theta_{0,k}(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} \in \Theta_k : \|D_k(\boldsymbol{\theta} - \boldsymbol{\theta}_k^*)\| \leq \mathbf{r}\}.$$

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

**Theorem B.1** (Concentration of the MLE, Spokoiny (2013)). *Let the conditions  $(ED_0)$ ,  $(ED_2)$ ,  $(\mathcal{L}_0)$ ,  $(\mathcal{I})$  and  $(\mathcal{L}_\mathbf{r})$  be fulfilled. If for each  $k = 1, \dots, K$  for the constants  $\mathbf{r}_{0,k} > 0$  and for the functions  $\mathbf{b}_k(\mathbf{r})$  from  $(\mathcal{L}_\mathbf{r})$  holds:*

$$\mathbf{b}_k(\mathbf{r})\mathbf{r} \geq 2 \left\{ \mathfrak{J}_{\text{qf}}(\mathbf{x}, \mathbb{B}_k) + 6\omega_k \nu_k \mathfrak{J}_k(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_{0,k})) \right\}, \quad \mathbf{r} > \mathbf{r}_{0,k} \quad (\text{B.1})$$

where the functions  $\mathfrak{Z}_k(\mathbf{x})$  and  $\mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathcal{B}_k)$  are defined in (B.3) and (B.4) respectively, then it holds for all  $k = 1, \dots, K$

$$\mathbb{P} \left( \tilde{\boldsymbol{\theta}}_k \notin \Theta_{0,k}(\mathbf{r}_{0,k}) \right) \leq 3e^{-\mathbf{x}}.$$

The constants  $\omega_k, \nu_k$  and  $\mathbf{a}_k$  come from the imposed conditions  $(\mathbf{ED}_0) - (\mathcal{I})$  (from Section 5). In the case 5.3  $\mathbf{r}_{0,k} \geq \mathbf{C}\sqrt{p_k + \mathbf{x}}$ .

**Theorem B.2** (Wilks approximation, Spokoiny (2013)). *Under the conditions of Theorem B.1 for some  $\mathbf{r}_{0,k} > 0$  s.t. (B.1) is fulfilled, it holds for each  $k = 1, \dots, K$  with probability  $\geq 1 - 5e^{-\mathbf{x}}$*

$$\begin{aligned} \left| 2 \left\{ L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) \right\} - \|\boldsymbol{\xi}_k\|^2 \right| &\leq \Delta_{k, \mathbb{W}^2}(\mathbf{r}_{0,k}, \mathbf{x}), \\ \left| \sqrt{2 \left\{ L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) \right\}} - \|\boldsymbol{\xi}_k\| \right| &\leq \Delta_{k, \mathbb{W}}(\mathbf{r}_{0,k}, \mathbf{x}) \end{aligned}$$

for

$$\Delta_{k, \mathbb{W}}(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} 3\mathbf{r} \{ \delta(\mathbf{r}) + 6\nu_k \mathfrak{Z}_k(\mathbf{x}) \omega_k \}, \quad (\text{B.2})$$

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

$$\mathfrak{Z}_k(\mathbf{x}) \stackrel{\text{def}}{=} 2\sqrt{p_k} + \sqrt{2\mathbf{x}} + 4p_k(\mathbf{x}\mathbf{g}_k^{-2} + 1)\mathbf{g}_k^{-1}. \quad (\text{B.3})$$

In the case 5.3 it holds for  $\mathbf{r} \leq \mathbf{r}_{0,k}$ :

$$\Delta_{k, \mathbb{W}}(\mathbf{r}, \mathbf{x}) \leq \mathbf{C} \frac{p_k + \mathbf{x}}{\sqrt{n}}, \quad \Delta_{k, \mathbb{W}^2}(\mathbf{r}, \mathbf{x}) \leq \mathbf{C} \sqrt{\frac{(p_k + \mathbf{x})^3}{n}}.$$

The constants  $\mathbf{g}_k$  and  $\delta_k(\mathbf{r})$  come from the imposed conditions  $(\mathbf{ED}_0)$ ,  $(\mathcal{L}_0)$  (from Section 5). The function  $\mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathcal{B}_k)$ , defined in (B.4), corresponds to the quantile function of deviations of the approximating random value  $\|\boldsymbol{\xi}_k\|$  (see Theorem B.3 below).

The following theorem characterizes the tail behaviour of the approximating terms  $\|\boldsymbol{\xi}_k\|^2$ . It means that with bounded exponential moments of the vectors  $\boldsymbol{\xi}_k$  (conditions  $(\mathbf{ED}_0)$ ,  $(\mathcal{I})$ ) its squared Euclidean norms  $\|\boldsymbol{\xi}_k\|^2$  have three regimes of deviations: sub-Gaussian, Poissonian and large-deviations' zone.

**Theorem B.3** (Deviation bound for a random quadratic form, Spokoiny (2012b)). *Let condition  $(\mathbf{ED}_0)$  be fulfilled, then for  $\mathbf{g}_k \geq \sqrt{2 \text{tr}(\mathcal{B}_k^2)}$  it holds for each  $k = 1, \dots, K$ :*

$$\mathbb{P} \left( \|\boldsymbol{\xi}_k\|^2 \geq \mathfrak{Z}_{\text{qf}}^2(\mathbf{x}, \mathcal{B}_k) \right) \leq 2e^{-\mathbf{x}} + 8.4e^{-\mathbf{x}_{c,k}},$$

where  $\mathcal{B}_k^2 \stackrel{\text{def}}{=} D_k^{-1} V_k^2 D_k^{-1}$ ,  $\lambda(\mathcal{B}_k)$  is a maximum eigenvalue of  $\mathcal{B}_k^2$ ,

$$\mathfrak{Z}_{\text{qf}}^2(\mathbf{x}, \mathcal{B}_k) \stackrel{\text{def}}{=} \begin{cases} \text{tr}(\mathcal{B}_k^2) + \sqrt{8 \text{tr}(\mathcal{B}_k^4)} \mathbf{x}, & \mathbf{x} \leq \sqrt{2 \text{tr}(\mathcal{B}_k^4) / \{18\lambda(\mathcal{B}_k)\}}, \\ \text{tr}(\mathcal{B}_k^2) + 6\mathbf{x}\lambda(\mathcal{B}_k), & \sqrt{2 \text{tr}(\mathcal{B}_k^4) / \{18\lambda(\mathcal{B}_k)\}} < \mathbf{x} \leq \mathbf{x}_{c,k}, \\ |\mathbf{z}_{c,k} + 2(\mathbf{x} - \mathbf{x}_{c,k})/\mathbf{g}_{c,k}|^2 \lambda(\mathcal{B}_k), & \mathbf{x} > \mathbf{x}_{c,k}, \end{cases} \quad (\text{B.4})$$

$$\begin{aligned} 2\mathbf{x}_{c,k} &\stackrel{\text{def}}{=} 2\mathbf{x}_{c,k}(\mathcal{B}_k) \stackrel{\text{def}}{=} \mu_c \mathbf{z}_{c,k}^2 + \log \det(\mathbf{I}_{p_k} - \mu_c \mathcal{B}_k^2 / \lambda(\mathcal{B}_k)), \\ \mathbf{z}_{c,k}^2 &\stackrel{\text{def}}{=} \{\mathbf{g}_k^2 / \mu_c^2 - \text{tr}(\mathcal{B}_k^2) / \mu_c\} / \lambda(\mathcal{B}_k), \\ \mathbf{g}_{c,k} &\stackrel{\text{def}}{=} \sqrt{\mathbf{g}_k^2 - \mu_c \text{tr}(\mathcal{B}_k^2)} / \sqrt{\lambda(\mathcal{B}_k)}, \\ \mu_c &\stackrel{\text{def}}{=} 2/3. \end{aligned} \quad (\text{B.5})$$

The matrices  $V_k^2$  come from condition **(ED<sub>0</sub>)** and can be defined as

$$V_k^2 \stackrel{\text{def}}{=} \text{Var} \{\nabla_{\boldsymbol{\theta}} L_k(\boldsymbol{\theta}_k^*)\}.$$

By condition **(I)**  $\text{tr}(\mathcal{B}_k^2) \leq \mathfrak{a}_k^2 p_k$ ,  $\text{tr}(\mathcal{B}_k^4) \leq \mathfrak{a}_k^4 p_k$  and  $\lambda(\mathcal{B}_k) \leq \mathfrak{a}_k^2$ . In the case 5.3  $\mathbf{g}_k = \mathbf{C}\sqrt{n}$ , hence  $\mathbf{x}_{c,k} = \mathbf{C}n$ , and for  $\mathbf{x} \leq \mathbf{x}_{c,k}$  it holds:

$$\mathfrak{Z}_{\text{qf}}^2(\mathbf{x}, \mathcal{B}_k) \leq \mathfrak{a}_k^2(p_k + 6\mathbf{x}). \quad (\text{B.6})$$

## B.2 Finite sample theory for the bootstrap world

Introduce for each  $k = 1, \dots, K$  the bootstrap score vectors at the point  $\boldsymbol{\theta} \in \Theta_k$ :

$$\begin{aligned} \boldsymbol{\xi}_k^\circ(\boldsymbol{\theta}) &\stackrel{\text{def}}{=} D_k^{-1} \nabla_{\boldsymbol{\theta}} \zeta_k^\circ(\boldsymbol{\theta}) \\ &= \sum_{i=1}^n D_k^{-1} \nabla_{\boldsymbol{\theta}} \ell_{i,k}(\boldsymbol{\theta})(u_i - 1). \end{aligned}$$

**Theorem B.4** (Bootstrap Wilks approximation, Spokoiny and Zhilova (2014)). *Under the conditions of Theorems B.1 and B.5 for each  $k = 1, \dots, K$  and some  $\mathbf{r}_{0,k}^2 \geq 0$  s.t. (B.1) and (B.9) are fulfilled, it holds for each  $k$  with  $\mathbb{P}$ -probability  $\geq 1 - 5e^{-\mathbf{x}}$*

$$\begin{aligned} \mathbb{P}^\circ \left( \left| \sup_{\boldsymbol{\theta} \in \Theta_k} 2 \left\{ L_k^\circ(\boldsymbol{\theta}) - L_k^\circ(\tilde{\boldsymbol{\theta}}_k) \right\} - \|\boldsymbol{\xi}_k^\circ(\tilde{\boldsymbol{\theta}}_k)\|^2 \right| \leq \Delta_{k, \text{W}^2}^\circ(\mathbf{r}_{0,k}, \mathbf{x}) \right) &\geq 1 - 4e^{-\mathbf{x}}, \\ \mathbb{P}^\circ \left( \left| \sqrt{\sup_{\boldsymbol{\theta} \in \Theta_k} 2 \left\{ L_k^\circ(\boldsymbol{\theta}) - L_k^\circ(\tilde{\boldsymbol{\theta}}_k) \right\}} - \|\boldsymbol{\xi}_k^\circ(\tilde{\boldsymbol{\theta}}_k)\| \right| \leq \Delta_{k, \text{W}}^\circ(\mathbf{r}_{0,k}, \mathbf{x}) \right) &\geq 1 - 4e^{-\mathbf{x}}. \end{aligned}$$

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

$$\begin{aligned}\Delta_{k,W}^\circ(\mathbf{r}, \mathbf{x}) &\stackrel{\text{def}}{=} 2\Delta_{k,W}(\mathbf{r}, \mathbf{x}) + 36\nu_k \mathbf{r} \omega_{1,k}(\mathbf{r}, \mathbf{x}) \mathfrak{Z}_k(\mathbf{x}), \\ \Delta_{k,W^2}^\circ(\mathbf{r}, \mathbf{x}) &\stackrel{\text{def}}{=} \frac{1}{18} \{12\mathbf{r} \Delta_{k,W}^\circ(\mathbf{r}, \mathbf{x}) + \Delta_{k,W}^\circ(\mathbf{r}, \mathbf{x})^2\}.\end{aligned}\tag{B.7}$$

$\Delta_{k,W}(\mathbf{r}, \mathbf{x})$  and  $\mathfrak{Z}_k(\mathbf{x})$  are defined in (B.2) and (B.3) respectively and

$$\omega_{1,k}(\mathbf{r}, \mathbf{x}) = \omega_{1,k} \stackrel{\text{def}}{=} \frac{\mathbf{C}_{m,k}(\mathbf{r})}{\sqrt{n}} + 2\omega_k \nu_k \sqrt{2\mathbf{x}},\tag{B.8}$$

where  $\mathbf{C}_{m,k}(\mathbf{r}), \omega_k, \nu_k$  come from the imposed conditions  $(\mathcal{L}_{0m}), (\mathbf{ED}_2)$  and  $(\mathbf{ED}_0)$ . For the case 5.3 and  $\mathbf{r} \leq \mathbf{r}_{0,k}$  it holds:

$$\Delta_{k,W}^\circ(\mathbf{r}, \mathbf{x}) \leq \mathbf{C} \frac{p_k + \mathbf{x}}{\sqrt{n}} \sqrt{\mathbf{x}}, \quad \Delta_{k,W^2}^\circ(\mathbf{r}, \mathbf{x}) \leq \mathbf{C} \sqrt{\frac{(p_k + \mathbf{x})^3}{n}} \sqrt{\mathbf{x}}.$$

and  $\omega_{1,k}(\mathbf{r}) \leq \mathbf{C}\mathbf{r}/n + \mathbf{C}\sqrt{\mathbf{x}/n}$ .

**Theorem B.5** (Concentration of the bootstrap MLE, Spokoiny and Zhilova (2014)).  
Let the conditions of Theorems B.1 and B.7,  $(\mathcal{L}_{0m})$  and  $(\mathbf{ED}_{2m})$  be fulfilled. If the following holds for each  $k = 1, \dots, K$ ,  $\omega_{1,k}(\mathbf{r}, \mathbf{x})$  defined in (B.8) and the  $\mathbb{P}$ -random matrices  $\mathcal{B}_k^2 \stackrel{\text{def}}{=} D_k^{-1} \text{Var}^\circ \{ \nabla_{\boldsymbol{\theta}} L_k^\circ(\boldsymbol{\theta}_k^*) \} D_k^{-1}$ :

$$\begin{aligned}\mathbf{b}_k(\mathbf{r})\mathbf{r} &\geq 2 \{ \mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathcal{B}_k) + \mathfrak{Z}_{\text{qf}}(\mathbf{x}, \mathcal{B}_k) + 6\nu_k \mathfrak{Z}_k(\mathbf{x}) \omega_{1,k}(\mathbf{r}_{0,k}) \mathbf{r}_{0,k} \} \\ &\quad + 12\nu_k (\omega_k + \omega_{1,k}(\mathbf{r}, \mathbf{x})) \mathfrak{Z}_k(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_{0,k})) \quad \text{for } \mathbf{r} > \mathbf{r}_{0,k},\end{aligned}\tag{B.9}$$

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

$$\mathbb{P}^\circ \left( \tilde{\boldsymbol{\theta}}_k^\circ \notin \Theta_{0,k}(\mathbf{r}_{0,k}) \right) \leq 3e^{-\mathbf{x}}.$$

Lemma B.6 below is implied straightforwardly by Lemma B.7 in Spokoiny and Zhilova (2014).

**Lemma B.6.** Let the conditions of  $(\mathbf{Eb}), (\mathcal{L}_{0m})$  and  $(\mathbf{ED}_{2m})$  be fulfilled, then for each  $k = 1, \dots, K$  it holds for  $\mathbf{r} \leq \mathbf{r}_{0,k}$  with  $\mathbb{P}$ -probability  $\geq 1 - e^{-\mathbf{x}}$

$$\mathbb{P}^\circ \left( \sup_{\boldsymbol{\theta} \in \Theta_{0,k}(\mathbf{r})} \|\boldsymbol{\xi}_k^\circ(\boldsymbol{\theta}) - \boldsymbol{\xi}_k^\circ(\boldsymbol{\theta}_k^*)\| \leq \Delta_{\boldsymbol{\xi},k}^\circ(\mathbf{r}, \mathbf{x}) \right) \geq 1 - e^{-\mathbf{x}},$$

where

$$\Delta_{\boldsymbol{\xi},k}^\circ(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} 6\nu_k \mathfrak{Z}_k(\mathbf{x}) \omega_{1,k}(\mathbf{r}, \mathbf{x}) \mathbf{r}$$

In the case 5.3 it holds for the bounding term

$$\Delta_{\boldsymbol{\xi}}^\circ(\mathbf{r}_{0,k}, \mathbf{x}) \leq \mathbf{C} \frac{p_k + \mathbf{x}}{\sqrt{n}} \sqrt{\mathbf{x}}.$$

**Theorem B.7** (Deviation bound for the bootstrap quadratic form, Spokoiny and Zhilova (2014)). *Let conditions  $(\mathbf{Eb})$ ,  $(\mathcal{I})$ ,  $(\widehat{SD}_1)$ ,  $(\mathcal{I}_B)$  be fulfilled, then for each  $k = 1, \dots, K$  and  $\mathbf{g}_k \geq \sqrt{2 \operatorname{tr}(\mathcal{B}_k^2)}$  it holds:*

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

where

$$\mathcal{B}_k^2 \stackrel{\text{def}}{=} D_k^{-1} \mathcal{V}^2(\boldsymbol{\theta}_k^*) D_k^{-1}, \quad \mathcal{V}_k^2(\boldsymbol{\theta}_k^*) \stackrel{\text{def}}{=} \operatorname{Var}^\circ \nabla_{\boldsymbol{\theta}} L_k^\circ(\boldsymbol{\theta}_k^*),$$

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

$$\begin{aligned} \mathfrak{z}_{\text{qf}}^2(\mathbf{x}, \mathcal{B}_k) &\leq \mathbf{a}_k^{\circ 2}(p_k + 6\mathbf{x}) \\ \text{for } \mathbf{a}_k^{\circ 2} &\stackrel{\text{def}}{=} (1 + \delta_{\mathcal{V},k}^2(\mathbf{x}))(\mathbf{a}_k^2 + \mathbf{a}_{B,k}^2) \end{aligned}$$

and  $\delta_{\mathcal{V},k}^2(\mathbf{x})$  defined in (C.1) (see Section C.1 on Bernstein matrix inequalities).

### B.3 Simultaneous square-root Wilks approximations

The statements below follow from the results from Sections B.1 and B.2 by probability union bound.

**Lemma B.8** (Simultaneous concentration bounds).

1. Let conditions of Theorem B.1 be fulfilled and (B.1) hold for each  $k = 1, \dots, K$  with  $\mathbf{x} = \mathbf{x}_1 + \log(K)$  for some  $\mathbf{x}_1 > 0$ , then

$$\mathbb{P} \left( \bigcup_{k=1}^K \left\{ \tilde{\boldsymbol{\theta}}_k \notin \Theta_{0,k}(\mathbf{r}_{0,k}) \right\} \right) \leq 3e^{-\mathbf{x}_1}.$$

2. Let conditions of Theorem B.5 be fulfilled and (B.9) hold for each  $k = 1, \dots, K$  with  $\mathbf{x} = \mathbf{x}_1 + \log(K)$  for some  $\mathbf{x}_1 > 0$ , then it holds with  $\mathbb{P}$ -probability  $\geq 1 - 3e^{-\mathbf{x}_1}$

$$\mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \tilde{\boldsymbol{\theta}}_k^\circ \notin \Theta_{0,k}(\mathbf{r}_{0,k}) \right\} \right) \leq 3e^{-\mathbf{x}_1}.$$

**Lemma B.9** (Simultaneous Wilks approximations).

1. Let the conditions of part 1 of Lemma B.8 be fulfilled for some  $\mathbf{r}_{0,k} > 0$  and  $\mathbf{x} = \mathbf{x}_1 + \log(K)$ , then it holds

$$\mathbb{P} \left( \bigcap_{k=1}^K \left\{ \left| 2 \left\{ L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) \right\} - \|\boldsymbol{\xi}_k\|^2 \right| \leq \Delta_{k, \mathbb{W}^2}(\mathbf{r}_{0,k}, \mathbf{x}_1 + \log(K)) \right\} \right) \geq 1 - 5e^{-\mathbf{x}_1},$$

$$\mathbb{P} \left( \bigcap_{k=1}^K \left\{ \left| \sqrt{2 \left\{ L_k(\tilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) \right\}} - \|\boldsymbol{\xi}_k\| \right| \leq \Delta_{k, \mathbb{W}}(\mathbf{r}_{0,k}, \mathbf{x}_1 + \log(K)) \right\} \right) \geq 1 - 5e^{-\mathbf{x}_1}.$$

2. Let the conditions of parts 1,2 of Lemma B.8 be fulfilled for some  $\mathbf{r}_{0,k} > 0$  and  $\mathbf{x} = \mathbf{x}_1 + \log(K)$ , then it holds with  $\mathbb{P}$ -probability  $\geq 1 - 5e^{-\mathbf{x}_1}$

$$\mathbb{P}^\circ \left( \bigcap_{k=1}^K \left\{ \left| \sup_{\boldsymbol{\theta} \in \Theta_k} 2 \left\{ L_k^\circ(\boldsymbol{\theta}) - L_k^\circ(\tilde{\boldsymbol{\theta}}_k) \right\} - \|\boldsymbol{\xi}_k^\circ(\tilde{\boldsymbol{\theta}}_k)\|^2 \right| \leq \Delta_{k,W^2}^\circ(\mathbf{r}_{0,k}, \mathbf{x}_1 + \log(K)) \right\} \right) \geq 1 - 4e^{-\mathbf{x}_1},$$

$$\mathbb{P}^\circ \left( \bigcap_{k=1}^K \left\{ \left| \sqrt{\sup_{\boldsymbol{\theta} \in \Theta_k} 2 \left\{ L_k^\circ(\boldsymbol{\theta}) - L_k^\circ(\tilde{\boldsymbol{\theta}}_k) \right\}} - \|\boldsymbol{\xi}_k^\circ(\tilde{\boldsymbol{\theta}}_k)\| \right| \leq \Delta_{k,W}^\circ(\mathbf{r}_{0,k}, \mathbf{x}_1 + \log(K)) \right\} \right) \geq 1 - 4e^{-\mathbf{x}_1}.$$

**Lemma B.10.** Let the conditions of Lemma B.6 be fulfilled, then it holds with  $\mathbb{P}$ -probability  $\geq 1 - e^{-\mathbf{x}}$

$$\mathbb{P}^\circ \left( \bigcap_{k=1}^K \left\{ \sup_{\substack{\boldsymbol{\theta} \in \Theta_{0,k}(\mathbf{r}), \\ \mathbf{r} \leq \mathbf{r}_{0,k}}} \|\boldsymbol{\xi}_k^\circ(\boldsymbol{\theta}) - \boldsymbol{\xi}_k^\circ(\boldsymbol{\theta}_k^*)\| \leq \Delta_{\xi,k}^\circ(\mathbf{r}, \mathbf{x} + \log(K)) \right\} \right) \geq 1 - e^{-\mathbf{x}}.$$

## C Proofs of the main results

Before proving the statements from Section 3.2 we formulate below the Bernstein matrix inequality, which is necessary for the further proofs.

### C.1 Bernstein matrix inequality

Here we restate the Theorem 1.4 by Tropp (2012) for the random  $p_{\text{sum}} \times p_{\text{sum}}$  matrix  $\widehat{\mathcal{V}}^2 \stackrel{\text{def}}{=} \text{Var}^\circ(\nabla_{\boldsymbol{\theta}} L_1^\circ(\boldsymbol{\theta}_1^*)^\top, \dots, \nabla_{\boldsymbol{\theta}} L_K^\circ(\boldsymbol{\theta}_K^*)^\top)^\top$  from the bootstrap world. Matrix  $\widehat{\mathcal{V}}^2$  equals to the sum of independent matrices  $\text{Var}^\circ(\nabla_{\boldsymbol{\theta}} \ell_{i,1}(\boldsymbol{\theta}_1^*)^\top u_i, \dots, \nabla_{\boldsymbol{\theta}} \ell_{i,K}(\boldsymbol{\theta}_K^*)^\top u_i)^\top$ . Let us denote

$$\mathbf{g}_i \stackrel{\text{def}}{=} \left( \nabla_{\boldsymbol{\theta}} \ell_{i,1}(\boldsymbol{\theta}_1^*)^\top, \dots, \nabla_{\boldsymbol{\theta}} \ell_{i,K}(\boldsymbol{\theta}_K^*)^\top \right)^\top \in \mathbb{R}^{p_{\text{sum}}},$$

$$\widehat{H}^2 \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} \left\{ \mathbf{g}_i \mathbf{g}_i^\top \right\},$$

$$\widehat{v}_i \stackrel{\text{def}}{=} \widehat{H}^{-1} \left\{ \mathbf{g}_i \mathbf{g}_i^\top - \mathbb{E} \left[ \mathbf{g}_i \mathbf{g}_i^\top \right] \right\} \widehat{H}^{-1},$$

then

$$\widehat{H}^2 = \mathbb{E} \widehat{\mathcal{V}}^2, \quad \sum_{i=1}^n \widehat{v}_i^2 = \widehat{H}^{-1} \widehat{\mathcal{V}}^2 \widehat{H}^{-1} - \mathbf{I}_{p_{\text{sum}}}.$$

Define also the deterministic scalar value

$$\widehat{\varkappa}_v^2 \stackrel{\text{def}}{=} \left\| \sum_{i=1}^n \mathbb{E} \widehat{v}_i^4 \right\|.$$



**Theorem C.1** (Bernstein inequality for  $\widehat{\mathcal{V}}^2$ ). *Let the condition  $(\widehat{SD}_1)$  be fulfilled, then it holds with probability  $\geq 1 - e^{-\mathbf{x}}$ :*

$$\|\widehat{H}^{-1}\widehat{\mathcal{V}}^2\widehat{H}^{-1} - \mathbf{I}_{p_{\text{sum}}}\| \leq \delta_{\widehat{\mathcal{V}}}^2(\mathbf{x}),$$

where the error term is defined as

$$\delta_{\widehat{\mathcal{V}}}^2(\mathbf{x}) \stackrel{\text{def}}{=} \sqrt{2\widehat{\mathcal{Z}}_v^2 \{\log(p_{\text{sum}}) + \mathbf{x}\}} + \frac{2}{3}\delta_{v^*}^2 \{\log(p_{\text{sum}}) + \mathbf{x}\} \quad (\text{C.1})$$

and is proportional to  $\sqrt{\{\log(p_{\text{sum}}) + \mathbf{x}\}/n}$  in the case 5.3.

We omit here the proof of Theorem C.1, since it follows straightforwardly from Theorem 1.4 by Tropp (2012), and is already given in Spokoiny and Zhilova (2014).

## C.2 Bootstrap validity for the case of one parametric model

Here we state the results on bootstrap validity from Spokoiny and Zhilova (2014), they will be used for some of the further proofs.

**Theorem C.2.** *Let the conditions of Section 5 be fulfilled, then it holds for each  $k = 1, \dots, K$ ,  $z_k \geq \max\{2, \sqrt{p_k}\} + \mathbf{C}(p_k + \mathbf{x})/\sqrt{n}$  with probability  $\geq 1 - 12e^{-\mathbf{x}}$ :*

$$\left| \mathbb{P} \left( L_k(\widetilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) > z_k^2/2 \right) - \mathbb{P}^\circ \left( L_k^\circ(\widetilde{\boldsymbol{\theta}}_k^\circ) - L_k^\circ(\widetilde{\boldsymbol{\theta}}_k) > z_k^2/2 \right) \right| \leq \Delta_{\text{full}, k}.$$

The error term  $\Delta_{\text{full}, k} \leq \mathbf{C}\{(p_k + \mathbf{x})^3/n\}^{1/8}$  in the case of i.i.d. model; see Section 5.3.

**Theorem C.3** (Validity of the bootstrap under a small modeling bias). *Assume the conditions of Theorem C.2. Then for  $\alpha \leq 1 - 8e^{-\mathbf{x}}$ , it holds*

$$\left| \mathbb{P} \left( L_k(\widetilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) > (\mathfrak{z}_k^\circ(\alpha))^2/2 \right) - \alpha \right| \leq \Delta_{\mathfrak{z}, \text{full}, k}.$$

The error term  $\Delta_{\mathfrak{z}, \text{full}, k} \leq \mathbf{C}\{(p_k + \mathbf{x})^3/n\}^{1/8}$  in the case of i.i.d. model; see Section 5.3.

**Theorem C.4** (Performance of the bootstrap for a large modeling bias). *Under the conditions of Section 5 except for  $(\widehat{\text{SmB}})$  it holds for  $z_k \geq \max\{2, \sqrt{p_k}\} + \mathbf{C}(p_k + \mathbf{x})/\sqrt{n}$  with probability  $\geq 1 - 14e^{-\mathbf{x}}$*

$$1. \quad \mathbb{P} \left( L_k(\widetilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) > z_k^2/2 \right) \leq \mathbb{P}^\circ \left( L_k^\circ(\widetilde{\boldsymbol{\theta}}_k^\circ) - L_k^\circ(\widetilde{\boldsymbol{\theta}}_k) > z_k^2/2 \right) + \Delta_{\text{b}, \text{full}, k}.$$

$$2. \quad \mathfrak{z}_k^\circ(\alpha) \geq \mathfrak{z}_k(\alpha + \Delta_{\text{b}, \text{full}, k}) \\ + \sqrt{\text{tr}\{D_k^{-1}H_k^2D_k^{-1}\}} - \sqrt{\text{tr}\{D_k^{-1}(H_k^2 - B_k^2)D_k^{-1}\}} - \Delta_{\text{qf}, 1, k},$$

$$\mathfrak{z}_k^\circ(\alpha) \leq \mathfrak{z}_k(\alpha - \Delta_{\text{b}, \text{full}, k}) \\ + \sqrt{\text{tr}\{D_k^{-1}H_k^2D_k^{-1}\}} - \sqrt{\text{tr}\{D_k^{-1}(H_k^2 - B_k^2)D_k^{-1}\}} + \Delta_{\text{qf}, 2, k}.$$

The term  $\Delta_{\text{b,full},k} \leq \mathbf{C}\{(p_k + \mathbf{x})^3/n\}^{1/8}$  in the case of i.i.d. model; see Section 5.3. The positive values  $\Delta_{\text{qf},1,k}, \Delta_{\text{qf},2,k}$  are bounded from above with  $(\mathbf{a}_k^2 + \mathbf{a}_{B,k}^2)(\sqrt{8\mathbf{x}p_k} + 6\mathbf{x})$  for the constants  $\mathbf{a}_k^2 > 0, \mathbf{a}_{B,k}^2 \geq 0$  from conditions  $(\mathcal{I}), (\mathcal{I}_B)$ .

### C.3 Proof of Theorem 3.1

**Lemma C.5** (Closeness of  $\mathcal{L}(\|\boldsymbol{\xi}_1\|, \dots, \|\boldsymbol{\xi}_K\|)$  and  $\mathcal{L}^\circ(\|\boldsymbol{\xi}_1^\circ\|, \dots, \|\boldsymbol{\xi}_K^\circ\|)$ ). *If the conditions  $(\mathbf{ED}_0), (\mathcal{I}), (\widehat{\mathbf{SmB}}), (\mathcal{I}_B), (\widehat{\mathbf{SD}}_1)$  and  $(\mathbf{Eb})$  are fulfilled, then it holds with probability  $\geq 1 - 6e^{-\mathbf{x}}$  for all  $\delta_{z_k} \geq 0$  and  $z_k \geq \sqrt{p_k} + \Delta_\varepsilon$  s.t.  $\mathbf{C} \max_{1 \leq k \leq K} \{n^{-1/2}, \delta_{z_k}\} \leq \Delta_\varepsilon \leq \mathbf{C} \min_{1 \leq k \leq K} \{1/z_k\}$  ( $\Delta_\varepsilon$  is given in (A.3)):*

$$\begin{aligned} \mathbb{P} \left( \bigcup_{k=1}^K \{\|\boldsymbol{\xi}_k\| > z_k\} \right) - \mathbb{P}^\circ \left( \bigcup_{k=1}^K \{\|\boldsymbol{\xi}_k^\circ\| > z_k - \delta_{z_k}\} \right) &\geq -\Delta_{\ell_2}, \\ \mathbb{P} \left( \bigcup_{k=1}^K \{\|\boldsymbol{\xi}_k\| > z_k\} \right) - \mathbb{P}^\circ \left( \bigcup_{k=1}^K \{\|\boldsymbol{\xi}_k^\circ\| > z_k + \delta_{z_k}\} \right) &\leq \Delta_{\ell_2}. \end{aligned}$$

for the deterministic nonnegative value

$$\Delta_{\ell_2} \leq 25\mathbf{C} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \{(\widehat{\mathbf{a}}^2 + \widehat{\mathbf{a}}_B^2) (1 + \delta_{\widehat{\mathbf{v}}}^2(\mathbf{x}))\}^{3/8}.$$

A more explicit bound on  $\Delta_{\ell_2}$  is given in Proposition A.1, see also Remark A.1.

*Proof of Lemma C.5.* The statement follows from Proposition A.1 and Theorem C.1. Let us take  $\phi_k := \boldsymbol{\xi}_k$  and  $\psi_k := \boldsymbol{\xi}_k^\circ$ . Define similarly to  $\Phi$  in (A.5)

$$\Xi \stackrel{\text{def}}{=} \left( \boldsymbol{\xi}_1^\top, \dots, \boldsymbol{\xi}_K^\top \right)^\top \quad \Xi^\circ \stackrel{\text{def}}{=} \left( \boldsymbol{\xi}_1^{\circ\top}, \dots, \boldsymbol{\xi}_K^{\circ\top} \right)^\top. \quad (\text{C.2})$$

Condition (A.4) rewrites for (C.2) as

$$\|\text{Var} \Xi - \text{Var}^\circ \Xi^\circ\|_{\max} \leq \delta_\Sigma^2$$

for some  $\delta_\Sigma^2 \geq 0$ . Denote

$$\begin{aligned} \widehat{D}^2 &\stackrel{\text{def}}{=} \text{diag} \{D_1^2, \dots, D_K^2\}, \\ \widehat{V}^2 &\stackrel{\text{def}}{=} \text{Var} \left( \nabla_{\boldsymbol{\theta}} L_1(\boldsymbol{\theta}_1^*)^\top, \dots, \nabla_{\boldsymbol{\theta}} L_K(\boldsymbol{\theta}_K^*)^\top \right)^\top. \end{aligned}$$

$\widehat{D}^2$  is a block-diagonal matrix and  $\widehat{V}^2$  is a block matrix. Both of them are symmetric, positive definite and have the dimension  $p_{\text{sum}} \times p_{\text{sum}}$ . Let also

$$\begin{aligned} \widehat{\mathcal{V}}^2 &\stackrel{\text{def}}{=} \text{Var}^\circ \left( \nabla_{\boldsymbol{\theta}} L_1^\circ(\boldsymbol{\theta}_1^*)^\top, \dots, \nabla_{\boldsymbol{\theta}} L_K^\circ(\boldsymbol{\theta}_K^*)^\top \right)^\top, \\ \mathbf{g}_i &\stackrel{\text{def}}{=} \left( \nabla_{\boldsymbol{\theta}} \ell_{i,1}(\boldsymbol{\theta}_1^*)^\top, \dots, \nabla_{\boldsymbol{\theta}} \ell_{i,K}(\boldsymbol{\theta}_K^*)^\top \right)^\top \in \mathbb{R}^{p_{\text{sum}}}, \\ \widehat{H}^2 &\stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} \{ \mathbf{g}_i \mathbf{g}_i^\top \}, \quad \widehat{B}^2 \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbb{E} \{ \mathbf{g}_i \} \mathbb{E} \{ \mathbf{g}_i \}^\top. \end{aligned}$$

It holds

$$\begin{aligned}\text{Var } \Xi &= \widehat{D}^{-1} \widehat{V}^2 \widehat{D}^{-1}, \quad \text{Var}^\circ \Xi^\circ = \widehat{D}^{-1} \widehat{V}^2 \widehat{D}^{-1}, \\ \widehat{H}^2 &= \mathbb{E} \widehat{V}^2, \quad \widehat{V}^2 = \widehat{H}^2 - \widehat{B}^2.\end{aligned}$$

Therefore

$$\begin{aligned}\|\text{Var } \Xi - \text{Var}^\circ \Xi^\circ\|_{\max} &= \|\widehat{D}^{-1}(\widehat{V}^2 - \widehat{V}^2) \widehat{D}^{-1}\|_{\max} \\ &\leq \|\widehat{D}^{-1}(\widehat{H}^2 - \widehat{V}^2) \widehat{D}^{-1}\|_{\max} + \|\widehat{D}^{-1} \widehat{B}^2 \widehat{D}^{-1}\|_{\max} \\ &\leq \delta_{\widehat{V}}^2(\mathbf{x}) \|\widehat{D}^{-1} \widehat{H}^2 \widehat{D}^{-1}\| + \|\widehat{D}^{-1} \widehat{B}^2 \widehat{D}^{-1}\| \quad (\text{C.3})\end{aligned}$$

$$\leq \{\delta_{\widehat{V}}^2(\mathbf{x}) + \widehat{\delta}_{\text{smb}}^2\}(\widehat{\mathbf{a}}^2 + \widehat{\mathbf{a}}_B^2) =: \delta_{\Sigma}^2. \quad (\text{C.4})$$

Here inequality (C.3) follows from the matrix Bernstein inequality by [Tropp \(2012\)](#) (see Section C.1). Inequality (C.4) is implied by conditions  $(\mathcal{I}_B)$  and  $(\widehat{\text{SmB}})$ , and Cauchy-Schwarz inequality.

Condition (C1) of Proposition A.1 is fulfilled for the vectors  $\boldsymbol{\xi}_{i,k}$  and  $\boldsymbol{\xi}_{i,k}^\circ$  due to conditions  $(\text{ED}_0)$ ,  $(\mathcal{I})$  and  $(\widehat{\text{SD}}_1)$ ,  $(\text{Eb})$ ,  $(\widehat{\text{SmB}})$ ,  $(\mathcal{I}_B)$  for  $\mathbf{c}_\phi := \widehat{\mathbf{a}}$  and  $\mathbf{c}_\psi^2 := (\widehat{\mathbf{a}}^2 + \widehat{\mathbf{a}}_B^2) \left\{ \delta_{\widehat{V}}^2 + \max_{1 \leq i \leq n} \|\widehat{H}^{-1} \mathbb{E} [\mathbf{g}_i \mathbf{g}_i^\top] \widehat{H}^{-1}\|^2 \right\}$ .  $\square$

*Proof of Theorem 3.1.* Let us denote  $\mathbf{x}_2 \stackrel{\text{def}}{=} \mathbf{x} + \log(K)$ . It holds with probability  $\geq 1 - 12e^{-x}$

$$\begin{aligned}& \mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^\circ(\tilde{\boldsymbol{\theta}}_k)} - 2L_k^\circ(\tilde{\boldsymbol{\theta}}_k) > z_k \right\} \right) \\ & \stackrel{\text{L.B.9}}{\geq} \mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \|\boldsymbol{\xi}_k^\circ(\tilde{\boldsymbol{\theta}}_k)\| \geq z_k + \Delta_{\text{W},k}^\circ(\mathbf{r}_{0,k}, \mathbf{x}_2) \right\} \right) \\ & \stackrel{\text{L.B.10}}{\geq} \mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \|\boldsymbol{\xi}_k^\circ(\boldsymbol{\theta}_k^*)\| > z_k + \Delta_{\text{W},k}^\circ(\mathbf{r}_{0,k}, \mathbf{x}_2) + \Delta_{\boldsymbol{\xi},k}^\circ(\mathbf{r}_{0,k}, \mathbf{x}_2) \right\} \right) \\ & \stackrel{\text{L.C.5}}{\geq} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \|\boldsymbol{\xi}_k\| > z_k - \Delta_{\text{W},k}(\mathbf{r}_{0,k}, \mathbf{x}_2) \right\} \right) - \Delta_{\text{total}} \\ & \stackrel{\text{L.B.9}}{\geq} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k)} - 2L_k(\boldsymbol{\theta}_k^*) > z_k \right\} \right) - \Delta_{\text{total}},\end{aligned}$$

for

$$\Delta_{\text{total}} \stackrel{\text{def}}{=} \Delta_{\ell_2}, \quad (\text{C.5})$$

$$\begin{aligned}\delta_{z_k} &:= \Delta_{\text{W},k}(\mathbf{r}_{0,k}, \mathbf{x} + \log(K)) + \Delta_{\text{W},k}^\circ(\mathbf{r}_{0,k}, \mathbf{x} + \log(K)) \\ &\quad + \Delta_{\boldsymbol{\xi},k}^\circ(\mathbf{r}_{0,k}, \mathbf{x} + \log(K))\end{aligned} \quad (\text{C.6})$$

$$\leq \mathbf{c} \frac{pk + \mathbf{x} + \log(K)}{\sqrt{n}} \sqrt{\mathbf{x} + \log(K)} \quad \text{in the case 5.3.} \quad (\text{C.7})$$

Definition of  $\Delta_{\ell_2}$  is given in Proposition A.1, see also Remark A.1. The bound from Lemma C.5 says:

$$\Delta_{\ell_2} \leq 25\mathbf{C} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \{ (\widehat{\mathbf{a}}^2 + \widehat{\mathbf{a}}_B^2) (1 + \delta_{\widehat{\mathbf{v}}}^2(\mathbf{x})) \}^{3/8}.$$

For  $\delta_{z_k}$  bounded as in (C.7) the conditions  $\mathbf{C} \max_{1 \leq k \leq K} \{n^{-1/2}, \delta_{z_k}\} \leq \Delta_\varepsilon \leq \mathbf{C} \min_{1 \leq k \leq K} \{1/z_k\}$  are fulfilled.  $\square$

### C.4 Proof of Theorem 3.2

*Proof of Theorem 3.2.* For the pointwise quantile functions  $\mathfrak{z}_k(\alpha)$  and  $\mathfrak{z}_k^\circ(\alpha)$  it holds for each  $k = 1, \dots, K$  with dominating probability:

$$\begin{aligned} \mathfrak{z}_k^\circ(\alpha + \Delta_{\text{full}, k}) &\leq \mathfrak{z}_k(\alpha), \\ \mathfrak{z}_k^\circ(\alpha) &\geq \mathfrak{z}_k(\alpha + \Delta_{\text{full}, k}) - \varepsilon_k \end{aligned} \tag{C.8}$$

here  $\Delta_{\text{full}, k} \leq \{(p_k + \mathbf{x})^3 / \sqrt{n}\}^{1/8}$ , it comes from Theorem C.2, and  $\varepsilon_k \leq \mathbf{C}(p_k + \mathbf{x}) / \sqrt{n}$ ,

$$\begin{aligned} \varepsilon_k &\stackrel{\text{def}}{=} \begin{cases} 0, & \text{if c.d.f. of } L_k(\widetilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*) \text{ is continuous in } \mathfrak{z}_k(\alpha + \Delta_{\text{full}, k}); \\ \mathbf{C}(p_k + \mathbf{x}) / \sqrt{n} & \text{s.t. (C.9) is fulfilled, otherwise.} \end{cases} \\ \mathbb{P} \left( \sqrt{2\{L_k(\widetilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*)\}} > \mathfrak{z}_k(\alpha + \Delta_{\text{full}, k}) - \varepsilon_k \right) &\geq \alpha + \Delta_{\text{full}, k}. \end{aligned} \tag{C.9}$$

Indeed, due to Theorem C.2 and definition (1.5)

$$\begin{aligned} \mathbb{P}^\circ \left( \sqrt{2\{L_k^\circ(\widetilde{\boldsymbol{\theta}}_k^\circ) - L_k^\circ(\widetilde{\boldsymbol{\theta}}_k)\}} > \mathfrak{z}_k(\alpha) \right) \\ \leq \mathbb{P} \left( \sqrt{2\{L_k(\widetilde{\boldsymbol{\theta}}_k) - L_k(\boldsymbol{\theta}_k^*)\}} > \mathfrak{z}_k(\alpha) \right) + \Delta_{\text{full}, k} \leq \alpha + \Delta_{\text{full}, k}, \end{aligned}$$

therefore, by definition (2.3)  $\mathfrak{z}_k^\circ(\alpha + \Delta_{\text{full}, k}) \leq \mathfrak{z}_k(\alpha)$ . The lower bound is derived similarly.

If there exist the inverse functions  $\mathbf{c}^{-1}(\cdot)$  and  $\mathbf{c}^{\circ-1}(\cdot)$ , then it holds for  $\beta \in (0, 1)$ :

$$\begin{aligned} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\widetilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq \mathfrak{z}_k(\beta) \right\} \right) &\leq \mathbf{c}^{-1}(\beta), \\ \mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^\circ(\widetilde{\boldsymbol{\theta}}_k^\circ) - 2L_k^\circ(\widetilde{\boldsymbol{\theta}}_k)} \geq \mathfrak{z}_k^\circ(\beta) \right\} \right) &\leq \mathbf{c}^{\circ-1}(\beta). \end{aligned} \tag{C.10}$$

Therefore, it holds

$$\begin{aligned}
& \mathfrak{c}^{\circ-1}(\beta + \Delta_{\text{full, max}}) \\
& \geq \mathbb{P}^{\circ} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k)} - 2L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k) \geq \mathfrak{z}_k^{\circ}(\beta + \Delta_{\text{full, } k}) \right\} \right) \\
& \stackrel{\text{by (C.8)}}{\geq} \mathbb{P}^{\circ} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k)} - 2L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k) \geq \mathfrak{z}_k(\beta) \right\} \right) \\
& \stackrel{\text{by Th. 3.1}}{\geq} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k)} - 2L_k(\boldsymbol{\theta}_k^*) \geq \mathfrak{z}_k(\beta) \right\} \right) - \Delta_{\text{total}} \\
& \stackrel{\text{by L. C.6 and (C.10)}}{\geq} \mathfrak{c}^{-1}(\beta) - \Delta_{\text{total}} - \Delta_{\text{ac, LR}},
\end{aligned}$$

here  $\Delta_{\text{ac, LR}} \leq \Delta_{\text{total}}$  (by Lemma C.6) and

$$\begin{aligned}
\Delta_{\text{full, max}} & \stackrel{\text{def}}{=} \max_{1 \leq k \leq K} \Delta_{\text{full, } k} \\
& \leq \mathfrak{C}\{(p_{\text{max}} + \mathbf{x})^3/n\}^{1/8} \text{ in the case 5.3.}
\end{aligned} \tag{C.11}$$

Thus

$$\begin{aligned}
\mathfrak{c}^{\circ-1}(\beta + \Delta_{\text{full, max}}) & \geq \mathfrak{c}^{-1}(\beta) - \Delta_{\text{total}} - \Delta_{\text{ac, LR}}, \\
\mathfrak{c}^{\circ}(\alpha) & \leq \mathfrak{c}(\alpha + \Delta_{\text{total}} + \Delta_{\text{ac, LR}}) + \Delta_{\text{full, max}}.
\end{aligned} \tag{C.12}$$

Hence it holds

$$\begin{aligned}
& \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k)} - 2L_k(\boldsymbol{\theta}_k^*) \geq \mathfrak{z}_k^{\circ}(\beta) \right\} \right) \\
& \stackrel{\text{by (C.8)}}{\leq} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k)} - 2L_k(\boldsymbol{\theta}_k^*) \geq \mathfrak{z}_k(\beta + \Delta_{\text{full, } k}) - \varepsilon_k \right\} \right) \\
& \stackrel{\text{by L. C.6 and (C.10)}}{\leq} \mathfrak{c}^{-1}(\beta + \Delta_{\text{full, max}}) + \Delta_{\text{ac, LR}}.
\end{aligned}$$

Therefore, if  $\mathfrak{c}(\alpha) \geq \Delta_{\text{full, max}}$ , then

$$\mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k)} - 2L_k(\boldsymbol{\theta}_k^*) \geq \mathfrak{z}_k^{\circ}(\mathfrak{c}(\alpha) - \Delta_{\text{full, max}}) \right\} \right) \leq \alpha + \Delta_{\text{ac, LR}}.$$

And by (C.12) for  $\mathfrak{c}^{\circ}(\alpha) \geq 2\Delta_{\text{full, max}}$  it holds

$$\begin{aligned}
& \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k)} - 2L_k(\boldsymbol{\theta}_k^*) \geq \mathfrak{z}_k^{\circ}(\mathfrak{c}^{\circ}(\alpha) - 2\Delta_{\text{full, max}}) \right\} \right) - \alpha \\
& \leq \Delta_{\text{total}} + 2\Delta_{\text{ac, LR}}.
\end{aligned}$$

Similarly for the inverse direction:

$$\begin{aligned}
\mathbf{c}^{\circ-1}(\beta) &\leq \mathbb{P}^{\circ} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k^{\circ}) - 2L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k)} \geq \mathfrak{z}_k^{\circ}(\beta) \right\} - \varepsilon_{1,k} \right) \\
&\leq \mathbb{P}^{\circ} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k^{\circ}) - 2L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k)} \geq \mathfrak{z}_k(\beta + \Delta_{\text{full},k}) - \varepsilon_{1,k} - \varepsilon_k \right\} \right) \\
&\leq \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq \mathfrak{z}_k(\beta + \Delta_{\text{full},k}) \right\} \right) + \Delta_{\text{total}} + \Delta_{\text{ac,LR}} \\
&\leq \mathbf{c}^{-1}(\beta + \Delta_{\text{full,max}}) + \Delta_{\text{total}} + \Delta_{\text{ac,LR}},
\end{aligned}$$

where  $0 \leq \varepsilon_{1,k} \leq \mathbf{C}(p_k + \mathbf{x})/\sqrt{n}$ . This implies

$$\begin{aligned}
\mathbf{c}^{\circ-1}(\beta) &\leq \mathbf{c}^{-1}(\beta + \Delta_{\text{full,max}}) + \Delta_{\text{total}} + \Delta_{\text{ac,LR}}, \\
\mathbf{c}^{\circ}(\alpha) &\geq \mathbf{c}(\alpha - \Delta_{\text{total}} - \Delta_{\text{ac,LR}}) - \Delta_{\text{full,max}}.
\end{aligned} \tag{C.13}$$

$$\begin{aligned}
&\mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq \mathfrak{z}_k^{\circ}(\beta + \Delta_{\text{full},k}) \right\} \right) \\
&\stackrel{\text{by (C.8)}}{\geq} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq \mathfrak{z}_k(\beta) \right\} \right) \\
&\geq \mathbf{c}^{-1}(\beta) - \Delta_{\text{ac,LR}}.
\end{aligned}$$

$$\mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq \mathfrak{z}_k^{\circ}(\mathbf{c}(\alpha) + \Delta_{\text{full,max}}) \right\} \right) \geq \alpha - \Delta_{\text{ac,LR}}.$$

And by (C.13)

$$\begin{aligned}
&\mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq \mathfrak{z}_k^{\circ}(\mathbf{c}^{\circ}(\alpha) + 2\Delta_{\text{full,max}}) \right\} \right) - \alpha \\
&\geq -\Delta_{\text{total}} - 2\Delta_{\text{ac,LR}}.
\end{aligned}$$

for

$$\Delta_{\mathfrak{z}, \text{total}} \stackrel{\text{def}}{=} \Delta_{\text{total}} + 2\Delta_{\text{ac,LR}} \leq 3\Delta_{\text{total}}. \tag{C.14}$$

Conditions of Theorem 3.1 include  $z_k \geq C\sqrt{p_k}$ , therefore, it has to be checked that  $\mathfrak{z}_k^{\circ}(\alpha) \geq C\sqrt{p_k}$ . It holds by Theorem B.4, Proposition A.1, Lemmas B.6 and C.7 with probability  $\geq 1 - 12e^{-x}$ :

$$\begin{aligned}
&\mathbb{P}^{\circ} \left( \sqrt{2\{L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k^{\circ}) - L_k^{\circ}(\tilde{\boldsymbol{\theta}}_k)\}} > \mathbf{C}\sqrt{p_k - \sqrt{2xp_k}} + \mathbf{C}(p_k + \mathbf{x})/\sqrt{n} \right) \\
&\geq 1 - 8e^{-x},
\end{aligned}$$

Taking  $1 - 8e^{-x} \geq \alpha$ , we have

$$\mathfrak{z}_k^\circ(\alpha) \geq \mathbf{C} \sqrt{p_k - \sqrt{2xp_k}} + \mathbf{C}2(p_k + \mathbf{x})/\sqrt{n}.$$

Inequalities for  $\mathfrak{c}^\circ(\alpha)$  had been already derived in (C.12) and (C.13) with

$$\Delta_{\mathfrak{c}} \stackrel{\text{def}}{=} \Delta_{\text{total}} + \Delta_{\text{ac,LR}}. \quad (\text{C.15})$$

□

**Lemma C.6.** *Let the conditions from Section 5.1 be fulfilled, and the values  $z_k \geq \sqrt{p_k}$  and  $\delta_{z_k} \geq 0$  be s.t.  $\mathbf{C} \max_{1 \leq k \leq K} \{n^{-1/2}, \delta_{z_k}\} \leq \Delta_\varepsilon \leq \mathbf{C} \min_{1 \leq k \leq K} \{1/z_k\}$  ( $\Delta_\varepsilon$  is given in (A.3)), then it holds with probability  $\geq 1 - 12e^{-x}$*

$$\begin{aligned} & \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq z_k \right\} \right) \\ & - \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} \geq z_k + \delta_{z_k} \right\} \right) \leq \Delta_{\text{ac,LR}}, \end{aligned}$$

where

$$\Delta_{\text{ac,LR}} \leq 12.5\mathbf{C} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \hat{\mathbf{a}}^{3/4}.$$

*Proof of Lemma C.6.* This statement's proof is similar to the one of Theorem 3.1 (see Section C.3). Here instead of the bootstrap statistics we consider only the values from the  $\mathbf{Y}$ -world. Let us denote  $\mathbf{x}_2 \stackrel{\text{def}}{=} \mathbf{x} + \log(K)$ . It holds with probability  $\geq 1 - 12e^{-x}$

$$\begin{aligned} & \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} > z_k \right\} \right) \\ & \stackrel{\text{L. B.9}}{\leq} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \|\boldsymbol{\xi}_k\| > z_k - \Delta_{\text{w},k}(\mathbf{r}_{0,k}, \mathbf{x}_2) \right\} \right) \\ & \stackrel{\text{Pr. A.1}}{\leq} \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \|\boldsymbol{\xi}_k\| > z_k + \delta_{z_k} + \Delta_{\text{w},k}(\mathbf{r}_{0,k}, \mathbf{x}_2) \right\} \right) + \Delta_{\text{ac,LR}} \\ & \leq \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\boldsymbol{\theta}}_k) - 2L_k(\boldsymbol{\theta}_k^*)} > z_k + \delta_{z_k} \right\} \right) + \Delta_{\text{ac,LR}}, \end{aligned}$$

where

$$\Delta_{\text{ac,LR}} \leq 12.5\mathbf{C} (p_{\max}^3/n)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \hat{\mathbf{a}}^{3/4}.$$

Similarly to (C.5) and (C.6) the term  $\Delta_{\text{ac,LR}}$  is equal to  $\Delta_{\ell_2}$  from Proposition A.1 with  $\Delta_{\Sigma}^2 := 0$ ,  $\delta_{z_k} := \delta_{z_k} + 2\Delta_{\text{w},k}(\mathbf{r}_{0,k}, \mathbf{x} + \log(K))$ . □

**Lemma C.7** (Lower bound for deviations of a Gaussian quadratic form). *Let  $\phi \sim \mathcal{N}(0, \mathbf{I}_p)$  and  $\Sigma$  is any symmetric non-negative definite matrix, then it holds for any  $\mathbf{x} > 0$*

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

*Proof of Lemma C.7.* It is sufficient to consider w.l.o.g. only the case of diagonal matrix  $\Sigma$ , 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  $a_1 \geq \dots \geq a_p$ ;  $U\phi \sim \mathcal{N}(0, \mathbf{I}_p)$ .

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

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

therefore

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

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

### C.5 Proof of Theorem 3.3

*Proof of Theorem 3.3.* Let us denote  $\mathbf{x}_2 \stackrel{\text{def}}{=} \mathbf{x} + \log(K)$ . By Lemmas B.9, B.10 and C.5 it holds with probability  $\geq 1 - 12e^{-\mathbf{x}}$

$$\begin{aligned} &\mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k^\circ(\tilde{\theta}_k)} - 2L_k^\circ(\tilde{\theta}_k) > z_k \right\} \right) \\ &\geq \mathbb{P}^\circ \left( \bigcup_{k=1}^K \left\{ \|\xi_k^\circ(\theta_k^*)\| > z_k + \Delta_{\mathbf{w},k}^\circ(\mathbf{r}_{0,k}, \mathbf{x}_2) + \Delta_{\xi,k}^\circ(\mathbf{r}_{0,k}, \mathbf{x}_2) \right\} \right) \\ &\geq \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \|\tilde{\xi}_k\| > z_k - \Delta_{\mathbf{w},k}(\mathbf{r}_{0,k}, \mathbf{x}_2) \right\} \right) - \Delta_{\text{b, total}} \end{aligned} \tag{C.16}$$

$$\geq \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \|\xi_k\| > z_k - \Delta_{\mathbf{w},k}(\mathbf{r}_{0,k}, \mathbf{x}_2) \right\} \right) - \Delta_{\text{b, total}} \tag{C.17}$$

$$\geq \mathbb{P} \left( \bigcup_{k=1}^K \left\{ \sqrt{2L_k(\tilde{\theta}_k)} - 2L_k(\theta_k^*) > z_k \right\} \right) - \Delta_{\text{b, total}},$$



here  $\tilde{\boldsymbol{\xi}}_k \stackrel{\text{def}}{=} (D_k^{-1} H_k^2 D_k^{-1})^{1/2} (\text{Var } \boldsymbol{\xi}_k)^{-1/2} \boldsymbol{\xi}_k$ , and  $\Delta_{\text{b,total}}$  is given below. Using the same notations as in the proof of Lemma C.5, we have

$$\begin{aligned} \tilde{\boldsymbol{\Xi}} &\stackrel{\text{def}}{=} \left( \tilde{\boldsymbol{\xi}}_1^\top, \dots, \tilde{\boldsymbol{\xi}}_K^\top \right)^\top \\ &= (\widehat{D}^{-1} \widehat{H}^2 \widehat{D}^{-1})^{1/2} (\text{Var } \boldsymbol{\Xi})^{-1/2} \boldsymbol{\Xi}, \end{aligned}$$

and by Theorem C.1 and by conditions  $(\mathcal{I})$ ,  $(\mathcal{I}_B)$ , it holds with probability  $\geq 1 - e^{-x}$

$$\begin{aligned} \|\text{Var } \tilde{\boldsymbol{\Xi}} - \text{Var}^\circ \boldsymbol{\Xi}^\circ\|_{\max} &= \|\widehat{D}^{-1} (\widehat{H}^2 - \widehat{V}^2) \widehat{D}^{-1}\|_{\max} \\ &\leq \delta_{\widehat{V}}^2(\mathbf{x}) \|\widehat{D}^{-1} \widehat{H}^2 \widehat{D}^{-1}\|_{\max} \\ &\leq \delta_{\widehat{V}}^2(\mathbf{x}) (\widehat{\mathbf{a}}^2 + \widehat{\mathbf{a}}_B^2). \end{aligned}$$

Thus, inequality (C.16) follows from Proposition A.1 applied to the sets of vectors  $\boldsymbol{\xi}_1^\circ(\boldsymbol{\theta}_1^*), \dots, \boldsymbol{\xi}_K^\circ(\boldsymbol{\theta}_K^*)$  and  $\tilde{\boldsymbol{\xi}}_1, \dots, \tilde{\boldsymbol{\xi}}_K$ . The error term  $\Delta_{\text{b,total}}$  is equal to  $\Delta_{\text{total}}$  from Theorem C.3 (see (C.5), (C.6)) with  $\widehat{\delta}_{\text{smb}}^2 := 0$ , thus

$$\Delta_{\text{b,total}} \leq 25\text{C} \left( \frac{p_{\max}^3}{n} \right)^{1/8} \log^{9/8}(K) \log^{3/8}(np_{\text{sum}}) \{ (\widehat{\mathbf{a}}^2 + \widehat{\mathbf{a}}_B^2) (1 + \delta_{\widehat{V}}^2(\mathbf{x})) \}^{3/8}.$$

Inequality (C.17) is implied by definitions of  $\tilde{\boldsymbol{\xi}}_k$  and matrices  $H_k^2, V_k^2$ , indeed:

$$\begin{aligned} &\left\| (D_k^{-1} H_k^2 D_k^{-1})^{-1/2} \text{Var } \boldsymbol{\xi}_k (D_k^{-1} H_k^2 D_k^{-1})^{-1/2} \right\| \\ &\leq \left\| (D_k^{-1} H_k^2 D_k^{-1})^{1/2} (D_k H_k^{-2} V_k^2 H_k^{-2} D_k) (D_k^{-1} H_k^2 D_k^{-1})^{1/2} \right\| \\ &\leq 1, \end{aligned}$$

therefore,  $\|\tilde{\boldsymbol{\xi}}_k\| \geq \|\boldsymbol{\xi}_k\|$ .

The second inequality in the statement is proven similarly to (C.12). It implies together with Theorem C.4 the rest part of the statement having

$$\Delta_{\text{b,c}} \stackrel{\text{def}}{=} \Delta_{\text{b,total}} + \Delta_{\text{ac,LR}}. \quad (\text{C.18})$$

□

## References

Arlot, S., Blanchard, G., and Roquain, E. (2010). Some nonasymptotic results on resampling in high dimension, II: Multiple tests. *The Annals of Statistics*, 38(1):83–99.

- Benjamini, Y. (2010). Simultaneous and selective inference: current successes and future challenges. *Biometrical Journal*, 52(6):708–721.
- Beran, R. (1988). Balanced simultaneous confidence sets. *Journal of the American Statistical Association*, 83(403):679–686.
- Beran, R. (1990). Refining bootstrap simultaneous confidence sets. *Journal of the American Statistical Association*, 85(410):417–426.
- Bickel, P. J. and Rosenblatt, M. (1973). On some global measures of the deviations of density function estimates. *The Annals of Statistics*, pages 1071–1095.
- Bonferroni, C. E. (1936). *Teoria statistica delle classi e calcolo delle probabilita*. Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze.
- Cao, H. and Kosorok, M. R. (2011). Simultaneous critical values for t-tests in very high dimensions. *Bernoulli*, 17(1):347–394.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2013a). Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *The Annals of Statistics*, 41(6):2786–2819.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2013b). Supplement to “Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors”.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2014a). Anti-concentration and honest, adaptive confidence bands. *The Annals of Statistics*, 42(5):1787–1818.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2014b). Central limit theorems and bootstrap in high dimensions. *arXiv preprint arXiv:1412.3661*.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2014c). Comparison and anti-concentration bounds for maxima of Gaussian random vectors. *Probability Theory and Related Fields*, 162:47–70.
- Claeskens, G. and Van Keilegom, I. (2003). Bootstrap confidence bands for regression curves and their derivatives. *Annals of Statistics*, pages 1852–1884.
- Dickhaus, T. (2014). *Simultaneous Statistical Inference: With Applications in the Life Sciences*. Springer.
- Hall, P. (1991). On convergence rates of suprema. *Probability Theory and Related Fields*, 89(4):447–455.

- Hall, P. (1993). On Edgeworth expansion and bootstrap confidence bands in nonparametric curve estimation. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 291–304.
- Hall, P. and Horowitz, J. (2013). A simple bootstrap method for constructing nonparametric confidence bands for functions. *The Annals of Statistics*, 41(4):1892–1921.
- Hall, P. and Pittelkow, Y. (1990). Simultaneous bootstrap confidence bands in regression. *Journal of Statistical Computation and Simulation*, 37(1-2):99–113.
- Härdle, W. (1989). Asymptotic maximal deviation of M-smoothers. *Journal of Multivariate Analysis*, 29(2):163–179.
- Härdle, W. and Marron, J. (1991). Bootstrap simultaneous error bars for nonparametric regression. *The Annals of Statistics*, pages 778–796.
- He, X. (1997). Quantile curves without crossing. *The American Statistician*, 51(2):186–192.
- Johnston, G. J. (1982). Probabilities of maximal deviations for nonparametric regression function estimates. *Journal of Multivariate Analysis*, 12(3):402–414.
- Kim, K. I. and van de Wiel, M. A. (2008). Effects of dependence in high-dimensional multiple testing problems. *BMC bioinformatics*, 9(1):114.
- Koenker, R. (2005). *Quantile regression*. Number 38. Cambridge university press.
- Koenker, R. and Bassett Jr, G. (1978). Regression quantiles. *Econometrica: Journal of the Econometric Society*, pages 33–50.
- Leek, J. T. and Storey, J. D. (2008). A general framework for multiple testing dependence. *Proceedings of the National Academy of Sciences*, 105(48):18718–18723.
- Lindeberg, J. W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeitsrechnung. *Mathematische Zeitschrift*, 15(1):211–225.
- Liu, W. (2010). *Simultaneous inference in regression*. CRC Press.
- Liu, Y. and Wu, Y. (2011). Simultaneous multiple non-crossing quantile regression estimation using kernel constraints. *Journal of nonparametric statistics*, 23(2):415–437.
- Manly, B. F. (2006). *Randomization, bootstrap and Monte Carlo methods in biology*, volume 70. CRC Press.
- Miller, R. G. (1981). *Simultaneous statistical inference*. Springer.

- Nadaraya, E. A. (1964). On estimating regression. *Theory of Probability & Its Applications*, 9(1):141–142.
- Neumann, M. H. and Polzehl, J. (1998). Simultaneous bootstrap confidence bands in nonparametric regression. *Journal of Nonparametric Statistics*, 9(4):307–333.
- Qu, Z. (2008). Testing for structural change in regression quantiles. *Journal of Econometrics*, 146(1):170–184.
- Röllin, A. (2013). Stein’s method in high dimensions with applications. In *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, volume 49, pages 529–549. Institut Henri Poincaré.
- Romano, J. P. and Wolf, M. (2005). Exact and approximate stepdown methods for multiple hypothesis testing. *Journal of the American Statistical Association*, 100(469):94–108.
- Shao, J. and Tu, D. (1995). *The jackknife and bootstrap*. Springer.
- Šidák, Z. (1967). Rectangular confidence regions for the means of multivariate normal distributions. *Journal of the American Statistical Association*, 62(318):626–633.
- Slepian, D. (1962). The one-sided barrier problem for Gaussian noise. *Bell System Technical Journal*, 41(2):463–501.
- 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*.
- Spokoiny, V. and Willrich, N. (2015). Bootstrap tuning in ordered model selection. *In preparation*.
- Spokoiny, V. and Zhilova, M. (2014). Bootstrap confidence sets under model misspecification. *To appear in the Annals of Statistics*. *arXiv:1410.0347*.
- Spokoiny, V. G. and Zhilova, M. M. (2013). Uniform properties of the local maximum likelihood estimate. *Automation and Remote Control*, 74(10):1656–1669.
- Stein, C. M. (1981). Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics*, pages 1135–1151.

- Talagrand, M. (2003). *Spin glasses: a challenge for mathematicians: cavity and mean field models*, volume 46. Springer.
- Tropp, J. A. (2012). User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12(4):389–434.
- Wasserman, L. (2006). *All of nonparametric statistics*. Springer Science & Business Media.
- Watson, G. S. (1964). Smooth regression analysis. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 359–372.
- Westfall, P. H. (1993). *Resampling-based multiple testing: Examples and methods for p-value adjustment*, volume 279. John Wiley & Sons.
- Working, H. and Hotelling, H. (1929). Applications of the theory of error to the interpretation of trends. *Journal of the American Statistical Association*, 24(165A):73–85.

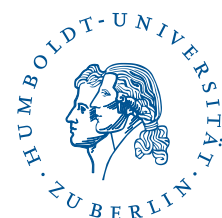
# SFB 649 Discussion Paper Series 2015

For a complete list of Discussion Papers published by the SFB 649, please visit <http://sfb649.wiwi.hu-berlin.de>.

- 001 "Pricing Kernel Modeling" by Denis Belomestny, Shujie Ma and Wolfgang Karl Härdle, January 2015.
- 002 "Estimating the Value of Urban Green Space: A hedonic Pricing Analysis of the Housing Market in Cologne, Germany" by Jens Kolbe and Henry Wüstemann, January 2015.
- 003 "Identifying Berlin's land value map using Adaptive Weights Smoothing" by Jens Kolbe, Rainer Schulz, Martin Wersing and Axel Werwatz, January 2015.
- 004 "Efficiency of Wind Power Production and its Determinants" by Simone Pieralli, Matthias Ritter and Martin Odening, January 2015.
- 005 "Distillation of News Flow into Analysis of Stock Reactions" by Junni L. Zhang, Wolfgang K. Härdle, Cathy Y. Chen and Elisabeth Bommers, January 2015.
- 006 "Cognitive Bubbles" by Ciril Bosch-Rosay, Thomas Meissner and Antoni Bosch-Domènech, February 2015.
- 007 "Stochastic Population Analysis: A Functional Data Approach" by Lei Fang and Wolfgang K. Härdle, February 2015.
- 008 "Nonparametric change-point analysis of volatility" by Markus Bibinger, Moritz Jirak and Mathias Vetter, February 2015.
- 009 "From Galloping Inflation to Price Stability in Steps: Israel 1985–2013" by Rafi Melnick and Till Strohsal, February 2015.
- 010 "Estimation of NAIRU with Inflation Expectation Data" by Wei Cui, Wolfgang K. Härdle and Weining Wang, February 2015.
- 011 "Competitors In Merger Control: Shall They Be Merely Heard Or Also Listened To?" by Thomas Giebe and Miyu Lee, February 2015.
- 012 "The Impact of Credit Default Swap Trading on Loan Syndication" by Daniel Streitz, March 2015.
- 013 "Pitfalls and Perils of Financial Innovation: The Use of CDS by Corporate Bond Funds" by Tim Adam and Andre Guettler, March 2015.
- 014 "Generalized Exogenous Processes in DSGE: A Bayesian Approach" by Alexander Meyer-Gohde and Daniel Neuhoff, March 2015.
- 015 "Structural Vector Autoregressions with Heteroskedasticity" by Helmut Lütkepohl and Aleksei Netšunajev, March 2015.
- 016 "Testing Missing at Random using Instrumental Variables" by Christoph Breunig, March 2015.
- 017 "Loss Potential and Disclosures Related to Credit Derivatives – A Cross-Country Comparison of Corporate Bond Funds under U.S. and German Regulation" by Dominika Paula Gałkiewicz, March 2015.
- 018 "Manager Characteristics and Credit Derivative Use by U.S. Corporate Bond Funds" by Dominika Paula Gałkiewicz, March 2015.
- 019 "Measuring Connectedness of Euro Area Sovereign Risk" by Rebekka Gätjen Melanie Schienle, April 2015.
- 020 "Is There an Asymmetric Impact of Housing on Output?" by Tsung-Hsien Michael Lee and Wenjuan Chen, April 2015.
- 021 "Characterizing the Financial Cycle: Evidence from a Frequency Domain Analysis" by Till Strohsal, Christian R. Proaño and Jürgen Wolters, April 2015.

**SFB 649, Spandauer Straße 1, D-10178 Berlin**  
**<http://sfb649.wiwi.hu-berlin.de>**

This research was supported by the Deutsche  
Forschungsgemeinschaft through the SFB 649 "Economic Risk".



## SFB 649 Discussion Paper Series 2015

For a complete list of Discussion Papers published by the SFB 649, please visit <http://sfb649.wiwi.hu-berlin.de>.

- 022 "Risk Related Brain Regions Detected with 3D Image FPCA" by Ying Chen, Wolfgang K. Härdle, He Qiang and Piotr Majer, April 2015.
- 023 "An Adaptive Approach to Forecasting Three Key Macroeconomic Variables for Transitional China" by Linlin Niu, Xiu Xu and Ying Chen, April 2015.
- 024 "How Do Financial Cycles Interact? Evidence from the US and the UK" by Till Strohsal, Christian R. Proaño, Jürgen Wolters, April 2015.
- 025 "Employment Polarization and Immigrant Employment Opportunities" by Hanna Wielandt, April 2015.
- 026 "Forecasting volatility of wind power production" by Zhiwei Shen and Matthias Ritter, May 2015.
- 027 "The Information Content of Monetary Statistics for the Great Recession: Evidence from Germany" by Wenjuan Chen and Dieter Nautz, May 2015.
- 028 "The Time-Varying Degree of Inflation Expectations Anchoring" by Till Strohsal, Rafi Melnick and Dieter Nautz, May 2015.
- 029 "Change point and trend analyses of annual expectile curves of tropical storms" by P.Burdejova, W.K.Härdle, P.Kokoszka and Q.Xiong, May 2015.
- 030 "Testing for Identification in SVAR-GARCH Models" by Helmut Luetkepohl and George Milunovich, June 2015.
- 031 " Simultaneous likelihood-based bootstrap confidence sets for a large number of models" by Mayya Zhilova, June 2015.

**SFB 649, Spandauer Straße 1, D-10178 Berlin**  
**<http://sfb649.wiwi.hu-berlin.de>**

This research was supported by the Deutsche  
Forschungsgemeinschaft through the SFB 649 "Economic Risk".

