

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

Preprint

ISSN 0946 – 8633

Critical dimension in profile semiparametric estimation

Andreas Andresen¹ Vladimir Spokoiny²

submitted: April 9, 2013

¹ Weierstrass Institute
Mohrenstr. 39
10117 Berlin
Germany
E-Mail: andreas.andresen@wias-berlin.de

² HU Berlin and Weierstrass Institute
Mohrenstr. 39
10117 Berlin
Germany
E-Mail: vladimir.spokoiny@wias-berlin.de

No. 1776
Berlin 2013



2010 *Mathematics Subject Classification.* Primary 62F10; secondary 62J12, 62F25, 62H12.

Key words and phrases. maximum likelihood, local quadratic bracketing, spread, local concentration.

The first author is supported by Research Units 1735 “Structural Inference in Statistics: Adaptation and Efficiency” Financial support by the German Research Foundation (DFG) through the Collaborative Research Center 649 “Economic Risk” is gratefully acknowledged .

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

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

Abstract

This paper revisits the classical inference results for profile quasi maximum likelihood estimators (profile MLE) in the semiparametric estimation problem. We mainly focus on two prominent theorems: the Wilks phenomenon and Fisher expansion for the profile MLE are stated in a new fashion allowing finite samples and model misspecification. The method of study is also essentially different from the usual analysis of the semiparametric problem based on the notion of the hardest parametric submodel. Instead we apply the local bracketing and the upper function devices from Spokoiny (2012). This novel approach particularly allows to address the important issue of the effective target and nuisance dimension and it does not involve any pilot estimator of the target parameter. The obtained nonasymptotic results are surprisingly sharp and yield the classical asymptotic statements including the asymptotic normality and efficiency of the profile MLE. The general results are specified to the important special cases of an i.i.d. sample.

1 Introduction

Many statistical tasks can be viewed as problems of semiparametric estimation when the unknown data distribution is described by a high or infinite dimensional parameter while the target is of low dimension. Typical examples are provided by functional estimation, estimation of a function at a point, or simply by estimating a given subvector of the parameter vector. The classical statistical theory provides a general solution to this problem: estimate the full parameter vector by the maximum likelihood method and project the obtained estimate onto the target subspace. This approach is known as *profile maximum likelihood* and it appears to be *semiparametrically efficient* under some mild regularity conditions. We refer to the papers Murphy and Van der Vaart (2000, 1999) and the book Kosorok (2005) for a detailed presentation of the modern state of the theory and further references. The famous Wilks result claims that the likelihood ratio test statistic in the semiparametric test problem is nearly chi-square with p degrees of freedom corresponding to the dimension of the target parameter. Various extensions of this result can be found e.g. in Fan et al. (2001); Fan and Huang (2005); Boucheron and Massart (2011); see also the references therein.

This study revisits the problem of profile semiparametric estimation and addresses some new issues. The most important difference between our approach and the classical theory is a nonasymptotic character of our study. A finite sample analysis is particularly challenging because most of notions, methods and tools in the classical theory are formulated in the asymptotic setup with growing sample size. Only few finite sample general results are available; see e.g. the recent paper Boucheron and Massart (2011). The results of this paper explicitly describes all “small” terms in the expansion of the log-likelihood. This helps to carefully treat the

question of applicability of the approach in different situations. A particularly important question is about the critical dimension of the target p and the full parameter dimension p^* for which the main results are still accurate. Another issue addressed in this paper is the model misspecification. In many practical problems, it is unrealistic to expect that the model assumptions are exactly fulfilled, even if some rich nonparametric models are used. This means that the true data distribution P does not belong to the considered parametric family. Applicability of the general semiparametric theory in such cases is questionable. An important feature of the new approach of Spokoiny (2012) is that it equally applies under a possible model misspecification.

The mentioned issues, especially the non-asymptotic character of study dictate to change entirely the tools and methods of analysis. We apply the recent bracketing approach of Spokoiny (2012) and demonstrate its power on the considered case of semiparametric estimation. Let \mathbf{Y} denote the observed random data, and P denote the data distribution. The parametric statistical model assumes that the unknown data distribution P belongs to a given parametric family (P_v) :

$$\mathbf{Y} \sim P = P_{v^*} \in (P_v, v \in \mathcal{Y}),$$

where \mathcal{Y} is some high dimensional or even infinite dimensional parameter space. This paper concentrates on a finite dimensional setting, however, an extension to a functional space is feasible and to be considered elsewhere. The maximum likelihood approach in the parametric estimation suggests to estimate the whole parameter vector v by maximizing the corresponding log-likelihood $\mathcal{L}(v) = \log \frac{dP_v}{d\mu_0}(\mathbf{Y})$ for some dominating measure μ_0 :

$$\tilde{v} \stackrel{\text{def}}{=} \operatorname{argmax}_{v \in \mathcal{Y}} \mathcal{L}(v).$$

Our study admits a model misspecification $P \notin (P_v, v \in \mathcal{Y})$. Equivalently, one can say that $\mathcal{L}(v)$ is the *quasi log-likelihood function* on \mathcal{Y} . The “target” value v^* of the parameter v can be defined by

$$v^* = \operatorname{argmax}_{v \in \mathcal{Y}} \mathbb{E} \mathcal{L}(v).$$

Under model misspecification, v^* defines the best parametric fit to P by the considered family. In the semiparametric framework, the target of analysis is only a low dimensional component θ of the whole parameter v . This means that the target of estimation is

$$\theta^* = \Pi_0 v^*,$$

for some mapping $\Pi_0 : \mathcal{Y} \rightarrow \mathbb{R}^p$, and $p \in \mathbb{N}$ stands for the dimension of the target.

The *profile maximum likelihood* approach defines the estimator of θ^* by projecting the obtained MLE \tilde{v} on the target space:

$$\tilde{\theta} = \Pi_0 \tilde{v}.$$

The Gauss-Markov Theorem claims the efficiency of such procedures for linear Gaussian models and linear mapping Π_0 , and the famous Fisher result extends it in the asymptotic sense to

the general situation under some regularity conditions. The Wilks phenomenon describes the limiting distribution of the likelihood ratio test statistic T :

$$T \stackrel{\text{def}}{=} \sup_{\mathbf{v} \in \mathcal{Y}} \mathcal{L}(\mathbf{v}) - \sup_{\substack{\mathbf{v} \in \mathcal{Y} \\ \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*}} \mathcal{L}(\mathbf{v}). \quad (1.1)$$

It appears that the distribution of this test statistic is nearly chi-square χ_p^2 as the samples size grows, Wilks (1938). In particular, this limiting behavior does not depend on the particular model structure and on the full dimension of the parameter \mathbf{v} , only the dimension of the target matters. The full parameter dimension can be even infinite under some upper bounds on its total entropy.

Below we consider a slightly different presentation of this estimator based on the partial optimization of the objective function $\mathcal{L}(\mathbf{v})$ for a fixed $\boldsymbol{\theta}$. Namely, define

$$\check{L}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \max_{\substack{\mathbf{v} \in \mathcal{Y} \\ \Pi_0 \mathbf{v} = \boldsymbol{\theta}}} \mathcal{L}(\mathbf{v}). \quad (1.2)$$

Then the profile MLE can be defined as the point of maximum of $\check{L}(\boldsymbol{\theta})$:

$$\tilde{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \check{L}(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \max_{\substack{\mathbf{v} \in \mathcal{Y} \\ \Pi_0 \mathbf{v} = \boldsymbol{\theta}}} \mathcal{L}(\mathbf{v}).$$

The test statistic T from (1.1) is also called the *semiparametric excess* and it can be defined as

$$\check{L}(\tilde{\boldsymbol{\theta}}) - \check{L}(\boldsymbol{\theta}^*) = \max_{\mathbf{v} \in \mathcal{Y}} \mathcal{L}(\mathbf{v}) - \max_{\substack{\mathbf{v} \in \mathcal{Y} \\ \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*}} \mathcal{L}(\mathbf{v}).$$

The Wilks result can be rewritten as

$$2\{\check{L}(\tilde{\boldsymbol{\theta}}) - \check{L}(\boldsymbol{\theta}^*)\} \xrightarrow{w} \chi_p^2.$$

The *local asymptotic normality* (LAN) approach by Le Cam leads to the most general setup in which the Wilks type results can be established. However, the classical theory of semiparametric estimation faces serious difficulties when the dimension of the nuisance parameter becomes large of infinite. The LAN property yields a local approximation of the log-likelihood of the full model by the log-likelihood of a linear Gaussian model, and this property is only validated in a root-n neighborhood of the true point. The non- and semiparametric cases require to consider larger neighborhoods where the LAN approach is not applicable any more. A proper extension of the Wilks result to the case of a growing or infinite nuisance dimension is quite challenging and involves special constructions like a pilot consistent estimator of the target, a hardest parametric submodel as well as some power tools of the empirical process theory; see Murphy and Van der Vaart (2000) or Kosorok (2005) for a comprehensive presentation.

The recent paper Spokoiny (2012) offers a new look at the classical LAN theory. The basic idea is to replace the local approximation by *local bracketing*. Instead of one approximating Gaussian log-likelihood, one builds two different quadratic processes such that the original log-likelihood

can be sandwiched between them up to a small error. It appears that the bracketing device can be applied for much larger neighborhoods than in the LAN approach. In this paper we show that the local bracketing approach of Spokoiny (2012) can be used for obtaining a version of the Wilks Theorem in a quite general semiparametric setup avoiding any special construction like “the hardest parametric submodel”.

Another important issue is that the new approach does not rely on any pilot estimator of the target. The usual assumption that a consistent pilot estimator is available can be even misleading in our setup because it separates local and global considerations. This paper attempts to figure out a list of condition ensuring global concentration and local expansion at the same time. This particularly allows to address the crucial question of the largest dimensionality or the nuisance parameter for which the Wilks result still holds. It appears that the profile semiparametric approach is validated under the constraint $p^{*3} \ll n$, where p^* is the full parameter dimension. It applies even if the dimension p of the target grows with the sample size under the mentioned constraint. The important *identifiability* issue is also addressed in a more careful way for the considered finite sample case.

For the further presentation we have to briefly outline the basic results from Spokoiny (2012). Introduce the log-likelihood ratio process

$$\mathcal{L}(\boldsymbol{v}, \boldsymbol{v}^*) = \mathcal{L}(\boldsymbol{v}) - \mathcal{L}(\boldsymbol{v}^*).$$

The key *bracketing* result of Spokoiny (2012) claims that $\mathcal{L}(\boldsymbol{v}, \boldsymbol{v}^*)$ can be sandwiched on a local elliptic set $\mathcal{Y}_o(\boldsymbol{r})$ around \boldsymbol{v}^* by two quadratic in \boldsymbol{v} processes $\mathbb{L}_\epsilon(\boldsymbol{v}, \boldsymbol{v}^*)$ and $\mathbb{L}_\epsilon(\boldsymbol{v}, \boldsymbol{v}^*)$:

$$\mathbb{L}_\epsilon(\boldsymbol{v}, \boldsymbol{v}^*) - \diamond_\epsilon(\boldsymbol{r}) \leq \mathcal{L}(\boldsymbol{v}, \boldsymbol{v}^*) \leq \mathbb{L}_\epsilon(\boldsymbol{v}, \boldsymbol{v}^*) + \diamond_\epsilon(\boldsymbol{r}), \quad \boldsymbol{v} \in \mathcal{Y}_o(\boldsymbol{r}), \quad (1.3)$$

where $\diamond_\epsilon(\boldsymbol{r}) > 0$ and $\diamond_\epsilon(\boldsymbol{r}) > 0$ are small terms. The value \boldsymbol{r} here can be viewed as the radius of the set $\mathcal{Y}_o(\boldsymbol{r})$ in the intrinsic semimetric corresponding to the process $L(\boldsymbol{\theta})$. See Section B for a precise formulation. This local result is accompanied with the deviation bound of the form

$$P(\tilde{\boldsymbol{v}} \in \mathcal{Y}_o(\boldsymbol{r})) \geq 1 - e^{-x},$$

where x grows almost linearly with \boldsymbol{r} . The bracketing result (1.3) yields a number of important and informative corollaries. One of them shows that the excess $\mathcal{L}(\tilde{\boldsymbol{v}}, \boldsymbol{v}^*)$ can be approximated by a quadratic form $\|\boldsymbol{\xi}\|^2/2$, where $\boldsymbol{\xi} \stackrel{\text{def}}{=} \mathcal{D}_0^{-1} \nabla \mathcal{L}(\boldsymbol{v}^*)$ is the normalized score while \mathcal{D}_0^2 approximates the total Fisher information matrix. Another important corollary of (1.3) is an expansion of the quasi MLE $\tilde{\boldsymbol{v}}$. The mentioned results can be written in the form

$$\begin{aligned} |2\mathcal{L}(\tilde{\boldsymbol{v}}, \boldsymbol{v}^*) - \|\boldsymbol{\xi}\|^2| &\leq 2\Delta_\epsilon, \\ \|\mathcal{D}_0(\tilde{\boldsymbol{v}} - \boldsymbol{v}^*) - \boldsymbol{\xi}\|^2 &\leq 2\Delta_\epsilon, \end{aligned} \quad (1.4)$$

where Δ_ϵ is a random term called the *spread* which is small with a large probability. In a typical situation with a correctly specified model, $\boldsymbol{\xi}$ is nearly standard normal and hence, $2\mathcal{L}(\tilde{\boldsymbol{v}}, \boldsymbol{v}^*)$ is nearly $\chi_{p^*}^2$, where p^* is the full parameter dimension, while the MLE $\tilde{\boldsymbol{v}}$ is asymptotically normal and efficient. The expansion (1.4) helps to build likelihood-based confidence sets for the

true parameter \boldsymbol{v}^* . Let χ_α be the $(1 - \alpha)$ -quantile of the chi-square distribution with p^* degrees of freedom. Set

$$\mathcal{E}(\alpha) \stackrel{\text{def}}{=} \{\boldsymbol{v} \in \mathcal{Y} : 2\mathcal{L}(\tilde{\boldsymbol{v}}, \boldsymbol{v}) \leq \chi_\alpha\}.$$

Then (1.4) ensures that the coverage probability $\mathbb{P}(\boldsymbol{v}^* \notin \mathcal{E}(\alpha))$ is close to α provided that Δ_ϵ is sufficiently small.

This paper aims at establishing a similar statements for the process $\check{L}(\boldsymbol{\theta})$ from (1.2). In particular, the Wilks result can be written as

$$\check{L}(\tilde{\boldsymbol{\theta}}) - \check{L}(\boldsymbol{\theta}^*) \cong \|\check{\boldsymbol{\xi}}\|^2/2,$$

where the random p -vector $\check{\boldsymbol{\xi}}$ satisfies $\mathbb{E}\check{\boldsymbol{\xi}} = 0$ and $\mathbb{E}\|\check{\boldsymbol{\xi}}\|^2 \cong p$. The deviation properties of $\|\check{\boldsymbol{\xi}}\|^2$ resemble the ones of a chi-square random variable with p degrees of freedom just as in the Wilks phenomenon. The expansion of the profile MLE $\tilde{\boldsymbol{\theta}}$ reads as

$$\check{D}_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \cong \check{\boldsymbol{\xi}}.$$

The symmetric matrix $\check{D}_0^2 \in \mathbb{R}^{p \times p}$ is usually called the influence matrix and it is the covariance of the efficient influence function; see Kosorok (2005).

Usually in the classical semiparametric setup, the vector \boldsymbol{v} is represented as $\boldsymbol{v} = (\boldsymbol{\theta}, \boldsymbol{\eta})$, where $\boldsymbol{\theta}$ is the target of analysis while $\boldsymbol{\eta}$ is the *nuisance parameter*. We refer to this situation as $(\boldsymbol{\theta}, \boldsymbol{\eta})$ -setup and our presentation follows this setting. An extension to the \boldsymbol{v} -setup with $\boldsymbol{\theta} = \Pi_0 \boldsymbol{v}$ is straightforward. Also for simplicity we only develop our results for the case that the full parameter space \mathcal{Y} is a subset of the Euclidean space of dimensionality p^* . An extension to an infinite dimensional parameter space is possible but involves a range of technical issues that have to be done elsewhere.

Section 2 introduces the objects and tools of the analysis and collects the main results including an extension of the Wilks Theorem, concentration properties of the profile estimator and the construction of confidence sets for the “true“ parameter $\boldsymbol{\theta}^*$. The concentration properties of the profile MLE are discussed in Section D.1. The appendix collects the conditions and the proofs of the main results.

2 Main results

This section presents our main results on the semiparametric profile estimator which include the Wilks expansion of the profile maximum likelihood and the Fisher expansion of the profile MLE $\tilde{\boldsymbol{\theta}}$. All the results are stated under the same list of conditions that can be found in Section A of the appendix. As already mentioned, our setup follows Spokoiny (2012). However, at one point there is an essential difference. The results of Spokoiny (2012) are stated for just one fixed finite sample. The same continues to hold for the results below. But we are also interested in understanding what happens if *the full dimension* p^* becomes large. For this we consider below an asymptotic setup with $p^* = p_n$, where n denotes the asymptotic parameter. It can

be viewed as the sample size with $n \rightarrow \infty$. We assume that all considered objects depend on n including the likelihood function, the full parameter set \mathcal{Y} and its dimension p^* , as well as all the constants in our conditions. The primary goal of our study is to fix the necessary and sufficient conditions on growth of p_n with n which ensures the Wilks and Fisher results.

Our result apply even if the target parameter θ is of growing dimension. The dimension p can be of order p^* . The case with a full dimensional target and low dimensional nuisance is also included.

2.1 The Wilks and Fisher expansion

This section states the key results in the semiparametric framework which heavily use the local bracketing idea of Spokoiny (2012). First we introduce the main elements of the bracketing device. This includes two $p^* \times p^*$ matrices \mathcal{V}_0^2 and \mathcal{D}_0^2 and two constants $\epsilon = (\delta, \rho)$. The matrix \mathcal{V}_0^2 describes the variability of the process $\mathcal{L}(\mathbf{v})$ around the true point \mathbf{v}^* :

$$\mathcal{V}_0^2 \stackrel{\text{def}}{=} \text{Var}\{\nabla \mathcal{L}(\mathbf{v}^*)\}. \quad (2.1)$$

The matrix \mathcal{D}_0^2 is defined similarly to the Fisher information matrix:

$$\mathcal{D}_0^2 \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}^*). \quad (2.2)$$

Here and in what follows we implicitly assume that the log-likelihood function $\mathcal{L}(\mathbf{v})$ is sufficiently smooth in \mathbf{v} , $\nabla \mathcal{L}(\mathbf{v})$ stands for the gradient and $\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v})$ for the Hessian of the expectation $\mathbb{E} \mathcal{L}$ at \mathbf{v} . It is worth mentioning that the matrices \mathcal{D}_0^2 and \mathcal{V}_0^2 coincide if the model $\mathbf{Y} \sim \mathbb{P}_{\mathbf{v}^*} \in (\mathbb{P}_{\mathbf{v}})$ is correctly specified and sufficiently regular; see e.g. Ibragimov and Khas'minskij (1981).

Now we switch to the (θ, η) -setup. Consider the block representation of the vector $\nabla \stackrel{\text{def}}{=} \nabla \mathcal{L}(\mathbf{v}^*)$ and of the matrices \mathcal{V}_0^2 from (2.1) and \mathcal{D}_0^2 from (2.2):

$$\nabla = \begin{pmatrix} \nabla_{\theta} \\ \nabla_{\eta} \end{pmatrix}, \quad \mathcal{D}_0^2 = \begin{pmatrix} D_0^2 & A_0 \\ A_0^\top & H_0^2 \end{pmatrix}, \quad \mathcal{V}_0^2 = \begin{pmatrix} V_0^2 & B_0 \\ B_0^\top & Q_0^2 \end{pmatrix}.$$

Define also the $p \times p$ matrix \check{D}_0^2 and p -vectors $\check{\nabla}_{\theta}$ and $\check{\xi}$ as

$$\begin{aligned} \check{D}_0^2 &= D_0^2 - A_0 H_0^{-2} A_0^\top, \\ \check{\nabla}_{\theta} &= \nabla_{\theta} - A_0 H_0^{-2} \nabla_{\eta}, \\ \check{\xi} &= \check{D}_0^{-1} \check{\nabla}_{\theta}. \end{aligned}$$

In what follows, by C we denote a generic fixed constant. For all results presented below we assume a sufficiently large value x to be fixed. It determines our level of overwhelming probability: a generic random set $\Omega(x)$ is of *dominating probability* if

$$\mathbb{P}(\Omega(x)) \geq 1 - C e^{-x}.$$

Similarly to p^* , the value x may depend on the asymptotic parameter n and grows to infinity with n . A particularly relevant choice is $x = x_n = C \log n$ for a fixed $C > 0$. We only require that x_n is not too large, more precisely, $x \leq x_c$; see (C.2) from Section C. In the i.i.d. setup x_c is of order $n^{1/2}$.

The other important value to be fixed is r_0 . This value determines the frontier between local and global consideration. In the local vicinity $\mathcal{I}_o(r_0)$ of radius r_0 we apply a very accurate local quadratic approximation of the log-likelihood process while outside of this vicinity a much more rough upper function device can be used; see Section B for more details. The general rule for the choice of r_0 is given by the condition $r_0^2 \geq C_0(p^* + x)$ for some specific constant C_0 that is independent of p^* and $n \in \mathbb{N}$ if $\nu_0/b > 0$ is independent of p^* and $n \in \mathbb{N}$ (see Section A). The quality of local quadratic approximation is measured by two functions $\delta(r)$ and $\omega(r)$ shown in local conditions $(\mathcal{E}\mathcal{D}_1)$, (\mathcal{L}_0) of Section A. More exactly, it can be described by the quantities τ_ϵ defined as

$$\tau_\epsilon \stackrel{\text{def}}{=} \delta(r_0) + 3\nu_0 a^2 \omega(r_0), \quad (2.3)$$

where the constants ν_0 and a are from conditions $(\mathcal{E}\mathcal{D}_1)$ and (\mathcal{I}) in Section A. The sub-index ϵ stands for the pair $\delta(r_0), \omega(r_0)$. Our results implicitly assume that τ_ϵ is small. We comment on typical behavior of τ_ϵ in Section 2.2 in context of i.i.d. models.

The first result can be viewed as an extension of the Wilks Theorem.

Theorem 2.1. *Assume $(\mathcal{E}\mathcal{D}_0)$, $(\mathcal{E}\mathcal{D}_1)$, (\mathcal{L}_0) , (\mathcal{I}) , $(\mathcal{E}\mathbf{x})$ and $(\mathcal{L}\mathbf{x})$ with $b(\mathbf{x}) \equiv b$; see Section A. Let also τ_ϵ from (2.3) fulfill $\tau_\epsilon \leq 1/2$. Then it holds on a random set $\Omega(x)$ of dominating probability*

$$|2\check{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\check{\boldsymbol{\xi}}\|^2| \leq C\tau_\epsilon(p^* + x), \quad (2.4)$$

Remark 2.1. In the case of the correct model specification with $\mathcal{D}_0^2 = \mathcal{V}_0^2$, the deviation properties of the quadratic form $\|\check{\boldsymbol{\xi}}\|^2 = \|\check{D}_0^{-1}\check{\nabla}_\theta\|^2$ are essentially the same as of a chi-square random variable with p degrees of freedom; see Theorem C.1 in the appendix. In the case of a possible model misspecification with $\mathcal{D}_0^2 \neq \mathcal{V}_0^2$, the behavior of the quadratic form $\|\check{\boldsymbol{\xi}}\|^2$ will depend on the characteristics of the matrix $B \stackrel{\text{def}}{=} \mathcal{D}_0^{-1}\mathcal{V}_0^2\mathcal{D}_0^{-1}$; see again Theorem C.1. Moreover, in the asymptotic setup the vector $\check{\boldsymbol{\xi}}$ is asymptotically standard normal; see Section 2.2 for the i.i.d. case.

Remark 2.2. The partial maximum likelihood process $\check{L}(\boldsymbol{\theta})$ can be used for defining the likelihood-based confidence sets of the form

$$\mathcal{E}(\mathfrak{z}) = \{\boldsymbol{\theta} : \check{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) \leq \mathfrak{z}\}$$

for some $\mathfrak{z} > 0$. The bound (2.4) helps to evaluate the coverage probability $IP(\boldsymbol{\theta}^* \notin \mathcal{E}(\mathfrak{z}))$ in terms of deviation probability for the quadratic form $\|\check{\boldsymbol{\xi}}\|^2$; cf. Corollary 3.2 in Spokoiny (2012).

The next result presents an expansion of the profile MLE $\tilde{\boldsymbol{\theta}}$.

Theorem 2.2. *Under the conditions of Theorem 2.1, it holds on a random set $\Omega(\mathbf{x})$ of dominating probability*

$$\|\check{D}_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\|^2 \leq C\tau_\epsilon(p^* + \mathbf{x}). \quad (2.5)$$

Remark 2.3. One can use the expansion (2.5) for describing the concentration probability for elliptic sets

$$\mathcal{A}(z) = \{\boldsymbol{\theta} : \|\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq z\};$$

cf. Corollary 3.5 in Spokoiny (2012).

In the next section the result (2.5) is used to show asymptotic normality and efficiency of the profile estimator in the i.i.d. setting and under the correct model specification.

2.2 The i.i.d. case and asymptotic efficiency

Here we briefly discuss the implications of our general results to the case with $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ where observations Y_i are i.i.d. from a measure P . The parametric assumption means $P = P_{\mathbf{v}^*} \in (P_{\mathbf{v}}, \mathbf{v} \in \mathcal{Y})$ for a given parametric family $(P_{\mathbf{v}})$, where \mathcal{Y} is a subset of the Euclidean space \mathbb{R}^{p^*} . We assume that $(P_{\mathbf{v}})$ obeys the regularity conditions listed in Section 5.1 of Spokoiny (2012). By $\ell(y, \mathbf{v})$ we denote the log-density of $P_{\mathbf{v}}$ w.r.t. some dominating measure μ_0 . For simplicity of comparison with the classical results we do not discuss the model misspecification issue, i.e. the parametric assumption is correct. However, an extension to the case of a misspecified model is straightforward. We utilize that $\mathcal{V}_0^2 = \mathcal{D}_0^2 = n\mathbb{F}$, $\omega(\mathbf{r}) = \omega^*\mathbf{r}/n^{1/2}$, $\delta(\mathbf{r}) = \delta^*\mathbf{r}/n^{1/2}$, and $\mathbf{g} = \mathbf{g}_1\sqrt{n}$; see Lemma 5.1 in Spokoiny (2012). Here \mathbb{F} is the Fisher information matrix of the family $(P_{\mathbf{v}})$ at the point \mathbf{v}^* , and ω^* , δ^* , and \mathbf{g}_1 are some positive constants.

It is shown in Spokoiny (2012) that the full parameter \mathbf{v}^* can be well estimated provided that p^*/n is sufficiently small. More precisely, the concentration property for the set $\mathcal{Y}_\circ(\mathbf{r})$ requires $\mathbf{r}^2 \geq C p^*$ for a fixed C , while the local bracketing device is validated up to the spread $\Delta_\epsilon(\mathbf{r})$ which is of order $p^*\delta(\mathbf{r}) \asymp p^*\mathbf{r}/n^{1/2} \asymp p^{*3/2}/n^{1/2}$. The range of applicability for the proposed approach can be informally defined by the rule “the spread is smaller than the value of the problem”, where the value of the problem is understood as the expected excess. If the full parameter \mathbf{v} is estimated, the value of the problem is of order p^* leading to the constraint “ p^*/n is small”. If the target parameter is of dimension p , then the value of the problem is also of order p leading to the constraint “ $p^{*3/2}/(n^{1/2}p)$ is small”.

Now we specify the results in the $(\boldsymbol{\theta}, \boldsymbol{\eta})$ semiparametric setup. To state the result we only need a version of the identifiability condition (\mathcal{I}) on the marginal distribution. Let \mathbb{F} be the Fisher information matrix of the family $(P_{\mathbf{v}})$ at the true point \mathbf{v}^* . Consider its block representation

$$\mathbb{F} = \begin{pmatrix} \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\theta}} & \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\eta}} \\ \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top & \mathbb{F}_{\boldsymbol{\eta}\boldsymbol{\eta}} \end{pmatrix}.$$

The required identifiability condition reads as follows:

(ι) There is a constant $\rho < 1$ such that

$$\|\mathbb{F}_{\theta\theta}^{-1/2}\mathbb{F}_{\theta\eta}\mathbb{F}_{\eta\eta}^{-1}\mathbb{F}_{\theta\eta}^\top\mathbb{F}_{\theta\theta}^{-1/2}\|_\infty \leq \rho. \quad (2.6)$$

Also define

$$\check{\mathbb{F}} \stackrel{\text{def}}{=} \mathbb{F}_{\theta\theta} - \mathbb{F}_{\theta\eta}\mathbb{F}_{\eta\eta}^{-1}\mathbb{F}_{\theta\eta}^\top.$$

The presented result admits that the full dimension p^* grows with the sample size but slower than $n^{1/3}$. The result is applicable even in the case when the target dimension also depends on the sample size.

Theorem 2.3. *Let Y_1, \dots, Y_n be i.i.d. $\mathbb{P}_{\mathbf{v}^*}$ and let (ed_0) , (ed_1) , (ℓ_0) , (eu) , and (ℓu) with $\mathfrak{b}(\mathfrak{u}) \equiv \mathfrak{b}$ of Spokoiny (2012) hold. In addition, assume (ι); see (2.6). Define for $\mathbf{x} = \mathbf{x}_n \leq n^{1/3}$*

$$\beta_n \stackrel{\text{def}}{=} (p^* + \mathbf{x}_n)^{3/2}/n^{1/2}.$$

It holds on the a set $\Omega(\mathbf{x}_n)$ of dominating probability:

$$\begin{aligned} \|(n\check{\mathbb{F}})^{1/2}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\|^2 &\leq \mathfrak{C}\beta_n, \\ |2\check{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\check{\boldsymbol{\xi}}\|^2| &\leq \mathfrak{C}\beta_n. \end{aligned}$$

Moreover, the p -vector $\check{\boldsymbol{\xi}} \stackrel{\text{def}}{=} \check{\mathbb{F}}^{-1/2}(\nabla_{\boldsymbol{\theta}} - \mathbb{F}_{\theta\eta}\mathbb{F}_{\eta\eta}^{-1}\nabla_{\boldsymbol{\eta}})$ is asymptotically standard normal as $n \rightarrow \infty$. This yields the asymptotic efficiency of the profile MLE $\tilde{\boldsymbol{\theta}}$.

2.3 Critical dimension

This section discusses the issue of a *critical dimension*. Namely we assume that the full dimension p^* grows with the sample size n and write $p^* = p_n$. Theorem 2.3 requires that $p_n = o(n^{1/3})$. Here we show that this condition is critical for the class of models satisfying the conditions of Section A. Namely, we present an example in which the behavior of the profile MLE $\tilde{\boldsymbol{\theta}}$ heavily depends on the value $\beta_n = \sqrt{p_n^3/n} \geq \beta > 0$. The conditions of Section A are satisfied such that if $\beta_n \rightarrow 0$ we derive asymptotic efficiency of $\tilde{\boldsymbol{\theta}}$. At the same time, if $\beta_n \geq \beta > 0$, then the MLE $\tilde{\boldsymbol{\theta}}$ is not anymore root- n consistent. For technical reasons we assume that $p_n/\sqrt{n} \rightarrow 0$ which is no restriction as it is easy to find examples that do not admit an efficient profile if $p_n^2/n \not\rightarrow 0$. Let a random vector $\mathbf{X} \in \mathbb{R}^{p_n}$ follow $\mathbf{X} \sim \mathcal{N}(\mathbf{v}^*, n^{-1}I_{p_n})$. Take for simplicity $\mathbf{v}^* = 0$ and let $\mathbb{P} = \mathbb{P}_0$ mean the distribution of \mathbf{X} . Introduce a special set $\mathfrak{S} \subset \mathbb{R}^{p_n}$ with

$$\begin{aligned} \mathfrak{S} \stackrel{\text{def}}{=} \left\{ \mathbf{v} = (v_1, \dots, v_{p_n}) : v_1 = \frac{z}{2}\sqrt{\beta_n/n}, z \in \mathbb{Z} \right\} \\ \cap \mathcal{Y}_o \left(\sqrt{2p_n/n} + \frac{1}{2}\sqrt{\beta_n/n} \right). \end{aligned} \quad (2.7)$$

We denote by \mathcal{S}_δ its δ -vicinity:

$$\mathcal{S}_\delta \stackrel{\text{def}}{=} \{\mathbf{v} : d(\mathbf{v}, \mathcal{S}) < \delta\},$$

where $d(\mathbf{v}, \mathcal{S})$ is the Euclidean distance from the point \mathbf{v} to the set \mathcal{S} . Also \mathcal{S}_δ^c stands for the complement of \mathcal{S}_δ . Below we fix $\delta = 1/n$. Consider a special parametric quasi log-likelihood ratio $\mathcal{L}(\mathbf{v}, 0)$ defined as

$$\mathcal{L}(\mathbf{v}, 0) = n\mathbf{X}^\top \mathbf{v} - n\|\mathbf{v}\|^2/2 + nf(\mathbf{v})\|\mathbf{v}\|^3.$$

Here $f : \mathbb{R} \mapsto \mathbb{R}$ is a smooth function with

$$f(\mathbf{v}) = \begin{cases} 1 & \mathbf{v} \in \mathcal{S}, \\ 0 & \mathbf{v} \in \mathcal{S}_\delta^c. \end{cases}$$

Below we consider the problem of estimating the first component $\theta \stackrel{\text{def}}{=} v_1 \in \mathbb{R}$. Since by assumption $p_n/\sqrt{n} \rightarrow 0$ it holds for n large enough and for any \mathbf{v} with $\|\mathbf{v}\|^2 \leq 4p_n/n + \beta_n/n$ that $n\|\mathbf{v}\|^2/2 \geq nf(\mathbf{v})\|\mathbf{v}\|^3$ and thus

$$\operatorname{argmax}_{\mathbf{v}} \mathbb{E}\mathcal{L}(\mathbf{v}) = \operatorname{argmin}_{\mathbf{v}} \{n\|\mathbf{v}\|^2/2 - nf(\mathbf{v})\|\mathbf{v}\|^3\} = 0.$$

It is easy to see that all conditions from Section A are satisfied with $\tau_\epsilon p_n \cong \beta_n^{1/2}$ and

$$\mathcal{D}_0^2 = \mathcal{V}_0^2 = nI_{p_n}.$$

Therefore, the results from Section 2.1 yield efficiency of the profile MLE $\tilde{\theta}$ if $p_n^3/n \rightarrow 0$. Moreover, it is straightforward to see that

$$\check{D}_0 = \sqrt{n}, \quad \check{\nabla}(\mathcal{L} - \mathbb{E}\mathcal{L}) = \nabla_{\theta}(\mathcal{L} - \mathbb{E}\mathcal{L}) = nX_1, \quad \text{and } \check{\xi} = \sqrt{n}X_1.$$

It follows similarly to Theorem 2.1 that if $\beta_n^2 = p_n^3/n \rightarrow 0$ then

$$\|\check{D}_0(\tilde{\theta} - \theta^*) - \check{\xi}\| = \sqrt{n}|\tilde{v}_1 - X_1| \rightarrow 0.$$

The next result shows that in the case when $\beta_n = \sqrt{p_n^3/n}$ is not small, the profile MLE $\tilde{\theta}$ is not root- n consistent.

Theorem 2.4. *Suppose that $\beta_n \rightarrow (6c)^2$ for some $c > 0$. Let also n be large enough to ensure*

$$\frac{2^{1/3} - 1}{2^{1/6}} \sqrt{p_n/n} \geq \frac{1}{2} (p_n/n)^{3/4}.$$

There exists a positive $\alpha > 0$ such that it holds with a probability exceeding α

$$\|\check{D}_0(\tilde{\theta} - \theta^*) - \check{\xi}\| \geq \frac{1}{6}\beta_n^{1/2} - \frac{1}{\sqrt{n}} \geq c - o_n(1).$$

If $\beta_n \rightarrow \infty$, then

$$\|\check{D}_0(\tilde{\theta} - \theta^*) - \check{\xi}\| \xrightarrow{\mathbb{P}} +\infty,$$

where $\xrightarrow{\mathbb{P}}$ means convergence in probability.

A Appendix

The appendix collects our conditions and proofs of the main results.

We adopt the conditions from Section 2 of Spokoiny (2012) with the obvious change of notations. The local conditions only describe the properties of the process $\mathcal{L}(\mathbf{v})$ for $\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}_0)$ with some fixed value \mathbf{r}_0 . The global conditions have to be fulfilled on the whole \mathcal{Y} . We start with the local conditions.

(\mathcal{ED}_0) There exists a constant $\nu_0 > 0$, a positive symmetric $p^* \times p^*$ matrix \mathcal{V}_0^2 satisfying $\text{Var}\{\nabla\zeta(\mathbf{v}^*)\} \leq \mathcal{V}_0^2$, and a constant $g > 0$ such that for all $|\mu| \leq g$

$$\sup_{\gamma \in \mathbb{R}^{p^*}} \log \mathbb{E} \exp \left\{ \mu \frac{\langle \nabla\zeta(\mathbf{v}^*), \gamma \rangle}{\|\mathcal{V}_0 \gamma\|} \right\} \leq \frac{\nu_0^2 \mu^2}{2}.$$

(\mathcal{ED}_1) For all $0 < \mathbf{r} < \mathbf{r}_0$, there exists a constant $\omega(\mathbf{r}) \leq 1/2$ such that for all $\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})$ and $|\mu| \leq g$

$$\sup_{\gamma \in \mathbb{R}^{p^*}} \log \mathbb{E} \exp \left\{ \mu \frac{\langle \gamma, \nabla\zeta(\mathbf{v}) - \nabla\zeta(\mathbf{v}^*) \rangle}{\omega(\mathbf{r}) \|\mathcal{V}_0 \gamma\|} \right\} \leq \frac{\nu_0^2 \mu^2}{2}.$$

(\mathcal{L}_0) There exists a symmetric $p^* \times p^*$ -matrix \mathcal{D}_0^2 such that it holds on the set $\mathcal{Y}_o(\mathbf{r}_0)$ for all $\mathbf{r} \leq \mathbf{r}_0$

$$\left| \frac{\nabla \mathbb{E} \mathcal{L}(\mathbf{v}, \mathbf{v}^*) - \mathcal{D}_0^2(\mathbf{v} - \mathbf{v}^*)}{\|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|} \right| \leq \delta(\mathbf{r}).$$

This condition together with the identity $\nabla \mathbb{E} \mathcal{L}(\mathbf{v}^*) = 0$ implies

$$\left| \frac{-2 \mathbb{E} \mathcal{L}(\mathbf{v}, \mathbf{v}^*)}{\|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2} - 1 \right| \leq \delta(\mathbf{r}).$$

The global conditions are:

(\mathcal{L}_r) For any $\mathbf{r} > \mathbf{r}_0$ there exists a value $\mathbf{b}(\mathbf{r}) > 0$, such that

$$\frac{-\mathbb{E} \mathcal{L}(\mathbf{v}, \mathbf{v}^*)}{\|\mathcal{V}_0(\mathbf{v} - \mathbf{v}^*)\|^2} \geq \mathbf{b}(\mathbf{r}).$$

(\mathcal{E}_r) For any $r \geq r_0$ there exists a constant $\nu_0 > 0$ and a constant $g(r) > 0$ such that

$$\sup_{\mathbf{v} \in \mathcal{T}_o(r)} \sup_{\mu \leq g(r)} \sup_{\gamma \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \mu \frac{\langle \nabla \zeta(\mathbf{v}), \gamma \rangle}{\|\mathcal{V}_0 \gamma\|} \right\} \leq \frac{\nu_0^2 \mu^2}{2}.$$

Our results are stated for $g(r) \equiv g > 0$, however, an extension to the case $g(r) \rightarrow 0$ can be made similarly to Spokoiny (2012).

Finally we specify the regularity conditions. We begin by representing the information and the covariance matrices in block form:

$$\mathcal{D}_0^2 = \begin{pmatrix} D_0^2 & A_0 \\ A_0^\top & H_0^2 \end{pmatrix}, \quad \mathcal{V}_0^2 = \begin{pmatrix} V_0^2 & B_0 \\ B_0^\top & Q_0^2 \end{pmatrix}.$$

The *identifiability conditions* in Spokoiny (2012) ensure that the matrix \mathcal{D}_0 is positive and satisfies $\alpha^2 \mathcal{D}_0^2 \geq \mathcal{V}_0^2$ for some $\alpha > 0$. Here we restate these conditions in the special block form which is specific for the $(\boldsymbol{\theta}, \boldsymbol{\eta})$ -setup.

(\mathcal{I}) There are constants $\alpha > 0$ and $\rho < 1$ such that

$$\alpha^2 D_0^2 \geq V_0^2, \quad \alpha^2 H_0^2 \geq Q_0^2, \quad \alpha^2 \mathcal{D}_0^2 \geq \mathcal{V}_0^2. \quad (\text{A.1})$$

and

$$\|D_0^{-1} A_0 H_0^{-2} A_0^\top D_0^{-1}\|_\infty \leq \rho. \quad (\text{A.2})$$

The quantity ρ bounds the angle between the target and nuisance subspaces in the tangent space. The regularity condition (\mathcal{I}) ensures that this angle is not too small and hence, the target and nuisance parameters are identifiable. In particular, the matrix $\check{\mathcal{D}}_0^2$ is well posed under \mathcal{I} .

The bounds in (A.1) are given with the same constant α only for simplifying the notation. One can show that the last bound on \mathcal{D}_0^2 follows from the first two and (A.2) with another constant α' depending on α and ρ only.

B Bracketing and upper function devices

This section briefly overviews the main constructions of Spokoiny (2012) including the bracketing bound and the upper function results. The bracketing bound describes the quality of quadratic approximation of the log-likelihood process $\mathcal{L}(\mathbf{v})$ in a local vicinity of the point \mathbf{v}^* , while the upper function method is used to show that the full MLE $\tilde{\mathbf{v}}$ belongs to this vicinity with a dominating probability. Given $r > 0$, define the local set

$$\mathcal{T}_o(r) \stackrel{\text{def}}{=} \{ \mathbf{v} : (\mathbf{v} - \mathbf{v}^*)^\top \mathcal{V}_0^2 (\mathbf{v} - \mathbf{v}^*) \leq r^2 \}.$$

For $\epsilon = (\delta, \varrho)$, define the bracketing quadratic processes $\mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*)$ and $\mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*)$:

$$\begin{aligned}\mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) &\stackrel{\text{def}}{=} (\mathbf{v} - \mathbf{v}^*)^\top \nabla \mathcal{L}(\mathbf{v}^*) - \|\mathcal{D}_\epsilon(\mathbf{v} - \mathbf{v}^*)\|^2/2, \\ \mathcal{D}_\epsilon^2 &\stackrel{\text{def}}{=} \mathcal{D}_0^2(1 - \delta) - \varrho \mathcal{V}_0^2,\end{aligned}\tag{B.1}$$

and accordingly for $\underline{\epsilon} = -\epsilon = (-\delta, -\varrho)$. The next result restates the local bracketing bound of Spokoiny (2012) in the semiparametric framework. The imposed conditions and the involved constants ν_0 , $\delta(\mathbf{r})$, and $\omega(\mathbf{r})$ are explained in Section A. The presented results implicitly assume that p^* is large, \mathbf{x} is large as well to ensure that $e^{-\mathbf{x}}$ is negligible. A proper choice is $\mathbf{x} = Cp^*$ for a fixed C .

Theorem B.1 (Spokoiny (2012), Theorem 3.1). *Assume $(\mathcal{E}\mathcal{D}_1)$ and (\mathcal{L}_0) . Let for some \mathbf{r} , the values $\varrho \geq 3\nu_0\omega(\mathbf{r})$ and $\delta \geq \delta(\mathbf{r})$ be such that $\mathcal{D}^2(1 - \delta) - \varrho\mathcal{V}_0^2 \geq 0$. Then*

$$\mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*) - \diamond_{\underline{\epsilon}}(\mathbf{r}) \leq \mathcal{L}(\mathbf{v}, \mathbf{v}^*) \leq \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) + \diamond_\epsilon(\mathbf{r}), \quad \mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r}),$$

where the random variables $\diamond_\epsilon(\mathbf{r}), \diamond_{\underline{\epsilon}}(\mathbf{r})$ fulfill on a random set $\Omega(\mathbf{x})$ of dominating probability

$$\diamond_\epsilon(\mathbf{r}) \leq C\varrho(p^* + \mathbf{x}), \quad \diamond_{\underline{\epsilon}}(\mathbf{r}) \leq C\varrho(p^* + \mathbf{x}).\tag{B.2}$$

In fact, Theorem 3.1 of Spokoiny (2012) states the following bound:

$$IP\{\varrho^{-1}\diamond_\epsilon(\mathbf{r}) \geq \mathfrak{z}(\mathbb{Q}, \mathbf{x})\} \leq \exp(-\mathbf{x}).$$

with $\mathbb{Q} = 2.4p^*$ and

$$\mathfrak{z}(\mathbb{Q}, \mathbf{x}) = \begin{cases} (1 + \sqrt{\mathbf{x} + \mathbb{Q}})^2 & \text{if } 1 + \sqrt{\mathbf{x} + \mathbb{Q}} < \frac{g}{\nu_0} \\ \left\{1 + \frac{\nu_0}{g}(\mathbf{x} + \mathbb{Q}) + \frac{g}{2\nu_0}\right\}^2 & \text{otherwise.} \end{cases}$$

Under the assumption that g is sufficiently large, that is, $g/\nu_0 \gg p^*$, we can apply $\mathfrak{z}(\mathbb{Q}, \mathbf{x}) \approx \mathbf{x} + \mathbb{Q} \leq C(p^* + \mathbf{x})$, and the result of Theorem B.1 follows.

The bracketing result of Theorem B.1 is local in the sense that it only applies for $\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})$. Following to the general approach of Spokoiny (2012) we accompany it with the large deviation bound on the concentration probability $IP(\tilde{\mathbf{v}} \in \mathcal{Y}_\circ(\mathbf{r}))$ when the local radius \mathbf{r} exceeds some level \mathbf{r}_0 which has to be sufficiently large, namely $\mathbf{r}_0^2 \geq Cp^*$. We adopt the upper function approach from Spokoiny (2012); cf. Corollary 4.4 therein. Again the constants $g(\mathbf{r})$ and $b(\mathbf{r})$ are introduced in Section A.

Theorem B.2 (Spokoiny (2012), Theorem 4.1). *Suppose $(\mathcal{E}\mathbf{r})$ and $(\mathcal{L}\mathbf{r})$ with $b(\mathbf{r}) \equiv b$. If for $\mathbf{r} \geq \mathbf{r}_0$, the following conditions are fulfilled:*

$$\begin{aligned}1 + \sqrt{\mathbf{x} + \mathbb{Q}} &\leq 3\nu_0^2 g(\mathbf{r})/b, \\ 6\nu_0\sqrt{\mathbf{x} + \mathbb{Q}} &\leq rb,\end{aligned}\tag{B.3}$$

then $\tilde{\mathbf{v}} \in \Upsilon_{\circ}(\mathbf{r}_0)$ on a random set $\Omega(\mathbf{x})$ of dominating probability. The same bound holds for the probability $\tilde{\mathbf{v}}_{\theta^*} \in \Upsilon_{\circ}(\mathbf{r}_0)$ where $\tilde{\mathbf{v}}_{\theta^*}$ maximizes $\mathcal{L}(\mathbf{v}, \mathbf{v}^*)$ subject to $\Pi_0 \mathbf{v} = \theta^*$:

$$\tilde{\mathbf{v}}_{\theta^*} \stackrel{\text{def}}{=} \underset{\substack{\mathbf{v} \in \mathcal{V} \\ \Pi_0 \mathbf{v} = \theta^*}}{\text{argmax}} \mathcal{L}(\mathbf{v}, \mathbf{v}^*).$$

Remark B.1. The condition (B.3) helps to understand which \mathbf{r}_0 ensures prescribed concentration properties of $\tilde{\mathbf{v}}$ and $\tilde{\mathbf{v}}_{\theta^*}$. Namely, if $g(\mathbf{r})$ is large enough, then (B.3) follows from the bound

$$\mathbf{r}_0 \geq 6\mathbf{b}^{-1}\nu_0\sqrt{\mathbf{x} + \mathbb{Q}}.$$

C Deviation bounds for quadratic forms

The following general result from Spokoiny (2013) helps to control the deviation for quadratic forms of type $\|\mathbb{B}\boldsymbol{\xi}\|^2$ for a given positive matrix \mathbb{B} and a random vector $\boldsymbol{\xi}$. It will be used several times in our proofs. Suppose that

$$\log \mathbb{E} \exp(\boldsymbol{\gamma}^\top \boldsymbol{\xi}) \leq \|\boldsymbol{\gamma}\|^2/2, \quad \boldsymbol{\gamma} \in \mathbb{R}^p, \|\boldsymbol{\gamma}\| \leq g. \quad (\text{C.1})$$

For a symmetric matrix \mathbb{B} , define

$$\mathbf{p} = \text{tr}(\mathbb{B}^2), \quad \mathbf{v}^2 = 2 \text{tr}(\mathbb{B}^4), \quad \lambda^* \stackrel{\text{def}}{=} \|\mathbb{B}^2\|_{\infty} \stackrel{\text{def}}{=} \lambda_{\max}(\mathbb{B}^2).$$

We suppose that $\lambda^* \leq 1$, otherwise we should replace everywhere \mathbb{B} with \mathbb{B}/λ^* .

Let g be shown in (C.1). Define ω_c by the equation

$$\frac{\omega_c(1 + \omega_c)}{(1 + \omega_c^2)^{1/2}} = g\mathbf{p}^{-1/2}.$$

Define also $\mu_c = \omega_c^2/(1 + \omega_c^2) \wedge 2/3$. Note that $\omega_c^2 \geq 2$ implies $\mu_c = 2/3$. Further define

$$\mathbf{y}_c^2 = (1 + \omega_c^2)\mathbf{p}, \quad 2\mathbf{x}_c = \mu_c\mathbf{y}_c^2 + \log \det\{I_p - \mu_c\mathbb{B}^2\}. \quad (\text{C.2})$$

Theorem C.1 (Spokoiny (2013)). *Let $\boldsymbol{\xi}$ fulfill (C.1) with $g^2 \geq 2\mathbf{p}$. Then we have for $\mathbf{x} \leq \mathbf{x}_c$ with \mathbf{x}_c from (C.2):*

$$\begin{aligned} \mathbb{P}(\|\mathbb{B}\boldsymbol{\xi}\|^2 \geq \mathfrak{z}(\mathbf{x}, \mathbb{B})) &\leq 2e^{-\mathbf{x}} + 8.4e^{-\mathbf{x}_c}, \\ \mathfrak{z}(\mathbf{x}, \mathbb{B}) &\stackrel{\text{def}}{=} \begin{cases} \mathbf{p} + 2\mathbf{v}\mathbf{x}^{1/2}, & \mathbf{x} \leq \mathbf{v}/18, \\ \mathbf{p} + 6\mathbf{x} & \mathbf{v}/18 < \mathbf{x} \leq \mathbf{x}_c. \end{cases} \end{aligned}$$

For $\mathbf{x} > \mathbf{x}_c$

$$\mathbb{P}(\|\mathbb{B}\boldsymbol{\xi}\|^2 \geq \mathfrak{z}_c(\mathbf{x}, \mathbb{B})) \leq 8.4e^{-\mathbf{x}}, \quad \mathfrak{z}_c(\mathbf{x}, \mathbb{B}) \stackrel{\text{def}}{=} |\mathbf{y}_c + 2(\mathbf{x} - \mathbf{x}_c)/g_c|^2.$$

It appears that the bound is slightly different in two zones separated by some specific value x_c from (C.2). It is large in typical situations as $x_c \cong g$ (it is of order \sqrt{n} in the i.i.d. case). For $x \leq x_c$, we obtain the same type of bounds as in the Gaussian case, for $x > x_c$ they are a bit worse.

D Proofs

This section collects the proofs of the results in chronological order.

D.1 Proof of Theorem 2.1

Define the $m \times m$ matrices H_ϵ^2 and $H_{\underline{\epsilon}}^2$ by

$$H_\epsilon^2 = H_0^2(1 - \delta) - \varrho Q_0^2, \quad H_{\underline{\epsilon}}^2 = H_0^2(1 + \delta) + \varrho Q_0^2;$$

cf. (B.1). Below we fix some constant r which is assumed to be large enough for ensuring the dominating probability for the concentration event $C_\epsilon(r)$ defined as

$$C_\epsilon(r) \stackrel{\text{def}}{=} \left\{ \|\mathcal{V}_0(\tilde{\mathbf{v}} - \mathbf{v}^*)\| \leq r, \|\mathcal{V}_0(\tilde{\mathbf{v}}_{\theta^*} - \mathbf{v}^*)\| \leq r, \right. \\ \left. \|\mathcal{V}_0 \mathcal{D}_\epsilon^{-2} \nabla\| \leq r, \quad \|\mathcal{Q}_0 H_\epsilon^{-2} \nabla_\eta\| \leq r \right\}. \quad (\text{D.1})$$

Note that the conditions $\|\mathcal{V}_0(\tilde{\mathbf{v}} - \mathbf{v}^*)\| \leq r$ and $\|\mathcal{V}_0(\tilde{\mathbf{v}}_{\theta^*} - \mathbf{v}^*)\| \leq r$ can be represented as $\{\tilde{\mathbf{v}} \in \mathcal{Y}_o(r)\}$ and $\{\tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_o(r)\}$. Similar representation holds for

$$\tilde{\mathbf{v}}_\epsilon \stackrel{\text{def}}{=} \mathcal{D}_\epsilon^{-2} \nabla = \underset{\mathbf{v}}{\text{argmin}} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*), \\ \tilde{\boldsymbol{\eta}}_\epsilon \stackrel{\text{def}}{=} H_\epsilon^{-2} \nabla_\eta = \underset{\substack{\mathbf{v} \in \mathcal{Y} \\ \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*}}{\text{argmin}} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*).$$

For instance, $\{\|\mathcal{V}_0 \mathcal{D}_\epsilon^{-2} \nabla\| \leq r\} = \{\tilde{\mathbf{v}}_\epsilon \in \mathcal{Y}_o(r)\}$. Later we show that a proper choice of r ensures a dominating probability of the random set $C_\epsilon(r)$; see Section D.1.

We first show that the bound (2.4) is fulfilled on the set $C_\epsilon(r)$ from (D.1) with

$$\Delta_\epsilon^+(\mathbf{r}) = \diamond_\epsilon(\mathbf{r}) + \diamond_{\underline{\epsilon}}(\mathbf{r}) + \frac{\tau_\epsilon}{1 - \tau_\epsilon} \|\mathcal{D}_0^{-1} \nabla\|^2 + \frac{\tau_\epsilon}{1 + \tau_\epsilon} \|H_0^{-1} \nabla_\eta\|^2, \quad (\text{D.2})$$

$$\Delta_\epsilon^-(\mathbf{r}) = \diamond_\epsilon(\mathbf{r}) + \diamond_{\underline{\epsilon}}(\mathbf{r}) + \frac{\tau_\epsilon}{1 + \tau_\epsilon} \|\mathcal{D}_0^{-1} \nabla\|^2 + \frac{\tau_\epsilon}{1 - \tau_\epsilon} \|H_0^{-1} \nabla_\eta\|^2. \quad (\text{D.3})$$

In analogy with Spokoiny (2012), the quantity $\Delta_\epsilon(\mathbf{r})$ with

$$\Delta_\epsilon(\mathbf{r}) = \Delta_\epsilon^+(\mathbf{r}) + \Delta_\epsilon^-(\mathbf{r}) \\ = 2\diamond_\epsilon(\mathbf{r}) + 2\diamond_{\underline{\epsilon}}(\mathbf{r}) + \frac{2\tau_\epsilon}{1 - \tau_\epsilon^2} \left(\|\mathcal{D}_0^{-1} \nabla\|^2 + \|H_0^{-1} \nabla_\eta\|^2 \right), \quad (\text{D.4})$$

can be called the *semiparametric spread*. It can be seen as a payment for the bracketing device. Below we show that $\Delta_\epsilon(\mathbf{r}) \leq C \tau_\epsilon(p^* + x)$ with a dominating probability.

We start with some technical results about the maximum of the upper and lower quadratic processes $\mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*)$ and $\mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*)$. Remind the notation $\nabla \stackrel{\text{def}}{=} \nabla \mathcal{L}(\mathbf{v}^*)$.

Lemma D.1. *It holds*

$$\sup_{\mathbf{v}} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) = \frac{1}{2} \|\mathcal{D}_\epsilon^{-1} \nabla\|^2, \quad \sup_{\mathbf{v}} \mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*) = \frac{1}{2} \|\mathcal{D}_{\underline{\epsilon}}^{-1} \nabla\|^2, \quad (\text{D.5})$$

where $\sup_{\mathbf{v}}$ means the maximum over all $\mathbf{v} \in \mathbb{R}^{p^*}$. Moreover, on the random set $\{\|\mathcal{V}_0 \mathcal{D}_0^{-2} \nabla\| \leq \mathbf{r}\}$ it holds

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*) = \sup_{\mathbf{v}} \mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*) = \frac{1}{2} \|\mathcal{D}_{\underline{\epsilon}}^{-1} \nabla\|^2.$$

Proof. The identity (D.5) directly follows by maximizing the quadratic expression

$$\mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) = (\mathbf{v} - \mathbf{v}^*)^\top \nabla - \|\mathcal{D}_\epsilon(\mathbf{v} - \mathbf{v}^*)\|^2/2,$$

with the maximum at $\mathbf{v} = \mathbf{v}^* + \mathcal{D}_\epsilon^{-2} \nabla$. Similarly, the maximum of $\mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*)$ is achieved at $\mathbf{v} = \mathbf{v}^* + \mathcal{D}_{\underline{\epsilon}}^{-2} \nabla \in C_\epsilon(\mathbf{r})$ which is within $\mathcal{Y}_o(\mathbf{r})$ under the condition

$$\|\mathcal{V}_0 \mathcal{D}_{\underline{\epsilon}}^{-2} \nabla\| \leq \|\mathcal{V}_0 \mathcal{D}_0^{-2} \nabla\| \leq \mathbf{r}.$$

This yields the claim. \square

The next lemma states similar results for the constrained maximum of \mathbb{L}_ϵ and $\mathbb{L}_{\underline{\epsilon}}$ subject to $\Pi_0 \mathbf{v} = \boldsymbol{\theta}^*$. The proof is the same as for Lemma D.1. Remember the notation $\nabla_\theta \stackrel{\text{def}}{=} \nabla_\theta \mathcal{L}(\mathbf{v}^*)$, $\nabla_\eta \stackrel{\text{def}}{=} \nabla_\eta \mathcal{L}(\mathbf{v}^*)$. We also use the block representation of \mathcal{D}_0^2 :

$$\mathcal{D}_0^2 = \begin{pmatrix} D_0^2 & A_0 \\ A_0^\top & H_0^2 \end{pmatrix}.$$

Lemma D.2. *It holds*

$$\sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) = \frac{1}{2} \|H_\epsilon^{-1} \nabla_\eta\|^2, \quad (\text{D.7})$$

Moreover, it holds on the random set $\{\|Q_0 H_0^{-2} \nabla_\eta\| \leq \mathbf{r}\}$

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}): \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*) = \sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*) = \frac{1}{2} \|H_{\underline{\epsilon}}^{-1} \nabla_\eta\|^2.$$

Further, define the process

$$\mathbb{L}(\mathbf{v}, \mathbf{v}^*) = (\mathbf{v} - \mathbf{v}^*)^\top \nabla - \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2/2.$$

Remember the definition of $\check{\nabla}_{\boldsymbol{\theta}}$ and \check{D}_0^2 :

$$\begin{aligned}\check{\nabla}_{\boldsymbol{\theta}} &\stackrel{\text{def}}{=} \nabla_{\boldsymbol{\theta}} - A_0 H_0^{-2} \nabla_{\boldsymbol{\eta}}, \\ \check{D}_0^2 &\stackrel{\text{def}}{=} D_0^2 - A_0 H_0^{-2} A_0^\top.\end{aligned}$$

Lemma D.3. *It holds on the random set $\{\|\mathcal{V}_0 \mathcal{D}_0^{-2} \nabla\| \leq r, \|Q_0 H_0^{-2} \nabla_{\boldsymbol{\eta}}\| \leq r\}$*

$$\begin{aligned}\sup_{\mathbf{v}} \mathbb{L}(\mathbf{v}, \mathbf{v}^*) &= \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \mathbb{L}(\mathbf{v}, \mathbf{v}^*) = \frac{1}{2} \|\mathcal{D}_0^{-1} \nabla\|^2, \\ \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}): \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}(\mathbf{v}, \mathbf{v}^*) &= \sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}(\mathbf{v}, \mathbf{v}^*) = \frac{1}{2} \|H_0^{-1} \nabla_{\boldsymbol{\eta}}\|^2, \\ \sup_{\mathbf{v}} \mathbb{L}(\mathbf{v}, \mathbf{v}^*) - \sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}(\mathbf{v}, \mathbf{v}^*) &= \frac{1}{2} \|\check{D}_0^{-1} \check{\nabla}_{\boldsymbol{\theta}}\|^2.\end{aligned}\tag{D.8}$$

Proof. First consider the adaptive cases with $A_0 = 0$ yielding $\check{D}_0^2 = D_0^2$ and $\check{\nabla}_{\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}}$. Then the process $\mathbb{L}(\mathbf{v}, \mathbf{v}^*)$ can be decomposed as

$$\begin{aligned}\mathbb{L}(\mathbf{v}, \mathbf{v}^*) &= (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla_{\boldsymbol{\theta}} - \frac{1}{2} \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \\ &\quad + (\boldsymbol{\eta} - \boldsymbol{\eta}^*)^\top \nabla_{\boldsymbol{\eta}} - \frac{1}{2} \|H_0(\boldsymbol{\eta} - \boldsymbol{\eta}^*)\|^2,\end{aligned}$$

and the partial optimization subject to $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ yields the results (D.6) and (D.7). Note that the constrained maximum is attained at $\boldsymbol{\eta} = \boldsymbol{\eta}^* + H_0^{-2} \nabla_{\boldsymbol{\eta}}$.

The general case can be reduced to the adaptive one by the change of variable. With $\boldsymbol{\gamma} \stackrel{\text{def}}{=} \boldsymbol{\eta} - \boldsymbol{\eta}^* + H_0^{-2} A_0^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)$, one can represent $\mathbb{L}(\mathbf{v}, \mathbf{v}^*)$ in the form

$$\mathbb{L}(\mathbf{v}, \mathbf{v}^*) = (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \check{\nabla}_{\boldsymbol{\theta}} - \|\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2 + \boldsymbol{\gamma}^\top \nabla_{\boldsymbol{\eta}} - \|H_0 \boldsymbol{\gamma}\|^2/2,$$

which corresponds to the decomposition in the adaptive case. \square

On the random set $\{\tilde{\mathbf{v}} \in \mathcal{Y}_o(\mathbf{r}), \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}_o(\mathbf{r})\}$, it holds

$$\check{\mathbb{L}}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \check{\mathbb{L}}(\tilde{\boldsymbol{\theta}}) - \check{\mathbb{L}}(\boldsymbol{\theta}^*) = \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \mathcal{L}(\mathbf{v}, \mathbf{v}^*) - \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}): \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathcal{L}(\mathbf{v}, \mathbf{v}^*).$$

Theorem B.1 implies

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*) - \diamond_{\underline{\epsilon}}(\mathbf{r}) \leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \mathcal{L}(\mathbf{v}, \mathbf{v}^*) \leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \mathbb{L}_{\epsilon}(\mathbf{v}, \mathbf{v}^*) + \diamond_{\epsilon}(\mathbf{r}).\tag{D.9}$$

The same bound applies to the maximum taken over the subset $\{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}) : \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*\}$. By Lemmas D.1 and D.2, on the random set $C_{\epsilon}(\mathbf{r})$, one can replace the sup of $\mathbb{L}_{\underline{\epsilon}}(\mathbf{v}, \mathbf{v}^*)$ over

$\mathcal{I}_\circ(\mathbf{r})$ by the sup over the whole vector space \mathbb{R}^{p^*} . Putting all the obtained bounds together yields

$$\begin{aligned}\check{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) &\geq \sup_{\mathbf{v}} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) - \sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) - \diamond_\epsilon(\mathbf{r}) - \diamond_\epsilon(\mathbf{r}), \\ \check{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) &\leq \sup_{\mathbf{v}} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) - \sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) + \diamond_\epsilon(\mathbf{r}) + \diamond_\epsilon(\mathbf{r}).\end{aligned}\tag{D.10}$$

Define

$$\begin{aligned}\square_\epsilon &\stackrel{\text{def}}{=} \sup_{\mathbf{v}} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) - \sup_{\mathbf{v}} \mathbb{L}(\mathbf{v}, \mathbf{v}^*), \\ \square_\epsilon &\stackrel{\text{def}}{=} \sup_{\mathbf{v}} \mathbb{L}(\mathbf{v}, \mathbf{v}^*) - \sup_{\mathbf{v}} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*).\end{aligned}$$

Lemmas D.1 implies

$$\begin{aligned}2\square_\epsilon &= \|\mathcal{D}_\epsilon^{-1}\nabla\|^2 - \|\mathcal{D}_0^{-1}\nabla\|^2, \\ 2\square_\epsilon &= \|\mathcal{D}_0^{-1}\nabla\|^2 - \|\mathcal{D}_\epsilon^{-1}\nabla\|^2.\end{aligned}$$

Define now

$$\begin{aligned}\alpha_\epsilon &\stackrel{\text{def}}{=} \|\mathcal{D}_0\mathcal{D}_\epsilon^{-2}\mathcal{D}_0 - I_{p^*}\|_\infty, \\ \alpha_\epsilon &\stackrel{\text{def}}{=} \|I_{p^*} - \mathcal{D}_0\mathcal{D}_\epsilon^{-2}\mathcal{D}_0\|_\infty.\end{aligned}$$

The regularity conditions $(\mathcal{I}) \quad \alpha^2\mathcal{D}_0^2 \geq \mathcal{V}_0^2$ implies for $\mathcal{D}_\epsilon^2 = \mathcal{D}_0^2(1 - \delta) - \varrho\mathcal{V}_0^2$

$$\begin{aligned}\mathcal{D}_0^2(1 - \tau_\epsilon) &\leq \mathcal{D}_\epsilon^2 \leq \mathcal{D}_0^2, \\ \mathcal{D}_0^2 &\leq \mathcal{D}_\epsilon^2 \leq \mathcal{D}_0^2(1 + \tau_\epsilon).\end{aligned}$$

with $\tau_\epsilon = \delta + \varrho\alpha^{-2}$ so that the quantities α_ϵ and α_ϵ satisfy

$$\alpha_\epsilon \leq \frac{1}{1 - \tau_\epsilon} - 1 = \frac{\tau_\epsilon}{1 - \tau_\epsilon}, \quad \alpha_\epsilon \leq 1 - \frac{1}{1 + \tau_\epsilon} = \frac{\tau_\epsilon}{1 + \tau_\epsilon}.$$

This yields

$$2\square_\epsilon \leq \alpha_\epsilon \|\mathcal{D}_0^{-1}\nabla\|^2, \quad 2\square_\epsilon \leq \alpha_\epsilon \|\mathcal{D}_0^{-1}\nabla\|^2.$$

Similarly by using the result of Lemma D.2

$$\begin{aligned}\square_{\epsilon,1} &\stackrel{\text{def}}{=} \sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) - \sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}(\mathbf{v}, \mathbf{v}^*) \\ &= \frac{1}{2} (\|H_\epsilon^{-1}\nabla_\eta\|^2 - \|H_0^{-1}\nabla_\eta\|^2) \leq \frac{\alpha_\epsilon}{2} \|H_0^{-1}\nabla_\eta\|^2, \\ \square_{\epsilon,1} &\stackrel{\text{def}}{=} \sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}(\mathbf{v}, \mathbf{v}^*) - \sup_{\mathbf{v}: \Pi_0 \mathbf{v} = \boldsymbol{\theta}^*} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) \\ &= \frac{1}{2} (\|H_0^{-1}\nabla_\eta\|^2 - \|H_\epsilon^{-1}\nabla_\eta\|^2) \leq \frac{\alpha_\epsilon}{2} \|H_0^{-1}\nabla_\eta\|^2.\end{aligned}$$

Further, (D.10) and (D.8) yield

$$\begin{aligned}\check{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) &\geq \frac{1}{2} \|\check{D}_0^{-1} \check{\nabla}_{\boldsymbol{\theta}}\|^2 - \square_{\underline{\epsilon}} - \square_{\epsilon,1} - \diamond_{\epsilon}(\mathbf{r}) - \diamond_{\underline{\epsilon}}(\mathbf{r}), \\ \check{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) &\leq \frac{1}{2} \|\check{D}_0^{-1} \check{\nabla}_{\boldsymbol{\theta}}\|^2 + \square_{\epsilon} + \square_{\epsilon,1} + \diamond_{\epsilon}(\mathbf{r}) + \diamond_{\underline{\epsilon}}(\mathbf{r}).\end{aligned}$$

The proof of (D.2) and (D.3) is completed.

The next step is to bound the spread $\Delta_{\epsilon}(\mathbf{r})$ from (D.4). The error terms $\diamond_{\epsilon}(\mathbf{r})$ and $\diamond_{\underline{\epsilon}}(\mathbf{r})$ follow the bound (B.2) of Theorem B.1 and they are of order $\varrho(p^* + \mathbf{x})$. Further we have to show that $\tau_{\epsilon} \|\mathcal{D}_0^{-1} \nabla\|^2$ is small relative to $\|\check{\boldsymbol{\xi}}\|^2$ and similarly for $\tau_{\epsilon} \|H_0^{-1} \nabla_{\boldsymbol{\eta}}\|^2$. Theorem C.1 provides a general deviation probability bound for such quadratic forms. In particular, for $B \stackrel{\text{def}}{=} \mathcal{D}_0^{-1} \mathcal{V}_0^2 \mathcal{D}_0^{-1}$ and $\mathbf{x} \leq \mathbf{x}_c$

$$\mathbb{P}(\|\mathcal{D}_0^{-1} \nabla\|^2 > \mathfrak{z}(\mathbf{x}, B)) \leq 2e^{-\mathbf{x}} + 8.4e^{-\mathbf{x}_c},$$

where $\mathfrak{z}(\mathbf{x}, B) \leq \text{tr}(B) + 6\mathbf{x}$ and the constant \mathbf{x}_c is large; see Section C for a precise formulation. Under the regularity condition (\mathcal{I}) it holds $\text{tr}(B) \leq \mathfrak{a}^2 p^*$. A similar bound holds for $\|H_0^{-1} \nabla_{\boldsymbol{\eta}}\|^2$. We conclude that the spread $\Delta_{\epsilon}(\mathbf{r})$ can be bounded with a probability of order $1 - e^{-\mathbf{x}}$ by $C\tau_{\epsilon}(p^* + \mathbf{x})$ for a fixed constant C .

Further we have to show that the random set $C_{\epsilon}(\mathbf{r})$ from (D.1) is of dominating probability if $\mathbf{r}^2 = C(p^* + \mathbf{x})$ for a proper constant C . By definition

$$C_{\epsilon}(\mathbf{r}) = \{\tilde{\mathbf{v}} \in \mathcal{Y}_o(\mathbf{r}), \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}_o(\mathbf{r}), \|\mathcal{V}_0 \mathcal{D}_{\underline{\epsilon}}^{-2} \nabla\| \leq \mathbf{r}, \|Q_0 H_{\underline{\epsilon}}^{-2} \nabla_{\boldsymbol{\eta}}\| \leq \mathbf{r}\}.$$

Theorem B.2 yields

$$\mathbb{P}\{\tilde{\mathbf{v}} \notin \mathcal{Y}_o(\mathbf{r})\} + \mathbb{P}\{\tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \notin \mathcal{Y}_o(\mathbf{r})\} \leq 2e^{-\mathbf{x}}.$$

To control the probability $\mathbb{P}(\|\mathcal{V}_0 \mathcal{D}_{\underline{\epsilon}}^{-2} \nabla\| > \mathbf{r})$ we apply Corollary C.1 with

$$B = \mathcal{D}_0^{-1} \mathcal{V}_0^2 \mathcal{D}_0^{-1}, \quad .$$

With the definitions from Section C

$$\begin{aligned}\mathbb{P}(\|\mathcal{V}_0 \mathcal{D}_{\underline{\epsilon}}^{-2} \nabla\| > \mathbf{r}) &\leq \mathbb{P}(\|\mathcal{D}_{\underline{\epsilon}}^{-1} \nabla\| \times \|\mathcal{V}_0 \mathcal{D}_{\underline{\epsilon}}^{-1}\|_{\infty} > \mathbf{r}) \\ &\leq \mathbb{P}(\|\mathcal{D}_0^{-1} \nabla\| \times \|\mathcal{V}_0 \mathcal{D}_0^{-1}\|_{\infty} \geq (1 - \tau_{\epsilon})\mathbf{r}) \\ &\leq \mathbb{P}(\|\mathcal{D}_0^{-1} \nabla\| \geq (1 - \tau_{\epsilon})\mathbf{r}/\mathfrak{a}) \\ &\leq \mathbb{P}\{\|\mathcal{D}_0^{-1} \nabla\|^2 \geq \mathfrak{z}(\mathbf{x}, B)\} \\ &< e^{-\mathbf{x}} + 8.4e^{-\mathbf{x}_c},\end{aligned}$$

provided that $\mathbf{r}^2 > \mathfrak{a}^4(1 - \tau_{\epsilon})^{-2}(p^* + 6\mathbf{x})$ and $\mathbf{x} \leq \mathbf{x}_c$. By similar arguments with $B_{\boldsymbol{\eta}} = H_0^{-1} Q_0^2 H_0^{-1}$ in place of B

$$\mathbb{P}(\|Q_0 H_{\underline{\epsilon}}^{-2} \nabla_{\boldsymbol{\eta}}\| > \mathbf{r}) < e^{-\mathbf{x}} + 8.4e^{-\mathbf{x}_c'}.$$

Putting the obtained bounds together shows that for $\mathbf{x} \leq \mathbf{x}_c$ and $r_0^2 \geq C_1(p^* + \mathbf{x})$, it holds

$$1 - \mathbb{P}(C_\epsilon(r_0)) \leq C_2 e^{-\mathbf{x}},$$

for some fixed constants C_1 and C_2 depending on τ_ϵ and α only. This completes the proof.

D.2 Proof of Theorem 2.2

We show that

$$2\Delta_\epsilon^*(\mathbf{r}) \stackrel{\text{def}}{=} \frac{2\Delta_\epsilon(\mathbf{r})}{1 - \tau_\epsilon} + \frac{\tau_\epsilon}{1 - \tau_\epsilon} \|\mathcal{D}_0^{-1}\nabla\|. \quad (\text{D.11})$$

First we derive the expansion for the whole parameter vector \mathbf{v} . On the set $C_\epsilon(\mathbf{r})$, the bracketing bound (D.9) and (D.5) imply

$$\begin{aligned} \mathcal{L}(\tilde{\mathbf{v}}, \mathbf{v}^*) &\geq \sup_{\mathbf{v}} \mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*) - \diamond_{\underline{\epsilon}}(\mathbf{r}) \\ &= \|\mathcal{D}_\epsilon^{-1}\nabla\|^2/2 - \diamond_{\underline{\epsilon}}(\mathbf{r}) \\ &\geq \|\mathcal{D}_\epsilon^{-1}\nabla\|^2/2 - \diamond_{\underline{\epsilon}}(\mathbf{r}) - \square_\epsilon - \square_{\underline{\epsilon}}. \end{aligned}$$

The bracketing bound (D.9) applied at $\tilde{\mathbf{v}}$ implies

$$\mathcal{L}(\tilde{\mathbf{v}}, \mathbf{v}^*) \leq \mathbb{L}_\epsilon(\tilde{\mathbf{v}}, \mathbf{v}^*) + \diamond_\epsilon(\mathbf{r}).$$

These two bounds together yield by the definition of $\mathbb{L}_\epsilon(\mathbf{v}, \mathbf{v}^*)$

$$(\tilde{\mathbf{v}} - \mathbf{v}^*)^\top \nabla - \frac{1}{2} \|\mathcal{D}_\epsilon(\tilde{\mathbf{v}} - \mathbf{v}^*)\|^2 \geq \frac{1}{2} \|\mathcal{D}_\epsilon^{-1}\nabla\|^2 - \diamond_\epsilon(\mathbf{r}) - \diamond_{\underline{\epsilon}}(\mathbf{r}) - \square_\epsilon - \square_{\underline{\epsilon}},$$

and thus

$$\|\mathcal{D}_\epsilon(\tilde{\mathbf{v}} - \mathbf{v}^*) - \mathcal{D}_\epsilon^{-1}\nabla\|^2 \leq 2\{\square_\epsilon + \square_{\underline{\epsilon}} + \diamond_\epsilon(\mathbf{r}) + \diamond_{\underline{\epsilon}}(\mathbf{r})\} \leq 2\Delta_\epsilon. \quad (\text{D.12})$$

The condition (\mathcal{I}) implies the inequality $\|\mathcal{D}_0^{-1}\mathcal{D}_\epsilon^2\mathcal{D}_0^{-1}\|_\infty \geq 1 - \tau_\epsilon$ and hence,

$$\|\mathcal{D}_0\mathcal{D}_\epsilon^{-2}\mathcal{D}_0\|_\infty \leq (1 - \tau_\epsilon)^{-1}.$$

This and (D.12) provide

$$\|\mathcal{D}_0(\tilde{\mathbf{v}} - \mathbf{v}^*) - \mathcal{D}_0\mathcal{D}_\epsilon^{-2}\nabla\|^2 \leq \frac{2\Delta_\epsilon}{1 - \tau_\epsilon}.$$

Similarly

$$\|\mathcal{D}_0\mathcal{D}_\epsilon^{-2}\nabla - \mathcal{D}_0^{-1}\nabla\| = \|(\mathcal{D}_0\mathcal{D}_\epsilon^{-2}\mathcal{D}_0 - I_{p^*})\mathcal{D}_0^{-1}\nabla\| \leq \frac{\tau_\epsilon}{1 - \tau_\epsilon} \|\mathcal{D}_0^{-1}\nabla\|.$$

Putting together the last two bounds yields

$$\|\mathcal{D}_0(\tilde{\mathbf{v}} - \mathbf{v}^*) - \mathcal{D}_0^{-1}\nabla\| \leq \sqrt{\frac{2\Delta_\epsilon}{1-\tau_\epsilon} + \frac{\tau_\epsilon}{1-\tau_\epsilon}} \|\mathcal{D}_0^{-1}\nabla\|.$$

It remains to note that for any $\mathbf{u} \in \mathbb{R}^p$, $\boldsymbol{\eta} \in \mathbb{R}^m$, and $\mathbf{w} = (\mathbf{u}, \boldsymbol{\eta}) \in \mathbb{R}^{p^*}$, it holds with $\boldsymbol{\gamma} \stackrel{\text{def}}{=} \boldsymbol{\eta} + H_0^{-2}A_0^\top \mathbf{u} \in \mathbb{R}^m$

$$\|\mathcal{D}_0\mathbf{w}\|^2 = \|\check{D}_0\mathbf{u}\|^2 + \|H_0\boldsymbol{\gamma}\|^2 \geq \|\check{D}_0\mathbf{u}\|^2. \quad (\text{D.13})$$

Also we use $\Pi_0\mathcal{D}_0^{-2}\nabla = \check{D}_0^{-2}\check{\nabla}$. This implies for $\mathbf{w} = \tilde{\mathbf{v}} - \mathbf{v}^*$ by (D.13)

$$\|\check{D}_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \check{D}_0^{-1}\check{\nabla}\| = \|\check{D}_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^* - \check{D}_0^{-2}\check{\nabla})\| \leq \|\mathcal{D}_0(\mathbf{w} - \mathcal{D}_0^{-2}\nabla)\|,$$

and the assertion (D.11) follows.

D.3 Proof of Theorem 2.3

Choose $\mathbf{x}_n \rightarrow \infty$ and $\mathbf{x}_n = o(n^{1/3})$, e.g. $\mathbf{x}_n = C \log(n)$. Then $\beta_n \rightarrow 0$ and $\mathbb{P}(\Omega(\mathbf{x}_n)) \rightarrow 1$. Moreover, in the i.i.d. setting $\mathbf{x}_c \cong \mathbf{g} \cong \sqrt{n}$ and thus $\mathbf{x}_n \leq \mathbf{x}_c$. Similarly for n large enough with $\mathbf{r}_0^2 = \mathbf{r}_0^2(\mathbf{x}_n) = C(p^* + \mathbf{x}_n)$

$$\tau_\epsilon \cong \mathbf{r}_0/\sqrt{n} \cong \sqrt{(p^* + \mathbf{x}_n)/n} < n^{-1/3} \leq 1/2.$$

Also the i.i.d. structure of the data yields

$$\check{D}_0^2 = n\check{\mathbb{F}}.$$

Now Theorems 2.1 and 2.2 can be applied yielding the first statement of the theorem. It remains to check asymptotic standard normality of the sum

$$\check{\boldsymbol{\xi}} = (n\check{\mathbb{F}})^{-1/2}\check{\nabla} = (n\check{\mathbb{F}})^{-1/2} \sum_{i=1}^n \check{\nabla}_i,$$

with

$$\check{\nabla}_i \stackrel{\text{def}}{=} \nabla_{\boldsymbol{\theta}} \ell(Y_i, \mathbf{v}) - \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\eta}} \mathbb{F}_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1} \nabla_{\boldsymbol{\eta}} \ell(Y_i, \mathbf{v}).$$

The result follows from the central limit theorem because $\text{Cov}(\check{\nabla}_i) = \check{\mathbb{F}}$ for all i .

D.4 Proof of Theorem 2.4

The first step of the proof shows that for n large enough, the MLE $\tilde{\mathbf{v}} \in \mathbb{R}^{p^*}$ belongs with probability close to one to the $\delta = 1/n$ vicinity \mathcal{S}_δ of the set \mathcal{S} from (2.7). The second step is to show that with a probability exceeding a fixed constant $\alpha > 0$, the profile MLE $\tilde{\boldsymbol{\theta}}$ differs

significantly from X_1 which is the profile MLE in the linear Gaussian model. The third step focuses on the case $\beta_n \rightarrow \infty$.

1. First we show that for n large enough, the MLE $\tilde{\mathbf{v}} \in \mathbb{R}^{p_n}$ lies in \mathcal{S}_δ with probability close to one. For this we check that the maximum of $\mathcal{L}(\mathbf{v})$ on \mathcal{S}_δ^c is smaller than a similar maximum on \mathcal{S} for “typical” values of \mathbf{X} and n large enough. Indeed, for any point $\mathbf{v} \in \mathcal{S}_\delta^c$

$$\begin{aligned} \mathcal{L}(\mathbf{v}, 0) &\leq \max_{\mathbf{v} \in \mathcal{S}_\delta^c} \mathcal{L}(\mathbf{v}, 0) = \max_{\mathbf{v} \in \mathcal{S}_\delta^c} \{n\mathbf{X}^\top \mathbf{v} - n\|\mathbf{v}\|^2/2\} \\ &\leq \max_{\mathbf{v} \in \mathbb{R}^{p_n}} \{n\mathbf{X}^\top \mathbf{v} - n\|\mathbf{v}\|^2/2\} = \frac{n}{2}\|\mathbf{X}\|^2. \end{aligned}$$

Further, introduce a random set of “typical” values \mathbf{X} :

$$C_1 \stackrel{\text{def}}{=} \left\{ \mathbf{X} : \frac{1}{2} \left(\frac{p_n}{n} \right)^{3/2} < \|\mathbf{X}\|^3 < \left(\frac{2p_n}{n} \right)^{3/2}, \text{ and } |X_1| \leq 1 \right\}.$$

It is straightforward to see that $P(\mathbf{X} \in C_1)$ is exponentially close to one for n large. Below we assume that $\mathbf{X} \in C_1$ and study the value $\mathcal{L}(\mathbf{v}, 0)$ for $\mathbf{v} \in \mathcal{S}$. Let also n be large enough to ensure that

$$\frac{2^{1/3} - 1}{2^{1/6}} \left(\frac{p_n}{n} \right)^{1/2} \geq \frac{1}{2} \left(\frac{p_n}{n} \right)^{3/4} = \frac{1}{2} \sqrt{\beta_n/n}. \quad (\text{D.14})$$

Introduce \mathbf{X}_s as the closest point in \mathcal{S} to \mathbf{X} with $|v_1| \geq |X_1|$. This point always exists by the definition of \mathcal{S} . Denote

$$\delta(\mathbf{X}) = \|\mathbf{X} - \mathbf{X}_s\| = |X_1 - v_1|.$$

By construction of \mathcal{S} , it holds $\delta(\mathbf{X}) \leq 0.5\sqrt{\beta_n/n}$ for $\mathbf{X} \in C_1$. For n satisfying (D.14) this also yields $[\|\mathbf{X}\| - \delta(\mathbf{X})]^3 \geq 1/2\|\mathbf{X}\|^3$. Now we have for $\mathbf{X} \in C_1$

$$\begin{aligned} \max_{\mathbf{v} \in \mathcal{S}} \mathcal{L}(\mathbf{v}, 0) &\geq \mathcal{L}(\mathbf{X}_s, 0) \\ &\geq n\|\mathbf{X}\|^2 - n|X_1|\delta(\mathbf{X}) - \frac{n}{2} \{ \|\mathbf{X}\|^2 - 2|X_1|\delta(\mathbf{X}) + \delta^2(\mathbf{X}) \} \\ &\quad + n \{ \|\mathbf{X}\|^2 - 2|X_1|\delta(\mathbf{X}) + \delta^2(\mathbf{X}) \}^{3/2} \\ &\geq \frac{n}{2} \|\mathbf{X}\|^2 - n\delta^2(\mathbf{X}) + n \{ \|\mathbf{X}\| - \delta(\mathbf{X}) \}^3 \\ &> \frac{n}{2} \|\mathbf{X}\|^2 - \frac{\beta_n}{4} + \frac{n}{2} \|\mathbf{X}\|^3 > \frac{n}{2} \|\mathbf{X}\|^2 = \max_{\mathbf{v} \in \mathcal{S}_\delta^c} \mathcal{L}(\mathbf{v}, 0). \end{aligned}$$

This implies $\tilde{\mathbf{v}} \in \mathcal{S}_\delta$.

2. Now we discuss the case when $\beta_n^2 = p_n^3/n \rightarrow (6c)^4$ for some $c \geq 0$ and show that the profile MLE $\tilde{\theta}$ deviates significantly from X_1 on a random set of positive probability. Define for each $n \in \mathbb{N}$

$$C_n \stackrel{\text{def}}{=} C_1 \cap \left\{ \|\mathbf{X} - \mathbf{X}_s\| \geq \frac{1}{6} \sqrt{\beta_n/n} \right\} = C_1 \cap \left\{ |X_1 - X_{s,1}| \geq \frac{1}{6} \sqrt{\beta_n/n} \right\}.$$

It is easy to see that $\mathbb{P}(C_n) \geq \alpha$ for some fixed $\alpha > 0$ and all n . It remains to note that on the set C_n it holds under (D.14)

$$\begin{aligned} \|\check{D}_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\| &= \sqrt{n}|\tilde{v}_1 - X_1| \\ &\geq \sqrt{n}|X_1 - X_{s,1}| - \sqrt{n}/n \\ &\geq \frac{1}{6}\beta_n^{1/2} - \frac{1}{\sqrt{n}} \rightarrow \begin{cases} \infty & p_n^3/n \rightarrow \infty, \\ c & p_n^3/n \rightarrow (6c)^4, \end{cases} \end{aligned}$$

which yields the claim.

3. Finally consider the case when $\beta_n \rightarrow \infty$. Fix any sequence c_n such that $c_n \rightarrow 0$ and $c_n^2\beta_n \rightarrow \infty$, e.g. $c_n = \beta_n^{-1/4}$. Consider the random set

$$C_n \stackrel{\text{def}}{=} C_1 \cap \left\{ \|\mathbf{X} - \mathbf{X}_s\| \geq \frac{c_n}{6} \sqrt{\beta_n/n} \right\} = C_1 \cap \left\{ |X_1 - X_{s,1}| \geq \frac{c_n}{6} \sqrt{\beta_n/n} \right\}.$$

Then $\mathbb{P}(C_n) \rightarrow 1$ and on C_n

$$\|\check{D}_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\| \geq \frac{c_n}{6}\beta_n^{1/2} - \frac{1}{\sqrt{n}} \rightarrow \infty,$$

as required.

References

- Boucheron, S. and Massart, P. (2011). A high-dimensional Wilks phenomenon. *Probability Theory and Related Fields*, 150:405–433. 10.1007/s00440-010-0278-7.
- Fan, J. and Huang, T. (2005). Profile likelihood inferences on semiparametric varying-coefficient partially linear models. *Bernoulli*, 11(6):1031–1057.
- Fan, J., Zhang, C., and Zhang, J. (2001). Generalized likelihood ratio statistics and Wilks phenomenon. *Ann. Stat.*, 29(1):153–193.
- Ibragimov, I. and Khas'minskij, R. (1981). *Statistical estimation. Asymptotic theory. Transl. from the Russian by Samuel Kotz*. New York - Heidelberg -Berlin: Springer-Verlag .
- Kosorok, M. (2005). *Introduction to Empirical Processes and Semiparametric Inference*. Springer in Statistics.
- Murphy, S. A. and Van der Vaart, A. W. (1999). Observed information in semi-parametric models. *Bernoulli*, 5(3):381–412.
- Murphy, S. A. and Van der Vaart, A. W. (2000). On profile likelihood. *Journal of the American Statistical Association*, 95(450):449–465.
- Spokoiny, V. (2012). Parametric estimation. Finite sample theory. *Ann. Statist.*, 40(6):2877–2909. arXiv:1111.3029.

Spokoiny, V. (2013). Sharp deviation bounds for quadratic forms. Manuscript. arXiv:1302.1699.

Wilks, S. (1938). The large-sample distribution of the likelihood ratio for testing composite hypotheses. *Ann. Math. Statist.*, Volume 9(1):60–62.