

Critical dimension in profile semiparametric estimation

Andreas Andresen*

*Weierstrass-Institute
Mohrenstr. 39
10117 Berlin, Germany*
e-mail: andresen@wias-berlin.de

and

Vladimir Spokoiny†

*Weierstrass Institute and HU Berlin
Moscow Institute of Physics and Technology
Mohrenstr. 39
10117 Berlin, Germany*
e-mail: spokoiny@wias-berlin.de

Abstract: This paper revisits the classical inference results for profile quasi maximum likelihood estimators (profile MLE) in semiparametric models. 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 derive finite sample deviation bounds for the linear approximation error for the gradient of the loglikelihood. This novel approach particularly allows to address the impact of the effective target and nuisance dimension on the accuracy of the results. 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 for the important special case of an i.i.d. sample and the analysis is exemplified with a single index model.

MSC 2010 subject classifications: Primary 62F10; secondary 62J12, 62F25, 62H12.

Keywords and phrases: Profile maximum likelihood, local linear approximation, spread, local concentration.

Received June 2014.

*The author is supported by Research Units 1735 “Structural Inference in Statistics: Adaptation and Efficiency”.

†The author is partially supported by Laboratory for Structural Methods of Data Analysis in Predictive Modeling, MIPT, RF government grant, ag. 11.G34.31.0073. Financial support by the German Research Foundation (DFG) through the CRC 649 “Economic Risk” is gratefully acknowledged.

Contents

1	Introduction	3078
2	Main results	3083
2.1	Conditions	3083
2.2	Wilks and Fisher expansions	3089
2.3	Large deviation bounds	3092
2.4	The i.i.d. case	3093
2.5	Critical dimension	3095
2.6	Infinite dimensional nuisance	3098
2.7	Sieve approach	3098
2.8	Bias constraints and efficiency	3099
2.9	Application to single index model	3103
A	Deviation bounds for quadratic forms	3106
B	Proofs	3107
C	A bound for the norm of a random process	3121
	Aknowledgements	3123
	References	3123

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 [22, 23] and the book [18] 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 [9, 8, 6]; 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 notions, methods and tools in the classical theory are formulated in the asymptotic setup with growing sample size. Only few general finite sample results are available; see e.g. the recent paper [6]. 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 the applicability of the approach in different situations. A particularly

important question concerns 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 \mathbb{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 presented approach is that it equally applies under a possible model misspecification.

Let \mathbf{Y} denote the observed random data, and \mathbb{P} denote the data distribution. The parametric statistical model assumes that the unknown data distribution \mathbb{P} belongs to a given parametric family $(\mathbb{P}_{\mathbf{v}})$:

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

where \mathcal{Y} is some high dimensional or even infinite dimensional parameter space.

The maximum likelihood approach in the parametric estimation suggests to estimate the whole parameter vector $\mathbf{v} \in \mathcal{Y}$ by maximizing the corresponding log-likelihood $\mathcal{L}(\mathbf{v}) = \log \frac{d\mathbb{P}_{\mathbf{v}}}{d\mu_0}(\mathbf{Y})$ for some dominating measure μ_0 :

$$\tilde{\mathbf{v}} \stackrel{\text{def}}{=} \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \mathcal{L}(\mathbf{v}). \quad (1.1)$$

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

$$\mathbf{v}^* = \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \mathbb{E} \mathcal{L}(\mathbf{v}). \quad (1.2)$$

Under model misspecification, \mathbf{v}^* defines the best parametric fit of the considered family to \mathbb{P} .

In the semiparametric framework, the target of analysis is only a low dimensional component $\boldsymbol{\theta} \in \mathbb{R}^p$ of the whole parameter \mathbf{v} . This means that the target of estimation is

$$\boldsymbol{\theta}^* = \Pi_{\boldsymbol{\theta}} \mathbf{v}^*,$$

for some mapping $\Pi_{\boldsymbol{\theta}} : \mathcal{Y} \rightarrow \mathbb{R}^p$, and $p \in \mathbb{N}$ stands for the dimension of the target. Often the vector \mathbf{v} is represented as $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta})$, 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.

Define the *profile likelihood*

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

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

$$\tilde{\theta} \stackrel{\text{def}}{=} \Pi_{\theta} \tilde{\mathbf{v}} = \Pi_{\theta} \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \mathcal{L}(\mathbf{v}) = \operatorname{argmax}_{\theta \in \Theta} \max_{\substack{\mathbf{v} \in \mathcal{Y} \\ \Pi_{\theta} \mathbf{v} = \theta}} \mathcal{L}(\mathbf{v}) \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta \in \Theta} \check{\mathcal{L}}(\theta). \quad (1.3)$$

The Gauss-Markov Theorem claims the efficiency of such procedures for linear Gaussian models and a linear mapping Π_{θ} , 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 which is also called the *semiparametric excess*: $T \stackrel{\text{def}}{=} 2\{\check{\mathcal{L}}(\tilde{\theta}) - \check{\mathcal{L}}(\theta^*)\}$. It appears that the distribution of this test statistic is nearly chi-square distributed with $p \in \mathbb{N}$ degrees of freedom as the sample size grows, [33]:

$$T \stackrel{\text{def}}{=} 2\left\{\max_{\mathbf{v} \in \mathcal{Y}} \mathcal{L}(\mathbf{v}) - \max_{\substack{\mathbf{v} \in \mathcal{Y} \\ \Pi_{\theta} \mathbf{v} = \theta^*}} \mathcal{L}(\mathbf{v})\right\} = 2\{\check{\mathcal{L}}(\tilde{\theta}) - \check{\mathcal{L}}(\theta^*)\} \xrightarrow{w} \chi_p^2.$$

In particular, the limit distribution 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.

The *local asymptotic normality* (LAN) approach by Le Cam leads to the most general setup in which the Wilks and Fisher type results can be established. However, the classical theory of semiparametric estimation faces serious difficulties when the dimension of the nuisance parameter becomes large or infinite. The LAN property yields a strong 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 and Fisher 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 [22] or [18] for a comprehensive presentation.

The recent paper [29] offers a new look at the classical LAN theory. The key steps are a local quadratic bracketing for the log-likelihood process and some concentration results for its stochastic component. The results can be stated for finite samples and do not involve any asymptotic consideration. It is also shown that many corollaries of the LAN property like Fisher and Wilks expansions only rely on these two facts. The bracketing idea of [29] is to build two different quadratic processes such that the original log-likelihood can be sandwiched between them up to a small error. This paper offers another approach based on the local linear approximation of the gradient of the log-likelihood process. This allows to improve the error term of the Fisher and Wilks expansion by a factor $\sqrt{p^*}$.

For the further presentation we briefly outline the basic steps of the analysis. Introduce for $\mathbf{v} \in \mathcal{Y}$ and $\mathbf{v}^* \in \mathcal{Y}$ as defined in (1.2), the log-likelihood ratio process

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

An important step of our approach is a deviation bound for the MLE $\tilde{\mathbf{v}} \in \mathcal{Y}$ from (1.1). Given some $\mathbf{x} > 0$, we define a radius $\mathbf{r}_0 = \mathbf{r}_0(\mathbf{x}) > 0$ that ensures that

$$\mathbb{P}(\tilde{\mathbf{v}} \in \mathcal{I}_o(\mathbf{r})) \geq 1 - e^{-\mathbf{x}}, \tag{1.4}$$

where $\mathcal{I}_o(\mathbf{r})$ is a ball of radius $\mathbf{r} > 0$ in the intrinsic semi-metric corresponding to the process $\mathcal{L}(\mathbf{v})$. We give conditions that ensure that the value $\mathbf{r}_0^2(\mathbf{x})$ grows almost linearly with \mathbf{x} . See Section 2.3 for a precise formulation. The second key step is to bound for $\mathbf{r} > 0$ the approximation error

$$\sup_{\mathbf{v} \in \mathcal{I}_o(\mathbf{r})} \|\check{\mathfrak{y}}(\mathbf{v})\| \stackrel{\text{def}}{=} \sup_{\mathbf{v} \in \mathcal{I}_o(\mathbf{r})} \left\| \check{D}^{-1} \{ \check{\nabla} \mathcal{L}(\mathbf{v}) - \check{\nabla} \mathcal{L}(\mathbf{v}^*) + \check{D}^2 (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \} \right\|, \tag{1.5}$$

where $\check{D}^{-2} = \Pi_{\boldsymbol{\theta}} \mathcal{D}^{-2} \Pi_{\boldsymbol{\theta}}^T \in \mathbb{R}^{p \times p}$ with the full information matrix $\mathcal{D}^2 = -\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}^*)$ and where the projected gradient $\check{\nabla} \mathcal{L}$ is defined below in the next section. Section B.2.1 provides the following bound on a set of probability of at least $1 - e^{-\mathbf{x}}$:

$$\sup_{\mathbf{v} \in \mathcal{I}_o(\mathbf{r})} \|\check{\mathfrak{y}}(\mathbf{v})\| \leq \check{\diamond}(\mathbf{r}, \mathbf{x}),$$

where $\check{\diamond}(\mathbf{r}, \mathbf{x})$ is a small error. In combination with the deviation bound (1.4) and the identity $\nabla \mathcal{L}(\tilde{\mathbf{v}}) = 0$, this allows to derive the following Fisher and Wilks type expansions: with probability greater $1 - 2e^{-\mathbf{x}}$

$$\|\check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\| \leq \check{\diamond}(\mathbf{r}_0, \mathbf{x}), \tag{1.6}$$

$$|\check{L}(\tilde{\boldsymbol{\theta}}) - \check{L}(\boldsymbol{\theta}^*) - \|\check{\boldsymbol{\xi}}\|^2/2| \leq \mathbf{C} \sqrt{p + \mathbf{x}} \check{\diamond}(\mathbf{r}_0, \mathbf{x}). \tag{1.7}$$

In the case of correctly specified i.i.d models \check{D}^2 is the covariance matrix of the efficient influence function; see [18]. The random vector $\check{\boldsymbol{\xi}}$ satisfies $\mathbb{E} \check{\boldsymbol{\xi}} = 0$ and $\mathbb{E} \|\check{\boldsymbol{\xi}}\|^2 \asymp p$ and $\mathbf{C} > 0$ is a constant independent of $\mathbf{x} > 0$ and full dimension p^* . The precise definitions of the random p -vector $\check{\boldsymbol{\xi}}$ is also given in the next section. Moreover, general deviation bounds for quadratic forms from [29] apply to $\|\check{\boldsymbol{\xi}}\|^2$ (see Section A for details). In the case of a correct model specification the tails of $\|\check{\boldsymbol{\xi}}\|^2$ behave like those of a chi-square random variable with p degrees of freedom, and the result (1.6) can be viewed as an extension of the Wilks phenomenon. Under general identifiability conditions, the radius \mathbf{r}_0 can be fixed by $\mathbf{r}_0^2 = \mathbf{C}_1(p^* + \mathbf{x})$ for a fixed constant \mathbf{C}_1 to ensure the concentration property (1.4).

With this choice of \mathbf{r} , in the important i.i.d. case, the error term $\check{\Delta}(\mathbf{r}_0, \mathbf{x})$ can be bounded by $\mathfrak{C}(p^* + \mathbf{x})/\sqrt{n}$. The results (1.6) and (1.7) are nonasymptotic and hold true even under model misspecification.

It is important to grasp the implications of (1.6) and (1.7). The central contribution of this work is to bound the term in (1.5). It appears in this or a similar form also in the asymptotic approaches (see [22]) but is shown to be a zero sequence in the sample size under certain complexity and smoothness assumptions on the set of scores $\{\check{\nabla}\mathcal{L}(\mathbf{v}), \mathbf{v} \in \mathcal{Y}\}$. We manage to quantify for finite samples upper bounds for this term as functions of the radius \mathbf{r}_0 and the full dimension p^* . This allows for example to address the error when constructing confidence sets. For this assume that the quantiles of $\|\check{\xi}\|$ are available or that they can be given up to small error based on the Berry Esseen theorem (see [4]) or Edgeworth expansions (see [13]). Then (1.6) and (1.7) allow the construction of “approximate” confidence sets that address the finite sample error term (1.5), see Remark 2.13. The obtained sets are more conservative, i.e. larger than the asymptotic ones, but guarantee that the desired confidence level is attained. The possible error made when neglecting the error term (1.5) is illustrated in an example in Remark 2.20. Note however that on this level the contribution is rather theoretical: as in the case of the asymptotic results in [22], crucial objects as the matrix \check{D}^2 are unknown and would have to be estimated as well. An honest real data application of these results, where all model specific constants are unknown, is not possible yet and would be well beyond the scope of this work.

The proposed approach does not assume that the profile is consistent but gives conditions that ensure the right concentration behavior. Simply assuming that the profile is consistent can be even misleading in our setup because this would separate local and global considerations. This paper attempts to figure out a list of conditions ensuring global concentration and local expansion at the same time. This particularly allows to address the crucial question of the largest dimension of the nuisance parameter for which the Wilks and Fisher expansions still hold. In the smooth semiparametric problem with a fixed dimension of the target parameter, both Fisher and Wilks results apply up to an error $p^*/n^{1/2}$. In particular, we obtain that the error term in the Fisher expansion can be by a factor \sqrt{p} smaller than the similar error term in the Wilks Theorem. This ratio $p^*/n^{1/2}$ is the critical bound for the quality of the Fisher and Wilks expansions under the imposed conditions which is confirmed by a specific counter-example. It is of interest to compare our statements with the existing literature on the growing parameter asymptotics. We particularly mention [19, 20, 21] and a series of papers by S. Portnoy, see e.g. [26, 25, 27]. The typical dimensional asymptotic is $p^* = o(n^{1/2})$, which corresponds to our results. For some particular special problems and examples the condition on the parameter dimension can be relaxed to $p = o(n^{3/2})$; see [25]. However, the results are mainly limited to linear or generalized linear regression with independent observations and heavily use the model structure. To the contrary, our results apply in a rather general situation and deliver some useful information even in the case when the model is misspecified.

We begin by developing the results for the case that the full parameter space Υ is a subset of the Euclidean space of dimension $p^* \in \mathbb{N}$. In Section 2.6 we will exemplify how to extend our approach to the case when \mathbf{v} is a functional parameter using the so called sieve approach; see e.g. [28]. The present paper combines the sieve approximation idea and the finite sample Fisher and Wilks results under a possibly misspecified model.

The paper is organized as follows. Section 2.1 contains the conditions that we impose for the approach. Section 2.2 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 $\boldsymbol{\theta}^*$. Section 2.3 explains how to control the large deviations of $\tilde{\mathbf{v}}$ from (1.1) and how to improve the accuracy of the main results. Section 2.4 explains how the results translate to the case of i.i.d. samples and how the approach allows to obtain asymptotic efficiency of the profile estimator in this setting. Section 2.5 presents an example that shows that the ratio $p^{*2}/n \rightarrow 0$ is critical to obtain the Wilks phenomenon and the Fisher expansion on the class of models that satisfy the conditions of Section 2.1. Section 2.6 discusses how the results can be extended to the case with the infinite full dimension via the sieve approach. We present further conditions on the correlation structure of the full gradient $\nabla \mathcal{L}(\mathbf{v}^*) \in \mathcal{X}$ to also treat the bias. Section 2.9 briefly outlines how the approach can be employed to derive the main results in the context of single index modeling and which ratio of full dimension to sample size is sufficient in that context. The details of this section can be found in [1]. The appendix collects 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 $\check{L}(\tilde{\boldsymbol{\theta}}) - \check{L}(\boldsymbol{\theta}^*) \in \mathbb{R}$ and the Fisher expansion of the profile MLE $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^p$.

Most of the results are stated in a finite sample setup for just one fixed sample. As we are also interested in understanding what happens if *the full dimension* p^* becomes large we also consider a specification of the general finite sample results to an asymptotic setup with $p^* = p_n$, where n denotes the asymptotic parameter, e.g. the sample size with $n \rightarrow \infty$. Our results apply also if the target parameter $\boldsymbol{\theta} \in \mathbb{R}^p$ is also of growing dimension. The dimension p can be of order p^* . Even the case with a full dimensional target and low dimensional nuisance is included.

2.1. Conditions

This section collects the conditions imposed on the model. Let the full dimension of the problem be finite, i.e. $p^* < \infty$. Our conditions involve the symmetric positive definite information matrix $\mathcal{D}^2 \in \mathbb{R}^{p^* \times p^*}$ and a central point $\mathbf{v}^\circ \in \mathbb{R}^{p^*}$. In typical situations for $p^* < \infty$, one can set $\mathbf{v}^\circ = \mathbf{v}^*$ where \mathbf{v}^* is the “true

point" from (1.2). The matrix $\mathcal{D}^2 \in \mathbb{R}^{p^* \times p^*}$ can be defined as follows:

$$\mathcal{D}^2 = -\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}^\circ).$$

It is worth mentioning that $-\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}^\circ) = \text{Cov}(\nabla \mathcal{L}(\mathbf{v}^*))$ if the model $\mathbf{Y} \sim \mathbb{P}_{\mathbf{v}^*} \in (\mathbb{P}_{\mathbf{v}})$ is correctly specified and sufficiently regular; see e.g. [15].

Remark 2.1. This is not the only possible choice for \mathcal{D}^2 and \mathbf{v}° . Another candidate is to use $\mathcal{D}^2 = -\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}^*)$ while $\mathbf{v}^\circ = \mathbf{v}_m^*$, where \mathbf{v}_m^* approximates \mathbf{v}^* with growing $m \in \mathbb{N}$, as in [1]. In general there is no restriction for the choice of \mathcal{D}^2 , as long as the following list of conditions can be satisfied. The same holds for the matrix $\mathcal{V}^2 \in \mathbb{R}^{p^* \times p^*}$ that we introduce below.

In the context of semiparametric estimation, it is convenient to represent the information matrix in block form:

$$\mathcal{D}^2 = \begin{pmatrix} D^2 & A \\ A^\top & H^2 \end{pmatrix}.$$

First we state an *identifiability condition*.

(\mathcal{I}) It holds for some $\rho < 1$

$$\|H^{-1}A^\top D^{-1}\| \leq \rho.$$

Remark 2.2. The condition (\mathcal{I}) allows to define the important $p \times p$ efficient information matrix \check{D}^2 which is defined as the inverse of the $\boldsymbol{\theta}$ -block of the inverse of the full dimensional matrix \mathcal{D}^2 . The exact formula is given by

$$\check{D}^2 \stackrel{\text{def}}{=} (\Pi_{\boldsymbol{\theta}} \mathcal{D}^{-2} \Pi_{\boldsymbol{\theta}}^\top)^{-1} = D^2 - AH^{-2}A^\top,$$

and (\mathcal{I}) ensures that the matrix \check{D}^2 is well posed, see for instance [5], Chapter 2.4.

Using the matrix \mathcal{D}^2 and the central point $\mathbf{v}^\circ \in \mathbb{R}^{p^*}$, we define the local set $\mathcal{Y}_o(\mathbf{r}) \subset \mathcal{Y} \subseteq \mathbb{R}^{p^*}$ with some $\mathbf{r} \geq 0$:

$$\mathcal{Y}_o(\mathbf{r}) \stackrel{\text{def}}{=} \{\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathcal{Y} : \|\mathcal{D}(\mathbf{v} - \mathbf{v}^\circ)\| \leq \mathbf{r}\}. \quad (2.1)$$

Remark 2.3. For readers familiar with [29] we remark that the use of \mathcal{D} instead of \mathcal{V} in the above definition has no deeper reason but is a choice of convenience.

We introduce $\tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}$, which maximizes $\mathcal{L}(\mathbf{v}, \mathbf{v}^*)$ subject to $\Pi_{\boldsymbol{\theta}} \mathbf{v} = \boldsymbol{\theta}^*$:

$$\tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \stackrel{\text{def}}{=} (\boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*}) \stackrel{\text{def}}{=} \underset{\substack{\mathbf{v} \in \mathcal{Y} \\ \Pi_{\boldsymbol{\theta}} \mathbf{v} = \boldsymbol{\theta}^*}}{\text{argmax}} \mathcal{L}(\mathbf{v}, \mathbf{v}^*),$$

and remember the definition of the radius $\mathbf{r}_0 > 0$

$$\mathbf{r}_0(\mathbf{x}) \stackrel{\text{def}}{=} \inf_{\mathbf{r} > 0} \{\mathbb{P}(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}_o(\mathbf{r})) \geq 1 - e^{-\mathbf{x}}\}, \quad (2.2)$$

which we set to infinity if $\tilde{\mathbf{v}} = \{ \}$ or $\tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} = \{ \}$. Under the conditions $(\mathcal{L}\mathbf{r})$ and $(\mathcal{E}\mathbf{r})$ Theorem 2.3 in Section 2.3 states that $\mathbf{r}_0 = \mathbf{r}_0(\mathbf{x}) \approx \mathbf{c}\sqrt{\mathbf{x} + \mathbf{p}^*} > 0$.

Here and in what follows we implicitly assume that the log-likelihood function $\mathcal{L}(\mathbf{v}): \mathbb{R}^{p^*} \rightarrow \mathbb{R}$ is sufficiently smooth in $\mathbf{v} \in \mathbb{R}^{p^*}$, $\nabla \mathcal{L}(\mathbf{v}) \in \mathbb{R}^{p^*}$ stands for the gradient and $\nabla^2 \mathbb{E}\mathcal{L}(\mathbf{v}) \in \mathbb{R}^{p^* \times p^*}$ for the Hessian of the expectation $\mathbb{E}\mathcal{L} : \mathbb{R}^{p^*} \rightarrow \mathbb{R}$ at $\mathbf{v} \in \mathbb{R}^{p^*}$. By smooth enough we mean that all appearing derivatives exist and that we can interchange $\nabla \mathbb{E}\mathcal{L}(\mathbf{v}) = \mathbb{E}\nabla \mathcal{L}(\mathbf{v})$ on $\mathcal{Y}_\circ(\mathbf{r}_0)$, where $\mathbf{r}_0 > 0$ is defined in Equation (2.2) and $\mathcal{Y}_\circ(\mathbf{r})$ in equation (2.1). The following two conditions further quantify the smoothness properties on $\mathcal{Y}_\circ(\mathbf{r})$ of the expected log-likelihood $\mathbb{E}\mathcal{L}(\mathbf{v})$ and of the stochastic component $\zeta(\mathbf{v}) = \mathcal{L}(\mathbf{v}) - \mathbb{E}\mathcal{L}(\mathbf{v})$.

($\check{\mathcal{L}}_0$) For each $\mathbf{r} \leq 4\mathbf{r}_0$, there is a constant $\delta(\mathbf{r})$ such that it holds on the set $\mathcal{Y}_\circ(\mathbf{r})$:

$$\begin{aligned} \|D^{-1}D^2(\mathbf{v})D^{-1} - I_p\| &\leq \check{\delta}(\mathbf{r}), \\ \|D^{-1}(A(\mathbf{v}) - A)H^{-1}\| &\leq \check{\delta}(\mathbf{r}), \\ \|D^{-1}AH^{-1}(I_m - H^{-1}H^2(\mathbf{v})H^{-1})\| &\leq \check{\delta}(\mathbf{r}), \end{aligned}$$

where

$$\mathcal{D}(\mathbf{v})^2 \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}\mathcal{L}(\mathbf{v}), \quad \mathcal{D}(\mathbf{v}) = \begin{pmatrix} D^2(\mathbf{v}) & A(\mathbf{v}) \\ A^\top(\mathbf{v}) & H^2(\mathbf{v}) \end{pmatrix}.$$

Remark 2.4. This condition describes the local smoothness properties of the function $\mathbb{E}\mathcal{L}(\mathbf{v})$. In particular, it allows to bound the error of local linear approximation of the gradient $\check{\nabla}_\theta \mathbb{E}\mathcal{L}(\mathbf{v})$ where

$$\check{\nabla}_\theta = \nabla_\theta - AH^{-2}\nabla_\eta.$$

Under condition ($\check{\mathcal{L}}_0$) it follows from the second order Taylor expansion for any $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_\circ(\mathbf{r})$ (see Lemma B.1)

$$\|\check{D}^{-1}(\check{\nabla} \mathbb{E}\mathcal{L}(\mathbf{v}) - \check{\nabla} \mathbb{E}\mathcal{L}(\mathbf{v}^*)) + \check{D}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{c}\check{\delta}(\mathbf{r})\mathbf{r}. \tag{2.3}$$

In the proofs we actually only need the inequality (2.3) which in some cases can be weaker than ($\check{\mathcal{L}}_0$). For readers familiar with the classical theory (for instance [22]) we remark that ($\check{\mathcal{L}}_0$) is related to the condition that $t \mapsto l(t, \boldsymbol{\eta}_t(\boldsymbol{\theta}, \boldsymbol{\eta}))$ is twice continuously differentiable where $\check{\delta}(\mathbf{r})$ quantifies how smooth the second derivative is. We impose such a qualified smoothness in order to give finite sample deviation bounds as a function of the radius of the local set $\mathcal{Y}_\circ(\mathbf{r})$.

The next condition concerns the regularity of the stochastic component $\zeta(\mathbf{v}) \stackrel{\text{def}}{=} \mathcal{L}(\mathbf{v}) - \mathbb{E}\mathcal{L}(\mathbf{v})$.

($\check{\mathcal{E}}\mathcal{D}_1$) $\zeta(\mathbf{v}) \rightarrow \zeta(\mathbf{v}')$ as $\mathbf{v} \rightarrow \mathbf{v}'$. Further for all $0 < \mathbf{r} < 4\mathbf{r}_0$, there exists a constant $\omega \leq 1/2$ such that for all $|\mu| \leq \check{\mathfrak{g}}$ and $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_\circ(\mathbf{r})$

$$\sup_{\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_\circ(\mathbf{r})} \sup_{\|\boldsymbol{\gamma}\| \leq 1} \log \mathbb{E} \exp \left\{ \frac{\mu \boldsymbol{\gamma}^\top \check{D}^{-1} \{ \check{\nabla}_{\boldsymbol{\theta}} \zeta(\mathbf{v}) - \check{\nabla}_{\boldsymbol{\theta}} \zeta(\mathbf{v}') \}}{\check{\omega} \|\mathcal{D}(\mathbf{v} - \mathbf{v}')\|} \right\} \leq \frac{\check{\nu}_1^2 \mu^2}{2}.$$

Remark 2.5. The above condition is strongly related to the assumption of Donsker- and and Glivenko-Cantelli properties in [22] in order to ensure that the error in the local linear approximation of $\check{\nabla} \mathcal{L}(\mathbf{v}) - \check{\nabla} \mathcal{L}(\mathbf{v}^*)$ disappears. We replace these conditions with the more specific assumption ($\check{\mathcal{E}}\mathcal{D}_1$), which in combination with the entropy of $\mathcal{Y}_\circ(\mathbf{r})$ yields the desired error bounds. Note that in linear models or regressions with bounded regressors this condition is automatically satisfied. In the single index example this condition becomes a condition on the smoothness of the employed basis functions $\mathbf{e}_k : \mathbb{R} \rightarrow \mathbb{R}$ and a sub exponential moment bound on the additive noise $\varepsilon \in \mathbb{R}$, see condition (**Cond** $_\varepsilon$) in Section 2.9.

The above conditions suffice for our main results. But we include another condition that allows to control the deviation behavior of $\|\check{D}^{-1} \check{\nabla} \zeta(\mathbf{v}^*)\|$.

($\check{\mathcal{E}}\mathcal{D}_0$) There exist a matrix $\check{V}^2 \in \mathbb{R}^{p \times p}$, constants $\nu_0 > 0$ and $\check{\mathfrak{g}} > 0$ such that for all $|\mu| \leq \check{\mathfrak{g}}$

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \mu \frac{\langle \check{\nabla}_{\boldsymbol{\theta}} \zeta(\mathbf{v}^\circ), \boldsymbol{\gamma} \rangle}{\|\check{V} \boldsymbol{\gamma}\|} \right\} \leq \frac{\check{\nu}_0^2 \mu^2}{2}.$$

Remark 2.6. One possible and natural choice for the matrices $\check{V}^2 \in \mathbb{R}^{p \times p}$ and $\mathcal{V}^2 \in \mathbb{R}^{p^* \times p^*}$ (see ($\mathcal{E}\mathcal{D}_0$) below) is

$$\mathcal{V}^2 \stackrel{\text{def}}{=} \text{Var}\{\nabla \mathcal{L}(\mathbf{v}^\circ)\}, \quad \check{V}^2 = \text{Cov}(\check{\nabla}_{\boldsymbol{\theta}} \zeta(\mathbf{v}^\circ)),$$

but also other matrices could be used as long as ($\check{\mathcal{E}}\mathcal{D}_0$) or ($\mathcal{E}\mathcal{D}_0$) can be satisfied.

In many situations the following, stronger conditions, are easier to check and allow a further improvement of the results of Theorem 2.2 with the help of Proposition 2.4:

(\mathcal{L}_0) For each $\mathbf{r} \leq \mathbf{r}_0$, there is a constant $\delta(\mathbf{r})$ such that it holds on the set $\mathcal{Y}_\circ(\mathbf{r})$:

$$\|\mathcal{D}^{-1} \{ \nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}) \} \mathcal{D}^{-1} - \mathbb{I}_{p^*}\| \leq \delta(\mathbf{r}).$$

($\mathcal{E}\mathcal{D}_1$) There exists a constant $\omega \leq 1/2$, such that for all $|\mu| \leq \mathfrak{g}$ and all $0 < \mathbf{r} < \mathbf{r}_0$

$$\sup_{\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_\circ(\mathbf{r})} \sup_{\|\boldsymbol{\gamma}\|=1} \log \mathbb{E} \exp \left\{ \frac{\mu \boldsymbol{\gamma}^\top \mathcal{D}^{-1} \{ \nabla \zeta(\mathbf{v}) - \nabla \zeta(\mathbf{v}') \}}{\omega \|\mathcal{D}(\mathbf{v} - \mathbf{v}')\|} \right\} \leq \frac{\nu_1^2 \mu^2}{2}.$$

($\mathcal{E}\mathcal{D}_0$) There exist a matrix $\mathcal{V}^2 \in \mathbb{R}^{p^* \times p^*}$, constants $\nu_0 > 0$ and $\mathbf{g} > 0$ such that for all $|\mu| \leq \mathbf{g}$

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

The following lemma shows, that these conditions imply the weaker ones from above:

Lemma 2.1. *Assume (I). Then $(\mathcal{E}\mathcal{D}_1)$ implies $(\check{\mathcal{E}}\mathcal{D}_1)$, (\mathcal{L}_0) implies $(\check{\mathcal{L}}_0)$, and $(\mathcal{E}\mathcal{D}_0)$ implies $(\check{\mathcal{E}}\mathcal{D}_0)$ with*

$$\check{\mathbf{g}} = \frac{\sqrt{1 - \rho^2}}{(1 + \rho)\sqrt{1 + \rho^2}} \mathbf{g}, \check{\nu}_i = \frac{(1 + \rho)\sqrt{1 + \rho^2}}{\sqrt{1 - \rho^2}} \nu_i, \check{\delta}(\mathbf{r}) = \delta(\mathbf{r}), \text{ and } \check{\omega} = \omega.$$

Remark 2.7. Note that with $(\check{\mathcal{L}}_0)$, $(\check{\mathcal{E}}\mathcal{D}_0)$ and $(\check{\mathcal{E}}\mathcal{D}_1)$ the smoothness and moment conditions do not have to be satisfied for the full gradient $\nabla \mathcal{L}(\cdot)$ but only for the projected one $(\nabla_{\boldsymbol{\theta}} + A H^{-1} \nabla_{\boldsymbol{\eta}}) \mathcal{L}(\cdot)$. This can make a tremendous difference to (\mathcal{L}_0) , $(\mathcal{E}\mathcal{D}_0)$ and $(\mathcal{E}\mathcal{D}_1)$ if $A(\cdot) \in \mathbb{R}^{p \times m}$ is small while $\nabla_{\boldsymbol{\eta}} \mathcal{L}(\cdot)$ is rather rough or possesses bad moment properties. In that case $(\mathcal{E}\mathcal{D}_0)$ and $(\mathcal{E}\mathcal{D}_1)$ might not be satisfied or $\check{\delta}(\mathbf{r})$, $\check{\omega}$ and $\check{\nu}_1$ would be considerably smaller than their counterparts $\delta(\mathbf{r})$, ω and ν_1 . This is particularly obvious if $A(\cdot) \equiv 0$.

Finally we present two conditions that allow a specific approach to determine a radius $\mathbf{r}_0(\mathbf{x}) > 0$ such that $\mathbb{P}(\tilde{\mathbf{v}} \in \mathcal{Y}(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}$ (see Section 2.3). These conditions have to be satisfied on the whole set $\mathcal{Y} \subseteq \mathbb{R}^{p^*}$. Note, however, that the conditions $(\mathcal{L}\mathbf{r})$ and $(\mathcal{E}\mathbf{r})$ can be substituted with any other set of conditions that allow to determine a value \mathbf{r}_0 ensuring $\mathbb{P}(\tilde{\mathbf{v}} \in \mathcal{Y}(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}$.

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

$$\frac{-\mathbb{E}\mathcal{L}(\mathbf{v}, \mathbf{v}^\circ)}{\|\mathcal{D}(\mathbf{v} - \mathbf{v}^\circ)\|^2} \geq \mathbf{b}(\mathbf{r}), \quad \mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r}).$$

($\mathcal{E}\mathbf{r}$) For any $\mathbf{r} \geq \mathbf{r}_0$ there exists a constant $\mathbf{g}(\mathbf{r}) > 0$ such that

$$\sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \sup_{\mu \leq \mathbf{g}(\mathbf{r})} \sup_{\gamma \in \mathbb{R}^{p^*}} \log \mathbb{E} \exp \left\{ \mu \frac{\langle \nabla \zeta(\mathbf{v}), \gamma \rangle}{\|\mathcal{D}\gamma\|} \right\} \leq \frac{\nu_{\mathbf{r}}^2 \mu^2}{2}.$$

Remark 2.8. These two conditions serve a qualified apriori concentration result for the full estimator $\tilde{\mathbf{v}}$, of the type $\mathbb{P}\{\tilde{\mathbf{v}} \in \mathcal{Y}_\circ(\mathbf{r}_0(\mathbf{x}))\} \geq 1 - e^{-\mathbf{x}}$. Condition $(\mathcal{L}\mathbf{r})$ is satisfied for many estimators that employ some least square functional as we do for the single index model in Section 2.9. In a more general setting it could be combined with yet another even rougher a priori consistency result $\mathbb{P}(\tilde{\mathbf{v}} \in U(\mathbf{v}^*))$ for some open neighborhood $U(\mathbf{v}^*) \subset \mathcal{Y}$. Then $(\mathcal{L}\mathbf{r})$ is automatically satisfied as smooth functions are quadratic around their maximum, in this case $\mathbb{E}\mathcal{L}$ around \mathbf{v}^* . Further the condition can be relaxed to $-\mathbb{E}\mathcal{L}(\mathbf{v}, \mathbf{v}^\circ)$ growing with

super linear speed in the distance $\|\mathcal{D}(\boldsymbol{v} - \boldsymbol{v}^*)\|$, see Theorem 2.1 in [30]. In this case the calculations become technically more involved which is why we focus on $(\mathcal{L}\mathbf{r})$ for the sake of readability. $(\mathcal{E}\mathbf{r})$ is a global exponential moment condition and ensures that the norm of the stochastic component $\nabla\zeta(\boldsymbol{v}) \in \mathbb{R}^{p^*}$ is bounded with high probability. For example in the least square setting with additive noise this is satisfied with $\mathbf{g}(\mathbf{r}) = \infty$ if the additive noise is sub Gaussian.

Remark 2.9. We briefly comment how restrictive the imposed conditions are. Our conditions on the regularity and smoothness of the log-likelihood process $\mathcal{L}(\boldsymbol{v})$ in terms of the second or even third derivative are stronger than usually required; cf. Chapters 1, 2 in [15]. But we aim not only for vanishing approximation error terms but for expressions that reveal the interplay of full dimension, smoothness of the functional \mathcal{L} and moments of the score. A quantification seems unavoidable of “how much smoother than twice differentiable” the function $\mathbb{E}\mathcal{L}(\cdot)$ is (i.e. condition (\mathcal{L}_0)), and of “how much smoother than once differentiable and well bounded in exponential moment terms” is $\nabla\zeta(\cdot)$ (i.e. condition $(\mathcal{E}\mathcal{D}_1)$). Note further, that we do not require that $\mathcal{L}(\boldsymbol{v})$ is the true log-likelihood. It comes from a parametric family chosen by a statistician. For typical examples, such a family possesses the required regularity. In particular, [29], Section 5.1, considered in details the i.i.d. case and presented some mild sufficient conditions on the parametric family which imply the above general conditions.

Concerning moments the conditions $(\check{\mathcal{E}}\mathcal{D}_1)$, $(\check{\mathcal{E}}\mathcal{D}_0)$, $(\mathcal{E}\mathcal{D}_1)$, $(\mathcal{E}\mathcal{D}_0)$ and $(\mathcal{E}\mathbf{r})$ require sub exponential moments of the observations (errors). Usually one only assumes finite second or third moments of the errors; cf. [15], Chapter 2. Our condition is a bit more restrictive but it allows to obtain finite sample bounds of the kind that with some small $\epsilon > 0$

$$\mathbb{P} \left\{ \|\check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\| \geq \epsilon(p^* + \mathbf{x}) \right\} \geq e^{-\mathbf{x}},$$

i.e. the bounds depend linearly on the exponent \mathbf{x} . Without comparable moment bounds these results do not seem to be attainable in such a general setting. Consider for instance the simple model

$$y = \sqrt{v^*} + \varepsilon \in \mathbb{R}, \quad \tilde{v} = \operatorname{argmax}_{v \in \mathbb{R}} (y - \sqrt{v})^2 / 2,$$

with $v^* \neq 0$, $\sqrt{x} \stackrel{\text{def}}{=} \operatorname{sign}(x)\sqrt{|x|}$, $\mathbb{E}\varepsilon = 0$ and $\operatorname{Cov}(\varepsilon) = 1$. Then up to the exponential moments all conditions from above are met with $\check{D}^2 = \mathcal{D}^2 = \frac{1}{4v^*}$ and $\check{\boldsymbol{\xi}} = \varepsilon$. We find

$$|\check{D}(\tilde{v} - v^*) - \check{\boldsymbol{\xi}}| = \left| \frac{1}{2\sqrt{v^*}} (y^2 - v^*) - \varepsilon \right| = \left| \left(\frac{2\sqrt{v^*} + \varepsilon}{2\sqrt{v^*}} - 1 \right) \varepsilon \right| = \frac{\varepsilon^2}{2\sqrt{v^*}}.$$

Now if $\log \mathbb{E}[\exp(\lambda\varepsilon)] < \lambda^2/2$ we can derive

$$\mathbb{P}(|\check{D}(\tilde{v} - v^*) - \check{\boldsymbol{\xi}}| \geq 8\sqrt{v^*\mathbf{x}}) \leq e^{-\mathbf{x}},$$

while obviously without comparable moment criteria such a result – a linear relation between the exponent on the right hand side and the bound on the left hand side – could not be attained.

To list some settings in which the conditions can be satisfied we name the regression and generalized regression models; cf. [10, 11] or [17]. [29], Section 5.2, argued that $(\mathcal{E}D_1)$ is automatically fulfilled for a generalized linear model, while $(\mathcal{E}D_0)$ requires that regression errors have to fulfill some exponential moments conditions. If this condition is too restrictive and a more stable (robust) estimation procedure is desirable, one can apply the LAD-type contrast leading to median regression. [29], Section 5.3, showed for the case of linear median regression that all the required conditions are fulfilled automatically if the sample size n exceeds Cp^* for a fixed constant C . [31] applied this approach for local polynomial quantile regression. [34] applied the approach to the problem of regression with Gaussian process where the unknown parameters enter in the likelihood in a rather complicated way. Further in this work we show how to satisfy them in a general i.i.d. setting and in the single index model, see Sections 2.4 and 2.9.

Remark 2.10. Another indication that the conditions are not too strong is served by an example in [2], where the error term $\check{\diamond}$ in our main result 2.2 is increased by a factor $\sqrt{p^*}$ if the condition (\mathcal{L}_0) is slightly relaxed to read

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

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

which appears in [29] under the label (\mathcal{L}_0) .

2.2. Wilks and Fisher expansions

This section states the main results in a finite dimensional framework.

First we introduce the main elements of the approach. Let the *information matrix* $\mathcal{D}^2 \in \mathbb{R}^{p^* \times p^*}$ be from the condition in Section 2.1, For the semiparametric $(\boldsymbol{\theta}, \boldsymbol{\eta})$ -setup, we consider the block representation of the vector $\nabla \stackrel{\text{def}}{=} \nabla \mathcal{L}(\mathbf{v}^*)$ and of the matrix \mathcal{D}^2

$$\nabla = \begin{pmatrix} \nabla_{\boldsymbol{\theta}} \\ \nabla_{\boldsymbol{\eta}} \end{pmatrix}, \quad \mathcal{D}^2 = \begin{pmatrix} D^2 & A \\ A^\top & H^2 \end{pmatrix}.$$

We repeat also the definition of the $p \times p$ matrix \check{D}^2

$$\check{D}^2 = D^2 - AH^{-2}A^\top,$$

and p -vectors $\check{\nabla}_{\boldsymbol{\theta}}$ and $\check{\boldsymbol{\xi}} \in \mathbb{R}^p$

$$\check{\nabla}_{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \nabla_{\boldsymbol{\theta}} \zeta(\mathbf{v}^*) - AH^{-2} \nabla_{\boldsymbol{\eta}} \zeta(\mathbf{v}^*), \quad \check{\boldsymbol{\xi}} \stackrel{\text{def}}{=} \check{D}^{-1} \check{\nabla}_{\boldsymbol{\theta}}.$$

The random variable $\check{\nabla}_{\theta} \in \mathbb{R}^p$ is related to the efficient influence function in semiparametric estimation and the matrix $\check{D}^2 \in \mathbb{R}^{p \times p}$ equals its covariance in the case of correct specification.

Remark 2.11. It seems worthy to point out that $\check{D}^{-2}\check{\nabla}_{\theta} = \Pi_{\theta}\mathcal{D}^{-2}\nabla$, see again [5], Chapter 2.4.

Define the *semiparametric spread* $\check{\diamond}(\mathbf{r}, \mathbf{x}) > 0$ as

$$\check{\diamond}(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} 4 \left(\frac{4}{(1-\rho^2)^2} \check{\delta}(4\mathbf{r}) + 6\nu_1 \check{\omega} \mathfrak{z}(\mathbf{x}, 2p^* + 2p) \right) \mathbf{r}, \quad (2.4)$$

where $\check{\delta}(\mathbf{r})$ is shown in the condition $(\check{\mathcal{L}}_0)$ and the constants $\check{\omega}$, ν_1 are from condition $(\check{\mathcal{E}}\mathcal{D}_1)$ in Section 2.1. The value $\mathfrak{z}(\mathbf{x}, 2p^* + 2p)$ is related to the entropy of the unit ball in a \mathbb{R}^{p^*+p} -dimensional Euclidean space with $\mathbf{g} > 0$ from $(\check{\mathcal{E}}\mathcal{D}_1)$ it is defined as

$$\mathfrak{z}(\mathbf{x}, \mathbb{Q}) \stackrel{\text{def}}{=} \begin{cases} \sqrt{2(\mathbf{x} + \mathbb{Q})} & \text{if } \sqrt{2(\mathbf{x} + \mathbb{Q})} \leq \mathbf{g}, \\ \mathbf{g}^{-1}(\mathbf{x} + \mathbb{Q}) + \mathbf{g}/2 & \text{otherwise,} \end{cases} \quad (2.5)$$

and one can apply $\mathfrak{z}(\mathbf{x}, p^*) \cong \sqrt{\mathbf{x} + p^*}$ for moderate choice of $\mathbf{x} > 0$; see Appendix C. The *semiparametric spread* $\check{\diamond}(\mathbf{r}, \mathbf{x})$ measures the quality of a linear approximation to $\check{\nabla}\mathcal{L}(\mathbf{v}) - \check{\nabla}\mathcal{L}(\mathbf{v}^*)$ in the local vicinity the local vicinity $\mathcal{Y}_o(\mathbf{r}) = \{\mathbf{v} \in \mathcal{Y} : \|\mathcal{D}(\mathbf{v} - \mathbf{v}^o)\| \leq \mathbf{r}\}$. Our results become accurate if $\check{\diamond}(\mathbf{r}_0, \mathbf{x})$ is small. The spread will be evaluated in the i.i.d. case in Section 2.4 below.

Theorem 2.2. Assume $(\check{\mathcal{E}}\mathcal{D}_1)$, $(\check{\mathcal{L}}_0)$, and (\mathcal{I}) with a central point $\mathbf{v}^o = \mathbf{v}^*$ and some matrix \mathcal{D}^2 . Further assume that the sets of maximizers $\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*}$ are not empty. Then it holds on a set $\Omega(\mathbf{x}) \subseteq \Omega$ of probability greater $1 - 2e^{-\mathbf{x}}$ for the profile MLE $\tilde{\theta}$ from (1.3)

$$\|\check{D}(\tilde{\theta} - \theta^*) - \check{\xi}\| \leq \check{\diamond}(\mathbf{r}_0, \mathbf{x}), \quad (2.6)$$

$$|2\check{L}(\tilde{\theta}, \theta^*) - \|\check{\xi}\|^2| \leq 4 \left(\|\check{\xi}\| + \check{\diamond}(\mathbf{r}_0, \mathbf{x}) \right) \check{\diamond}(\mathbf{r}_0, \mathbf{x}) + \check{\diamond}(\mathbf{r}_0, \mathbf{x})^2, \quad (2.7)$$

where the spread $\check{\diamond}(\mathbf{r}_0, \mathbf{x})$ is defined in (2.4) and where $\mathbf{r}_0 > 0$ is defined in (2.2).

Remark 2.12. The Wilks expansion claims that the profile maximum likelihood $\check{L}(\tilde{\theta}, \theta^*) \stackrel{\text{def}}{=} \check{L}(\tilde{\theta}) - \check{L}(\theta^*)$ can be approximated by a quadratic form $\|\check{\xi}\|^2/2$ with $\check{\xi} = \check{D}^{-1}\check{\nabla}_{\theta}$. In the correctly specified i.i.d setting the vector $\check{\xi}$ is asymptotically standard normal and the quadratic form $\|\check{\xi}\|^2 = \|\check{D}^{-1}\check{\nabla}_{\theta}\|^2$ weakly converges to a chi-square random variable with $p \in \mathbb{N}$ degrees of freedom, which follows from the central limit theorem and the fact that then $\text{Cov}(\check{\xi}) = I_p$. In the general case, the behavior of the quadratic form $\|\check{\xi}\|^2$ depends on the characteristics of the matrix $\check{B} \stackrel{\text{def}}{=} \check{D}^{-1}\check{V}^2\check{D}^{-1}$ where $\check{V}^2 \in \mathbb{R}^{p \times p}$ is from $(\check{\mathcal{E}}\mathcal{D}_0)$

and in many cases equals $\check{V}^2 = \text{Cov}(\check{\nabla}_{\boldsymbol{\theta}})$. More precisely one can find an upper quantile function $\mathfrak{z}(\mathbf{x}, \check{B})$ of this quadratic form ensuring

$$\mathbb{P}(\|\check{\boldsymbol{\xi}}\| > \mathfrak{z}(\mathbf{x}, \check{B})) \leq 2e^{-\mathbf{x}};$$

see Proposition A.1. One can use the bound $\mathfrak{z}^2(\mathbf{x}, \check{B}) \leq \mathfrak{C}(p + \mathbf{x})$ in most situations. We call $\check{B} \in \mathbb{R}^{p \times p}$ *semiparametric misspecification matrix* as it is related to the misspecification matrix introduced in [14]. \check{B} is equal to the identity matrix if a correctly specified log likelihood is used.

Remark 2.13. One can use the expansion (2.6) for the construction of elliptic confidence sets

$$\mathcal{E}(\mathfrak{z}) = \{\boldsymbol{\theta} : \|\check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})\| \leq \mathfrak{z}\};$$

for some $\mathfrak{z}(\mathbf{x}) > 0$. More precisely let $q_\alpha > 0$ be the α -level quantile of $\|\check{\boldsymbol{\xi}}\|$. Then we find with the triangular inequality and (2.6)

$$\begin{aligned} \mathbb{P}\left\{\boldsymbol{\theta}^* \notin \mathcal{E}\left(q_\alpha + \check{\diamond}(\mathbf{r}_0, \mathbf{x})\right)\right\} &= \mathbb{P}\left\{\|\check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \geq q_\alpha + \check{\diamond}(\mathbf{r}_0, \mathbf{x})\right\} \\ &\leq \mathbb{P}\left\{\|\check{\boldsymbol{\xi}}\| \geq q_\alpha\right\} + 2e^{-\mathbf{x}} = 1 - (\alpha - 2e^{-\mathbf{x}}), \end{aligned}$$

and

$$\begin{aligned} \mathbb{P}\left\{\boldsymbol{\theta}^* \in \mathcal{E}\left(q_\alpha - \check{\diamond}(\mathbf{r}_0, \mathbf{x})\right)\right\} &= \mathbb{P}\left\{\|\check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq q_\alpha - \check{\diamond}(\mathbf{r}_0, \mathbf{x})\right\} \\ &\leq \mathbb{P}\left\{\|\check{\boldsymbol{\xi}}\| \leq q_\alpha\right\} + 2e^{-\mathbf{x}} = \alpha + 2e^{-\mathbf{x}}. \end{aligned}$$

So up to $\check{\diamond}(\mathbf{r}_0, \mathbf{x})$ and $2e^{-\mathbf{x}}$ the set $\mathcal{E}(q_\alpha)$ serves as a confidence set. The choice of \mathbf{x} determines the trade off between the closeness of $q_\alpha \pm \check{\diamond}(\mathbf{r}_0, \mathbf{x})$ to q_α and the probability level $\alpha + 2e^{-\mathbf{x}}$ to α .

Remark 2.14. The profile maximum likelihood process $\check{L}(\boldsymbol{\theta})$ can be used for defining the likelihood-based confidence sets of the form

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

The bound (2.7) helps to evaluate the coverage probability $\mathbb{P}(\boldsymbol{\theta}^* \notin \mathcal{E}(\mathfrak{z}))$ in terms of deviation probability for the quadratic form $\|\check{\boldsymbol{\xi}}\|^2$ and in term; cf. Corollary 3.2 in [29].

Remark 2.15. In the classical finite dimensional case, a usual choice for the central point \mathbf{v}° is $\mathbf{v}^\circ = \mathbf{v}^* = \text{argmax}_{\mathbf{v} \in \mathcal{Y}} \mathbb{E}\mathcal{L}(\mathbf{v})$ and one can define the matrix \mathcal{D}^2 as $\mathcal{D}^2 = -\nabla^2 \mathbb{E}\mathcal{L}(\mathbf{v}^*)$. However, for the sieve semiparametric problem in Section 2.6, we use another definition related to the infinite dimensional model.

2.3. Large deviation bounds

In this section we want to present a way to determine a value $\mathbf{r}_0 > 0$ such that the full MLE $\tilde{\mathbf{v}} \in \mathbb{R}^{p^*}$ belongs to the local vicinity $\mathcal{Y}_o(\mathbf{r}_0) \subset \mathbb{R}^{p^*}$ with high probability. As a first step we adopt the upper function approach from [29]; cf. Theorem 4.2 therein. It is important to note that Corollary 4.4 is one particular approach which could be replaced by any other proper technique. For instance, in the model with i.i.d. observations, Theorem 5.3 of [15] might serve as a tool. The required conditions can be substantially weakened to upper and lower bounds on the Hellinger distance between models for distinct parameters. We follow the general way of [29] because it allows to address possible model misspecification and finite samples.

A close look at the proof of Theorem 4.2 of [29] shows that it actually yields the following modified version:

Theorem 2.3 ([29], Theorem 4.2). *Suppose $(\mathcal{E}\mathbf{r})$ and $(\mathcal{L}\mathbf{r})$ with $\mathbf{b}(\mathbf{r}) \equiv \mathbf{b}$. Further define the following random set*

$$\mathcal{Y}(K) \stackrel{\text{def}}{=} \{\mathbf{v} \in \mathcal{Y} : \mathcal{L}(\mathbf{v}, \mathbf{v}^*) \geq -K\}.$$

If for a fixed \mathbf{r}_0 and any $\mathbf{r} \geq \mathbf{r}_0$, the following conditions are fulfilled:

$$\begin{aligned} 1 + \sqrt{\mathbf{x} + 2p^*} &\leq 3\nu_{\mathbf{r}}^2 \mathbf{g}(\mathbf{r})/\mathbf{b}, \\ 6\nu_{\mathbf{r}} \sqrt{\mathbf{x} + 2p^* + \frac{\mathbf{b}}{9\nu_{\mathbf{r}}^2} K} &\leq \mathbf{r}\mathbf{b}, \end{aligned} \quad (2.8)$$

then

$$\mathbb{P}(\mathcal{Y}(K) \subseteq \mathcal{Y}_o(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}. \quad (2.9)$$

Remark 2.16. Note that this Theorem also ensures that the maximum of $\mathcal{L} : \mathbb{R}^{p^*} \rightarrow \mathbb{R}$ is actually attained. Clearly $\mathbf{v}^* \in \mathcal{Y}(0)$ such that it is nonempty. Further

$$\mathbb{P}(\mathcal{Y}(0) \subseteq \mathcal{Y}_o(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}},$$

such that $\mathcal{Y}(0) \subseteq \mathcal{Y}_o(\mathbf{r}_0) \subset \mathbb{R}^{p^*}$ is compact and thus \mathcal{L} attains its maximum on $\mathcal{Y}(0)$, which will be the global maximum $\tilde{\mathbf{v}}$. The same holds for $\tilde{\mathbf{v}}_{\theta^*} \in \mathbb{R}^{p^*}$.

Remark 2.17. The condition (2.8) helps to understand which $\mathbf{r}_0 > 0$ ensures prescribed concentration properties of $\tilde{\mathbf{v}} \in \mathbb{R}^{p^*}$ and $\tilde{\mathbf{v}}_{\theta^*} \in \mathbb{R}^{p^*}$ because by definition both are in the set $\mathcal{Y}(0)$. Consequently, if $\mathbf{g}(\mathbf{r}) > 0$ is large enough, (2.8) follows from the bound

$$\mathbf{r}_0 \geq 6\mathbf{b}^{-1}\nu_{\mathbf{r}}\sqrt{\mathbf{x} + p^*}. \quad (2.10)$$

The upper function approach in Theorem 2.3 of showing the consistency for an M-estimator can be rather rough and the bound (2.10) could lead to quite large values of $\mathbf{r}_0 > 0$. As the obtained value $\mathbf{r}_0 > 0$ enters into the error term $\check{\diamond}(\mathbf{r}_0, \mathbf{x}) > 0$ of Theorem 2.2 it is desirable to obtain a general refined bound for $\mathbf{r}_1 \leq \mathbf{r}_0$ that still ensures that $\mathbb{P}(\tilde{\mathbf{v}} \in \mathcal{Y}_o(\mathbf{r}_1)) \geq 1 - \mathbf{C}e^{-\mathbf{x}}$ with a small constant $\mathbf{C} > 0$. Such an improvement is possible as the following proposition shows. Define the *parametric spread*:

$$\diamond(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \{\delta(\mathbf{r}) + 6\nu_1 \mathfrak{z}(\mathbf{x}, 4p^*)\omega\} \mathbf{r}, \tag{2.11}$$

where $\mathfrak{z}(\mathbf{x}, \mathbb{Q})$ is defined in (2.5). Further with $\mathcal{V}^2 \in \mathbb{R}^{p^* \times p^*}$ from condition $(\mathcal{E}\mathcal{D}_0)$ introduce the *misspecification matrix* $B \in \mathbb{R}^{p^* \times p^*}$ given by the famous sandwich formula; see [14]:

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

In the case of correct model specification with $\mathcal{D}^2 = \mathcal{V}^2$, the *sandwich matrix* B becomes the identity: $B = \mathbb{I}_{p^*}$. Theorem A.1 tells us that

$$\mathbb{P} \{ \|\mathcal{D}^{-1}\nabla\mathcal{L}(\mathbf{v}^*)\| \geq \mathfrak{z}(\mathbf{x}, B) \} \leq 2e^{-\mathbf{x}},$$

where $\mathfrak{z}(\mathbf{x}, B) \leq \mathbf{C}\sqrt{\text{tr}(B^2) + \mathbf{x}}$ for moderate choice of $\mathbf{x} > 0$, see (A.2).

Proposition 2.4. *Assume the conditions of Theorem 2.2 and additionally assume $(\mathcal{E}\mathcal{D}_1)$, (\mathcal{L}_0) and $(\mathcal{E}\mathcal{D}_0)$ with $\mathcal{V}^2 \in \mathbb{R}^{p^* \times p^*}$. Let $\mathbf{r}_0 > 0$ be such that (2.9) holds and define the radius*

$$\mathbf{r}_1 \stackrel{\text{def}}{=} \mathfrak{z}(\mathbf{x}, B) + \diamond(\mathbf{r}_0, \mathbf{x}) \wedge \mathbf{r}_0 \leq \mathbf{r}_0.$$

Then the result of Theorem 2.2 applies with the error term $\check{\diamond}(\mathbf{r}_1, \mathbf{x})$ in place of $\check{\diamond}(\mathbf{r}_0, \mathbf{x})$ and with probability greater $1 - 5e^{-\mathbf{x}}$.

2.4. The i.i.d. case

In this section we want to illustrate the results for the case of a smooth i.i.d. model. This means that given i.i.d. $(\mathbf{Y}_1, \dots, \mathbf{Y}_n) \in \otimes_{i=1}^n \mathcal{Y}$ we use

$$\mathcal{L}(\mathbb{Y}, \mathbf{v}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{Y}_i, \mathbf{v}), \quad \mathbb{E}_{\mathbb{P}}\mathcal{L}(\mathbf{v}) = \mathbb{E}_{\mathbb{P}_{\mathbf{Y}}}\ell(\mathbf{Y}_1, \mathbf{v}),$$

where $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a suitable functional. As above we omit the data in the following and write $\ell_i(\mathbf{v}) \stackrel{\text{def}}{=} \ell(\mathbf{Y}_i, \mathbf{v})$. Note that

$$\mathbf{v}^* \stackrel{\text{def}}{=} \underset{\mathbf{v} \in \mathcal{Y}}{\text{argmax}} \mathbb{E}\mathcal{L}(\mathbf{v}) = \underset{\mathbf{v} \in \mathcal{Y}}{\text{argmax}} \mathbb{E}\ell(\mathbf{v}),$$

$$\begin{aligned} \mathcal{D}^2 &\stackrel{\text{def}}{=} \nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}^*) = n d^2 \stackrel{\text{def}}{=} n \nabla^2 \mathbb{E} \ell(\mathbf{v}^*), \\ \mathcal{V}^2 &\stackrel{\text{def}}{=} \text{Cov}(\nabla \zeta(\mathbf{v}^*)) = n v^2 \stackrel{\text{def}}{=} n \text{Cov}(\nabla(\ell - \mathbb{E} \ell)(\mathbf{v}^*)). \end{aligned}$$

To check the conditions of section 2.1 in principle we only have to assume that they are met with \mathcal{L}, \mathcal{D} replaced by ℓ, d with some $\nu_0^*, \omega_1^*, \delta(\mathbf{r}) = \delta^* \mathbf{r}, \mathbf{b}(\mathbf{r}) = \mathbf{b}^*$ and $\mathbf{g} = \mathbf{g}_1$. Under these conditions, one can easily check the conditions in Section 2.1 for the full log-likelihood $\mathcal{L}(\mathbf{v}) = \sum_{i=1}^n \ell(y_i, \mathbf{v})$ with $\omega = \omega_1 n^{-1/2}$, $\delta(\mathbf{r}) = \delta^* \mathbf{r} n^{-1/2}$, $\mathbf{b}(\mathbf{r}) = \mathbf{b}^*$, and $\mathbf{g} = \mathbf{g}_1 n^{1/2}$; cf. Lemma 5.1 in [29]. To gain a bit more intuition let us consider the following stronger sufficient list of assumptions:

- (ℓ_0) The matrix valued function $\nabla^2 \mathbb{E}[\ell(\cdot)] : \mathcal{Y} \rightarrow \mathbb{R}^{p^* \times p^*}$ is locally Lipschitz continuous with Lipschitz constant δ^* in an open neighborhood $U \ni \mathbf{v}^*$.
- (ed_1) There are constants $\nu_0^*, \mathbf{g}^* > 0$ and an open neighborhood $U \ni \mathbf{v}^*$ such that for all $\mathbf{v} \in U$ the random matrix valued function $\nabla^2(\ell - \mathbb{E} \ell)(\cdot, \mathbf{Y}) \mathcal{Y} \rightarrow \mathbb{R}^{p^* \times p^*}$ satisfies for all $|\lambda| \leq \mathbf{g}^*$

$$\sup_{\substack{\gamma_1, \gamma_2 \in \mathbb{R}^{p^*} \\ \|\gamma_1\| = \|\gamma_2\| = 1}} \log \mathbb{E} \exp \left\{ \lambda \gamma_1^\top d^{-1} \nabla^2(\ell - \mathbb{E} \ell)(\mathbf{v}) d^{-1} \gamma_2 \right\} \leq \nu_0^* \lambda^2 / 2.$$

- (ed_0) The random vector valued function $\nabla(\ell - \mathbb{E} \ell)(\cdot, \mathbf{Y}) \mathcal{Y} \rightarrow \mathbb{R}^{p^* \times p^*}$ satisfies for all $|\lambda| \leq \mathbf{g}^*$ and all $\mathbf{v} \in \mathcal{Y}$

$$\sup_{\substack{\gamma \in \mathbb{R}^{p^*} \\ \|\gamma\| = 1}} \log \mathbb{E} \exp \left\{ \lambda \gamma^\top d^{-1} \nabla(\ell - \mathbb{E} \ell)(\mathbf{v}) \right\} \leq \nu_0^* \lambda^2 / 2.$$

- (ℓ_r) There is a constant $\mathbf{b}^* > 0$ such that

$$\mathbb{E} [\ell(\mathbf{v}) - \ell(\mathbf{v}^*)] \geq \mathbf{b}^* \|d(\mathbf{v} - \mathbf{v}^*)\|^2.$$

- (ι) There is a constant $c_d > 0$ such that the matrix $d^2 \stackrel{\text{def}}{=} \nabla^2 \mathbb{E} \ell(\mathbf{v}^*)$ satisfies $\gamma^\top d^2 \gamma \geq c_d \|\gamma\|^2$ for all $\gamma \in \mathbb{R}^{p^*}$.

Lemma 2.5. Assume that $n \in \mathbb{N}$ is large enough to ensure that the local neighborhood $U \subset \mathcal{Y}$ of \mathbf{v}^* from conditions (ℓ_0) and (ed_1) satisfies

$$\begin{aligned} \mathcal{Y}_\circ(\mathbf{r}^*) &\stackrel{\text{def}}{=} \{ \mathbf{v} \in \mathcal{Y} : \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}^* \} \\ &= \frac{1}{\sqrt{n}} \{ \mathbf{v} \in \mathcal{Y} : \|d(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}^* \} \subseteq U. \end{aligned}$$

Then the conditions (ℓ_0), (ed_1), (ed_0), (ℓ_r) and (ι) imply (\mathcal{L}_0), ($\mathcal{E} \mathcal{D}_1$), ($\mathcal{E} \mathcal{D}_0$), ($\mathcal{E} \mathcal{D}_r$), (\mathcal{L}_0) and (\mathcal{I}) with $\delta(\mathbf{r}) = \frac{\delta^*}{\sqrt{nc_d^*}} \mathbf{r}$, $\omega = \frac{1}{\sqrt{n}}$, $\mathbf{g} = \sqrt{n} \mathbf{g}^*$, $\nu_1 = \nu_0 = \nu_0^*$, $\mathbf{g}(\mathbf{r}) = \sqrt{n} \mathbf{g}^*$, $\mathbf{b} = \mathbf{b}^*$ for all $\mathbf{r} \leq \mathbf{r}^*$. Further $\rho^2 \geq 1 - \frac{c_d}{\|d_\theta^2\| \|\mathbf{v}\| h^2}$ where $d_\theta^2 = \Pi_\theta^\top d \Pi_\theta \in \mathbb{R}^{p \times p}$ and $h^2 = \Pi_\eta^\top d \Pi_\eta \in \mathbb{R}^{m \times m}$.

Remark 2.18. To keep things simple we do not elaborate on how to check $(\check{\mathcal{L}}_0)$, $(\check{\mathcal{E}}\mathcal{D}_1)$, $(\check{\mathcal{E}}\mathcal{D}_0)$ but refer to Lemma 2.1.

Noting that $\mathcal{L}(\tilde{\mathbf{v}}, \mathbf{v}^*) \geq 0$ and $\mathcal{L}(\tilde{\mathbf{v}}_{\theta^*}, \mathbf{v}^*) \geq 0$ Theorem 2.3 yields that

$$\mathbb{P}(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \Upsilon_{\circ}(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}, \quad \text{with } \mathbf{r}_0(\mathbf{x}) = 6 \frac{\nu_0^*}{\mathbf{b}^*} \sqrt{2p^* + \mathbf{x}}.$$

Theorem 2.2 applies with $\mathcal{D}^2 = n\nabla^2 \mathbb{E} \ell(\mathbf{v}^*)$ and $\mathbf{v}^{\circ} = \mathbf{v}^*$. We immediately obtain the following result.

Corollary 2.6. Let Y_1, \dots, Y_n be i.i.d. and let the conditions (ℓ_0) , (ed_1) , (ed_0) , $(\ell_{\mathbf{r}})$ and (ι) be met. Assume that $\mathbf{r}_0(\mathbf{x}) = 6 \frac{\nu_0^*}{\mathbf{b}^*} \sqrt{2p^* + \mathbf{x}} \leq \mathbf{r}^*$. Then we get the Fisher and Wilks results of Theorem 2.2 for $\mathbf{x} \ll \sqrt{n} \mathbf{g}^*$ with

$$\check{\diamond}(\mathbf{r}_0, \mathbf{x}) \leq \frac{36\nu_0^*}{\sqrt{n}\mathbf{b}^*} \left(\frac{4}{(1-\rho^2)^2} \frac{\delta^* \nu_0}{c_d^3 \mathbf{b}^*} (\mathbf{x} + 2p^*) + \nu_0 \mathfrak{z}(\mathbf{x}, 2p^* + 2p) \sqrt{\mathbf{x} + 2p^*} \right).$$

Remark 2.19. The definition of $\mathfrak{z}(\mathbf{x}, 2p^* + 2p)$ in (C.2) implies for moderate values of $\mathbf{x} > 0$ that

$$\check{\diamond}(\mathbf{r}_0, \mathbf{x}) \leq \mathbf{C}_{\diamond}(\mathbf{x} + p^*)/\sqrt{n},$$

with some fixed constant \mathbf{C}_{\diamond} . The Fisher result (2.6) is meaningful if $\check{\diamond}(\mathbf{r}_0, \mathbf{x})$ is small yielding the constraint $p^* \ll n^{1/2}$. If the target dimension p is fixed, the same condition is sufficient for the Wilks expansion in (2.7). However, if the target dimension p is of order p^* , the constraint for the Wilks theorem becomes $p^* = o(n^{1/3})$. See [2] for an example that shows, that this difference actually occurs in certain examples.

2.5. Critical dimension

This section discusses the issue of *critical parameter dimensions* when the full dimension p^* grows with the sample size n . We write $p^* = p_n$. The results of Theorem 2.2 refined by Proposition 2.4 are accurate if the spread function $\diamond(\mathbf{r}, \mathbf{x})$ from (2.11) fulfills $\diamond(\mathbf{r}_0, \mathbf{x}) \leq \mathfrak{z}(\mathbf{x}, B)$ and $\check{\diamond}(\mathbf{r}_1, \mathbf{x})$ is small, with $\mathbf{r}_1 = 2\mathfrak{z}(\mathbf{x}, B)$. Usually $\mathfrak{z}(\mathbf{x}, B) \leq \mathbf{C}\sqrt{\mathbf{x} + p^*}$ leading to

$$\check{\diamond}(\mathbf{r}_1, \mathbf{x}) \asymp \check{\delta}(\mathbf{r}_1)\mathbf{r}_1 + \check{\omega}\mathbf{r}_1^2 \quad \text{is small for } \mathbf{r}_1^2 \asymp p^*. \tag{2.12}$$

The critical size of p^* then depends on the exact bounds for $\check{\delta}(\cdot), \check{\omega}$. If $\check{\delta}(\mathbf{r})/\mathbf{r} \asymp \check{\omega} \asymp 1/\sqrt{n}$ (as in Corollary 2.6) the condition (2.12) reads “ $\check{\diamond}(\mathbf{r}_1, \mathbf{x}) \asymp p^*/\sqrt{n}$ is small”. This means that one needs that “ p^{*2}/n is small” to obtain an accurate non asymptotic version of the Wilks phenomenon and the Fisher Theorem. Similar conclusions were obtained by Portnoy in series of papers on growing dimension in generalized linear models and for natural exponential families, see

e.g. [26, 25, 27]. Our results are non-asymptotic and apply to general statistical models under the conditions of Section 2.1. The following example shows that the constraint “ p^*/n is small” is critical.

Consider single observation model

$$\mathbf{Y} = f(\mathbf{v}) + \varepsilon,$$

$$f(\mathbf{v}) = f(\theta, \boldsymbol{\eta}) \stackrel{\text{def}}{=} \begin{pmatrix} \theta \\ \boldsymbol{\eta}_1 \\ \vdots \\ \boldsymbol{\eta}_{p_n-1} \end{pmatrix} + \begin{pmatrix} \|\boldsymbol{\eta}\|^2 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^{p_n},$$

with $\varepsilon \sim \mathcal{N}(0, \frac{1}{n}I_{p_n})$ and $\mathbf{v} = (\theta, \boldsymbol{\eta}) \in \mathbb{R} \times \mathbb{R}^{p_n-1}$. This model is equivalent to the i.i.d. observations in the same model with the errors $\varepsilon_i \sim \mathcal{N}(0, I_{p_n})$. Assume that the parameter of interest is $\theta \in \mathbb{R}$ and that the true point satisfies $\mathbf{v}^* = 0 \in \mathbb{R}^{p^*}$.

Proposition 2.7. *Under $p_n/\sqrt{n} \rightarrow 0$, the Fisher expansion is accurate and the profile MLE asymptotically standard normal. If $p_n/\sqrt{n} \not\rightarrow 0$ the profile MLE in the above model is not root- n consistent. For $\sqrt{n} = o(p_n)$ the root- n bias tends to infinity almost surely. Finally, the Wilks phenomenon occurs if and only if $p_n = o(\sqrt{n})$.*

Remark 2.20. The above example can also be used to illustrate the difference between a finite sample approach and using asymptotic normality for the construction of confidence sets. For fixed dimension the profile MLE is asymptotically standard normal, i.e. with $q_\alpha > 0$ denoting the α -level quantile of a chi-square distribution with one degree of freedom

$$\mathbb{P}\left(\theta^* \in \left\{|\tilde{\theta} - \theta|^2 \leq q_\alpha/n\right\}\right) \rightarrow \alpha. \quad (2.13)$$

But the proof of Proposition 2.7 gives

$$|\tilde{\theta} - \theta^*| = |\varepsilon_\theta - \|\boldsymbol{\varepsilon}_\boldsymbol{\eta}\|^2|,$$

where $n\|\boldsymbol{\varepsilon}_\boldsymbol{\eta}\|^2 \sim \chi_{p_n-1}^2$ and $\varepsilon_\theta \sim \mathcal{N}(0, 1/n)$. It is known that the median of a chi-square distribution converges to its number of degrees of freedom when the degrees of freedom tend to infinity. This means that for any $0 < \epsilon < 1$ the set

$$C \stackrel{\text{def}}{=} \{n\|\boldsymbol{\varepsilon}_\boldsymbol{\eta}\|^2 \geq (1 - \epsilon)p_n\},$$

is of probability greater 1/2 for $n, p_n \in \mathbb{N}$ large enough. Let $f_{\chi_{p_n-1}^2} : [0, \infty) \rightarrow \mathbb{R}$ denote the Lebesgue density of a $\chi_{p_n-1}^2$ random variable. We can use the independence of $\|\boldsymbol{\varepsilon}_\boldsymbol{\eta}\|$ and ε_θ and Fubini's Theorem to estimate

$$\mathbb{P}\left(\theta^* \in \left\{|\tilde{\theta} - \theta|^2 \leq q_\alpha/n\right\}\right) = \int_0^\infty \mathbb{P}\left(|\varepsilon_\theta - z/n|^2 \leq q_\alpha/n\right) f_{\chi_{p_n-1}^2}(z) dz$$

$$\begin{aligned}
 &= \int_0^\infty \left[\Phi \left(\frac{z}{\sqrt{n}} + \sqrt{q_\alpha} \right) - \Phi \left(\frac{z}{\sqrt{n}} - \sqrt{q_\alpha} \right) \right] f_{\chi_m^2}(z) dz \\
 &< \frac{1}{2} \left[\alpha + \Phi \left(\frac{(1-\epsilon)p_n}{\sqrt{n}} + \sqrt{q_\alpha} \right) - \Phi \left((1-\epsilon) \frac{p_n}{\sqrt{n}} - \sqrt{q_\alpha} \right) \right],
 \end{aligned}$$

where $\Phi : \mathbb{R} \rightarrow [0, 1]$ denotes the distribution function of a standard normal random variable. If p_n/\sqrt{n} is significantly larger than 0, the value

$$\Phi \left(\frac{(1-\epsilon)p_n}{\sqrt{n}} + \sqrt{q_\alpha} \right) - \Phi \left(\frac{(1-\epsilon)p_n}{\sqrt{n}} - \sqrt{q_\alpha} \right),$$

is distinctively smaller α . For example for $\alpha = 0.95$ and $(1-\epsilon)p_n/\sqrt{n} = 11/12$ we get

$$\mathbb{P} \left(\theta^* \in \left\{ |\tilde{\theta} - \theta|^2 \leq q_{0.95}/n \right\} \right) < 0.9.$$

In other words the asymptotic confidence statement in (2.13) is way off in the finite sample case because the error term in the local linear approximation is not addressed. This is exactly where a large full dimension has an impact on the behavior of the estimator. Our results in Theorem 2.2 quantify the size of these terms for a large set of models and give a guideline for how to correct confidence sets to address this effect. The price are more conservative sets, but their coverage property is ensured.

Remark 2.21. There is an interesting connection of the condition $p^*/\sqrt{n} \rightarrow 0$ with the general theory on semiparametric M-estimators. In the common asymptotic approach to semiparametric M estimators one assumes a priori consistency of the estimator $\tilde{\nu} = (\tilde{\theta}, \tilde{\eta})$, more precisely in the case that the functional $\check{\nabla} \mathcal{L}(\cdot)$ is smooth enough one assumes that $\|\tilde{\theta} - \theta^*\| = o_{\mathbb{P}}(1)$ and $\|\tilde{\eta} - \eta^*\| = O_{\mathbb{P}}(n^{-1/4})$, see [18] Section 21.1.4. On the other hand the results of Theorem 2.2 are accurate if $\check{\diamond}(\mathbf{r}_0, \mathbf{x})$ is small. As explained above this means in the i.i.d. setting that $\check{\diamond}(\mathbf{r}_0, \mathbf{x}) = o(1)$. Neglecting the contribution of $\|\tilde{\theta} - \theta^*\|$ to \mathbf{r}_0 this can be ensured if

$$\check{\diamond}(\mathbf{r}_0, \mathbf{x}) \leq \mathbf{C}(p^* + \mathbf{r}_0^2)/\sqrt{n} \leq o(1) + \mathbf{C}\sqrt{n}\|\tilde{\eta} - \eta^*\|^2 \rightarrow 0,$$

i.e. if $\|\tilde{\eta} - \eta^*\| = o(n^{-1/4})$. But consider the radius $\mathbf{r}_1 > 0$ from Proposition 2.4. It is of order $\sqrt{p^* + m}$ if $\check{\diamond}(\mathbf{r}_0) = O(\sqrt{p^* + \mathbf{x}})$. In that case in the i.i.d. setting the constraint on the a priori deviation bound becomes $\check{\diamond}(\mathbf{r}_0, \mathbf{x}) = O(\sqrt{p^* + \mathbf{x}})$ which can be ensured if

$$\check{\diamond}(\mathbf{r}_0, \mathbf{x}) \leq \mathbf{C}(p^* + \mathbf{r}_0^2)/\sqrt{n} \leq o(1) + \mathbf{C}\sqrt{n}\|\tilde{\eta} - \eta^*\|^2 = O(\sqrt{p^* + \mathbf{x}}),$$

which means if $p^* + \mathbf{x} = o(\sqrt{n})$ that $\|\tilde{\eta} - \eta^*\| = o(n^{-1/8})$, which is a considerably weaker constraint. These bounds only concern the finite dimensional case. In the infinite dimensional setting, treated in Section 2.8 we have to impose conditions

that ensure that the bias induced by the sieve approach is small enough. [3] serve such conditions for the Hilbertspace setting. One of these conditions reads that $\|H(\boldsymbol{\eta}^* - \Pi_m \boldsymbol{\eta}^*)\|^2 \leq Cm$, i.e. the true nuisance component $\boldsymbol{\eta}^* \in \mathcal{X}$ is well approximated by its projection into the span of the first $m \in \mathbb{N}$ basis elements $(\mathbf{e}_k) \subset \mathcal{X}$. If we represent with some $\alpha > 0$

$$\boldsymbol{\eta}^* = \sum_{k=1}^m \eta_k^* \mathbf{e}_k, \quad \sum_{k=1}^m \eta_k^{*2} k^{2\alpha} < \infty,$$

we obtain the constraint $m \leq n^{1/(2\alpha+1)}$, which means that we need $\alpha > 1/2$ to get $m = o(n^{1/2})$ and in that case $\boldsymbol{\eta}^* \in \mathcal{X}$ is nonparametrically estimable with rate $n^{-1/4}$.

Remark 2.22. Concerning the difference between the critical dimensions in the Wilks (2.7) and the Fisher expansion (2.6) we remark that [2] presents an example where $p = p^*/2$, where all conditions of Section 2.1 are met with $\check{\delta}(\mathbf{r})/\mathbf{r} \asymp \check{\omega} \asymp 1/\sqrt{n}$ and where the Wilks phenomenon occurs iff $p^{*3}/n \rightarrow 0$ while for the Fisher expansion p^{*2}/n suffices.

2.6. Infinite dimensional nuisance

This section discusses how the approach can be extended to the infinite dimensional case. First the basic idea of projecting the infinite dimensional problem down to a finite dimensional one is explained. Then we prove under bias constraints that the projected *sieve* estimator is nearly normal and efficient. To avoid further technical distractions (or obstacles) we present the case of a separable Hilbert space. The ideas can be modified to treat the case when the nuisance parameter belongs to a Banach space.

2.7. Sieve approach

Consider the $(\boldsymbol{\theta}, \mathbf{f})$ -setup with $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^p$ and $\mathbf{f} \in \mathcal{X}$, where \mathcal{X} is an infinite dimensional separable Hilbert space. The target parameter $\boldsymbol{\theta}^*$ can be defined as

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \sup_{\mathbf{f} \in \mathcal{X}} \mathbb{E} \mathcal{L}(\boldsymbol{\theta}, \mathbf{f}). \quad (2.14)$$

As the Hilbert space \mathcal{X} is assumed to be separable it possesses a countable orthonormal basis $\{\mathbf{e}_1, \mathbf{e}_2, \dots\} \subset \mathcal{X}$. Any vector $\mathbf{f} \in \mathcal{X}$ admits a unique decomposition in the form

$$\mathbf{f} = \sum_{j=1}^{\infty} \eta_j \mathbf{e}_j,$$

where $\eta_j = \langle \mathbf{f}, \mathbf{e}_j \rangle$ is the usual Fourier coefficient. In the *sieve* approach one assumes that for any $m \in \mathbb{N}$ a finite set $\mathbf{e}_1, \dots, \mathbf{e}_m$ of elements in \mathcal{X} is fixed

and the vector \mathbf{f} can be approximated by a finite linear combination $\mathbf{f}_m(\boldsymbol{\eta})$ of the \mathbf{e}_j 's:

$$\mathbf{f}_m(\boldsymbol{\eta}) \stackrel{\text{def}}{=} \sum_{j=1}^m \eta_j \mathbf{e}_j.$$

We denote the parameter by $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^p \times l^2$. In the following we will need to quantify the accuracy of approximating \mathbf{f} by \mathbf{f}_m as m grows; see condition (**bias**) below.

Let $\mathcal{L}(\boldsymbol{\theta}, \mathbf{f})$ be the log-likelihood in the original model. Define by abuse of notation

$$\begin{aligned} \mathcal{L}(\mathbf{v}) &\stackrel{\text{def}}{=} \mathcal{L}\left(\boldsymbol{\theta}, \sum_{j=1}^{\infty} \eta_j \mathbf{e}_j\right) \\ \mathbf{v}^* &\stackrel{\text{def}}{=} \operatorname{argmax}_{(\boldsymbol{\theta}, \boldsymbol{\eta}) \in l^2} \mathbb{E} \left[\mathcal{L}\left(\boldsymbol{\theta}, \sum_{k=1}^{\infty} \eta_k \mathbf{e}_k\right) \right], \end{aligned}$$

and the m -dimensional sieve approximation $\mathcal{L}_m(\mathbf{v})$ of $\mathcal{L}(\mathbf{v})$ by

$$\begin{aligned} \mathcal{L}_m(\boldsymbol{\theta}, \boldsymbol{\eta}) &\stackrel{\text{def}}{=} \mathcal{L}(\boldsymbol{\theta}, \mathbf{f}_m(\boldsymbol{\eta})), \\ (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathcal{Y}_m &\stackrel{\text{def}}{=} \{\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^{p^*} : (\boldsymbol{\theta}, \mathbf{f}_m(\boldsymbol{\eta})) \in \mathcal{Y}\}. \end{aligned}$$

The corresponding sieve profile estimator $\tilde{\boldsymbol{\theta}}_m$ and its target $\boldsymbol{\theta}_m^*$ for this parametric m -submodel are defined in the usual way:

$$\begin{aligned} \tilde{\boldsymbol{\theta}}_m &\stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \tilde{\mathbf{v}}_m \stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}_m} \mathcal{L}_m(\boldsymbol{\theta}, \boldsymbol{\eta}), \\ \boldsymbol{\theta}_m^* &\stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \mathbf{v}_m^* \stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}_m} \mathbb{E} \mathcal{L}_m(\boldsymbol{\theta}, \boldsymbol{\eta}). \end{aligned} \tag{2.15}$$

The question we are interested in can be formulated as follows: is $\tilde{\boldsymbol{\theta}}_m$ a good (efficient) estimator of $\boldsymbol{\theta}^*$ from (2.14) under a proper choice of m ?

2.8. Bias constraints and efficiency

The parametric results obtained in Section 2 claim that $\tilde{\boldsymbol{\theta}}_m \in \mathbb{R}^p$ estimates well $\boldsymbol{\theta}_m^* \in \mathbb{R}^p$ if the spread $\check{\diamond}(\mathbf{r}_0, \mathbf{x}) > 0$ is small. More precisely we have the following: Define for fixed $\mathbf{x} > 0$ the value $\mathbf{r}_0 > 0$ by

$$\begin{aligned} \mathbf{r}_0(\mathbf{x}) &\stackrel{\text{def}}{=} \inf_{\mathbf{r} \geq 0} \{\mathbb{P}\{\tilde{\mathbf{v}}_m, \tilde{\mathbf{v}}_{\boldsymbol{\theta}_m^*, m} \in \mathcal{Y}_{0,m}(\mathbf{r})\} \geq 1 - e^{-\mathbf{x}}\}, \\ \mathcal{Y}_{0,m}(\mathbf{r}) &\stackrel{\text{def}}{=} \{\mathbf{v} \in \mathcal{Y}_m, \|\mathcal{D}_m(\mathbf{v} - \mathbf{v}_m^*)\| \leq \mathbf{r}\}. \end{aligned}$$

where $\mathbf{v}_m^* = (\boldsymbol{\theta}_m^*, \boldsymbol{\eta}_m^*) = \operatorname{argmax}_{\mathbf{v}} \mathbb{E} \mathcal{L}_m(\mathbf{v})$ and

$$\tilde{\mathbf{v}}_{\boldsymbol{\theta}, m} \stackrel{\text{def}}{=} \operatorname{argmax}_{\substack{\mathbf{v} \in \mathcal{Y}_m \\ \Pi_{\boldsymbol{\theta}} \mathbf{v} = \boldsymbol{\theta}}} \mathcal{L}_m(\mathbf{v}, \mathbf{v}^*).$$

Further the matrix \check{D}_m^2 is defined as

$$\check{D}_m^2(\mathbf{v}_m^*) \stackrel{\text{def}}{=} (\Pi_{\boldsymbol{\theta}} \mathcal{D}_m^{-2} \Pi_{\boldsymbol{\theta}}^\top)^{-1} \in \mathbb{R}^{p \times p}, \quad \mathcal{D}_m^2 \stackrel{\text{def}}{=} \nabla_{p+m}^2 \mathbb{E}[\mathcal{L}(\mathbf{v}_m^*)] \in \mathbb{R}^{p^* \times p^*},$$

i.e. the derivatives of $\mathbb{E}[\mathcal{L}]$ are only taken with respect to the first $p + m \in \mathbb{N}$ coordinates of $\mathbf{v} \in l^2$ and the Hessian is evaluated in $\mathbf{v}_m^* \in \mathbb{R}^{p^*}$. Applying Theorem 2.2 to $\tilde{\boldsymbol{\theta}}_m$ from (2.15) we find that with probability greater $1 - 2e^{-x}$

$$\|\check{D}_m(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*) - \check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| \leq \check{\diamond}(\mathbf{r}_0, \mathbf{x}), \quad (2.16)$$

The result (2.16) involves two kinds of bias, one that concerns the difference $\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*$ and the other the difference between $\check{D}_m \in \mathbb{R}^{p \times p}$ and $\check{D} \in \mathbb{R}^{p \times p}$ where

$$\check{D}^2 \stackrel{\text{def}}{=} (\Pi_{\boldsymbol{\theta}} \nabla^2 \mathbb{E}[\mathcal{L}(\mathbf{v}^*)]^{-1} \Pi_{\boldsymbol{\theta}}^\top)^{-1} \in \mathbb{R}^{p \times p},$$

i.e. the derivatives of $\mathbb{E}[\mathcal{L}]$ are taken with respect to all coordinates of $\mathbf{v} \in l^2$ and the Hessian is calculated in the “true point” $\mathbf{v}^* \in l^2$. The second bias – i.e. bounds for $\|I - \check{D}_m^{-1}(\mathbf{v}_m^*) \check{D}^2(\mathbf{v}^*) \check{D}_m^{-1}(\mathbf{v}_m^*)\|$ – will be neglected for now, as only the operator $\check{D}_m^2(\mathbf{v}_m^*) \in \mathbb{R}^{p \times p}$ is available in practice. We will come back to it, when we derive efficiency for the sieve profile estimator $\tilde{\boldsymbol{\theta}}_m \in \mathbb{R}^p$.

Remark 2.23. To be more precise we assume that $\mathbb{E} \mathcal{L} : \mathcal{Y} \rightarrow \mathbb{R}$ is Fréchet differentiable and that each element of the gradient $\langle \nabla \mathbb{E} \mathcal{L}, \mathbf{e}_k \rangle$ again is Fréchet differentiable aswell. We denote the resulting operator by $\mathcal{D}^2 = \nabla^2 \mathbb{E}[\mathcal{L}(\mathbf{v}^*)] : \overline{\operatorname{span}} \mathcal{Y} \rightarrow \overline{\operatorname{span}} \mathcal{Y}$.

For the first type of bias we impose the following condition:

(bias) There exists a function $\alpha : \mathbb{N} \rightarrow \mathbb{R}_+$ such that

$$\|\check{D}_m(\mathbf{v}_m^*)(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*)\| \leq \alpha(m), \quad \alpha(m) \rightarrow 0, \quad \text{as } m \rightarrow \infty.$$

Remark 2.24. For now we focus on the result 2.2 and thus we do not elaborate on approximation theory. But [3] presents conditions on the structure of $\mathcal{D} : l^2 \rightarrow l^2$ and on the sequence $\boldsymbol{\eta}^* \in l^2$ that yield **(bias)**.

We represent

$$\mathcal{D}_m^2(\mathbf{v}_m^*) = \begin{pmatrix} D^2(\mathbf{v}_m^*) & A_m^\top(\mathbf{v}_m^*) \\ A_m(\mathbf{v}_m^*) & H_m^2(\mathbf{v}_m^*) \end{pmatrix} \in \mathbb{R}^{(p+m) \times (p+m)}.$$

With Theorem 2.2 and **(bias)** we directly get the following corollary:

Corollary 2.8. Assume (bias) and that the conditions $(\check{\mathcal{D}})$, $(\check{\mathcal{D}}_1)$ and $(\check{\mathcal{L}}_0)$ from Section 2.1 are satisfied for all $m \geq m_0$ for some $m_0 \in \mathbb{N}$ and with $\mathcal{D}^2 = \nabla_{p+m}^2 \mathbb{E} \mathcal{L}_m(\mathbf{v}_m^*) \in \mathbb{R}^{p^* \times p^*}$, $\mathcal{V}^2 = \text{Cov}[\nabla_{p+m} \mathcal{L}_m(\mathbf{v}_m^*)] \in \mathbb{R}^{p^* \times p^*}$ and $\mathbf{v}^\circ = \mathbf{v}_m^* \in \mathbb{R}^{p^*}$. Assume that $\tilde{\mathbf{v}}_m \neq \{\}$ and $\tilde{\mathbf{v}}_{\theta_m^*} \neq \{\}$. Choose $\mathbf{r}_0(\mathbf{x}) > 0$ such that $\mathbb{P}(\tilde{\mathbf{v}}_m, \tilde{\mathbf{v}}_{\theta_m^*, m} \in \mathcal{Y}_{0,m}(\mathbf{r}_0(\mathbf{x}))) \geq 1 - e^{-x}$. Then it holds for any $m \geq m_0$ with probability greater $1 - 2e^{-x}$

$$\|\check{D}_m(\mathbf{v}_m^*)(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| \leq \check{\diamond}(\mathbf{r}_0, \mathbf{x}) + \alpha(m),$$

where

$$\check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*) \stackrel{\text{def}}{=} \check{D}_m^{-1}(\nabla_{\boldsymbol{\theta}} - A_m H_m^{-1} \nabla_{\boldsymbol{\eta}}) \mathcal{L}_m(\mathbf{v}_m^*).$$

Define

$$\check{L}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \max_{\boldsymbol{\eta} \in \mathbb{R}^m} \mathcal{L}_m(\boldsymbol{\theta}, \boldsymbol{\eta}),$$

where it is important to note that the maximization is restricted to the finite dimensional space \mathbb{R}^m . As above abbreviate $\check{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} \check{L}(\boldsymbol{\theta}) - \check{L}(\boldsymbol{\theta}^*)$. For the bias in the Wilks result a bit more work is needed. We can show the following:

Theorem 2.9. Assume the same as in Corollary 2.8. Pick a radius $0 < \mathbf{r}_0^\circ$ such that

$$\mathbb{P}(\{\tilde{\mathbf{v}}_m, \tilde{\mathbf{v}}_{\theta_m^*, m}, \tilde{\mathbf{v}}_{\theta^*, m} \in \mathcal{Y}_{0,m}(\mathbf{r}_0^\circ)\}) > 1 - e^{-x},$$

Then we get with probability greater $1 - 2e^{-x}$

$$\begin{aligned} & |2\check{L}(\tilde{\boldsymbol{\theta}}_m, \boldsymbol{\theta}^*) - \|\check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\|^2| \\ & \leq 8 \left(\|\check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| + \check{\diamond}(\mathbf{r}_0^\circ, \mathbf{x}) \right) \check{\diamond}(2(1 + \rho)\mathbf{r}_0^\circ, \mathbf{x}) + \check{\diamond}(\mathbf{r}_0^\circ, \mathbf{x})^2 \\ & \quad + \alpha(m) \left(2\|\check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| + \alpha(m) + 2\check{\diamond}(2(1 + \rho)\mathbf{r}_0^\circ, \mathbf{x}) \right) \end{aligned}$$

Remark 2.25. With condition $(\check{\mathcal{D}}_0)$ we can use Theorem A.1 to obtain

$$\mathbb{P} \left(\|\check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| \geq \mathfrak{z}(\mathbf{x}, \check{B}) \right) \leq e^{-x}.$$

Remark 2.26. The radius $\mathbf{r}_0^\circ \in \mathbb{R}$ can be determined again using the tools of Section 2.3. Clearly Theorem 2.3 can be applied to find some $\mathbf{r}_0 \leq \mathbf{r}_0^\circ$ such that

$$\mathbb{P}(\tilde{\mathbf{v}}_m, \tilde{\mathbf{v}}_{\theta_m^*, m} \in \mathcal{Y}_{0,m}(\mathbf{r}_0)) > 1 - e^{-x}.$$

Further note that by the mean value theorem

$$\begin{aligned} \mathcal{L}_m(\boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\theta^*, m}) - \mathcal{L}_m(\mathbf{v}_m^*) & \geq \mathcal{L}_m(\boldsymbol{\theta}^*, \boldsymbol{\eta}_m^*) - \mathcal{L}_m(\mathbf{v}_m^*) \\ & \geq -(1 + \rho)\alpha(m) \sup_{\mathbf{v} \in \mathcal{Y}_\circ((1+\rho)\alpha(m))} \|D^{-1} \nabla_{\boldsymbol{\theta}} \mathcal{L}_m(\mathbf{v})\|. \end{aligned}$$

With condition $(\check{\mathcal{D}}_0)$ and $(\check{\mathcal{D}}_1)$ the right hand side can be bounded by some constant $-\alpha(m)\mathcal{C}(p^* + \mathbf{x}) \in \mathbb{R}$ with probability greater $1 - 2e^{-\mathbf{x}}$ using the tools of Section 2.3. Combining this with Theorem 2.3 gives that with

$$\mathbf{r}_0^\circ = 6\mathbf{b}^{-1}\nu_{\mathbf{r}}\sqrt{\mathbf{x} + \log(4) + p^* + \frac{\mathbf{b}}{9\nu_{\mathbf{r}}^2}\alpha(m)\mathcal{C}(p^* + \mathbf{x})},$$

it holds that

$$\begin{aligned} \mathbb{P}(\{\tilde{\mathbf{v}}_m, \tilde{\mathbf{v}}_{\theta^*, m} \in \mathcal{Y}_{0,m}(\mathbf{r}_0)\} \cap \{\tilde{\mathbf{v}}_{\theta^*, m} \in \mathcal{Y}_{0,m}(\mathbf{r}_0^\circ)\}) &> 1 - 4e^{-\mathbf{x} - \log(4)} \\ &= 1 - e^{-\mathbf{x}}. \end{aligned}$$

This means that $\mathbf{r}_0^\circ \approx \mathbf{r}_0$ as long as $\alpha(m) \rightarrow 0$.

Now we want to show how this approach allows to prove the classical weak convergence statements for the sieve profile ME and efficiency of the sieve profile MLE $\hat{\boldsymbol{\theta}}_m \in \mathbb{R}^p$. From this point on we focus on the i.i.d. model in which n denotes the sample size and the functional is of the form $\mathcal{L} = \sum_{i=1}^n \ell(\boldsymbol{\theta}, \boldsymbol{\eta}, \mathbf{Y}_i)$. As in Section 2.4 this gives that $\mathcal{D}_m^2 = n\check{d}_m$, $\check{D}_m^2 = n\check{d}_m$ and $\check{D}^2 = n\check{d}$. As the efficient covariance is derived for the score evaluated in the true full target $\mathbf{v}^* \in l^2$ we need further assumptions on the bias:

(bias') With $\|\cdot\|$ denoting the spectral norm and with some function $\beta(m) \rightarrow 0$ as $m \rightarrow \infty$

$$\begin{aligned} \|I - \check{D}_m(\mathbf{v}^*)^{-1}\check{D}(\mathbf{v}^*)^2\check{D}_m(\mathbf{v}^*)^{-1}\| &\leq \beta(m), \\ \|I - \check{D}_m(\mathbf{v}_m^*)^{-1}\check{D}_m(\mathbf{v}^*)^2\check{D}_m(\mathbf{v}_m^*)^{-1}\| &\leq \beta(m). \end{aligned}$$

Remark 2.27. This paper focuses on the result 2.2 and thus we do not elaborate on approximation theory. But [3] presents conditions on the structure of $\mathcal{D} : l^2 \rightarrow l^2$ and on the sequence $\boldsymbol{\eta}^* \in l^2$ that yield **(bias')**.

Further we need convergence of the covariance of the weighted score. For this define

$$\begin{aligned} \check{v}_{m,\mathcal{D}}^2(\mathbf{v}_m^*) &\stackrel{\text{def}}{=} \text{Cov}(\nabla_{\boldsymbol{\theta}}\ell_1(\mathbf{v}_m^*) - A_m H_m^{-2}\nabla_{\boldsymbol{\eta}}\ell_1(\mathbf{v}_m^*)), \\ \check{v}^2(\mathbf{v}^*) &\stackrel{\text{def}}{=} \text{Cov}(\nabla_{\boldsymbol{\theta}}\ell_1(\mathbf{v}^*) - A H^{-2}\nabla_{\boldsymbol{\eta}}\ell_1(\mathbf{v}^*)). \end{aligned}$$

(bias'') As $m \rightarrow \infty$ with $\|\cdot\|$ denoting the spectral norm

$$\|\check{D}_m^{-1}(\mathbf{v}_m^*)\check{V}_{m,\mathcal{D}}^2(\mathbf{v}_m^*)\check{D}_m^{-1}(\mathbf{v}_m^*) - \check{d}^{-1}\check{v}^2\check{d}^{-1}\| \rightarrow 0.$$

Remark 2.28. This is a condition on how the covariance operator of $\nabla_{p+m}\mathcal{L}(\mathbf{v}) \in \mathbb{R}^{p+m}$ is affected when it is evaluated in $\mathbf{v}_m^* \in \mathbb{R}^{p+m}$ instead of $\mathbf{v}^* \in l^2$. In the single-index example we get **(bias'')** due to the smoothness of the functional.

Corollary 2.8 and Theorem 2.9 allow to derive the following corollary which yields the asymptotic efficiency of $\tilde{\boldsymbol{\theta}}_m$ and the classical Wilks phenomenon.

Corollary 2.10. *Assume that we have iid observations from $\mathbb{P} = \mathbb{P}_{\boldsymbol{\theta}^*, \boldsymbol{\eta}^*}$ and that for some $m_0 \in \mathbb{N}$ any $m \geq m_0$ the conditions of Theorem 2.8 and the condition $(\check{\mathcal{E}}\mathcal{D}_0)$ are satisfied with $\mathbf{v}^\circ = \mathbf{v}_m^*$. Further let the conditions (bias') and (bias'') be satisfied. Assume that for any $\mathbf{r} > 0$ that $\check{\delta}_n(\mathbf{r}) \rightarrow 0$ as $n \in \mathbb{N}$ tends to infinity, that $\check{\omega}_n \rightarrow 0$ and that $\mathbf{r}_0(\mathbf{x}) < \infty$ for any $\mathbf{x} > 0$, $m, n \in \mathbb{N}$, where $\mathbf{r}_0(\mathbf{x})$ is chosen such that $\mathbb{P}(\tilde{\mathbf{v}}_m, \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*, m}, \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*, m} \in \mathcal{Y}_\circ(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}$. Then there is a sequence $m_n \rightarrow \infty$ such that as $n \rightarrow \infty$*

$$\begin{aligned} n\check{d}(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}} &\xrightarrow{\mathbb{P}} 0, \\ n\check{d}(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) &\xrightarrow{w} \mathcal{N}(0, \check{d}^{-1}\check{v}^2\check{d}^{-1}), \\ 2\check{L}(\tilde{\boldsymbol{\theta}}_m, \boldsymbol{\theta}^*) &\xrightarrow{w} \mathcal{L}(\|\check{\boldsymbol{\xi}}_\infty\|), \quad \check{\boldsymbol{\xi}}_\infty \sim \mathcal{N}(0, \check{d}^{-1}\check{v}^2\check{d}^{-1}). \end{aligned}$$

Remark 2.29. On this level of generality we can not specify the right choice of $m_n \in \mathbb{N}$ that ensures the convergence. But in [1] it is shown that it equals the optimal choice for a series estimator of the nuisance component $\boldsymbol{\eta}^* \in l^2$ for know $\boldsymbol{\theta}^*$ - as pointed out in [24] the best choice is $m = n^{1/(2\alpha+1)}$, with $\alpha > 1/2$ quantifying the "smoothness" of $\boldsymbol{\eta}^*$ - is admissible.

Remark 2.30. For the case of the profile MLE $\ell(\boldsymbol{\theta}, \boldsymbol{\eta}, \mathbf{Y}_i)$ is the log-likelihood for a single observation. In that case assume that the linear operator $\mathbb{F}_{\mathbf{v}^*}^2 \stackrel{\text{def}}{=} \text{Cov}\{\nabla\ell(\mathbf{v}^*)\} : l^2 \rightarrow \text{Im}(\mathbb{F}_{\mathbf{v}^*}^2)$ is invertible and that $\nabla\ell(\mathbf{v}^*) \in \text{Im}(\mathbb{F}_{\mathbf{v}^*}^2)$. It is known from the convolution theorem (see [32], Theorem 3.11.2 p. 414, setting $\varkappa(\mathbb{P}_{\mathbf{v}}) = \boldsymbol{\theta}$) that the asymptotically optimal variance for regular estimators is given by the inverse of the partial information matrix

$$\check{\mathbb{F}}_{\mathbf{v}^*} = \left(\Pi_{\boldsymbol{\theta}} \text{Cov}\{\nabla\ell(\mathbf{v}^*)\}^{-1} \Pi_{\boldsymbol{\theta}}^\top \right)^{-1},$$

where as above $\Pi_{\boldsymbol{\theta}}$ is the orthogonal projection onto the $\boldsymbol{\theta}$ -components, and $\Pi_{\boldsymbol{\theta}}^\top$ its adjoint operator. In the case of correct specification we have that $\check{v}^2 = \check{d}^{-1} = \check{\mathbb{F}}_{\boldsymbol{\theta}, \boldsymbol{\eta}}$, such that

$$\check{d}^{-1}\check{v}^2\check{d}^{-1} = I_p.$$

In that case Corollary 2.10 yields the efficiency of the sieve profile MLE and we recover the Wilks phenomenon for that estimator.

2.9. Application to single index model

We illustrate how the results from Section 2 and the last statement can be derived for Single Index modeling. We focus on the complete set of assumptions

that allow to apply the results from above. For a detailed treatment of this model see [1]. Consider the following model

$$\mathbf{Y}_i = f(\mathbf{X}_i^\top \boldsymbol{\theta}^*) + \varepsilon_i, \quad i = 1, \dots, n,$$

for some $f : \mathbb{R} \rightarrow \mathbb{R}$ and $\boldsymbol{\theta}^* \in S_1^{p,+} \subset \mathbb{R}^p$, i.i.d errors $\varepsilon_i \in \mathbb{R}$ with $\mathbb{E}\varepsilon_i = 0$ and $\text{Var}(\varepsilon_i) = \sigma^2$ and i.i.d random variables $\mathbf{X}_i \in \mathbb{R}^p$ with distribution denoted by $\mathbb{P}^{\mathbf{X}}$. The single-index model is widely applied in statistics. For example in econometric studies it serves as a compromise between too restrictive parametric models and flexible but hardly estimable purely nonparametric models. Usually the statistical inference focuses on estimating the index vector $\boldsymbol{\theta}^*$. A lot of research has already been done in this field. For instance, [7] show the asymptotic efficiency of the general semiparametric maximum-likelihood estimator for particular examples and in [12] the right choice of bandwidth for the nonparametric estimation of the link function is analyzed.

To ensure identifiability of $\boldsymbol{\theta}^* \in \mathbb{R}^p$ we assume that it lies in the half sphere $S_1^{p,+} \stackrel{\text{def}}{=} \{\boldsymbol{\theta} \in \mathbb{R}^p : \|\boldsymbol{\theta}\| = 1, \theta_1 > 0\} \subset \mathbb{R}^p$. For simplicity we assume that the support of the $\mathbf{X}_i \in \mathbb{R}^p$ is contained in the ball of radius $s_{\mathbf{X}} > 0$. This allows to approximate $f \in \{f : [-s_{\mathbf{X}}, s_{\mathbf{X}}] \mapsto \mathbb{R}\}$ by an orthonormal C^2 -Daubechies-wavelet basis, i.e. for a suitable function $\mathbf{e}_0 \stackrel{\text{def}}{=} \psi : [-s_{\mathbf{X}}, s_{\mathbf{X}}] \mapsto \mathbb{R}$ we set for $k = (2^{j_k} - 1)13 + r_k$ with $j_k \in \mathbb{N}_0$ and $r_k \in \{0, \dots, (2^{j_k})13 - 1\}$

$$\mathbf{e}_k(t) = 2^{j_k/2} \psi(2^{j_k}(t - 2r_k s_{\mathbf{X}})), \quad k \in \mathbb{N}.$$

Our aim is to analyze the properties of the profile MLE

$$\tilde{\boldsymbol{\theta}}_m \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta}}{\text{argmax}} \max_{\boldsymbol{\eta} \in \mathbb{R}^m} \mathcal{L}_m(\boldsymbol{\theta}, \boldsymbol{\eta}),$$

where

$$\mathcal{L}_m(\boldsymbol{\theta}, \boldsymbol{\eta}) \stackrel{\text{def}}{=} -\frac{1}{2} \sum_{i=1}^n \left| \mathbf{Y}_i - \sum_{k=0}^m \eta_k \mathbf{e}_k(\mathbf{X}_i^\top \boldsymbol{\theta}) \right|^2.$$

[16] analyzed a very similar estimator in a more general setting based on a kernel estimation of $\mathbb{E}[\mathbf{Y} \mid f(\boldsymbol{\theta}^\top \mathbf{X})]$ instead of using a parametric sieve approximation $\sum_{k=0}^m \eta_k \mathbf{e}_k$. He showed \sqrt{n} -consistency and asymptotic normality of the proposed estimator.

To apply the technique presented above we need a list of assumptions denoted by (\mathcal{A}) :

(Cond_X) The measure $\mathbb{P}^{\mathbf{X}}$ is absolutely continuous with respect to the Lebesgue measure. The Lebesgue density $d_{\mathbf{X}} : \mathbb{R}^p \rightarrow \mathbb{R}$ of $\mathbb{P}^{\mathbf{X}}$ is only positive on the ball $B_{s_{\mathbf{x}}+h}(0) \subset \mathbb{R}^p$ with some small $h > 0$ and Lipschitz continuous on $B_{s_{\mathbf{x}}}(0) \subset \mathbb{R}^p$ with Lipschitz constant $L_{d_{\mathbf{X}}} \in \mathbb{R}_+$. Also the density $d_{\mathbf{X}} : \mathbb{R}^p \rightarrow \mathbb{R}$ of the regressors satisfies $c_{d_{\mathbf{X}}} \leq d_{\mathbf{X}} \leq C_{d_{\mathbf{X}}}$ on $B_{s_{\mathbf{x}}}(0) \subset \mathbb{R}^p$

for constants $0 < c_{d_{\mathbf{X}}} \leq C_{d_{\mathbf{X}}} < \infty$. Further we assume that for any $\boldsymbol{\theta} \perp \boldsymbol{\theta}^*$ with $\|\boldsymbol{\theta}\| = 1$ we have $\text{Var}(\mathbf{X}^\top \boldsymbol{\theta} | \mathbf{X}^\top \boldsymbol{\theta}^*) > \sigma_{\mathbf{X} | \boldsymbol{\theta}^*}^2$ for some constant $\sigma_{\mathbf{X} | \boldsymbol{\theta}^*}^2 > 0$ that does not depend on $\mathbf{X}^\top \boldsymbol{\theta}^* \in \mathbb{R}$.

(Cond_f) For some $\boldsymbol{\eta}^* \in l^2$

$$f = f_{\boldsymbol{\eta}^*} = \sum_{k=1}^{\infty} \eta_k^* \mathbf{e}_k,$$

where with some $\alpha > 2$ and a constant $C_{\|\boldsymbol{\eta}^*\|} > 0$

$$\sum_{l=0}^{\infty} l^{2\alpha} \eta_l^{*2} \leq C_{\|\boldsymbol{\eta}^*\|}^2 < \infty.$$

(Cond_{X $\boldsymbol{\theta}^*$}) It holds true that $\mathbb{P}(|f'_{\boldsymbol{\eta}^*}(\mathbf{X}^\top \boldsymbol{\theta}^*)| > c_{f'_{\boldsymbol{\eta}^*}}) > c_{\mathbb{P}_{f'}}$ for some $c_{f'_{\boldsymbol{\eta}^*}}, c_{\mathbb{P}_{f'}} > 0$.

(Cond _{ε}) The errors $(\varepsilon_i) \in \mathbb{R}$ are i.i.d. with $\mathbb{E}[\varepsilon_i] = 0$, $\text{Cov}(\varepsilon_i) = \sigma^2$ and satisfy for all $|\mu| \leq \tilde{g}$ for some $\tilde{g} > 0$ and some $\tilde{\nu}_{\boldsymbol{\varepsilon}} > 0$

$$\log \mathbb{E}[\exp\{\mu \varepsilon_1\}] \leq \tilde{\nu}_{\boldsymbol{\varepsilon}}^2 \mu^2 / 2.$$

(Cond _{\mathcal{Y}}) $\mathcal{Y} \subseteq \mathcal{Y}_o(\sqrt{n} \mathbf{r}^\circ) \subset \mathbb{R}^{p+m}$ with $\mathbf{r}^\circ \in \mathbb{R}$, i.e. $d_{\mathcal{Y}} \stackrel{\text{def}}{=} \text{diam}(\mathcal{Y}) < \infty$.

If these conditions denoted by (\mathcal{A}) are met we can proof the following results:

Proposition 2.11. Assume (\mathcal{A}) with $\alpha = 2 + \epsilon$ for some $\epsilon > 0$ with $p^{*5}/n \rightarrow 0$ but $p^{*5+2\epsilon}/n \rightarrow \infty$. If $n \in \mathbb{N}$ is large enough it holds with probability greater $1 - 4e^{-x} - \exp\{-m^3\} - \exp\{-nc_{(\mathcal{Q})}/4\}$

$$\|\check{D}_m(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*) - \check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| \leq C_{\diamond} \frac{(p^* + \mathbf{x})^{5/2}}{\sqrt{n}},$$

$$|2\check{L}(\tilde{\boldsymbol{\theta}}_m, \boldsymbol{\theta}_m^*) - \|\check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\|^2| \leq \left(\sqrt{p + \mathbf{x}} + C_{\diamond} \frac{(p^* + \mathbf{x})^{5/2}}{\sqrt{n}} \right) C_{\diamond} \frac{(p^* + \mathbf{x})^{5/2}}{\sqrt{n}}.$$

where $c_{(\mathcal{Q})} > 0$. Further as $n \rightarrow \infty$

$$\begin{aligned} \|\check{D}(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| &\xrightarrow{\mathbb{P}} 0, \\ \check{D}(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) &\xrightarrow{w} \mathcal{N}(0, \sigma^2 \mathbb{I}_p), \\ 2\check{L}(\tilde{\boldsymbol{\theta}}_m, \boldsymbol{\theta}^*) &\xrightarrow{w} \chi_p^2. \end{aligned}$$

Remark 2.31. The constant $C_{\diamond} > 0$ is a polynomial of $\|\psi\|_{\infty}, \|\psi'\|_{\infty}, \|\psi''\|_{\infty}$ and $s_{\mathbf{X}}$ that is independent of \mathbf{x}, n, p^* . The constant $c_{(\mathcal{Q})} > 0$ is related to \mathbf{b} from $(\mathcal{L}_{\mathbf{r}})$ and also does not depend on \mathbf{x}, n, p^* .

Remark 2.32. The necessary size of $n \in \mathbb{N}$ is determined by the size of $p^{*5/2}/\sqrt{n} \rightarrow 0$ and $m^{-2\alpha-1}n \rightarrow 0$. In the proof of Proposition 2.11 we impose conditions on $n \in \mathbb{N}$ of the kind

$$p^{*5/2}/\sqrt{n} \leq C_1^{-1}, \quad m^{-2\alpha-1}n \leq C_2^{-1},$$

for certain constants $C_1, C_2 > 0$ that are polynomials of $\|\psi\|_\infty, \|\psi'\|_\infty, \|\psi''\|_\infty, C_{\|f^*\|}$ and $s_{\mathbf{X}}$. These constants enter into the bound for C_\diamond .

For details see [1].

Appendix A: Deviation bounds for quadratic forms

The following general results from the supplement of [29] help to control the deviation for quadratic forms of type $\|B\xi\|^2$ for a given positive matrix B – i.e. $BB^\top > 0$ – and a random vector ξ . It will be used several times in our proofs. Suppose that

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

Remark A.1. In the setting of Section 2 we have either $\xi = \check{V}^{-1}\check{\nabla}_\theta \zeta(\mathbf{v}^*)$ and $B = \check{D}^{-1}\check{V}$ or $\xi = \mathcal{V}^{-1}\nabla \zeta(\mathbf{v}^*)$ and $B = \mathcal{D}^{-1}\mathcal{V}$.

For a matrix B , define

$$\mathfrak{p} = \text{tr}(BB^\top), \quad \mathfrak{v}^2 = 2 \text{tr}(BB^\top BB^\top), \quad \lambda_B \stackrel{\text{def}}{=} \|BB^\top\| \stackrel{\text{def}}{=} \lambda_{\max}(BB^\top).$$

For ease of presentation, suppose that $\mathfrak{g}^2 \geq 2\mathfrak{p}_B$. The other case only changes the constants in the inequalities. Define $\mu_c = 2/3$ and

$$\begin{aligned} \mathfrak{g}_c &\stackrel{\text{def}}{=} \sqrt{\mathfrak{g}^2 - \mu_c \mathfrak{p}_B}, \\ 2(\mathbf{x}_c + 2) &\stackrel{\text{def}}{=} (\mathfrak{g}^2/\mu_c - \mathfrak{p}_B)/\lambda_B + \log \det(\mathbb{I}_p - \mu_c B/\lambda_B). \end{aligned} \quad (\text{A.1})$$

Proposition A.1. Let (ED_0) hold with $\nu_0 = 1$ and $\mathfrak{g}^2 \geq 2\mathfrak{p}_B$. Then for each $\mathbf{x} > 0$

$$\mathbb{P}(\|B\xi\| \geq \mathfrak{z}(\mathbf{x}, BB^\top)) \leq 2e^{-\mathbf{x}},$$

where $\mathfrak{z}(\mathbf{x}, BB^\top)$ is defined by

$$\begin{aligned} \mathfrak{z}^2(\mathbf{x}, BB^\top) & \hspace{15em} (\text{A.2}) \\ & \stackrel{\text{def}}{=} \begin{cases} \mathfrak{p}_B + 2\mathfrak{v}_B(\mathbf{x} + 1)^{1/2}, & \mathbf{x} + 1 \leq \mathfrak{v}_B/(18\lambda_B), \\ \mathfrak{p}_B + 6\lambda_B(\mathbf{x} + 1), & \mathfrak{v}_B/(18\lambda_B) < \mathbf{x} + 1 \leq \mathbf{x}_c + 2, \\ |y_c + 2\lambda_B(\mathbf{x} - \mathbf{x}_c + 1)/\mathfrak{g}_c|^2, & \mathbf{x} > \mathbf{x}_c + 1, \end{cases} \end{aligned}$$

with $y_c^2 \leq \mathfrak{p}_B + 6\lambda_B(\mathbf{x}_c + 2)$.

Depending on the value \mathbf{x} , we observe three types of tail behavior of the quadratic form $\|B\xi\|^2$. The sub-Gaussian regime for $\mathbf{x} + 1 \leq v_B/(18\lambda_B)$ and the Poissonian regime for $\mathbf{x} \leq \mathbf{x}_c + 1$ are similar to the case of a Gaussian quadratic form. The value \mathbf{x}_c from (A.1) is of order \mathbf{g}^2 . In all our results we suppose that \mathbf{g}^2 and hence, \mathbf{x}_c is sufficiently large and the quadratic form $\|\xi\|^2$ can be bounded with a dominating probability by $p_B + 6\lambda_B(\mathbf{x} + 1)$ for a proper \mathbf{x} . We refer to the supplement of [29] for the proof of this and related results, further discussion and references.

Appendix B: Proofs

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

B.1. Proof of Lemma 2.1

Proof. Take any $\gamma \in \mathbb{R}^p$ with $\|\gamma\| = 1$ then

$$\begin{aligned} \gamma^\top \check{D}^{-1} \check{\nabla}_\theta \zeta(\mathbf{v}) &= \gamma^\top \left(\check{D}^{-1} D \quad \check{D}^{-1} A H^{-1} \right) \begin{pmatrix} D^{-1} & 0 \\ 0 & H^{-1} \end{pmatrix} \nabla \zeta(\mathbf{v}) \\ &\stackrel{\text{def}}{=} \hat{\gamma}^\top \mathcal{D}^{-1} \nabla \zeta(\mathbf{v}), \end{aligned}$$

where

$$\|\hat{\gamma}\| \leq \left\| \left(\check{D}^{-1} D \quad \check{D}^{-1} A H^{-1} \right) \right\| \left\| \begin{pmatrix} D^{-1} & 0 \\ 0 & H^{-1} \end{pmatrix} \mathcal{D} \right\| \leq \frac{(1 + \rho)\sqrt{1 + \rho^2}}{\sqrt{1 - \rho^2}}.$$

This gives that $(\mathcal{E}\mathcal{D}_1)$ implies $(\check{\mathcal{E}}\mathcal{D}_1)$ and $(\mathcal{E}\mathcal{D}_0)$ implies $(\check{\mathcal{E}}\mathcal{D}_0)$ with

$$\check{\mathbf{g}} = \frac{\sqrt{1 - \rho^2}}{(1 + \rho)\sqrt{1 + \rho^2}} \mathbf{g}, \quad \check{\nu} = \frac{(1 + \rho)\sqrt{1 + \rho^2}}{\sqrt{1 - \rho^2}} \nu.$$

Further for any $\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})$

$$\begin{aligned} \|I_p - D^{-1} D^2(\mathbf{v}) D^{-1}\| &= \|D^{-1} (D^2 - D^2(\mathbf{v})) D^{-1}\| \\ &= \|D^{-1} \Pi_\theta (\mathcal{D}^2 - \mathcal{D}^2(\mathbf{v})) \Pi_\theta^\top D^{-1}\| \\ &= \|D^{-1} \Pi_\theta \mathcal{D} (I_{p^*} - \mathcal{D}^{-1} \mathcal{D}^2(\mathbf{v}) \mathcal{D}^{-1}) \mathcal{D} \Pi_\theta^\top D^{-1}\| \\ &\leq \|D^{-1} \Pi_\theta \mathcal{D}\|^2 \|I_{p^*} - \mathcal{D}^{-1} \mathcal{D}^2(\mathbf{v}) \mathcal{D}^{-1}\| = \delta(\mathbf{r}). \end{aligned}$$

Also

$$\begin{aligned} \|D^{-1} (A(\mathbf{v}) - A) H^{-1}\| &= \|D^{-1} \Pi_\theta (\mathcal{D}^2(\mathbf{v}) - \mathcal{D}^2) \Pi_\theta^\top H^{-1}\| \\ &= \|D^{-1} \Pi_\theta \mathcal{D} (\mathcal{D}^{-1} \mathcal{D}^2(\mathbf{v}) \mathcal{D}^{-1} -) \mathcal{D} \Pi_\theta^\top H^{-1}\| \end{aligned}$$

$$\begin{aligned} &\leq \|D^{-1} \Pi_{\theta} \mathcal{D}\| \|H^{-1} \Pi_{\eta} \mathcal{D}\| \|I_{p^*} - \mathcal{D}^{-1} \mathcal{D}^2(\mathbf{v}) \mathcal{D}^{-1}\| \\ &\leq \delta(\mathbf{r}). \end{aligned}$$

With the same arguments

$$\|D^{-1} A H^{-1} (I_m - H^{-1} H^2(\mathbf{v}) H^{-1})\| \leq \rho \delta(\mathbf{r}).$$

□

B.2. Proof of Theorem 2.2

Remember the semiparametric spread

$$\check{\diamond}(\mathbf{x}, \mathbf{x}) \stackrel{\text{def}}{=} 4 \left(\frac{4}{(1 - \rho^2)^2} \check{\delta}(4\mathbf{r}) + 6\nu_1 \check{\omega}_3(\mathbf{x}, 2p^* + 2p) \right) \mathbf{r}.$$

For $\zeta(\mathbf{v}) = \mathcal{L}(\mathbf{v}) - \mathbb{E}\mathcal{L}(\mathbf{v})$ define the semiparametric normalized stochastic gradient gap

$$\check{\mathfrak{Y}}(\mathbf{v}) = \check{D}^{-1} \left(\check{\nabla}_{\theta} \zeta(\mathbf{v}) - \check{\nabla}_{\theta} \zeta(\mathbf{v}^*) \right). \tag{B.1}$$

Fix the radius $\mathbf{r}_0(\mathbf{x}) > 0$ that ensures $\mathbb{P}\{\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_o(\mathbf{r}_0)\} \geq 1 - e^{-\mathbf{x}}$. Define $C(\mathbf{r}_0, \mathbf{x}) \subseteq \Omega$ as

$$C(\mathbf{r}_0, \mathbf{x}) \stackrel{\text{def}}{=} \{\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_o(\mathbf{r}_0)\} \cap \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_o(4\mathbf{r}_0)} \|\check{\mathfrak{Y}}(\mathbf{v})\| \leq 6\nu_1 \check{\omega}_3(\mathbf{x}, \mathbb{Q}) 4\mathbf{r}_0 \right\}.$$

In the following we will derive statements that hold true on this set $C(\mathbf{r}_0, \mathbf{x}) \subseteq \Omega$ which is of probability greater $1 - 2e^{-\mathbf{x}}$ because it follows right away from the definition of $\mathbf{r}_0 > 0$ that

$$\mathbb{P}\{\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \notin \mathcal{Y}_o(\mathbf{r}_0)\} \leq e^{-\mathbf{x}},$$

and by Theorem C.1 which is applicable because $(\check{\mathfrak{E}}\mathcal{D}_1)$ implies (C.1) with $\|\cdot\|_{\mathfrak{y}} = \|\mathcal{D}(\cdot)\|$

$$\mathbb{P} \left(\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}_0)} \|\check{\mathfrak{Y}}(\mathbf{v})\| \leq 6\nu_1 \check{\omega}_3(\mathbf{x}, 2p^* + 2p) \mathbf{r}_0 \right) \geq 1 - e^{-\mathbf{x}}.$$

B.2.1. Proof of claim on $C(\mathbf{r}_0, \mathbf{x}) \subseteq \Omega$

Before we prove the claim we prove the following useful lemma:

Lemma B.1. *Assume that the condition $(\check{\mathcal{L}}_0)$ is fulfilled. Then*

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \left\| \check{D}^{-1} \left(\check{\nabla} \mathbb{E}\mathcal{L}(\mathbf{v}) - \check{\nabla} \mathbb{E}\mathcal{L}(\mathbf{v}^*) \right) + \check{D}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right\| \leq \frac{4}{(1 - \rho^2)^2} \mathbf{r} \check{\delta}(\mathbf{r}).$$

Proof. We have with Taylor expansion and some $\widehat{\boldsymbol{v}} \in \Upsilon_o(\boldsymbol{r})$

$$\begin{aligned} \nabla \mathbb{E} \mathcal{L}(\boldsymbol{v}) - \nabla \mathbb{E} \mathcal{L}(\boldsymbol{v}^*) &= \nabla^2 \mathbb{E} \mathcal{L}(\widehat{\boldsymbol{v}})(\boldsymbol{v} - \boldsymbol{v}^*) \\ &\stackrel{\text{def}}{=} -\mathcal{D}^2(\widehat{\boldsymbol{v}})(\boldsymbol{v} - \boldsymbol{v}^*) \\ &= - \begin{pmatrix} D^2(\widehat{\boldsymbol{v}}) & A(\widehat{\boldsymbol{v}}) \\ A^\top(\widehat{\boldsymbol{v}}) & H^2(\widehat{\boldsymbol{v}}) \end{pmatrix} (\boldsymbol{v} - \boldsymbol{v}^*). \end{aligned}$$

This gives

$$\begin{aligned} &-\check{D}^{-1} \left(\check{\nabla} \mathbb{E} \mathcal{L}(\boldsymbol{v}) - \check{\nabla} \mathbb{E} \mathcal{L}(\boldsymbol{v}^*) \right) \\ &= \check{D}^{-1} \left(D^2(\widehat{\boldsymbol{v}}) - AH^{-2}A^\top(\widehat{\boldsymbol{v}}) \quad A(\widehat{\boldsymbol{v}}) - AH^{-2}H^2(\widehat{\boldsymbol{v}}) \right) (\boldsymbol{v} - \boldsymbol{v}^*) \\ &= \check{D}^{-1} \left(D^2(\widehat{\boldsymbol{v}}) - AH^{-2}A^\top(\widehat{\boldsymbol{v}}) \right) \check{D}^{-1} \check{D}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \\ &\quad + \left(\check{D}^{-1}A(\widehat{\boldsymbol{v}}) - \check{D}^{-1}AH^{-2}H^2(\widehat{\boldsymbol{v}}) \right) (\boldsymbol{\eta} - \boldsymbol{\eta}^*). \end{aligned}$$

We estimate separately using $(\check{\mathcal{L}}_0)$ and (\mathcal{I})

$$\begin{aligned} &\| \check{D}^{-1} \left(D^2(\widehat{\boldsymbol{v}}) - AH^{-2}A^\top(\widehat{\boldsymbol{v}}) \right) \check{D}^{-1} - I_p \| \\ &= \left\| \check{D}^{-1} \left(D^2(\widehat{\boldsymbol{v}}) - D^2 - \{AH^{-2}(A^\top(\widehat{\boldsymbol{v}}) - A^\top)\} \right) \check{D}^{-1} \right\| \\ &\leq \| \check{D}^{-1}D \|^2 (\|D^{-1}D^2(\widehat{\boldsymbol{v}})D^{-1} - I_p\| \\ &\quad + \|D^{-1}AH^{-1}\| \|D^{-1}(A(\widehat{\boldsymbol{v}}) - A)H^{-1}\|) \\ &\leq \frac{1 + \rho}{1 - \rho^2} \check{\delta}_0(\boldsymbol{r}), \end{aligned}$$

and

$$\begin{aligned} &\left\| \left(\check{D}^{-1}A(\widehat{\boldsymbol{v}}) - \check{D}^{-1}AH^{-2}H^2(\widehat{\boldsymbol{v}}) \right) (\boldsymbol{\eta} - \boldsymbol{\eta}^*) \right\| \\ &\leq \left\| \check{D}^{-1}A(\widehat{\boldsymbol{v}})H^{-1} - \check{D}^{-1}AH^{-2}H^2(\widehat{\boldsymbol{v}})H^{-1} \right\| \|H(\boldsymbol{\eta} - \boldsymbol{\eta}^*)\| \\ &\leq \| \check{D}^{-1}D \| \{ \|D^{-1}(A(\widehat{\boldsymbol{v}}) - A)H^{-1}\| \\ &\quad + \|D^{-1}AH^{-1}(I_m - H^{-1}H^2(\widehat{\boldsymbol{v}})H^{-1})\| \} \|H(\boldsymbol{\eta} - \boldsymbol{\eta}^*)\| \\ &\leq \frac{2}{\sqrt{1 - \rho^2}} \check{\delta}_0(\boldsymbol{r}) \|H(\boldsymbol{\eta} - \boldsymbol{\eta}^*)\|. \end{aligned}$$

Further

$$\| \check{D}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \| \vee \| H(\boldsymbol{\eta} - \boldsymbol{\eta}^*) \| \leq \frac{1}{\sqrt{1 - \rho^2}} \frac{1}{\sqrt{1 - \rho^2}} \| \mathcal{D}(\boldsymbol{v} - \boldsymbol{v}^*) \| \leq \frac{1}{1 - \rho^2} \boldsymbol{r}.$$

Together this gives that

$$\begin{aligned} & \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \left\| \check{D}^{-1} \left(\check{\nabla} \mathbb{E} \mathcal{L}(\mathbf{v}) - \check{\nabla} \mathbb{E} \mathcal{L}(\mathbf{v}^*) \right) + \check{D}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right\| \\ & \leq \left(\frac{1 + \rho}{1 - \rho^2} + \frac{2}{\sqrt{1 - \rho^2}} \right) \frac{1}{1 - \rho^2} \mathbf{r} \check{\delta}(\mathbf{r}) \\ & \leq \frac{4}{(1 - \rho^2)^2} \mathbf{r} \check{\delta}(\mathbf{r}). \end{aligned}$$

□

The next Lemma already completes the proof of (2.6) and (2.7) on $C(\mathbf{r}_0, \mathbf{x}) \subset \Omega$:

Lemma B.2. *Assume that the condition $(\check{\mathcal{L}}_0)$ is fulfilled. Then on the set $C(\mathbf{r}_0, \mathbf{x}) \subset \Omega$ the approximations (2.6) and (2.7) are valid.*

Proof. Using $\check{\nabla}_{\boldsymbol{\theta}} \mathcal{L}(\tilde{\mathbf{v}}) = 0$, that by assumption $\check{\nabla} \mathbb{E} \mathcal{L} = \mathbb{E} \check{\nabla} \mathcal{L}$ and the triangular inequality we find

$$\begin{aligned} \|\check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\| &= \left\| \check{D}^{-1} \left\{ \check{\nabla} \mathcal{L}(\tilde{\mathbf{v}}) - \check{\nabla} \mathcal{L}(\mathbf{v}^*) \right\} + \check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \right\| \\ &\leq \left\| \check{D}^{-1} \left(\check{\nabla} \mathbb{E} \mathcal{L}(\tilde{\mathbf{v}}) - \check{\nabla} \mathbb{E} \mathcal{L}(\mathbf{v}^*) \right) + \check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \right\| \\ &\quad + \left\| \check{D}^{-1} \left\{ \check{\nabla}_{\boldsymbol{\theta}} \zeta(\tilde{\mathbf{v}}) - \check{\nabla}_{\boldsymbol{\theta}} \zeta(\mathbf{v}^*) \right\} \right\|. \end{aligned}$$

Note that by condition $(\check{\mathcal{L}}_0)$ we get with Lemma B.1 as we assume that $\tilde{\mathbf{v}} \in \mathcal{Y}_\circ(\mathbf{r}_0)$

$$\left\| \check{D}^{-1} \left(\check{\nabla} \mathbb{E} \mathcal{L}(\tilde{\mathbf{v}}) - \check{\nabla} \mathbb{E} \mathcal{L}(\mathbf{v}^*) \right) + \check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \right\| \leq \frac{4}{(1 - \rho^2)^2} \mathbf{r}_0 \check{\delta}(\mathbf{r}_0).$$

For the remainder we use that on $C(\mathbf{r}_0, \mathbf{x}) \subset \Omega$

$$\left\| \check{D}^{-1} \left\{ \check{\nabla}_{\boldsymbol{\theta}} \zeta(\tilde{\mathbf{v}}) - \check{\nabla}_{\boldsymbol{\theta}} \zeta(\mathbf{v}^*) \right\} \right\| \leq \sup_{\mathbf{v} \in \mathcal{Y}_\circ(4\mathbf{r}_0)} \|\check{\mathfrak{Y}}(\mathbf{v})\| \leq 6\check{\nu}_1 \check{\omega}_3(\mathbf{x}, \mathbb{Q}) 4\mathbf{r}_0.$$

This gives (2.6) on $C(\mathbf{r}_0, \mathbf{x}) \subset \Omega$. For (2.7) we will first show that on $C(\mathbf{r}_0, \mathbf{x}) \subset \Omega$

$$\begin{aligned} & \left| \check{L}(\tilde{\boldsymbol{\theta}}) - \check{L}(\boldsymbol{\theta}^*) - \left(\check{\nabla} \zeta(\mathbf{v}^*)(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \|\check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2/2 \right) \right| \quad (\text{B.2}) \\ & \leq \left(\|\check{D}^{-1} \check{\nabla}\| + \check{\diamond}(\mathbf{r}_0, \mathbf{x}) \right) \check{\diamond}(\mathbf{r}_0, \mathbf{x}). \end{aligned}$$

To show this we use some ideas of the proof of Theorem 1 of [22], that is we define

$$l : \mathbb{R}^p \times \mathcal{Y} \rightarrow \mathbb{R}, \quad (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\eta}) \mapsto \mathcal{L}(\boldsymbol{\theta}_1, \boldsymbol{\eta} + H^{-2} A^\top (\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1)). \quad (\text{B.3})$$

Note that

$$\begin{aligned} \nabla_{\theta_1} l(\theta_1, \theta_2, \eta) &= \check{\nabla}_{\theta} \mathcal{L}(\theta_1, \eta + H^{-2} A^\top (\theta_2 - \theta_1)), \\ \text{i.e. } \nabla_{\theta_1} l(\theta^*, \theta^*, \eta^*) &= \check{\nabla} \zeta(\mathbf{v}^*). \end{aligned}$$

Remark B.1. If the model was correctly specified and $\tilde{\mathcal{L}}$ the true log likelihood $\nabla_{\theta_1} l(\theta^*, \theta^*, \eta^*)$ would be equal to $\sum_{i=1}^n \tilde{\psi}_{\mathbb{P}}(\mathbf{Y}_i)$, with $\tilde{\psi}_{\mathbb{P}}$ the efficient influence function.

We can represent:

$$\check{L}(\tilde{\theta}) - \check{L}(\theta^*) = l(\tilde{\theta}, \tilde{\theta}, \tilde{\eta}) - l(\theta^*, \theta^*, \tilde{\eta}_{\theta^*}), \quad \tilde{\eta}_{\theta^*} \stackrel{\text{def}}{=} \Pi_{\eta} \underset{\substack{\mathbf{v} \in \mathcal{Y}, \\ \Pi_{\theta} \mathbf{v} = \theta^*}}{\text{argmax}} \mathcal{L}(\mathbf{v}).$$

This allows to bound from above

$$\begin{aligned} \check{L}(\tilde{\theta}) - \check{L}(\theta^*) &\leq l(\tilde{\theta}, \tilde{\theta}, \tilde{\eta}) - l(\theta^*, \tilde{\theta}, \tilde{\eta}) \\ &= \nabla_{\theta_1} l(\theta^*, \theta^*, \eta^*)(\tilde{\theta} - \theta^*) - \|\check{D}(\tilde{\theta} - \theta^*)\|^2/2 + \check{\alpha}(\tilde{\theta}, \theta^*), \end{aligned}$$

where

$$\begin{aligned} \check{\alpha}(\theta_1, \theta_2) &\stackrel{\text{def}}{=} l(\theta_1, \tilde{\theta}, \tilde{\eta}) - l(\theta_2, \tilde{\theta}, \tilde{\eta}) - \nabla_{\theta_1} l(\theta^*, \theta^*, \eta^*)(\theta_1 - \theta_2) \\ &\quad + \|\check{D}(\theta_1 - \theta_2)\|^2/2. \end{aligned}$$

We will show

$$\check{\alpha}(\tilde{\theta}, \theta^*) \leq \left(\|\check{D}^{-1} \check{\nabla}\| + \check{\diamond}(\mathbf{r}_0, \mathbf{x}) \right) \check{\diamond}(\mathbf{r}_0, \mathbf{x}), \tag{B.4}$$

which gives the upper bound of (B.2). Note that $\check{\alpha}(\theta^*, \theta^*) = 0$ such that we get with Taylor expansion

$$\check{\alpha}(\tilde{\theta}, \theta^*) \leq \|\check{D}(\tilde{\theta} - \theta^*)\| \sup_{\theta \in \Pi_{\theta} \mathcal{Y}_c(\mathbf{r}_0)} |\check{D}^{-1} \nabla_{\theta_1} \check{\alpha}(\theta, \theta^*)|.$$

We find

$$\begin{aligned} \nabla_{\theta_1} \check{\alpha}(\theta, \theta^*) &= \nabla_{\theta_1} l(\theta, \tilde{\theta}, \tilde{\eta}) - \nabla_{\theta_1} l(\theta^*, \theta^*, \eta^*) + \check{D}(\theta - \theta^*) \\ &= \check{\nabla} \zeta(\mathbf{v}^\circ) - \check{\nabla} \zeta(\mathbf{v}^*) + \mathbb{E} \left[\check{\nabla} \mathcal{L}(\mathbf{v}^\circ) - \check{\nabla} \mathcal{L}(\mathbf{v}^*) \right] + \check{D}(\theta - \theta^*), \end{aligned}$$

where

$$\begin{aligned} \mathbf{v}^\circ &\stackrel{\text{def}}{=} (\theta, \tilde{\eta} + H^{-2} A^\top (\tilde{\theta} - \theta)), \\ \|\mathcal{D}(\mathbf{v}^\circ - \mathbf{v}^*)\| &\leq \|D(\theta - \theta^*)\| + \|H(\tilde{\eta} - \eta^*)\| + \rho \|D(\tilde{\theta} - \theta)\| \\ &\leq 2(1 + \rho) \mathbf{r}_0 < 4\mathbf{r}_0. \end{aligned}$$

Using Lemma B.1 and the definition of $C(\mathbf{r}_0, \mathbf{x})$ we can bound

$$\sup_{\boldsymbol{\theta} \in \Pi_{\boldsymbol{\theta}} \tilde{\mathcal{Y}}_o(\mathbf{r}_0)} |\check{D}^{-1} \nabla_{\boldsymbol{\theta}_1} \check{\alpha}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)| \leq \check{\diamond}(\mathbf{r}_0, \mathbf{x}).$$

Using (2.6) we find on $C(\mathbf{r}_0, \mathbf{x})$

$$\|\check{D}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq \|\check{D}^{-1} \check{\nabla}\| + \check{\diamond}(\mathbf{r}_0, \mathbf{x}).$$

This gives (B.4). Similarly we can bound from below:

$$\check{L}(\tilde{\boldsymbol{\theta}}) - \check{L}(\boldsymbol{\theta}^*) \geq l(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*}) - l(\boldsymbol{\theta}^*, \boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*}),$$

and repeat the same arguments using that $\tilde{\boldsymbol{v}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}_o(\mathbf{r}_0)$ on $C(\mathbf{r}_0, \mathbf{x}) \subset \Omega$ to obtain the lower bound of (B.2). Plugging (2.6) into (B.2) this gives

$$\begin{aligned} |2\check{L}(\tilde{\boldsymbol{\theta}}) - 2\check{L}(\boldsymbol{\theta}^*) - \|\check{D}^{-1} \check{\nabla} \zeta(\mathbf{v}^*)\|^2| &\leq 4 \left(\|\check{D}^{-1} \check{\nabla}\| + \check{\diamond}(\mathbf{r}_0, \mathbf{x}) \right) \check{\diamond}(\mathbf{r}_0, \mathbf{x}) \\ &\quad + \check{\diamond}(\mathbf{r}_0, \mathbf{x})^2. \end{aligned}$$

□

B.3. Proof of Proposition 2.4

We start with an auxiliary result. Define the *parametric gradient gap*

$$\mathfrak{Y}(\mathbf{v}) = \mathcal{D}^{-1} \left(\nabla \zeta(\mathbf{v}) - \nabla \zeta(\mathbf{v}^*) \right).$$

Lemma B.3. *Assume that the condition (\mathcal{L}_0) is fulfilled. Then for $0 \leq \mathbf{r}$ on the set*

$$\mathcal{M}(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \{ \tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}_o(\mathbf{r}) \} \cap \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\mathfrak{Y}(\mathbf{v})\| \leq 6\nu_1 \omega_{\mathfrak{Z}}(\mathbf{x}, 4p^*) \mathbf{r} \right\}, \quad (\text{B.5})$$

we have

$$\begin{aligned} \|\mathcal{D}(\tilde{\mathbf{v}} - \mathbf{v}^*) - \mathcal{D}^{-1} \nabla \mathcal{L}(\mathbf{v}^*)\| &\leq \diamond(\mathbf{r}, \mathbf{x}), \\ \|H(\tilde{\mathbf{v}} - \mathbf{v}^*) - H^{-1} \nabla_{\boldsymbol{\eta}} \mathcal{L}(\mathbf{v}^*)\| &\leq \diamond(\mathbf{r}, \mathbf{x}). \end{aligned}$$

Proof. Since $\nabla \mathcal{L}(\tilde{\mathbf{v}}) = 0$ we find with the triangular inequality

$$\begin{aligned} \|\mathcal{D}(\tilde{\mathbf{v}} - \mathbf{v}^*) - \mathcal{D}^{-1} \nabla \mathcal{L}(\mathbf{v}^*)\| &\leq \|\mathcal{D}^{-1} \left(\nabla \zeta(\tilde{\mathbf{v}}) - \nabla \zeta(\mathbf{v}^*) \right)\| \\ &\quad + \|\mathcal{D}^{-1} \mathbb{E} \nabla \mathcal{L}(\mathbf{v}) - \mathcal{D}^{-1} \mathbb{E} \nabla \mathcal{L}(\mathbf{v}^*) + \mathcal{D}(\tilde{\mathbf{v}} - \mathbf{v}^*)\|. \end{aligned}$$

In section 2.1 we assume that $\mathcal{L} : \mathbb{R}^{p^*} \rightarrow \mathbb{R}$ is smooth enough such that we can interchange $\nabla \mathbb{E} \mathcal{L}(\mathbf{v}) = \mathbb{E} \nabla \mathcal{L}(\mathbf{v})$ on $\mathcal{Y}_o(\mathbf{r}_0)$. This gives by condition (\mathcal{L}_0) and

Taylor expansion

$$\begin{aligned} & \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{x})} \|\mathcal{D}^{-1} \mathbb{E} \nabla \mathcal{L}(\mathbf{v}) - \mathcal{D}^{-1} \mathbb{E} \nabla \mathcal{L}(\mathbf{v}^*) + \mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \\ & \leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{x})} \|\mathcal{D}^{-1} \nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}) \mathcal{D}^{-1} + \mathbb{I}_{p^*}\| \mathbf{r} \leq \delta(\mathbf{x}) \mathbf{r}. \end{aligned}$$

For the remainder we use the definition of $\mathcal{M}(\mathbf{x}, \mathbf{x})$ in (B.5). This gives the first claim. For the second claim we repeat the same arguments with the restriction to the set $\mathcal{Y}_{o, \boldsymbol{\theta}^*}(\mathbf{x}) \stackrel{\text{def}}{=} \{(\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathcal{Y}_o(\mathbf{x}) : \boldsymbol{\theta} = \boldsymbol{\theta}^*\}$. We bound on $\mathcal{Y}_{o, \boldsymbol{\theta}^*}(\mathbf{x})$

$$\begin{aligned} & \|H^{-1} \{ \nabla_{\boldsymbol{\eta}} \mathcal{L}(\mathbf{v}) - \nabla_{\boldsymbol{\eta}} \mathcal{L}(\mathbf{v}^*) + H^2(\boldsymbol{\eta} - \boldsymbol{\eta}^*) \}\| \leq \|H^{-1} \{ \nabla_{\boldsymbol{\eta}} \mathcal{L}(\mathbf{v}) - \nabla_{\boldsymbol{\eta}} \mathcal{L}(\mathbf{v}^*) \}\| \\ & + \|H^{-1} \{ \nabla_{\boldsymbol{\eta}} \mathbb{E} \mathcal{L}(\mathbf{v}) - \nabla_{\boldsymbol{\eta}} \mathbb{E} \mathcal{L}(\mathbf{v}^*) + H^2(\boldsymbol{\eta} - \boldsymbol{\eta}^*) \}\|. \end{aligned}$$

Take any $\boldsymbol{\gamma} \in \mathbb{R}^m$ with $\|\boldsymbol{\gamma}\| = 1$ then

$$\boldsymbol{\gamma}^\top H^{-1} \nabla_{\boldsymbol{\eta}} \mathcal{L}(\mathbf{v}) = (0, H^{-1} \boldsymbol{\gamma})^\top \nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v}) = (0, H^{-1} \boldsymbol{\gamma})^\top \mathcal{D} \mathcal{D}^{-1} \nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v}).$$

Now note that $\|\mathcal{D}(0, H^{-1} \boldsymbol{\gamma})\|^2 = \|\boldsymbol{\gamma}\|^2 = 1$ such that

$$\begin{aligned} \|H^{-1} \{ \nabla_{\boldsymbol{\eta}} \zeta(\mathbf{v}) - \nabla_{\boldsymbol{\eta}} \zeta(\mathbf{v}^*) \}\| &= \sup_{\substack{\boldsymbol{\gamma} \in \mathbb{R}^m \\ \|\boldsymbol{\gamma}\|=1}} \boldsymbol{\gamma}^\top H^{-1} \{ \nabla_{\boldsymbol{\eta}} \zeta(\mathbf{v}) - \nabla_{\boldsymbol{\eta}} \zeta(\mathbf{v}^*) \} \\ &\leq \sup_{\substack{\boldsymbol{\gamma} \in \mathbb{R}^{p^*} \\ \|\boldsymbol{\gamma}\|=1}} \boldsymbol{\gamma}^\top \mathcal{D}^{-1} \{ \nabla_{\mathbf{v}} \zeta(\mathbf{v}) - \nabla_{\mathbf{v}} \zeta(\mathbf{v}^*) \} \\ &= \|\mathcal{D}^{-1} \{ \nabla_{\mathbf{v}} \zeta(\mathbf{v}) - \nabla_{\mathbf{v}} \zeta(\mathbf{v}^*) \}\| \\ &\leq 6\nu_1 \omega_3(\mathbf{x}, 4p^*) \mathbf{x}. \end{aligned}$$

As above we find with Taylor expansion

$$\begin{aligned} & \sup_{\mathbf{v} \in \mathcal{Y}_{o, \boldsymbol{\theta}^*}(\mathbf{x})} \|H^{-1} \{ \nabla_{\boldsymbol{\eta}} \mathbb{E} \mathcal{L}(\mathbf{v}) - \nabla_{\boldsymbol{\eta}} \mathbb{E} \mathcal{L}(\mathbf{v}^*) + H^2(\boldsymbol{\eta} - \boldsymbol{\eta}^*) \}\| \\ & \leq \sup_{\mathbf{v} \in \mathcal{Y}_{o, \boldsymbol{\theta}^*}(\mathbf{x})} \|H^{-1} H^2(\mathbf{v}) H^{-1} - \mathbb{I}_m\| \mathbf{r}. \end{aligned}$$

We can bound using $\|\mathcal{D}(0, H^{-1} \boldsymbol{\gamma})\|^2 = \|\boldsymbol{\gamma}\|^2$ and (\mathcal{L}_0)

$$\begin{aligned} \|H^{-1} H^2(\mathbf{v}) H^{-1} - \mathbb{I}_m\| &= \sup_{\substack{\boldsymbol{\gamma} \in \mathbb{R}^m \\ \|\boldsymbol{\gamma}\|=1}} (H^{-1} \boldsymbol{\gamma})^\top \{ H^2(\mathbf{v}) - H^2 \} H^{-1} \boldsymbol{\gamma} \\ &= \sup_{\substack{\boldsymbol{\gamma} \in \mathbb{R}^m \\ \|\boldsymbol{\gamma}\|=1}} (0, H^{-1} \boldsymbol{\gamma})^\top \{ \mathcal{D}^2(\mathbf{v}) - \mathcal{D}^2 \} (0, H^{-1} \boldsymbol{\gamma}) \end{aligned}$$

$$\leq \sup_{\substack{\gamma \in \mathbb{R}^p \\ \|\gamma\|=1}} \gamma \{ \mathcal{D}^{-1} \mathcal{D}^2(\mathbf{v}) \mathcal{D}^{-1} - \mathbb{I}_{p^*} \} \gamma \leq \delta(\mathbf{x}).$$

This gives the claim. □

Now we can proof Proposition 2.4. Define

$$\begin{aligned} C'(\mathbf{r}_0, \mathbf{x}) \stackrel{\text{def}}{=} & \left\{ \tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_o(\mathbf{r}_0), \right. & \text{(B.6)} \\ & \left. \|\mathcal{D}^{-1} \nabla\| \leq \mathfrak{z}(\mathbf{x}, B), \ \|H^{-1} \nabla_{\boldsymbol{\eta}}\| \leq \mathfrak{z}(\mathbf{x}, B), \right\} \\ & \cap \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}_0)} \|\mathcal{Y}(\mathbf{v})\| \leq 6\nu_1 \omega \mathfrak{z}(\mathbf{x}, 4p^*) \mathbf{r}_0 \right\} \\ & \cap \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_o(4\mathbf{r}_1)} \|\check{\mathcal{Y}}(\mathbf{v})\| \leq 6\nu_1 \check{\omega} \mathfrak{z}(\mathbf{x}, 2p^* + 2p) 4\mathbf{r}_1 \right\}. \end{aligned}$$

The desired result occurs on this set. First we show that $\mathbb{P}(C'(\mathbf{r}_0, \mathbf{x})) \geq 1 - 5e^{-x}$. Lemma B.4 yields

$$\|H^{-1} \nabla_{\boldsymbol{\eta}}\|^2 \leq \|\mathcal{D}^{-1} \nabla\|^2,$$

which implies that

$$\{\|\mathcal{D}^{-1} \nabla\| \leq \mathfrak{z}(\mathbf{x}, B)\} \subseteq \{\|H^{-1} \nabla_{\boldsymbol{\eta}}\| \leq \mathfrak{z}(\mathbf{x}, B)\}.$$

To control the probability $\mathbb{P}(\|\mathcal{D}^{-1} \nabla\| > \mathfrak{z}(\mathbf{x}, B))$ we apply Proposition A.1 with

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

We obtain

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

By Theorem C.1 with $p = p^*$ we have

$$\mathbb{P} \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}_0)} \|\mathcal{Y}(\mathbf{v})\| \leq 6\nu_1 \omega \mathfrak{z}(\mathbf{x}, 4p^*) \mathbf{r}_0 \right\} \geq 1 - e^{-x}.$$

This gives that $\mathbb{P}(C'(\mathbf{r}_0, \mathbf{x})) \geq 1 - 5e^{-x}$. Lemma B.3 gives that on the set $C'(\mathbf{r}_0, \mathbf{x})$ from (B.6) we have

$$\|\mathcal{D}(\tilde{\mathbf{v}} - \mathbf{v}^*) - \mathcal{D}^{-1} \nabla\| \leq \diamond(\mathbf{r}_0, \mathbf{x}).$$

With the triangular inequality this gives

$$\|\mathcal{D}(\tilde{\mathbf{v}} - \mathbf{v}^*)\| \leq \|\mathcal{D}^{-1} \nabla\| + \diamond(\mathbf{r}_0, \mathbf{x}).$$

Now on $C'(\mathbf{r}_0, \mathbf{x})$ we have $\|\mathcal{D}^{-1}\nabla\| \leq \mathfrak{z}(\mathbf{x}, B)$, which implies

$$\|\mathcal{D}(\tilde{\mathbf{v}} - \mathbf{v}^*)\| \leq \mathfrak{z}(\mathbf{x}, B) + \diamond(\mathbf{r}_0, \mathbf{x}).$$

The same can be done for $\|\mathcal{D}(\tilde{\mathbf{v}}_{\theta^*} - \mathbf{v}^*)\|$ which gives

$$\begin{aligned} C'(\mathbf{r}_0, \mathbf{x}) &\subseteq \{\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_o(\mathbf{r}_1)\} \cap \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}_1)} \|\check{\mathfrak{Y}}(\mathbf{v})\| \leq 6\nu_1 \check{\mathfrak{z}}(\mathbf{x}, p^*) \mathbf{r}_1 \right\} \\ &= C(\mathbf{r}_1, \mathbf{x}) \subset \Omega. \end{aligned}$$

Now the claim follows as in the proof of Theorem 2.2 with $\mathbf{r}_0 > 0$ replaced with $\mathbf{r}_1 > 0$.

Lemma B.4. Let $\mathcal{D} \in \mathbb{R}^{(p+p) \times (p+p)}$ be invertible and

$$\mathcal{D}^2 = \begin{pmatrix} D^2 & A \\ A^\top & H^2 \end{pmatrix} \in \mathbb{R}^{(p+p) \times (p+p)}, \quad D \in \mathbb{R}^{p \times p}, H \in \mathbb{R}^{m \times m} \text{ invertible.}$$

Then for any $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^{p+m}$ we have $\|H^{-1}\boldsymbol{\eta}\| \vee \|D^{-1}\boldsymbol{\theta}\| \leq \|\mathcal{D}^{-1}\mathbf{v}\|$.

Proof. With $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^{p+m}$

$$\|D^{-1}\boldsymbol{\theta}\| = \|D^{-1}\Pi_{\boldsymbol{\theta}}\mathcal{D}^{-1}\mathbf{v}\| \leq \|D^{-1}\Pi_{\boldsymbol{\theta}}\mathcal{D}\| \|\mathcal{D}^{-1}\mathbf{v}\| = \|\mathcal{D}^{-1}\mathbf{v}\|,$$

because

$$\|D^{-1}\Pi_{\boldsymbol{\theta}}\mathcal{D}\|^2 = \sup_{\|\boldsymbol{\gamma}\|=1} \boldsymbol{\gamma}^\top D^{-1}\Pi_{\boldsymbol{\theta}}\mathcal{D}^2\Pi_{\boldsymbol{\theta}}^\top D^{-1}\boldsymbol{\gamma} = \|\boldsymbol{\gamma}\| = 1.$$

The same argument works for $\|H^{-1}\boldsymbol{\eta}\|$. □

B.4. Proof of Lemma 2.5

Proof. First note that due to (l) we have

$$\|\mathcal{D}^{-1}\| = \frac{1}{\sqrt{n}} \|d^{-1}\| \leq \frac{1}{\sqrt{nc_d}}. \tag{B.7}$$

Now we prove the implications.

(\mathcal{L}_0) As by assumption $\mathcal{Y}_o(\mathbf{r}^*) \subset U$ we simply estimate using (B.7) and (ℓ_0) for any $\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}^*)$

$$\begin{aligned} \|\mathbb{I} - \mathcal{D}^{-1}\nabla^2\mathbb{E}\mathcal{L}(\mathbf{v})\mathcal{D}^{-1}\| &\leq \frac{1}{nc_d^2} \|\mathcal{D}^2 - \nabla^2\mathbb{E}\mathcal{L}(\mathbf{v})\| \\ &= \frac{1}{c_d^2} \|\nabla^2\mathbb{E}\ell(\mathbf{v}^*) - \nabla^2\mathbb{E}\ell(\mathbf{v})\| \leq \frac{\delta^*}{\sqrt{nc_d^3}} \mathbf{r}. \end{aligned}$$

($\mathcal{E}\mathcal{D}_1$) Abbreviate $\zeta_i = (\ell_i - \mathbb{E}\ell_i)$ and $\zeta = (\mathcal{L} - \mathbb{E}\mathcal{L})$. Take any $\gamma \in \mathbb{R}^{p^*}$ and $\mathbf{v}, \mathbf{v}' \in \Upsilon_\circ(\mathbf{x}^*) \subset U$ and use the mean value theorem to find some $\hat{\mathbf{v}} \in \text{conv}(\mathbf{v}, \mathbf{v}') \subset U$

$$\begin{aligned} & \log \mathbb{E} \exp \left\{ \frac{\mu \gamma^\top \mathcal{D}^{-1} \{ \nabla \zeta(\mathbf{v}) - \nabla \zeta(\mathbf{v}') \}}{\omega \|\mathcal{D}(\mathbf{v} - \mathbf{v}')\|} \right\} \\ &= \log \mathbb{E} \exp \left\{ \frac{\mu}{\omega n} \gamma^\top d^{-1} \left\{ \sum_{i=1}^n \nabla^2 \zeta_i(\hat{\mathbf{v}}) \right\} d^{-1} \frac{d(\mathbf{v} - \mathbf{v}')}{\|d(\mathbf{v} - \mathbf{v}')\|} \right\}. \end{aligned}$$

Using independence and (ed_1) this gives with $\omega = \frac{1}{\sqrt{n}}$ and $|\mu| \leq \sqrt{n} \mathbf{g}_0^*$

$$\begin{aligned} & \log \mathbb{E} \exp \left\{ \frac{\mu \gamma^\top \mathcal{D}^{-1} \{ \nabla \zeta(\mathbf{v}) - \nabla \zeta(\mathbf{v}') \}}{\omega \|\mathcal{D}(\mathbf{v} - \mathbf{v}')\|} \right\} \\ & \leq \sum_{i=1}^n \sup_{\substack{\gamma \in \mathbb{R}^{p^*} \\ \|\gamma\|=1}} \log \mathbb{E} \exp \left\{ \frac{\mu}{\sqrt{n}} \gamma_1^\top d^{-1} \nabla^2 \zeta_i(\hat{\mathbf{v}}) d^{-1} \gamma_2 \right\} \leq \nu_0^* \mu^2 / 2. \end{aligned}$$

Further (\mathcal{I}) is a consequence of Lemma B.5 and (ι). The other claims can be shown with the same argument or follow trivially from the setting. \square

Lemma B.5. For a positive definite symmetric matrix

$$\mathcal{D}^2 = \begin{pmatrix} D^2 & A \\ A^\top & H^2 \end{pmatrix},$$

with $c_{\mathcal{D}} \|\mathbf{v}\|^2 \leq \mathbf{v}^\top \mathcal{D} \mathbf{v}$ for some $c_{\mathcal{D}} > 0$ we have that

$$\|D^{-1} A H^{-2} A^\top D^{-1}\| =: \rho^2 \leq 1 - \frac{c_{\mathcal{D}}}{\|D\|^2 \wedge \|H\|^2}.$$

Proof. For any $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^{p+m}$ we have

$$\begin{aligned} \mathbf{v}^\top \mathcal{D}^2 \mathbf{v} &= (\boldsymbol{\theta}^\top, \boldsymbol{\eta}^\top) \begin{pmatrix} D^2 & A \\ A^\top & H^2 \end{pmatrix} \begin{pmatrix} \boldsymbol{\theta} \\ \boldsymbol{\eta} \end{pmatrix} \\ &= (\boldsymbol{\theta}^\top D^\top, \boldsymbol{\eta}^\top H^\top) \begin{pmatrix} I_p & D^{-1} A H^{-1} \\ H^{-1} A^\top D^{-1} & I_m \end{pmatrix} \begin{pmatrix} D \boldsymbol{\theta} \\ H \boldsymbol{\eta} \end{pmatrix} \\ &= \|D \boldsymbol{\theta}\|^2 + \|H \boldsymbol{\eta}\|^2 + 2 \langle H \boldsymbol{\eta}, H^{-1} A^\top D^{-1} \boldsymbol{\theta} \rangle. \end{aligned}$$

Minimized with respect to $\boldsymbol{\eta}$, i.e. with $H \boldsymbol{\eta} = -H^{-1} A^\top D^{-1} D \boldsymbol{\theta}$ we find

$$\mathbf{v}^\top \mathcal{D}^2 \mathbf{v} = \|D \boldsymbol{\theta}\|^2 - \|H^{-1} A^\top D^{-1} D \boldsymbol{\theta}\|^2 = (D \boldsymbol{\theta})^\top (I_p - D^{-1} A H^{-2} A^\top D^{-1}) D \boldsymbol{\theta},$$

which gets minimal – i.e. equal to $(1 - \rho^2)\|D\boldsymbol{\theta}\|$ - if

$$D^{-1}AH^{-2}A^{\top}D^{-1}D\boldsymbol{\theta} = \|D^{-1}AH^{-2}A^{\top}D^{-1}\|D\boldsymbol{\theta} = \rho^2D\boldsymbol{\theta},$$

i.e. if $D\boldsymbol{\theta} \in \mathbb{R}^p$ is a maximal eigenvalue of $D^{-1}AH^{-2}A^{\top}D^{-1} \in \mathbb{R}^{p \times p}$. With the assumption $c_{\mathcal{D}}\|\mathbf{v}\|^2 \leq \mathbf{v}^{\top}\mathcal{D}\mathbf{v}$ this gives

$$c_{\mathcal{D}}\|\mathbf{v}\|^2 \leq \mathbf{v}^{\top}\mathcal{D}^2\mathbf{v} = (1 - \rho^2)\|D\boldsymbol{\theta}\|^2, \quad \|\mathbf{v}\|^2 = \|\boldsymbol{\theta}\|^2 + \|H^{-2}A^{\top}\boldsymbol{\theta}\|^2,$$

such that

$$\rho^2 \leq 1 - c_{\mathcal{D}}\frac{\|\boldsymbol{\theta}\|^2}{\|D\boldsymbol{\theta}\|^2} \leq 1 - \frac{c_{\mathcal{D}}}{\|D\|^2}.$$

With analogous arguments we can obtain

$$\rho^2 \leq 1 - c_{\mathcal{D}}\frac{\|\boldsymbol{\eta}\|^2}{\|H\boldsymbol{\eta}\|^2} \leq 1 - \frac{c_{\mathcal{D}}}{\|H\|^2},$$

which completes the proof. □

B.5. Proof of Proposition 2.7

The profile MLE can be calculated easily

$$\tilde{\boldsymbol{\theta}} = \Pi_{\boldsymbol{\theta}}f^{-1}(\mathbf{Y}) = \Pi_{\boldsymbol{\theta}}f^{-1}(f(\mathbf{v}^*) + \boldsymbol{\varepsilon}_i) = \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}_{\boldsymbol{\theta}} - \|\boldsymbol{\varepsilon}_{\boldsymbol{\eta}}\|^2,$$

where $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_{\boldsymbol{\theta}}, \boldsymbol{\varepsilon}_{\boldsymbol{\eta}}) \in \mathbb{R} \times \mathbb{R}^{p_n-1}$. It is straight forward to show, that the conditions of Section 2.1 are satisfied with $\mathcal{D}^2 = n\mathbb{E}[\nabla f \nabla f^{\top}(\mathbf{v}^*)] = Id_{p^*}$, $\check{D}^2 = n$ and $\check{\boldsymbol{\xi}} = \sqrt{n}\boldsymbol{\varepsilon}_{\boldsymbol{\theta}}$. But we immediately see that

$$\sqrt{n}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \sqrt{n}\boldsymbol{\varepsilon}_{\boldsymbol{\theta}} = -\sqrt{n}\|\boldsymbol{\varepsilon}_{\boldsymbol{\eta}}\|^2 \sim \frac{-\chi_{p_n-1}^2}{\sqrt{n}}.$$

This means that if $p_n = O(n^{1/2})$ the estimator is not root-n consistent. For $\sqrt{n} = o(p_n)$ the root-n bias goes to infinity almost surely. Clearly if $p_n = o(n^{1/2})$ the Fisher expansion is accurate.

Concerning the Wilks phenomenon note that $\mathcal{L}(\tilde{\mathbf{v}}) = 0$. On the other hand

$$\begin{aligned} -\max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}^*, \boldsymbol{\eta}) &= n \min_{\lambda \in \mathbb{R}} \left\{ \left(y_{\boldsymbol{\theta}} - \lambda^2 \|\mathbf{Y}_{\boldsymbol{\eta}}\|^2 \right)^2 + (1 - \lambda)^2 \|\mathbf{Y}_{\boldsymbol{\eta}}\|^2 \right\} \\ &= n \min_{\lambda \in \mathbb{R}} \left\{ \left(\boldsymbol{\varepsilon}_{\boldsymbol{\theta}} - \lambda^2 \|\boldsymbol{\varepsilon}_{\boldsymbol{\eta}}\|^2 \right)^2 + (1 - \lambda)^2 \|\boldsymbol{\varepsilon}_{\boldsymbol{\eta}}\|^2 \right\} \\ &= n \min_{\lambda \in \mathbb{R}} \left\{ \boldsymbol{\varepsilon}_{\boldsymbol{\theta}}^2 + \|\boldsymbol{\varepsilon}_{\boldsymbol{\eta}}\|^2 \left(\lambda^4 \|\boldsymbol{\varepsilon}_{\boldsymbol{\eta}}\|^2 - \lambda^2 \boldsymbol{\varepsilon}_{\boldsymbol{\theta}} + (1 - \lambda)^2 \right) \right\}, \end{aligned}$$

where $\mathbf{Y} = (y_\theta, \mathbf{Y}_\eta) \in \mathbb{R} \times \mathbb{R}^{p_n-1}$ and $\boldsymbol{\varepsilon} = (\varepsilon_\theta, \boldsymbol{\varepsilon}_\eta) \in \mathbb{R} \times \mathbb{R}^{p_n-1}$. Now clearly $\|\boldsymbol{\varepsilon}_\eta\|^2 = O(p_n/n) \rightarrow 0$ a.s. and $\varepsilon_\theta \rightarrow 0$ a.s. such that the sequence of minimizers satisfies $\lambda_n \rightarrow 1$ a.s.. This gives for any $\tau > 0$ and $n \geq n_\tau \in \mathbb{N}$ large enough

$$-\max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}^*, \boldsymbol{\eta}) \geq n\varepsilon_\theta^2 + (1-\tau)n\|\boldsymbol{\varepsilon}_\eta\|^4 - (1+\tau)n\varepsilon_\theta\|\boldsymbol{\varepsilon}_\eta\|^2. \quad (\text{B.8})$$

Further we get setting $\lambda = 1$

$$-\max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}^*, \boldsymbol{\eta}) \leq n\varepsilon_\theta^2 + n\|\boldsymbol{\varepsilon}_\eta\|^4 - n\varepsilon_\theta\|\boldsymbol{\varepsilon}_\eta\|^2. \quad (\text{B.9})$$

As $\check{L}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\theta}}^*) = -\max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}^*, \boldsymbol{\eta})$ the inequalities (B.8) and (B.9) combine to

$$\begin{aligned} n\varepsilon_\theta^2 + n(1-\tau)\|\boldsymbol{\varepsilon}_\eta\|^4 - (1+\tau)n\varepsilon_\theta\|\boldsymbol{\varepsilon}_\eta\|^2 &\leq \check{L}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\theta}}^*) \\ &\leq n\varepsilon_\theta^2 + n\|\boldsymbol{\varepsilon}_\eta\|^4 - n\varepsilon_\theta\|\boldsymbol{\varepsilon}_\eta\|^2. \end{aligned}$$

This gives the Wilks phenomenon if $p_n^2/n \rightarrow 0$. Now if $p_n^2/n \rightarrow \infty$ the right hand side in (B.8) diverges since with $\tau = 1/2$

$$\begin{aligned} n\varepsilon_\theta^2 + \frac{n}{2}\|\boldsymbol{\varepsilon}_\eta\|^4 - 2n\varepsilon_\theta\|\boldsymbol{\varepsilon}_\eta\|^2 &\sim \chi_1^2 * (\chi_{p_n-1}^4/2n) * \{-\mathcal{N}(0,1)(2\chi_{p_n-1}^2/\sqrt{n})\} \\ &\xrightarrow{w} \delta_\infty. \end{aligned}$$

If $p_n^2/n \rightarrow \mathbb{C}$ then $\check{L}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\theta}}^*)$ can not converge to a χ^2 -distribution with one degree of freedom as one can let $\tau > 0$ tend 0. This completes the proof.

B.6. Proof of Theorem 2.9

Remember the definition

$$\tilde{\boldsymbol{\nu}}_{\boldsymbol{\theta}_m^*, m} = (\boldsymbol{\theta}_m^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}_m^*}) \stackrel{\text{def}}{=} \operatorname{argmax}_{\substack{\boldsymbol{v} \in \mathcal{Y} \\ \Pi_{\boldsymbol{\theta}} \boldsymbol{v} = \boldsymbol{\theta}_m^*}} \mathcal{L}_m(\boldsymbol{v}), \quad \tilde{\boldsymbol{\nu}}_{\boldsymbol{\theta}^*, m} = (\boldsymbol{\theta}_m^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*}) \stackrel{\text{def}}{=} \operatorname{argmax}_{\substack{\boldsymbol{v} \in \mathcal{Y} \\ \Pi_{\boldsymbol{\theta}} \boldsymbol{v} = \boldsymbol{\theta}^*}} \mathcal{L}_m(\boldsymbol{v}).$$

Define for some $0 < \mathbf{r}_0^\circ$

$$\begin{aligned} \mathcal{A}(\mathbf{x}, \mathbf{r}_0^\circ) &\stackrel{\text{def}}{=} \{\tilde{\boldsymbol{\nu}}_m, \tilde{\boldsymbol{\nu}}_{\boldsymbol{\theta}_m^*, m}, \tilde{\boldsymbol{\nu}}_{\boldsymbol{\theta}^*, m} \in \mathcal{Y}_{0,m}(\mathbf{r}_0^\circ)\} \\ &\cap \left\{ \sup_{\boldsymbol{v} \in \mathcal{Y}_\circ(4\mathbf{r}_0^\circ)} \|\check{\boldsymbol{y}}(\boldsymbol{v})\| \leq 6\nu_1 \check{\omega}_3(\mathbf{x}, 2p^* + 2p) 4\mathbf{r}_0^\circ \right\} \subset \Omega, \end{aligned}$$

with $\check{\boldsymbol{y}}(\boldsymbol{v}) \in \mathbb{R}^{p^*}$ from (B.1).

We prove this claim in a similar fashion as in Section B.2.1. With the function $l_m : \mathbb{R}^p \times \mathcal{Y} \rightarrow \mathbb{R}$ defined as in (B.3) with \mathcal{L} replaced by \mathcal{L}_m we can represent:

$$\check{L}_m(\boldsymbol{\theta}_m^*) - \check{L}_m(\boldsymbol{\theta}^*) = l_m(\boldsymbol{\theta}_m^*, \boldsymbol{\theta}_m^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}_m^*}) - l_m(\boldsymbol{\theta}^*, \boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*}),$$

Repeating the same arguments as in Section B.2.1 we obtain

$$\begin{aligned} \check{L}_m(\boldsymbol{\theta}_m^*) - \check{L}_m(\boldsymbol{\theta}^*) &\leq l_m(\boldsymbol{\theta}_m^*, \boldsymbol{\theta}_m^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}_m^*}) - l_m(\boldsymbol{\theta}^*, \boldsymbol{\theta}_m^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}_m^*}) \\ &= \check{\nabla}_{\boldsymbol{\theta}} \mathcal{L}_m(\mathbf{v}^*)(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*) - \|\check{D}_m(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*)\|^2/2 \\ &\quad + \check{\alpha}_m^*(\boldsymbol{\theta}_m^*, \boldsymbol{\theta}^*), \end{aligned}$$

where $\check{\alpha}_m^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \in \mathbb{R}$ is defined as

$$\begin{aligned} \check{\alpha}_m^*(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) &\stackrel{\text{def}}{=} l(\boldsymbol{\theta}_1, \boldsymbol{\theta}_m^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}_m^*}) - l(\boldsymbol{\theta}_2, \boldsymbol{\theta}_m^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}_m^*}) \\ &\quad - \nabla_{\boldsymbol{\theta}_1} l(\boldsymbol{\theta}^*, \boldsymbol{\theta}^*, \boldsymbol{\eta}^*)(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2) - \|\check{D}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)\|^2/2. \end{aligned}$$

and satisfies

$$\begin{aligned} \check{\alpha}_m^*(\boldsymbol{\theta}_m^*, \boldsymbol{\theta}^*) &\leq \|\check{D}_m(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*)\| \sup_{\boldsymbol{\theta} \in \Pi_{\boldsymbol{\theta}} \mathcal{Y}_\circ(4\mathbf{r}_0^\circ)} |\check{D}_m^{-1} \nabla_{\boldsymbol{\theta}_1} \check{\alpha}_m(\boldsymbol{\theta}, \boldsymbol{\theta}^*)| \\ &\leq \alpha(m) \check{\diamond}(2(1+\rho)4\mathbf{r}_0^\circ, \mathbf{x}), \end{aligned}$$

since $\mathcal{A}(\mathbf{x}, \mathbf{r}_0^\circ) \subseteq \{\tilde{\boldsymbol{\nu}}_{\boldsymbol{\theta}_m^*}, \tilde{\boldsymbol{\nu}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}_\circ(\mathbf{r}_0^\circ)\}$. With similar arguments for the lower bound this gives

$$2|\check{L}_m(\boldsymbol{\theta}_m^*) - \check{L}_m(\boldsymbol{\theta}^*)| \leq \alpha(m) \left(2\|\check{D}^{-1} \check{\nabla} \mathcal{L}_m(\mathbf{v}^*)\| + \alpha(m) + 2\check{\diamond}(\mathbf{r}_0^\circ, \mathbf{x}) \right).$$

The claim follows because the result (2.7) of Theorem 2.2 occurs on $\mathcal{A}(\mathbf{x}, \mathbf{r}_0^\circ) \subseteq C(\mathbf{x}, \mathbf{r}_0^\circ) \subset \Omega$. It remains to note that the set $\mathcal{A}(\mathbf{x}, \mathbf{r}_0^\circ) \subset \Omega$ is of probability greater $1 - 2e^{-x}$ by the choice of $\mathbf{r}_0^\circ > 0$.

B.7. Proof of Corollary 2.10

We will only prove the asymptotic normality as the the proof the Wilks phenomenon is very similar. Define

$$\begin{aligned} \mathcal{V}_m^2(\mathbf{v}_m^*) &= \text{Cov}(\nabla_{p+m} \mathcal{L}_m(\mathbf{v}_m^*)), \quad B_m = \mathcal{D}_m^{-1} \mathcal{V}_m^2 \mathcal{D}_m^{-1}, \\ \check{\nabla}_{\boldsymbol{\theta}, m} &= \nabla_{\boldsymbol{\theta}} - A_m H_m^{-2} \nabla_{\boldsymbol{\eta}}, \quad \check{V}_m^2 = \text{Cov}(\check{\nabla}_{\boldsymbol{\theta}} \zeta(\mathbf{v}_m^*)), \quad \check{B}_m = \check{D}_m^{-1} \check{V}_m^2 \check{D}_m^{-1}. \end{aligned}$$

Remember $p^* = p + m \in \mathbb{N}$ and that the point $\mathbf{v}_m^* \in \mathbb{R}^p \times \mathbb{R}^m$ is defined by maximizing the expected log-likelihood for the sieved functional models \mathcal{L}_m and the operators $\mathcal{D}_m^2 \in \mathbb{R}^{p^* \times p^*}$, $\check{D}_m^2 \in \mathbb{R}^{p \times p}$ correspond to this point, i.e. we abbreviate $\mathcal{D}_m^2 \stackrel{\text{def}}{=} \mathcal{D}_m^2(\mathbf{v}_m^*)$, while $\mathcal{D}^2 = \mathcal{D}^2(\mathbf{v}^*)$ and $\check{D}_m^2 = \check{D}_m^2(\mathbf{v}_m^*)$, $\check{D}^2 = \check{D}^2(\mathbf{v}^*)$, where $\mathbf{v}^* = \text{argmax}_{\mathbf{v} \in \mathcal{Y}} \mathbb{E} \mathcal{L}(\mathbf{v})$, i.e. the true full maximizer.

We get with Theorem 2.2 applied to $\tilde{\boldsymbol{\theta}}_m$ from (2.15) that with probability greater $1 - 2e^{-x}$

$$\|\check{D}_m(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*) - \check{\xi}_m(\mathbf{v}_m^*)\| \leq \diamond(\mathbf{r}_0, \mathbf{x}). \quad (\text{B.10})$$

We write

$$\begin{aligned} & \check{D}(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*) \\ &= \check{D}_m(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*) - \check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*) + (\check{D}_m - \check{D})(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*) + \check{D}_m(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*). \end{aligned}$$

By (B.10) it suffices to bound $\|(\check{D}_m - \check{D})(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*)\|$ and $\|\check{D}_m(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*)\|$. With assumption (*bias*) we get

$$\|\check{D}_m(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*)\| \leq \alpha(m).$$

Further

$$\begin{aligned} & \|(\check{D}_m - \check{D})(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*)\| \\ & \leq \|(\check{D}_m - \check{D}_m(\mathbf{v}^*))(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*)\| + \|(\check{D}_m(\mathbf{v}^*) - \check{D})(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*)\| \\ & \leq \|\check{D}_m(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*)\| \left(\|\mathbb{I} - \check{D}_m^{-1} \check{D}_m^2(\mathbf{v}^*) \check{D}_m^{-1}\|^{1/2} \right. \\ & \quad \left. + \|\mathbb{I} - \check{D}_m(\mathbf{v}^*)^{-1} \check{D}^2(\mathbf{v}^*) \check{D}_m(\mathbf{v}^*)^{-1}\|^{1/2} \|\check{D}_m(\mathbf{v}^*) \check{D}_m^{-1}\| \right). \end{aligned}$$

Condition ($\check{\mathcal{D}}$) yields that $\mathbb{P}(\|\check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| \leq \mathfrak{z}(\mathbf{x}_n, \check{B}_m)) \geq 1 - 2e^{-\mathbf{x}_n}$ (see Section A). This gives with (B.10) that with probability greater $1 - 4e^{-\mathbf{x}_n}$

$$\|\check{D}_m(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*)\| \leq \|\check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| + \diamond(\mathbf{r}_0, \mathbf{x}) \leq \mathfrak{z}(\mathbf{x}, \check{B}_m) + \diamond(\mathbf{r}_0, \mathbf{x}),$$

where $\mathfrak{z}(\mathbf{x}, \check{B}_m) = O(\sqrt{p + \mathbf{x}})$. Combining these bounds gives with (*bias'*)

$$\|\check{D}(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| \leq \diamond(\mathbf{r}_0, \mathbf{x}) + \beta(m) \left(\mathfrak{z}(\mathbf{x}, \check{B}_m) + \diamond(\mathbf{r}_0, \mathbf{x}) \right) + \alpha(m),$$

where $\mathbf{r}_0(\mathbf{x})$ is chosen such that $\mathbb{P}(\tilde{\mathbf{v}}_n, \tilde{\mathbf{v}}_{\boldsymbol{\theta}_m^*, m} \in \mathcal{Y}_{0,m}(\mathbf{r}_0(\mathbf{x}))) \geq 1 - e^{-\mathbf{x}}$. By assumption $\mathbf{r}_0(\mathbf{x}) < \infty$ for any $\mathbf{x} > 0$, $m, n \in \mathbb{N}$. Remember that $\check{\diamond}(\mathbf{r}_0, \mathbf{x}_n) \approx \check{\delta}_n(\mathbf{r}_0) \mathbf{r}_0 + \check{\omega}_n \sqrt{\mathbf{x} + p + m_n} \mathbf{r}_0$ where by assumption $\check{\delta}_n(\mathbf{r}) \rightarrow 0$ for any $\mathbf{r} > 0$ and $\check{\omega}_n \rightarrow 0$. This implies that there exist sequences $(m_n) \subset \mathbb{N}$ with $m_n \rightarrow \infty$ and $\mathbf{x}_n \rightarrow \infty$ with

$$\diamond(\mathbf{r}_0, \mathbf{x}) + \beta(m) \left(\mathfrak{z}(\mathbf{x}, \check{B}_{m_n}) + \diamond(\mathbf{r}_0, \mathbf{x}) \right) + \alpha(m_n) \rightarrow 0 \quad (\text{B.11})$$

as $n \rightarrow \infty$. Fix such sequences $m_n \rightarrow \infty$ and $\mathbf{x}_n \rightarrow \infty$. Then we have due to (B.11) that for any $\epsilon > 0$ there exists an $n \in \mathbb{N}$ such that

$$\mathbb{P}(\|\check{D}(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)\| \geq \epsilon) \leq 4e^{-\mathbf{x}_n}.$$

As $\mathbf{x}_n \rightarrow \infty$ we get the claim by Slutsky's Lemma once we showed that $\check{\boldsymbol{\xi}}_m(\mathbf{v}_m^*)$ is asymptotically $\mathcal{N}(0, \check{d}^{-1} \check{v}^2 \check{d}^{-1})$ -distributed.

For this observe

$$\begin{aligned} \check{\xi}_m(\mathbf{v}_m^*) &= \check{D}_m^{-1}(\nabla_{\boldsymbol{\theta}} - A_m H_m^{-2} \nabla_{\boldsymbol{\eta}}) \mathcal{L}(\mathbf{v}_m^*) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\frac{1}{\sqrt{n}} \check{D}_m \right)^{-1} (\nabla_{\boldsymbol{\theta}} \ell_i(\mathbf{v}_m^*) - A_m H_m^{-2} \nabla_{\boldsymbol{\eta}} \ell_i(\mathbf{v}_m^*)) \\ &\stackrel{\text{def}}{=} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i. \end{aligned}$$

Due to assumptions (*bias''*) we have $\text{Cov}(\mathbf{X}_i) \rightarrow \check{d}^{-1} \check{v}^2 \check{d}^{-1} \in \mathbb{R}^{p \times p}$. Consequently

$$\check{\xi}_m(\mathbf{v}_m^*) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i,$$

where the random vectors \mathbf{X}_i are i.i.d. with zero mean and covariance tending to $\check{d}^{-1} \check{v}^2 \check{d}^{-1}$, such that by a slightly generalized central limit theorem

$$\check{\xi}_m(\mathbf{v}_m^*) \xrightarrow{w} \mathcal{N}(0, \check{d}^{-1} \check{v}^2 \check{d}^{-1}).$$

Appendix C: A bound for the norm of a random process

We want to derive for a random process $\mathcal{Y}(\mathbf{v}) \in \mathbb{R}^p$ and $\mathbf{v} \in \mathcal{Y}_o(\mathbf{r}) \subset \mathbb{R}^{p^*}$ a bound of the kind

$$\mathbb{P} \left(\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| \geq \mathbf{C}_3(\mathbf{x}, 2p^* + 2p)\mathbf{r} \right) \leq e^{-\mathbf{x}}.$$

In the following we elaborate how to extend the results of the supplement of [29] on empirical processes to this situation without substantial changes to the bounds.

For this let $\mathcal{Y}(\mathbf{v})$ be a smooth centered random vector process with values in \mathbb{R}^p . We aim at bounding the maximum of the norm $\|\mathcal{Y}(\mathbf{v})\|$ over a vicinity $\mathcal{Y}_o(\mathbf{r}) \stackrel{\text{def}}{=} \{\|\mathbf{v} - \mathbf{v}^*\|_{\mathcal{Y}} \leq \mathbf{r}\}$ of \mathbf{v}^* with some norm $\|\cdot\|_{\mathcal{Y}}$. Suppose that $\mathcal{Y}(\mathbf{v})$ satisfies for each $0 < \mathbf{r} < \mathbf{r}^*$ and for all pairs $\mathbf{v}, \mathbf{v}^\circ \in \mathcal{Y}_o(\mathbf{r}) = \{\mathbf{v} \in \mathcal{Y} : \|\mathbf{v} - \mathbf{v}^*\|_{\mathcal{Y}} \leq \mathbf{r}\} \subset \mathbb{R}^{p^*}$ and $|\lambda| \leq \mathbf{g}$

$$\sup_{\|\mathbf{u}\| \leq 1} \log \mathbb{E} \exp \left\{ \lambda \frac{\mathbf{u}^\top (\mathcal{Y}(\mathbf{v}) - \mathcal{Y}(\mathbf{v}^\circ))}{\omega \|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}}} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}. \tag{C.1}$$

Remark C.1. In the setting of Theorem 2.2 and Proposition 2.4 we have

$$\mathcal{Y}(\mathbf{v}) = \check{D}^{-1} \left(\check{\nabla} \zeta(\mathbf{v}) - \check{\nabla} \zeta(\mathbf{v}^*) \right), \quad \mathcal{Y}(\mathbf{v}) = \mathcal{D}^{-1} \left(\nabla \zeta(\mathbf{v}) - \nabla \zeta(\mathbf{v}^*) \right),$$

respectively and in both cases the norm becomes $\|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}} = \|\mathcal{D}(\mathbf{v} - \mathbf{v}^\circ)\|$ and condition (C.1) becomes $(\check{\mathcal{E}}\mathcal{D}_1)$ from Section 2.1.

Theorem C.1. *Let a random p -vector process $\mathcal{Y}(\mathbf{v})$ fulfill $\mathcal{Y}(\mathbf{v}^*) = 0$ and the condition (C.1) be satisfied. Then for each $\mathbf{r} > 0$, on a set of probability greater $1 - e^{-\mathbf{x}}$*

$$\sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| \leq 6\omega\nu_1\mathfrak{z}(\mathbf{x}, 2p^* + 2p)\mathbf{r},$$

where with $\mathfrak{g}_0 = \nu_0\mathfrak{g}$ and for some $\mathbb{Q} > 0$

$$\mathfrak{z}(\mathbf{x}, \mathbb{Q}) \stackrel{\text{def}}{=} \begin{cases} \sqrt{2(\mathbf{x} + \mathbb{Q})} & \text{if } \sqrt{2(\mathbf{x} + \mathbb{Q})} \leq \mathfrak{g}_0, \\ \mathfrak{g}_0^{-1}(\mathbf{x} + \mathbb{Q}) + \mathfrak{g}_0/2 & \text{otherwise.} \end{cases} \quad (\text{C.2})$$

Remark C.2. Note that the entropy of the original set is increased by adding $p \in \mathbb{N}$ as the supremum is taken over $\mathcal{Y}_\circ(\mathbf{r}) \times \mathcal{B}_\mathbf{r}(0) \subset \mathbb{R}^{p^*} \times \mathbb{R}^p$.

Proof. In what follows, we use the representation

$$\|\mathcal{Y}(\mathbf{v})\| = \sup_{\|\mathbf{u}\| \leq \mathbf{r}} \frac{1}{\mathbf{r}} \mathbf{u}^\top \mathcal{Y}(\mathbf{v}).$$

This implies

$$\sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| = 2 \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \sup_{\|\mathbf{u}\| \leq \mathbf{r}} \frac{1}{2\mathbf{r}} \mathbf{u}^\top \mathcal{Y}(\mathbf{v}).$$

Due to Lemma C.2 the process $\mathcal{U}(\mathbf{v}, \mathbf{u}) \stackrel{\text{def}}{=} \frac{1}{2\mathbf{r}} \mathbf{u}^\top \mathcal{Y}(\mathbf{v})$ satisfies the condition (C.4) as process on $\mathbb{R}^{p^*} \times \mathbb{R}^p$. This allows to apply Corollary 2.2 of the supplement of [29] to obtain the desired result. We get on a set of probability greater $1 - e^{-\mathbf{x}}$

$$\begin{aligned} \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| &\leq 2 \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \sup_{\|\mathbf{u}\| \leq \mathbf{r}} \left\{ \frac{1}{2\mathbf{r}} \mathbf{u}^\top \mathcal{Y}(\mathbf{v}) \right\} \\ &\leq 6\nu_1\mathbf{r}\mathfrak{z}(\mathbf{x}, \mathbb{Q}(\mathcal{Y}_\circ(\mathbf{r}) \times \mathcal{B}_\mathbf{r}(0))). \end{aligned}$$

The constant $\mathbb{Q}(\mathcal{Y}_\circ(\mathbf{r}) \times \mathcal{B}_\mathbf{r}(0)) > 0$ quantifies the complexity of the set $\mathcal{Y}_\circ(\mathbf{r}) \times \mathcal{B}_\mathbf{r}(0) \subset \mathbb{R}^{p^*} \times \mathbb{R}^p$. We point out that for compact $M \subset \mathbb{R}^{p^*}$ we have $\mathbb{Q}(M) = 2p^*$ (see Supplement of [29], Lemma 2.10). This gives $\mathbb{Q}(\mathcal{Y}_\circ(\mathbf{r}) \times \mathcal{B}_\mathbf{r}(0)) = 2p^* + 2p$. \square

Lemma C.2. *Suppose that $\mathcal{Y}(\mathbf{v})$ satisfies for each $\|\mathbf{u}\| \leq 1$*

$$\sup_{\mathbf{v} \in \mathcal{Y}_\circ} \log \mathbb{E} \exp \left\{ \frac{\lambda(\mathcal{Y}(\mathbf{v}) - \mathcal{Y}(\mathbf{v}^\circ))^\top \mathbf{u}}{\omega \|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}}} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathfrak{g}. \quad (\text{C.3})$$

Then for any $\mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^p$ with $\|\mathbf{u}_i\|_{\mathcal{Y}} \leq 2\mathbf{r}$ and $\|\mathbf{u}_i\| \leq \mathbf{r}$

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{2\mathbf{r}} \frac{(\mathcal{Y}(\mathbf{v}))^\top \mathbf{u}_1 - \mathcal{Y}(\mathbf{v}^\circ)^\top \mathbf{u}_2}{\omega \sqrt{\|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}}^2 + \|\mathbf{u}_1 - \mathbf{u}_2\|^2}} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathbf{g}. \quad (\text{C.4})$$

Proof. We simply plug in the definition to find for $\mathbf{v}, \mathbf{v}^\circ \in \mathcal{T}_\circ(\mathbf{r})$

$$\begin{aligned} & \log \mathbb{E} \exp \left\{ \frac{\lambda}{2\mathbf{r}} \frac{\mathbf{u}_1^\top \mathcal{Y}(\mathbf{v}) - \mathbf{u}_2^\top \mathcal{Y}(\mathbf{v}^\circ)}{\omega \sqrt{\|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}}^2 + \|\mathbf{u}_1 - \mathbf{u}_2\|^2}} \right\} \\ &= \log \mathbb{E} \exp \left\{ \frac{\lambda}{2\mathbf{r}} \frac{\mathbf{u}_1^\top (\mathcal{Y}(\mathbf{v}) - \mathcal{Y}(\mathbf{v}^\circ)) + (\mathbf{u}_1^\top - \mathbf{u}_2^\top) \mathcal{Y}(\mathbf{v}^\circ)}{\omega \sqrt{\|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}}^2 + \|\mathbf{u}_1 - \mathbf{u}_2\|^2}} \right\}. \end{aligned}$$

By the Hölder inequality and (C.3) we infer

$$\begin{aligned} & \log \mathbb{E} \exp \left\{ \frac{\lambda}{2\mathbf{r}} \frac{\mathbf{u}_1^\top (\mathcal{Y}(\mathbf{v}) - \mathcal{Y}(\mathbf{v}^\circ)) + (\mathbf{u}_1^\top - \mathbf{u}_2^\top) \mathcal{Y}(\mathbf{v}^\circ)}{\omega \sqrt{\|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}}^2 + \|\mathbf{u}_1 - \mathbf{u}_2\|^2}} \right\} \\ & \leq \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\mathbf{r}} \frac{\mathbf{u}_1^\top (\mathcal{Y}(\mathbf{v}) - \mathcal{Y}(\mathbf{v}^\circ))}{\omega \|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}}} \right\} \\ & \quad + \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\mathbf{r}} \frac{(\mathbf{u}_1^\top - \mathbf{u}_2^\top) \mathcal{Y}(\mathbf{v}^\circ)}{\omega \|\mathbf{u}_1 - \mathbf{u}_2\|} \right\} \\ & \leq \sup_{\|\mathbf{u}\| \leq 1} \frac{1}{2} \log \mathbb{E} \exp \left\{ \lambda \frac{\mathbf{u}^\top (\mathcal{Y}(\mathbf{v}) - \mathcal{Y}(\mathbf{v}^\circ))}{\omega \|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}}} \right\} \\ & \quad + \sup_{\|\mathbf{u}\| \leq 1} \frac{1}{2} \log \mathbb{E} \exp \left\{ \lambda \frac{\mathbf{u}^\top (\mathcal{Y}(\mathbf{v}^\circ) - \mathcal{Y}(\mathbf{v}^*))}{\omega \|\mathbf{v}^\circ - \mathbf{v}^*\|_{\mathcal{Y}}} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad \lambda \leq \mathbf{g}. \end{aligned}$$

□

Acknowledgements

We are grateful to the referees and editors for very helpful remarks and comments.

References

- [1] ANDRESEN, A., Finite sample analysis of profile m-estimation in the single index model. ArXiv:[1406.4052](#), 2014.
- [2] ANDRESEN, A., A note on critical dimensions in profile semiparametric estimation. ArXiv:[1410.4709](#), 2014.
- [3] ANDRESEN, A., A note on the bias of sieve profile estimation. ArXiv:[1406.4045](#), 2014.

- [4] BERRY, A., The accuracy of the gaussian approximation to the sum of independent variates. *Transactions of the American Mathematical Society*, 49(1):122–136, 1941. [MR0003498](#)
- [5] BICKEL, P. J., KLAASSEN, C. A. J., RITOV, Y., and WELLNER, J. A., *Efficient and Adaptive Estimation for Semiparametric Models*. Springer, 1998. [MR1623559](#)
- [6] BOUCHERON, S. and MASSART, P., A high-dimensional Wilks phenomenon. *Probability Theory and Related Fields*, 150:405–433, 2011. [10.1007/s00440-010-0278-7](#). [MR2824862](#)
- [7] DELECROIX, M., HAERDLE, W., and HRISTACHE, M., Efficient estimation in single-index regression. Technical report, SFB 373, Humboldt Univ. Berlin, 1997.
- [8] FAN, J. and HUANG, T., Profile likelihood inferences on semiparametric varying-coefficient partially linear models. *Bernoulli*, 11(6):1031–1057, 2005. [MR2189080](#)
- [9] FAN, J., ZHANG, C., and ZHANG, J., Generalized likelihood ratio statistics and Wilks phenomenon. *Ann. Stat.*, 29(1):153–193, 2001. [MR1833962](#)
- [10] GHOSAL, S., Asymptotic normality of posterior distributions in high-dimensional linear models. *Bernoulli*, 5(2):315–331, 1999. [MR1681701](#)
- [11] GHOSAL, S., Asymptotic normality of posterior distributions for exponential families when the number of parameters tends to infinity. *J. Multivariate Anal.*, 74(1):49–68, 2000. [MR1790613](#)
- [12] HAERDLE, W., HALL, P., and ICHIMURA, H., Optimal smoothing in single-index models. *Ann. Statist.*, 21:157–178, 1993. [MR1212171](#)
- [13] HALL, P., *The Bootstrap and Edgeworth Expansion*. Springer, 1992. [MR1145237](#)
- [14] HUBER, P. J., The behavior of maximum likelihood estimates under non-standard conditions. Proc. 5th Berkeley Symp. Math. Stat. Probab., Univ. Calif. 1965/1966, 1, 221–233 (1967). [MR0216620](#)
- [15] IBRAGIMOV, I. A. and KHAS’MINSKIJ, R. Z., *Statistical Estimation. Asymptotic Theory. Transl. from the Russian by Samuel Kotz*. Springer-Verlag, New York–Heidelberg–Berlin, 1981. [MR0620321](#)
- [16] ICHIMURA, H., Semiparametric least squares (sls) and weighted sls estimation of single-index models. *J. Econometrics*, 58:71–120, 1993. [MR1230981](#)
- [17] KIM, Y., The Bernstein-von Mises theorem for the proportional hazard model. *Ann. Statist.*, 34(4):1678–1700, 2006. [MR2283713](#)
- [18] KOSOROK, M. R., *Introduction to Empirical Processes and Semiparametric Inference*. Springer in Statistics, 2005. [MR2724368](#)
- [19] MAMMEN, E., Asymptotics with increasing dimension for robust regression with applications to the bootstrap. *Ann. Stat.*, 17(1):382–400, 1989. [MR0981457](#)
- [20] MAMMEN, E., Bootstrap and wild bootstrap for high dimensional linear models. *Ann. Stat.*, 21(1):255–285, 1993. [MR1212176](#)
- [21] MAMMEN, E., Empirical process of residuals for high-dimensional linear models. *Ann. Stat.*, 24(1):307–335, 1996. [MR1389892](#)

- [22] MURPHY, S. A. and VAN DER VAART, A. W., On profile likelihood. *Journal of the American Statistical Association*, 95(450):449–465, 2000. [MR1803168](#)
- [23] MURPHY, S. A. and VAN DER VAART, A. W., Observed information in semi-parametric models. *Bernoulli*, 5(3):381–412, 1999. [MR1693616](#)
- [24] NEWEY, W. K., Convergence rates and asymptotic normality for series estimators. *Journal of Econometrics*, 79(1):147–168, 1997.
- [25] PORTNOY, S., Asymptotic behavior of m estimators of p regression parameters when p^2/n is large: Ii normal approximation. *The Annals of Statistics*, 13(4):1403–1417, 1985. [MR0811499](#)
- [26] PORTNOY, S., Asymptotic behavior of M-estimators of p regression parameters when p^2/n is large. I. Consistency. *Ann. Stat.*, 12:1298–1309, 1984. [MR0760690](#)
- [27] PORTNOY, S., Asymptotic behavior of the empiric distribution of M-estimated residuals from a regression model with many parameters. *Ann. Stat.*, 14:1152–1170, 1986. [MR0856812](#)
- [28] SHEN, J., SHI, X., Sieve likelihood ratio inference on general parameter space. *Science in China*, 48(1):67–78, 2005. [MR2156616](#)
- [29] SPOKOINY, V., Parametric estimation. Finite sample theory. *Ann. Statist.*, 40(6):2877–2909, 2012. [MR3097963](#)
- [30] SPOKOINY, V., Bernstein–von Mises theorem for growing parameter dimension. Manuscript. ArXiv:[1302.3430](#), 2013.
- [31] SPOKOINY, V., WANG, W., and HÄRDLE, W., Local quantile regression (with rejoinder). *J. of Statistical Planning and Inference*, 143(7):1109–1129, 2013. ArXiv:[1208.5384](#). [MR3049611](#)
- [32] VAN DER VAART, A. W. and WELLNER, J., *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Series in Statistics. Springer, 1996. [MR1385671](#)
- [33] WILKS, S. S., The large-sample distribution of the likelihood ratio for testing composite hypotheses. *Ann. Math. Statist.*, 9(1):60–62, 1938.
- [34] ZAITSEV, A., BURNAEV, E., and SPOKOINY, V., Properties of the posterior distribution of a regression model based on gaussian random fields. *Automation and Remote Control*, 74(10):1645–1655, 2013. [MR3219856](#)