

SFB
823

Sieve maximum likelihood estimation in a semi-parametric regression model with errors in variables

Denis Belomestny, Egor Klochkov,
Vladimir Spokoiny

Nr. 15/2016

Discussion Paper



Sieve maximum likelihood estimation in a semi-parametric regression model with errors in variables

Denis Belomestny*

Duisburg-Essen University, Germany
and National Research University Higher School of Economics, Russian Federation
`denis.belomestny@uni-due.de`

Egor Klochkov

Humboldt University Berlin,
IRTG 1792
`eklochkov@gmail.com`

Vladimir Spokoiny

Weierstrass-Institute,
Humboldt University Berlin and IITP RAS,
Mohrenstr. 39, 10117 Berlin, Germany,
`spokoiny@wias-berlin.de`

March 29, 2016

Abstract

The paper deals with a semi-parametric regression problem under deterministic and regular design which is observed with errors. We first linearise the problem using a sieve approach and then apply the total penalised maximum likelihood estimator to the linearised model. Sufficient conditions for \sqrt{n} -consistency and efficiency under parametric assumption are derived and a possible misspecification bias under different smoothness assumptions on the design is analysed. The Monte Carlo simulations show the performance of the estimator with simulated data.

AMS 2000 Subject Classification: Primary 62F10. Secondary

Keywords: errors-in-variables model, regression, \sqrt{n} -consistency

*The author's research is supported by the Deutsche Forschungsgemeinschaft through the SFB 823 "Statistical modelling of nonlinear dynamic processes"

1 Introduction

Consider a classical regression model

$$Y_i = f(X_i) + \varepsilon_i, \quad i = 1, \dots, n, \quad (1.1)$$

which describes the relation between the response variable Y and the independent variable $X \in \mathbb{R}^d$. Here Y_i , $i = 1, \dots, n$, are independent observations of Y and ε_i , $i = 1, \dots, n$, stand for the measurement errors, which are assumed to be i.i.d. random variables. We consider the situation when the regressors X_i are observed with some errors as well. Measurement error (or errors-in-variables) models have been extensively studied over the last decades, see, e.g., the monographs of Schneeweiß and Mittag (1986), Fuller (2009), Cheng and Van Ness (1999), Wansbeek and Meijer (2000), and Carroll et al. (1995). In particular, the last book deals almost exclusively with nonlinear regression models. According to the usual interpretation of measurement error models, Y is an observable variable depending on the regressor vector X , here via a nonlinear regression function f . The variable X , however, is not directly observable. Instead a perturbed vector $Z \in \mathbb{R}^d$ is observed, which is related to X via

$$Z_i = X_i + \sigma \nu_i, \quad i = 1, \dots, n,$$

where σ^2 stands for the noise variance in the regressors Z_i , and ν_i are standardised random errors. The aim is to estimate the regression function f from the observations (Y_i, Z_i) . The model (1.1) obviously faces an identification problem: if the distribution of the errors-in-regressors ν_i is unknown, the function f cannot be recovered consistently. Some information about the distribution of the ν_i can be helpful in this context. Indeed, taking the expectation of (1.1) reveals that the conditional mean of Y_i is a convolution of the regression function f with some kernel corresponding to the distribution of the ν_i :

$$E[Y | X] = \int f(X + \mathbf{x}) \varphi_\sigma(\mathbf{x}) d\mathbf{x}$$

Here $\varphi(\mathbf{x})$ is the density of the ν_i 's and $\varphi_\sigma(\mathbf{x}) = \sigma^{-1} \varphi(\mathbf{x}/\sigma)$. To recover f , one has to make a deconvolution which leads to an ill-posed inverse problem (see, e.g. Fan et al. (1991)). Note that the classical assumption that (ν_i) are normally distributed is in some sense the worst-case and the function f can only be recovered with a $\log n$ accuracy in this case. We refer to Butucea and Taupin (2008) for the detailed discussion and the overview of the literature on this problem. In this paper we suggest a new approach to this problem which treats the unobserved regressors (X_i) as a high-dimensional nuisance parameter. Such a model is characterised by an increasing number of nuisance parameters

and a fixed number of parameters of interest. It is well known that in models with an increasing number of nuisance parameters usual estimation procedures may fail to be consistent. The main problem tackled here is the elimination of the nuisance parameters in such a way that the estimating procedure delivers \sqrt{n} -consistent estimators. While in linear models the total least squares approach gives, under some conditions on the design, consistent and efficient estimators for the parameter of interest, the situation is much more involved in the case of nonlinear models.

For simplicity we assume below that $d = 1$ and the function f is univariate. We apply the linear sieve approach which approximates the target function f as a linear combination of the fixed basis functions $\Psi(x) = (\psi_1(x), \dots, \psi_p(x))$:

$$f(x) \approx f(x, \boldsymbol{\theta}) = \Psi(x)^\top \boldsymbol{\theta} = \sum_{m=1}^p \theta_m \psi_m(x).$$

The function $f(x, \boldsymbol{\theta})$ is described by the vector $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^\top$ and the estimation problem for the function f is equivalent to estimation of $\boldsymbol{\theta}$. Suppose for a moment that the errors ε_i are i.i.d. standard normal and the errors ν_i are i.i.d. zero mean normal with the variance σ^2 . Then the log-likelihood for the full parameter $\mathbf{v} = (\boldsymbol{\theta}, \mathbf{X})$ with $\mathbf{X} = (X_1, \dots, X_n)$ reads as follows:

$$L(\mathbf{v}) = L(\boldsymbol{\theta}, \mathbf{X}) = -\frac{1}{2} \|\mathbf{Y} - \Psi(\mathbf{X})^\top \boldsymbol{\theta}\|^2 - \frac{1}{2\sigma^2} \|\mathbf{Z} - \mathbf{X}\|^2$$

Here $\Psi(\mathbf{X})$ is the $p \times n$ -matrix with entries $\psi_m(X_i)$ for $m = 1, \dots, p$ and $i = 1, \dots, n$. Even if $d = 1$, the parameter dimension is $p+n$ and it is larger than the sample size. So we need to reduce the complexity of estimation of the high-dimensional nuisance vector \mathbf{X} . Let $\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_n$ are n orthonormal vectors in \mathbb{R}^n , i.e., it holds $\boldsymbol{\phi}_i^\top \boldsymbol{\phi}_j = \delta_{ij}$, then we can always represent the vector \mathbf{X} in the form

$$\mathbf{X} = \eta_1 \boldsymbol{\phi}_1 + \dots + \eta_n \boldsymbol{\phi}_n$$

for an unknown vector $\boldsymbol{\eta} \in \mathbb{R}^n$. We assume that \mathbf{X} can be approximately spanned by first $q \ll n$ of $\boldsymbol{\phi}_j$, i.e.

$$\mathbf{X} \approx \eta_1 \boldsymbol{\phi}_1 + \dots + \eta_q \boldsymbol{\phi}_q = \Phi \boldsymbol{\eta}, \quad \boldsymbol{\eta} \in \mathbb{R}^q, \quad \Phi = (\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_q). \quad (1.2)$$

This leads to the following quasi-log-likelihood parametrized by $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta})$

$$L_\Phi(\mathbf{v}) = L_\Phi(\boldsymbol{\theta}, \boldsymbol{\eta}) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}, \Phi \boldsymbol{\eta}) = -\frac{1}{2} \|\mathbf{Y} - \Psi(\Phi \boldsymbol{\eta})^\top \boldsymbol{\theta}\|^2 - \frac{1}{2\sigma^2} \|\mathbf{Z} - \Phi \boldsymbol{\eta}\|^2, \quad (1.3)$$

with the total dimension $p+q$ and dimension of the target parameter p . Note, for example, that the equality $\mathbf{X} = \Phi \boldsymbol{\eta}$ with $\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_q$ being the first q vectors of Haar

basis means the aggregation of points, which is equivalent to *repeated measurements* model Fuller (2009). In general $\mathbf{X} \neq \Phi\boldsymbol{\eta}$ and we need some conditions that allow to control bias caused by dimension reduction.

2 Main results

Let us consider the following auxiliary semiparametric model

$$\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}, \quad \mathbf{Z} = \mathbf{X}^* + \sigma\boldsymbol{\nu}, \quad (2.1)$$

with $\mathbf{f}^* = (f(X_1^*), \dots, f(X_n^*))^\top$, where $\mathbf{X}^* = (X_1, \dots, X_n)$ is the true (unknown) design. Let us approximate f via $f(x, \boldsymbol{\theta}) = \sum_{m=1}^p \theta_m \psi_m(x)$ with $\boldsymbol{\theta} \in \mathbb{R}^p$. The model (2.1) can be viewed as a sieve approximation for the original model (1.1). In the next sections we study the properties of this sieve maximum likelihood approach.

Note that non-concavity of the log-likelihood (1.3) is a problem from both optimization and statistical points of view. In Section 3 we suggest to conduct an iterative optimization procedure starting from the *plug-in* estimator:

$$\begin{aligned} \tilde{\mathbf{v}}^{(\text{pl})} &= (\tilde{\boldsymbol{\theta}}^{(\text{pl})}, \tilde{\boldsymbol{\eta}}^{(\text{pl})}) \\ \tilde{\boldsymbol{\theta}}^{(\text{pl})} &= (\tilde{\boldsymbol{\Psi}}\tilde{\boldsymbol{\Psi}}^\top)^{-1}\tilde{\boldsymbol{\Psi}}\mathbf{Y}, \quad \tilde{\boldsymbol{\Psi}} = \boldsymbol{\Psi}(\Phi\tilde{\boldsymbol{\eta}}^{(\text{pl})}), \\ \tilde{\boldsymbol{\eta}}^{(\text{pl})} &= \Phi^\top\mathbf{Z}. \end{aligned}$$

This is typically a good initial estimator of the full parameter $\mathbf{v}^* = (\boldsymbol{\theta}^*, \boldsymbol{\eta}^*)$. However, in order to proceed with semiparametric theory, we need to have at least one stationary point of the likelihood, which is close to the true point with high probability, see Section 4.7. Since $\tilde{\mathbf{v}}^{(\text{pl})}$ is usually already close to \mathbf{v}^* , we suggest to define the estimator as the closest stationary point of the likelihood to plug-in $\tilde{\mathbf{v}}^{(\text{pl})}$:

$$\tilde{\mathbf{v}} = \arg \min_{\substack{\mathbf{v}=(\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^{p+q}: \\ \nabla L_\Phi(\mathbf{v})=0}} \|\tilde{\boldsymbol{\Psi}}^\top(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}^{(\text{pl})})\|^2 + \sigma^{-2}\|\boldsymbol{\eta} - \tilde{\boldsymbol{\eta}}^{(\text{pl})}\|^2, \quad (2.2)$$

and, similarly,

$$\tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*} = \arg \min_{\substack{\boldsymbol{\eta} \in \mathbb{R}^q: \\ \nabla_{\boldsymbol{\eta}} L_\Phi(\boldsymbol{\theta}^*, \boldsymbol{\eta})=0}} \|\boldsymbol{\eta} - \tilde{\boldsymbol{\eta}}^{(\text{pl})}\|.$$

2.1 Consistency and efficiency

Suppose that $\boldsymbol{\theta}^* \in \mathbb{R}^p$, $\boldsymbol{\eta}^* \in \mathbb{R}^q$ are two vectors and (ψ_m) , (ϕ_m) are two sieves such that the following conditions hold.

(i) Suppose that for each $i = 1, \dots, n$, and each $s = 0, 1, 2$ it holds

$$|f^{(s)}(X_i)| = |\Psi^{(s)}(X_i)^\top \boldsymbol{\theta}^*| = \left| \sum_{m=1}^p \theta_m^* \psi_m^{(s)}(X_i) \right| \leq \mathbf{c}_{f,s}.$$

(ii) Assume that the eigenvalues of the matrix

$$\frac{\Psi(\Phi \boldsymbol{\eta}^*) \Psi(\Phi \boldsymbol{\eta}^*)^\top}{n}$$

belong to the interval $[\mathbf{f}, \mathbf{F}]$ for fixed positive constants \mathbf{f} and \mathbf{F} ;

(iii) For each $s = 0, 1, 2$, there are some fixed constants $\mu_{1,s} \leq \mu_{2,s} \leq \dots \leq \mu_{p,s}$ such that each function $\psi_m(x)$ fulfills for all x

$$|\psi_m^{(s)}(x)| \leq \mu_{m,s}. \quad (2.3)$$

(iv) The orthogonal vectors $\boldsymbol{\phi}_j$ fulfil $\|\boldsymbol{\phi}_j\| = 1$ and

$$\max_{i=1, \dots, n} \sum_{j=1}^q \phi_{ij}^2 \leq \frac{\mathbf{v}_q^2}{n},$$

where \mathbf{v}_q , as we will see, typically depends on q .

(v) The error of approximating \mathbf{f}^* via $\Psi(\Phi \boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*$ is not large, in the sense that

$$\square \stackrel{\text{def}}{=} \|\mathbf{f}^* - \Psi(\Phi \boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\| \leq \sqrt{n}.$$

(vi) Errors $(\varepsilon_i), (\nu_i)$ are zero mean mutually independent and satisfy for some $0 < \nu_\varepsilon \leq 1$, $i = 1, \dots, n$, and each λ

$$\log \mathbb{E} \exp(\lambda \varepsilon_i) \leq \frac{\lambda^2 \nu_\varepsilon^2}{2}, \quad \log \mathbb{E} \exp(\lambda \nu_i) \leq \frac{\lambda^2 \nu_\varepsilon^2}{2},$$

which implies $\max\{\text{Var}(\boldsymbol{\varepsilon}), \text{Var}(\boldsymbol{\nu})\} \leq \nu_\varepsilon^2 \mathbf{I}_n$.

(vii) With \mathbf{f} from (ii) and $\mu_{m,s}$ from (iii), it holds for some $\mathbf{w}_{p,2} \geq \mathbf{w}_{p,1} \geq \mathbf{w}_{p,0} > 0$,

$$\mathbf{f}^{-1} \sum_{m=1}^p \mu_{m,s}^2 \leq \mathbf{w}_{p,s}^2, \quad s = 0, 1, 2, \quad p \in \mathbb{N}.$$

Moreover,

$$\begin{aligned} \sigma \mathbf{w}_{p,1} \sqrt{p+q+\mathbf{x}} &\lesssim \sqrt{n}, & \square &\lesssim \mathbf{w}_{p,1}^{-1} \sqrt{n}, \\ \square &\leq (\mathbf{w}_{p,2} \mathbf{v}_q)^{-1/2} \sqrt{n}, & \square &\leq \mathbf{v}_q^{-1} \sqrt{n}, \end{aligned}$$

where expression “ \lesssim ” means inequality up to a small constant depending on σ , $C_{f,1}$, $C_{f,2}$, \mathbf{f} and ν_ε only. It is worth to mention, that in typical situations $w_{p,s} \asymp p^{s+\frac{1}{2}}$ and the two last inequalities are likely to be fulfilled, while the first two require n to be large enough relatively to p and q .

Main notations The semiparametric *Fisher information* matrix reads as follows,

$$\check{D}_{\theta\theta}^2 \doteq D_{\theta\theta}^2 - A_{\theta\eta} D_{\eta\eta}^{-2} A_{\theta\eta}^\top, \quad (2.4)$$

where

$$\begin{aligned} D_{\theta\theta}^2 &\doteq \Psi(\Phi\eta^*) \Psi(\Phi\eta^*)^\top, \\ A_{\theta\eta} &\doteq \Psi(\Phi\eta^*) \text{diag}\{\Psi'(\Phi\eta^*)^\top \theta^*\} \Phi, \\ D_{\eta\eta}^2 &\doteq \sigma^{-2} I_q + \Phi^\top \text{diag}\{\Psi'(\Phi\eta^*)^\top \theta^*\}^2 \Phi. \end{aligned}$$

Set

$$\diamond(\mathbf{x}) \asymp \sigma \left(w_{p,1} + v_q + w_{p,2} v_q \sqrt{n^{-1}(p+q+\mathbf{x})} \right) (p+q+\mathbf{x}) n^{-1/2} + \square. \quad (2.5)$$

Theorem 2.1. *Under the above assumptions, there exists random vector $\check{\xi} \in \mathbb{R}^p$ satisfying for each $\mathbf{x} \geq 0$,*

$$\mathbb{P} \left(\|\check{\xi}\| > \sqrt{p} + \sqrt{2\mathbf{x}} \right) \leq 2e^{-\mathbf{x}},$$

such that the following statements hold.

- If $\diamond^2(\mathbf{x})/p \rightarrow 0$, then $\tilde{\theta}$ is root- n consistent and $\mathbb{E}\|\tilde{\theta} - \theta^*\|^2 \asymp p/n$.
- If $\diamond(\mathbf{x}) \rightarrow 0$, then $\tilde{\theta}$ is root- n asymptotically normal and the Fisher expansion holds

$$\|\check{D}_{\theta\theta}(\tilde{\theta} - \theta^*) - \check{\xi}\| \leq \diamond(\mathbf{x})$$

on a set with probability at least $1 - 6e^{-\mathbf{x}}$.

- If $p \diamond^2(\mathbf{x}) \rightarrow 0$, then the Wilks expansion holds

$$|L_\Phi(\tilde{\theta}, \tilde{\eta}) - L_\Phi(\theta^*, \tilde{\eta}_{\theta^*}) - \|\check{\xi}\|^2/2| \leq \sqrt{p} \diamond(\mathbf{x})$$

on a set with probability at least $1 - 6e^{-\mathbf{x}}$.

Discussion Note, that the estimator $\tilde{\boldsymbol{\theta}}$ depends on the matrix Φ through the linear subspace $\{\Phi\boldsymbol{\eta} : \boldsymbol{\eta} \in \mathbb{R}^q\}$. It is reasonable to expect that the condition (iv) is actually a condition for the span of columns of Φ . Indeed, if we have another orthonormal matrix $\tilde{\Phi}$, that generates the same subspace, then there is S such that

$$\tilde{\Phi} = \Phi S, \quad S^\top S = I.$$

Introduce

$$V(\Phi) \stackrel{\text{def}}{=} n \max_{i=1, \dots, n} \|\Phi_i\|^2,$$

where Φ_1, \dots, Φ_n are the rows of the matrix Φ . Then, the rows of matrix $\tilde{\Phi}$ are $\Phi_i S$ and we have

$$V(\tilde{\Phi}) = n \max_{i=1, \dots, n} \|\Phi_i S\|^2 = V(\Phi), \quad (2.6)$$

since S is orthogonal. Thus, the value $V(\Phi)$ remains the same for all orthogonal matrices Φ representing the same subspace. Now let us analyse the value \mathfrak{v}_q from the condition (iv). Since Φ is orthonormal, we have $\text{tr}(\Phi^\top \Phi) = q$. On the other hand,

$$n \max_{i \leq n} \|\Phi_i\|^2 \geq \sum_{i=1}^n \|\Phi_i\|^2 = \text{tr}(\Phi \Phi^\top),$$

implying

$$\mathfrak{v}_q \geq \sqrt{q}. \quad (2.7)$$

Below we show two examples with $\mathfrak{v}_q \asymp \sqrt{q}$.

Example 2.1. A possible representation of a “smooth” design \mathbf{X}^* can be given by

$$X_i^* = g(i/n), \quad g : [0, 1] \rightarrow \mathbb{R}, \quad i = 1, \dots, n, \quad (2.8)$$

with smoothness of the function $g(t)$ corresponding to that of \mathbf{X}^* . Suppose we are given an orthonormal basis $\{h_m(t)\}_{m=1}^\infty$ with $\langle h_i, h_j \rangle_{L^2[0,1]} = \delta_{ij}$. Consider the corresponding sieve approximation of the function g :

$$g(t) \approx \sum_{m=1}^q \eta_m h_m(t).$$

This corresponds to Φ in (1.2) defined as follows,

$$\phi_m = \sqrt{\frac{1}{n}} \left(h_m \left(\frac{1}{n} \right), \dots, h_m(1) \right), \quad \Phi = (\phi_1, \phi_2, \dots, \phi_q). \quad (2.9)$$

These vectors satisfy condition (iv), but fail to be orthogonal. Nevertheless, due to orthogonality of the basis $(h_m(t), m = 1, \dots, q)$ under mild conditions we can orthogonalize Φ without much loss in $V(\Phi)$.

Lemma 2.2. *Suppose, that each h_m is twice differentiable and*

$$\|h_m^{(s)}\|_\infty \leq A_s m^s, \quad s = 0, 1, 2.$$

If $A_0^2 q/n + (A_0 A_2 + A_1^2) q^3 / (6n^2) \leq \frac{1}{2}$, then the matrix Φ in (2.9) can be transformed into $\tilde{\Phi}$, such that $\tilde{\Phi}^\top \tilde{\Phi} = I_q$ and $V(\tilde{\Phi}) \leq 2V(\Phi)$.

Example 2.2. Another example is given by Haar basis, which might be relevant to image processing, when the function $g(t)$ in (2.8) is piecewise smooth. Suppose that $n = 2^k$ and consider the sequence of partitions $\mathcal{S}^{(0)}, \mathcal{S}^{(1)}, \dots, \mathcal{S}^{(k)}$ of the set of indices $\mathcal{N} = \{1, 2, \dots, n\}$. First, we have $\mathcal{S}^{(0)} = \{A_1^{(0)} = \mathcal{N}\}$ — the whole set. Further, $\mathcal{S}^{(l+1)}$ is made by splitting each $A_j^{(l)} \in \mathcal{S}^{(l)}$ into two equal parts $A_{2j-1}^{(l+1)}$ and $A_{2j}^{(l+1)}$. Denoting $\mathcal{S}^{(l)} = \{A_1^{(l)}, \dots, A_{2^l}^{(l)}\}$, we have

$$\begin{aligned} A_1^{(1)} &= \{1, \dots, n/2\}, & A_1^{(2)} &= \{1, \dots, n/4\}, & \dots \\ A_2^{(1)} &= \{n/2 + 1, \dots, n\}, & A_2^{(2)} &= \{n/4 + 1, \dots, n/2\}, & \dots \\ & & A_3^{(2)} &= \{n/2 + 1, \dots, 3n/4\}, \\ & & A_4^{(2)} &= \{3n/4 + 1, \dots, n\}, \end{aligned}$$

Continuing the sequence we get for $l \leq k-1$ and $1 \leq r \leq 2^l$ that $A_r^{(l)} = \{(r-1)n/2^l + 1, \dots, rn/2^l\}$. Note, that $\#A_r^{(l)} = n2^{-l}$. Now, denoting by $\phi(i)$ the i th component of the vector ϕ , introduce

$$\phi_1(i) = \frac{1}{\sqrt{n}}, \quad \phi_{2^{l-1}+r}(i) = \sqrt{\frac{2^{l-1}}{n}} \begin{cases} 1, & i \in A_{2^{r-1}}^{(l)}, \\ -1, & i \in A_{2^r}^{(l)}, \\ 0, & \text{otherwise} \end{cases}$$

where, obviously, each $j = 2, \dots, n$, is uniquely representable as $j = 2^{l-1} + r$ with $1 \leq l \leq k$ and $1 \leq r \leq 2^{l-1}$. By the construction, $(\phi_j)_{j=1}^n$ is an orthonormal basis in \mathbb{R}^n . Moreover, it's easy to see, that for $q = 2^l$ the span of vectors ϕ_1, \dots, ϕ_q is

$$\{\mathbf{X} \in \mathbb{R}^n : X_i \text{ is constant over each } A \in \mathcal{S}^{(l)}\},$$

and that the projector $\Phi\Phi^\top$ on that subspace, when applied to $\mathbf{X} \in \mathbb{R}^n$, averages values X_i over $i \in A$ for each $A \in \mathcal{S}^{(l)}$. By (2.6) we can consider any other orthonormal $\tilde{\Phi}$,

generating the same subspace. One of them is

$$\tilde{\Phi} = \sqrt{\frac{2^l}{n}} \begin{pmatrix} \overbrace{1 \ 1 \ \dots \ 1}^{n/2^l} & \overbrace{0 \ 0 \ \dots \ 0}^{n/2^l} & \dots & \overbrace{0 \ 0 \ \dots \ 0}^{n/2^l} \\ 0 \ 0 \ \dots \ 0 & \overbrace{1 \ 1 \ \dots \ 1}^{n/2^l} & \dots & 0 \ 0 \ \dots \ 0 \\ & & \dots & \\ 0 \ 0 \ \dots \ 0 & 0 \ 0 \ \dots \ 0 & \dots & \overbrace{1 \ 1 \ \dots \ 1}^{n/2^l} \end{pmatrix},$$

for which condition (v) is satisfied with $\nu_q = \sqrt{2^l} = \sqrt{q}$, which is exactly the lower bound (2.7). Note, that when n is not a power of 2, one can make the same sequence of partition dividing at each step $A_r^{(l-1)}$ into $A_{2r-1}^{(l)}$ and $A_{2r}^{(l)}$ such that $|\#A_j^{(l)} - \#A_{j+1}^{(l)}| \leq 1$ for each $l \leq \log_2 n$, $j < 2^l$.

2.2 Application to composite function estimation

Consider a regression model of the form

$$Y_i = f(g(i/n)) + \varepsilon_i, \quad Z_i = g(i/n) + \nu_i, \quad i = 1, \dots, n, \quad (2.10)$$

where f and g are two real valued functions. Here the problem is to estimate the function f from the observations $(Y_1, Z_1), \dots, (Y_n, Z_n)$. This problem is related to the problem of composite function estimation which recently got much attention in the literature (see, e.g. Juditsky et al. (2009) and references therein). However our main interest here lies in estimating the function f and not the whole composite function $f(g(\cdot))$. Let us apply finite-dimensional sieve approximations to f and g :

$$f(x) \approx f_p(x) = \sum_{m=1}^p \theta_m \psi_m(x), \quad g(x) \approx g_q(x) = \sum_{m=1}^q \eta_m \phi_m(x).$$

Thus, our problem is reduced to a parametric model (2.1) with parameter $(\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^{p+q}$ and the estimator $\tilde{f}_{p,q}$ of f is given by $\tilde{\boldsymbol{\theta}}$ from (2.2) corresponding to $\tilde{\Phi}$ given in (2.9):

$$\tilde{f}_{p,q}(x) = \Psi(x)^\top \tilde{\boldsymbol{\theta}} = \sum_{m=1}^p \tilde{\theta}_m \psi_m(x).$$

The performance of the estimator depends on the values p, q and on the accuracy of corresponding sieve approximation. One can describe the dependence of the bias of $\tilde{f}_{p,q}$ on the values of p and q through the so-called *Sobolev ellipsoids*

$$S_\psi(\beta, Q) = \left\{ f(x) = \sum_{m=1}^{\infty} \theta_m \psi_m(x) : \boldsymbol{\theta} \in \ell^2, \sum_{m=1}^{\infty} m^{2\beta} \theta_m^2 \leq Q^2 \right\}.$$

We refer to Tsybakov (2009) for the proof of the following result.

Lemma 2.3. *Let $f \in S_\psi(\beta, Q)$, then*

$$\|f - f_p\|_2 \leq Qp^{-\beta}. \quad (2.11)$$

Moreover, if $\beta > s + 1/2$ and each $\psi_m(x)$, is $s \geq 0$ times continuously differentiable with $\|\psi_m^{(s)}\|_\infty \leq Am^s$, $m = 1, 2, \dots$ for some $A > 0$, then

$$\|f_p^{(s)}\|_\infty \leq \frac{AQ}{\sqrt{2(\beta - s) - 1}}, \quad \|f^{(s)} - f_p^{(s)}\|_\infty \leq \frac{AQ}{\sqrt{2(\beta - s) - 1}} p^{-(\beta - s) + \frac{1}{2}}. \quad (2.12)$$

In order to apply Theorem 2.1, we need to check conditions (i) - (vi). Here we introduce a list of conditions for the nonparametric problem, which ensure those of the reduced parametric problem.

(a) The bases (ψ_m) and (ϕ_m) satisfy, for some $A > 0$ and each $s = 0, 1, 2$,

$$\|\psi_m^{(s)}\|_\infty \leq Am^s, \quad \|\phi_m^{(s)}\|_\infty \leq Am^s;$$

(b) $f \in S_\psi(\beta, Q)$ and $g \in S_\phi(\beta_X, Q)$;

(c) $\Psi(\mathbf{X}^*)\Psi(\mathbf{X}^*)^\top \geq (\mathbf{f}n)I_p$;

(d) for some \mathbf{C}_0 and each interval $I \subset \mathbb{R}$ it holds

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}\{X_i^* \in I\} \leq \mathbf{C}_0 \max\left(\mu(I), \frac{1}{n}\right),$$

where $\mu(I)$ is the length of the interval.

Discussion Let us remark on the assumption (c). One may argue that this condition is not exactly nonparametric, as it depends on the dimension p . In that case we suggest to consider another condition, that seems more interpretable. The following lemma states, that if the the points X_i^* are dense enough on the interval $[0, 1]$, then (c) is satisfied.

Lemma 2.4. *Let (ψ_m) satisfy (a). Suppose, that with some constant $\mathbf{c}_0 > 0$ and some small enough $\delta > 0$ it holds*

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}\{X_i^* \in I\} \geq \mathbf{c}_0(\mu(I) - \delta p^{-2})$$

for each interval $I \subset [0, 1]$. Then we have,

$$n^{-1}\Psi(\mathbf{X}^*)\Psi(\mathbf{X}^*)^\top \geq \mathbf{c}_0(1 - \alpha(\delta))I_p,$$

with $\alpha(\delta) \rightarrow 0$ as $\delta \rightarrow 0$.

Let us check other conditions of Theorem 2.1. Condition (v) follows from the following lemma.

Lemma 2.5. *Let (a) and (d) hold, and (b) be satisfied with $\beta, \beta_X > 3/2$. Then,*

$$\|\mathbf{f}^* - \Psi(\mathbf{X}^*)^\top \boldsymbol{\theta}^*\| \leq \mathbf{c}_0^{1/2} Q \sqrt{n} p^{-\beta} + \frac{\mathbf{c}_0^{1/2} A Q}{\sqrt{\beta - 3/2}} p^{-\beta+1},$$

and

$$\|\Phi \boldsymbol{\eta}^* - \mathbf{X}^*\| \leq Q \sqrt{n} q^{-\beta_X} + \frac{A Q}{\sqrt{\beta_X - 3/2}} q^{-\beta_X+1}.$$

Additionally, if $p, q \leq \sqrt{n}$, then it holds

$$\|\mathbf{f}^* - \Psi(\mathbf{X}^*)^\top \boldsymbol{\theta}^*\| \leq Q' \sqrt{n} p^{-\beta}, \quad \|\Phi \boldsymbol{\eta}^* - \mathbf{X}^*\| \leq Q' \sqrt{n} q^{-\beta_X}, \quad (2.13)$$

with $Q' = \mathbf{c}_0^{1/2} Q \left(1 + \frac{A}{\sqrt{\min\{\beta, \beta_X\} - 3/2}} \right)$ not depending on p, q and n .

Condition (i) follows from inequality (2.12), and the rest of the conditions can be checked straightforward. Manipulating with the error term $\diamond(\mathbf{x})$, given in (2.5), we can show the following result.

Proposition 2.6. *Let $\beta > 2.5$ and $\beta_X > 3/2$. Then, by choosing*

$$p \stackrel{\text{def}}{=} n^{\frac{1}{2\beta+1}}, \quad q \stackrel{\text{def}}{=} p^{\beta/\beta_X} = n^{\frac{\beta}{2\beta\beta_X + \beta_X}},$$

we get the standard non-parametric rate of convergence

$$\|f - \tilde{f}_{p,q}\|_2 \lesssim_{\mathcal{P}} n^{-\frac{\beta}{2\beta+1}}.$$

This rate is in fact optimal over the class $S_\psi(\beta, Q)$.

A stronger assumption, that f is Hölder smooth on $[0, 1]$, i.e. $f \in \Sigma(\beta, L)$ (see, e.g. Tsybakov (2009)) for some $\beta > 0$ and $\psi_m(x) = \cos(2\pi m x)$, yields $f \in S(\beta, L/\pi^\beta)$ and the corresponding inequality (2.11). However, one can get a better bound in the sup-norm

$$\|f^{(s)} - f_p^{(s)}\|_\infty \leq L' p^{-(\beta-s)} \ln p,$$

for some L' not depending on p and $\beta > s$ (compared to $\beta > s + 1/2$ in (2.12)). Thus, such an assumption helps to “relax” the conditions $\beta > 2.5$ and $\beta_X > 3/2$ in the Proposition 2.6. Assume that

(a') both basis are trigonometric: $\psi_m(x) = \phi_m(x) = \cos(2\pi m x)$;

(b') $f \in \Sigma(\beta, L)$ and $g \in \Sigma(\beta_X, L)$.

and the conditions (c), (d) remain the same. Then we have the following proposition.

Proposition 2.7. *Assume (a'), (b'), (c) and (d). Let $\beta > 2$ and $\beta_X \geq 3/2$. Then, by choosing*

$$p \stackrel{\text{def}}{=} n^{\frac{1}{2\beta+1}}, \quad q \stackrel{\text{def}}{=} p^{\beta/\beta_X} = n^{\frac{\beta}{2\beta_X + \beta_X}},$$

we get the standard non-parametric rate of convergence

$$\|f - \tilde{f}_{p,q}\|_2 \lesssim_{\mathbb{P}} n^{-\frac{\beta}{2\beta+1}}.$$

This rate is in fact optimal over the class $\Sigma(\beta, L)$.

3 Simulation results

Consider first a semiparametric regression model (2.10) with

$$f(x) = \sum_{k=1}^p \theta_k^* \cos(\pi k x)$$

and

$$g(x) = \sum_{k=1}^{\infty} \frac{\cos(\pi k x)}{k^{2\beta}} = (-1)^{\beta-1} (2\pi)^{2\beta} \frac{B_{2\beta}(x/2)}{2(2\beta)!}, \quad x \in [0, 1]$$

for some natural $\beta, p > 0$, where $B_m(x)$ is a Bernoulli polynomial of order m . Note that $g \in S_{\cos}(2\beta + 1/2, Q)$ for some $Q > 0$. The aim is to estimate the function f , i.e., the vector $\boldsymbol{\theta}^* = (\theta_1^*, \dots, \theta_p^*)$ based on the sample $(Y_1, Z_1), \dots, (Y_n, Z_n)$. We assume that the errors $\varepsilon_i, \nu_i, i = 1, \dots, n$, in (2.10) are normal i.i.d. with zero mean and variance 0.1. We consider the following sieve approximation for g

$$g(x) \approx g_q(x) = \sum_{m=1}^q \eta_m^* \cos(\pi m x)$$

with $\eta_m^* = m^{-2\beta}$. In order to compute the corresponding ML estimates for $\boldsymbol{\theta}^*$ and $\boldsymbol{\eta}^* = (\eta_1^*, \dots, \eta_q^*)$ we use the following strategy. First we compute a preliminary plug-in estimate $\boldsymbol{\eta}^{(0)}$ for $\boldsymbol{\eta}^*$ by applying the standard least-squares approach to the linear model

$$Z_i = g_q(i/n) + \nu_i, \quad i = 1, \dots, n.$$

A preliminary estimate $\boldsymbol{\theta}^{(0)}$ for $\boldsymbol{\theta}^*$ is then obtained from the model

$$Y_i = f(g_q^{(0)}(i/n)) + \varepsilon_i, \quad i = 1, \dots, n,$$

with $g_q^{(0)}(x) = \sum_{m=1}^q \eta_m^{(0)} \phi_m(x)$. Then we maximize the ML function

$$L_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) = -\|\mathbf{Y} - \Psi(\Phi\boldsymbol{\eta})^{\top} \boldsymbol{\theta}\|^2 - \|\mathbf{Z} - \Phi\boldsymbol{\eta}\|^2$$

iteratively over $\boldsymbol{\theta}$ and $\boldsymbol{\eta}$ starting at $(\boldsymbol{\theta}^{(0)}, \boldsymbol{\eta}^{(0)})$. To minimize $L_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta})$ over $\boldsymbol{\eta}$, we use the Levenberg-Marquardt algorithm (`minpack.lm` library in `R`). Such alternation maximization procedure usually converges rather quickly. In Figure 3.1 we show the boxplots of the weighted norm $(n/p)\|\boldsymbol{\theta}^* - \boldsymbol{\theta}^{(10)}\|^2$ based on 1000 independent samples each of length n from an EIV model with parameters $\boldsymbol{\theta}^* = (1, \dots, 1)$, $\beta = 2$, $q = 5$ as function of n (left) and p (right). A stable behaviour can be observed for increasing values of n and p supporting the results of Theorem 2.1. Next we turn to a nonparametric model

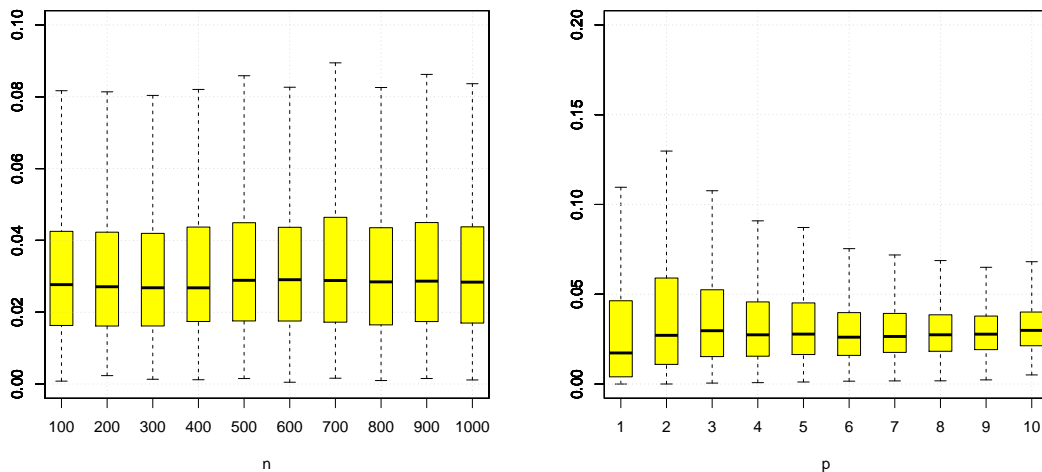


Figure 3.1: Left: the boxplots of the loss $(n/p)\|\boldsymbol{\theta}^* - \boldsymbol{\theta}^{(10)}\|^2$ as a function of the sample size n . Right: the boxplots of the loss $(n/p)\|\boldsymbol{\theta}^* - \boldsymbol{\theta}^{(10)}\|^2$ as a function of the dimension p .

with

$$f(x) = \sum_{k=1}^{\infty} \frac{\sin(\pi kx)}{k^{2\alpha+1}} = (-1)^{\alpha-1} (2\pi)^{2\alpha+1} \frac{B_{2\alpha+1}(x/2)}{2(2\alpha+1)!}, \quad x \in [0, 1],$$

where $\alpha = 3$. In Figure 3.2 we show the empirical error $\|\mathbf{Y} - \Psi(\Phi\boldsymbol{\eta}^{(10)})^{\top} \boldsymbol{\theta}^{(10)}\|^2/n$ as function of p and q for $n = 100$. Clear minima can be observed in each plot and their locations are in accordance with the results of Proposition 2.6. As also can be seen, the dependence of estimation error on the dimension q of the nuisance parameter $\boldsymbol{\eta}$ is much weaker than on p (at least in the case of large p and q). Figure 3.3 graphically shows 500 realisations of the estimate $\tilde{f}_{p,q}$ for $n = 100$ (left) and $n = 500$ (right) and suitably chosen p, q .

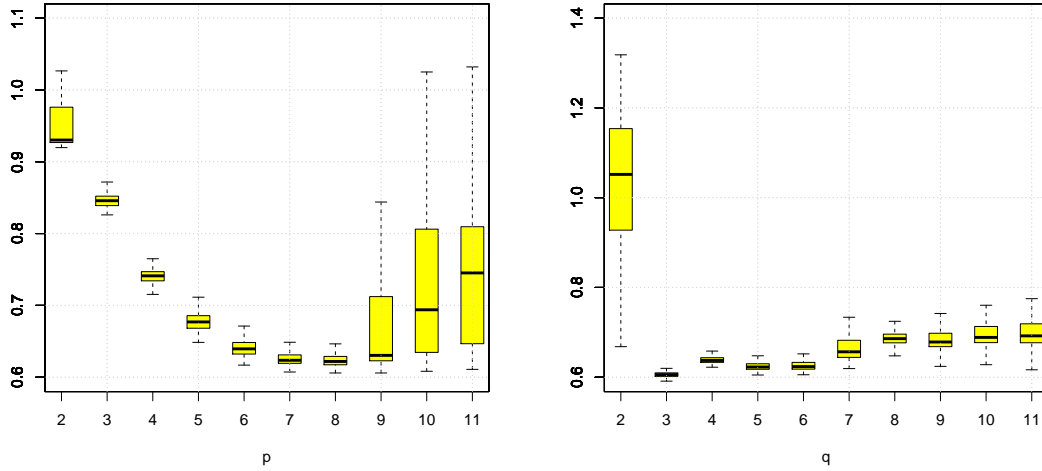


Figure 3.2: Left: the boxplots of the empirical loss $(1/n)\|\mathbf{f} - \Psi(\Phi\boldsymbol{\eta}^{(10)})^\top \boldsymbol{\theta}^{(10)}\|^2$ as a function of the dimension p . Right: the boxplots of the loss $(1/n)\|\mathbf{f} - \Psi(\Phi\boldsymbol{\eta}^{(10)})^\top \boldsymbol{\theta}^{(10)}\|^2$ as a function of the dimension q .

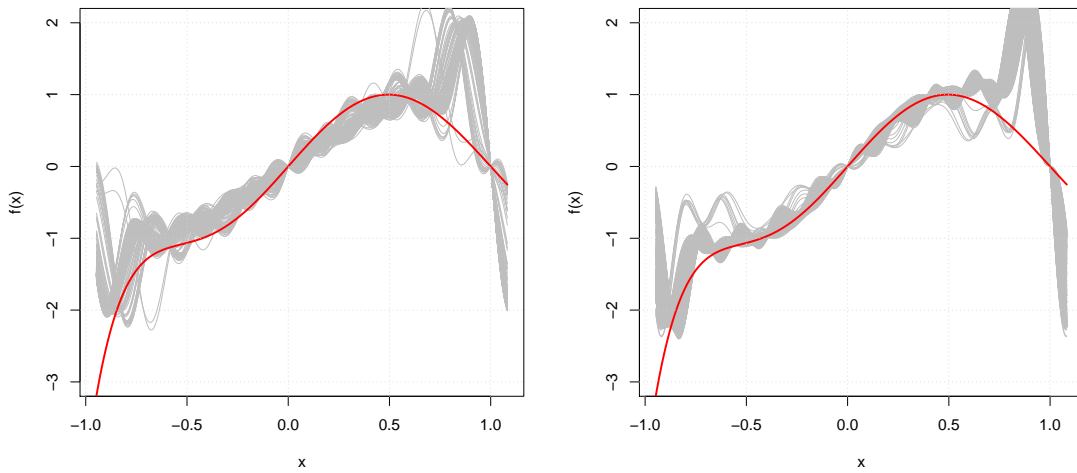


Figure 3.3: Function $f(x)$ (red) together with 500 estimates based on 500 independent samples each of the length $n = 100$ (left) and $n = 500$ (right) with properly chosen p and q .

4 Proofs

This section is organized as follows. First, we apply general finite sample theory imposed in Spokoiny (2012) and Andresen and Spokoiny (2013) to parametric problem (2.1),

which will lead us to the proof of Theorem 2.1. Then, we present proofs for the rest of statements, in the order in which they appear above.

4.1 Main objects

We start with describing main objects that will be helpful through the whole paper. The *stochastic term* $\zeta_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta})$ is easy to describe:

$$\zeta_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) = L_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) - \mathbb{E}L_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \boldsymbol{\varepsilon}^{\top} \Psi(\Phi \boldsymbol{\eta})^{\top} \boldsymbol{\theta} + \sigma^{-1} \boldsymbol{\nu}^{\top} \Phi \boldsymbol{\eta}.$$

Below we denote by $\Psi^{(s)}(\mathbf{X})$ the matrix with the entries $\psi_m^{(s)}(X_i)$ which is obtained by entry-wise differentiation of $\Psi(\mathbf{X})$. For a vector $\mathbf{a} = (a_i) \in \mathbb{R}^n$, the notation $\text{diag}(\mathbf{a})$ means a $n \times n$ diagonal matrix with the diagonal entries a_i . Then the *score* $\nabla \zeta_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta})$ can be written in the block-wise form as

$$\nabla \zeta_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \begin{pmatrix} \nabla_{\boldsymbol{\theta}} \zeta_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) \\ \nabla_{\boldsymbol{\eta}} \zeta_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) \end{pmatrix} = \begin{pmatrix} \Psi(\Phi \boldsymbol{\eta}) \boldsymbol{\varepsilon} \\ \sigma^{-1} \Phi^{\top} \boldsymbol{\nu} + \Phi^{\top} \text{diag}\{\Psi'(\Phi \boldsymbol{\eta})^{\top} \boldsymbol{\theta}\} \boldsymbol{\varepsilon} \end{pmatrix}.$$

Here $\text{diag}\{\Psi'(\mathbf{X})^{\top} \boldsymbol{\theta}\}$ means the diagonal matrix with the diagonal entries $\sum_m \theta_m \psi'_m(X_i)$. The *Hessian of the stochastic component* reads as

$$\nabla^2 \zeta_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \begin{pmatrix} 0 & \Psi'(\Phi \boldsymbol{\eta}) \text{diag}(\boldsymbol{\varepsilon}) \Phi \\ \Phi^{\top} \text{diag}(\boldsymbol{\varepsilon}) \Psi'(\Phi \boldsymbol{\eta})^{\top} & \Phi^{\top} \text{diag}\{\Psi''(\Phi \boldsymbol{\eta})^{\top} \boldsymbol{\theta}\} \text{diag}(\boldsymbol{\varepsilon}) \Phi \end{pmatrix}$$

The lower right block of this matrix is diagonal with the diagonal entries $\sum_m \theta_m \psi''_m(X_i) \varepsilon_i$ for $i = 1, \dots, n$ (not taking into account left and right multiplying by Φ).

Further, the expected log-likelihood reads (up to a constant term) as

$$-2\mathbb{E}L_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \|\mathbf{f}^* - \Psi(\Phi \boldsymbol{\eta})^{\top} \boldsymbol{\theta}\|^2 + \sigma^{-2} \|\Phi \boldsymbol{\eta} - \mathbf{X}^*\|^2.$$

The related matrix $D^2(\boldsymbol{\theta}, \boldsymbol{\eta})$ is

$$D^2(\mathbf{v}) = -\nabla^2 \mathbb{E}L_{\Phi}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \begin{pmatrix} D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2(\mathbf{v}) & A_{\boldsymbol{\theta}\boldsymbol{\eta}}(\mathbf{v}) \\ A_{\boldsymbol{\theta}\boldsymbol{\eta}}^{\top}(\mathbf{v}) & D_{\boldsymbol{\eta}\boldsymbol{\eta}}^2(\mathbf{v}) \end{pmatrix},$$

where

$$D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2(\mathbf{v}) = \Psi(\Phi \boldsymbol{\eta}) \Psi(\Phi \boldsymbol{\eta})^{\top},$$

$$A_{\boldsymbol{\theta}\boldsymbol{\eta}}(\mathbf{v}) = \Psi(\Phi \boldsymbol{\eta}) \text{diag}\{\Psi'(\Phi \boldsymbol{\eta})^{\top} \boldsymbol{\theta}\} \Phi + \Psi'(\Phi \boldsymbol{\eta}) \text{diag}\{\mathbf{f}^* - \Psi(\Phi \boldsymbol{\eta})^{\top} \boldsymbol{\theta}\} \Phi,$$

$$D_{\boldsymbol{\eta}\boldsymbol{\eta}}^2(\mathbf{v}) = \sigma^{-2} I_q + \Phi^{\top} \text{diag}\{\Psi'(\Phi \boldsymbol{\eta})^{\top} \boldsymbol{\theta}\}^2 \Phi + \Phi \text{diag}\{\Psi''(\Phi \boldsymbol{\eta})^{\top} \boldsymbol{\theta}\} \text{diag}\{\mathbf{f}^* - \Psi(\Phi \boldsymbol{\eta})^{\top} \boldsymbol{\theta}\} \Phi,$$

and define also

$$\begin{aligned} D_{\theta\theta}^2 &= \Psi(\Phi\boldsymbol{\eta}^*)\Psi(\Phi\boldsymbol{\eta}^*)^\top, \\ A_{\theta\eta} &= \Psi(\Phi\boldsymbol{\eta}^*) \text{diag}\{\Psi'(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\} \Phi, \\ D_{\eta\eta}^2 &= \sigma^{-2}I_q + \Phi^\top \text{diag}\{\Psi'(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\}^2 \Phi, \end{aligned}$$

which are the blocks of the matrix

$$D^2 = \begin{pmatrix} D_{\theta\theta}^2 & A_{\theta\eta} \\ A_{\theta\eta}^\top & D_{\eta\eta}^2 \end{pmatrix}.$$

The matrix D^2 is defined in a that way to be close to $D^2(\mathbf{v}^*)$, but is more convenient to work with.

Finally, we need to measure distance between \mathbf{v} 's. Motivated by the view of D^2 define

$$H_0^2 = \begin{pmatrix} D_{\theta\theta}^2 & 0 \\ 0 & \sigma^{-2}I_q \end{pmatrix},$$

and the *local vicinity* for parameter $(\boldsymbol{\theta}, \boldsymbol{\eta})$ among projection $(\boldsymbol{\theta}^*, \boldsymbol{\eta}^* = \Phi^\top \mathbf{X}^*)$

$$\mathcal{V}_\circ(\mathbf{r}) = \{\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) : \|H_0(\mathbf{v} - \mathbf{v}^*)\|^2 = \|D_{\theta\theta}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 + \sigma^{-2}\|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|^2 \leq \mathbf{r}^2\}.$$

4.2 Auxiliary lemmas

The first lemma describes the variability of the matrix $\Psi(\mathbf{X})$.

Lemma 4.1. *For any vector $\boldsymbol{\gamma} \in \mathbb{R}^p$ with $\|\boldsymbol{\gamma}\| \leq 1$, it holds with $w_{p,s}$ from (2.3)*

$$\|\Psi^{(s)}(\mathbf{X})^\top D_{\theta\theta}^{-1} \boldsymbol{\gamma}\|_\infty \leq \frac{w_{p,s}}{\sqrt{n}}, \quad s = 0, 1, 2$$

Proof. By Cauchy-Schwartz inequality

$$|\Psi^{(s)}(X_i)^\top D_{\theta\theta}^{-1} \boldsymbol{\gamma}|^2 \leq \|D_{\theta\theta}^{-1} \boldsymbol{\gamma}\|^2 \times \sum_{i=1}^p \psi_m^{(s)}(X_i) \leq \mathbf{f}^{-1} \sum_{i=1}^p \mu_{m,s}^2.$$

□

In addition, we need to bound the value $\Psi^{(s)}(\mathbf{X})^\top \boldsymbol{\theta}$ on the neighbourhood of $\boldsymbol{\theta}^*$.

Lemma 4.2. *For $s = 0, 1, 2$ and any $\boldsymbol{\theta} \in \mathbb{R}^p$ satisfying $\|D_{\theta\theta}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}$ and each $X \in \mathbb{R}$ it holds*

$$\begin{aligned} |\Psi^{(s)}(X)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)| &\leq \frac{w_{p,s} \mathbf{r}}{\sqrt{n}}, \\ |\Psi^{(s)}(X)^\top \boldsymbol{\theta}| &\leq \mathbf{c}_{f,s} + \frac{w_{p,s} \mathbf{r}}{\sqrt{n}}. \end{aligned}$$

Proof. Using previous lemma with $\gamma = \mathbf{r}^{-1}D_{\theta\theta}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$,

$$|\Psi^{(s)}(X)^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)| = |\Psi^{(s)}(X)^\top D_{\theta\theta}^{-1} \times D_{\theta\theta}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)| \leq \frac{\mathbf{r}w_{p,s}}{\sqrt{n}},$$

and the second follows by $|\Psi^{(s)}(X)^\top \boldsymbol{\theta}| \leq |\Psi^{(s)}(X)^\top \boldsymbol{\theta}^*| + |\Psi^{(s)}(X)^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)|$. \square

In particular, if $\mathbf{r} \leq \sqrt{n}/w_{p,s}$, then $|\Psi^{(s)}(X)^\top \boldsymbol{\theta}| \leq \mathbf{C}_{f,s} + 1$.

Lemma 4.3. *If $\sigma^{-1}\|\mathbf{X} - \mathbf{X}^*\| \leq \mathbf{r}$, then it holds for $s = 0, 1, 2$*

$$\|\{\Psi^{(s)}(\mathbf{X}) - \Psi^{(s)}(\mathbf{X}^*)\}^\top D_{\theta\theta}^{-1} \gamma\| \leq \frac{\sigma w_{p,s+1} \mathbf{r}}{\sqrt{n}},$$

$$\|\{\Psi^{(s)}(\mathbf{X}) - \Psi^{(s)}(\mathbf{X}^*)\}^\top \boldsymbol{\theta}^*\| \leq \sigma \mathbf{C}_{f,s+1} \mathbf{r}.$$

Proof. We have with some \mathbf{X}' on the line connecting \mathbf{X} , \mathbf{X}^*

$$\{\Psi^{(s)}(\mathbf{X}) - \Psi^{(s)}(\mathbf{X}^*)\}^\top D_{\theta\theta}^{-1} \gamma = (\mathbf{X} - \mathbf{X}^*)^\top \text{diag}\{\Psi^{(s+1)}(\mathbf{X}')^\top D_{\theta\theta}^{-1} \gamma\},$$

and the first inequality by Lemma 4.1.

By analogy, the second inequality follows from $|\Psi^{(s+1)}(\mathbf{X}')^\top \boldsymbol{\theta}^*| \leq \mathbf{C}_{f,s+1}$ and

$$\{\Psi^{(s)}(\mathbf{X}) - \Psi^{(s)}(\mathbf{X}^*)\}^\top \boldsymbol{\theta}^* = (\mathbf{X} - \mathbf{X}^*)^\top \text{diag}\{\Psi(\mathbf{X}')^\top \boldsymbol{\theta}^*\}.$$

\square

The next lemma checks identifiability condition (\mathcal{I}) of Andresen and Spokoiny (2013).

Lemma 4.4. *Define*

$$\rho = \sqrt{\frac{\mathbf{C}_{f,1}^2}{\sigma^{-2} + \mathbf{C}_{f,1}^2}}.$$

Then,

$$\|D_{\theta\theta}^{-1} A_{\theta\eta} D_{\eta\eta}^{-1}\| \leq \rho, \quad D^2 \geq (1 - \rho)H_0^2, \quad D^2 \leq (1 - \rho)^{-1}H_0^2.$$

Proof. Define, $\tilde{D}_{\eta\eta}^2 = D_{\eta\eta}^2 - \sigma^{-2}I_q$. Then,

$$\begin{aligned} \tilde{D}^2 &= \begin{pmatrix} D_{\theta\theta}^2 & A_{\theta\eta} \\ A_{\theta\eta}^\top & \tilde{D}_{\eta\eta}^2 \end{pmatrix} = \begin{pmatrix} \Psi(\Phi\boldsymbol{\eta}^*)\Psi(\Phi\boldsymbol{\eta}^*)^\top & \Psi(\Phi\boldsymbol{\eta}^*) \text{diag}\{\Psi'(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\} \Phi \\ \Phi^\top \text{diag}\{\Psi'(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\} \Psi(\Phi\boldsymbol{\eta}^*) & \Phi^\top \text{diag}\{\Psi'(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\}^2 \Phi \end{pmatrix} \\ &= \begin{pmatrix} \Psi(\Phi\boldsymbol{\eta}^*) \\ \Phi^\top \text{diag}\{\Psi'(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\} \end{pmatrix} \begin{pmatrix} \Psi(\Phi\boldsymbol{\eta}^*) \\ \Phi^\top \text{diag}\{\Psi'(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\} \end{pmatrix}^\top, \end{aligned}$$

meaning that \tilde{D}^2 is nonnegative and $\|D_{\theta\theta}^{-1}A_{\theta\eta}\tilde{D}_{\eta\eta}^{-1}\| \leq 1$. Thus,

$$\|D_{\theta\theta}^{-1}A_{\theta\eta}D_{\eta\eta}^{-1}\| \leq \|\tilde{D}_{\eta\eta}D_{\eta\eta}^{-1}\| \leq \|\text{diag}\{\Psi'(\Phi\eta^*)^\top\theta^*\}^2(\sigma^{-2} + \text{diag}\{\Psi'(\Phi\eta^*)^\top\theta^*\}^2)^{-1}\|^{1/2},$$

which is less than ρ since the function $x/(\sigma^{-2} + x)$ is non decreasing and $|\Psi'(X)^\top\theta^*| \leq C_{f,1}$ by the condition (i). So, by this inequality we have

$$D^2 \geq (1 - \rho) \begin{pmatrix} D_{\theta\theta}^2 & 0 \\ 0 & D_{\eta\eta}^2 \end{pmatrix} \geq (1 - \rho)H_0^2.$$

Further, using $D_{\eta\eta}^2 \leq (1 + \sigma^2 C_{f,1}^2)H_\sigma^2 = (1 - \rho^2)^{-1}H_\sigma^2$ we get

$$D^2 \leq (1 + \rho) \begin{pmatrix} D_{\theta\theta}^2 & 0 \\ 0 & D_{\eta\eta}^2 \end{pmatrix} \leq (1 + \rho)(1 - \rho^2)^{-1}H_0^2.$$

□

Lemma 4.5. *The vector $\mathbf{b} = D^{-1}\nabla L_\Phi(\mathbf{v}^*)$ satisfies $\|\mathbf{b}\| \leq \|\mathbf{f}^* - \Psi(\Phi\eta^*)^\top\theta^*\|$.*

Proof. The explicit expression for the vector is as follows

$$\begin{aligned} \nabla EL_\Phi(\mathbf{v}^*) &= \begin{pmatrix} \Psi(\Phi\eta^*)(\mathbf{f}^* - \Psi(\Phi\eta^*)^\top\theta^*) \\ \Phi^\top \text{diag}\{\Psi'(\Phi\eta^*)^\top\theta^*\}(\mathbf{f}^* - \Psi(\Phi\eta^*)^\top\theta^*) \end{pmatrix} \\ &= \begin{pmatrix} \Psi(\Phi\eta^*) \\ \Phi^\top \text{diag}\{\Psi'(\Phi\eta^*)^\top\theta^*\} \end{pmatrix} (\mathbf{f}^* - \Psi(\Phi\eta^*)^\top\theta^*). \end{aligned}$$

Define $A = \Psi(\Phi\eta^*)$, $B = \Phi^\top \text{diag}\{\Psi'(\Phi\eta^*)^\top\theta^*\}$. Then,

$$D^2 = \begin{pmatrix} AA^\top & AB^\top \\ AB^\top & \sigma^{-2}I_p + BB^\top \end{pmatrix} \geq \begin{pmatrix} A \\ B \end{pmatrix} \begin{pmatrix} A \\ B \end{pmatrix}^\top.$$

Hence,

$$\|D^{-1}\nabla EL_\Phi(\mathbf{v}^*)\| \leq \left\| D^{-1} \begin{pmatrix} A \\ B \end{pmatrix} \right\| \times \|\mathbf{f}^* - \Psi(\Phi\eta^*)^\top\theta^*\| \leq \|\mathbf{f}^* - \Psi(\Phi\eta^*)^\top\theta^*\|.$$

□

4.3 Local perturbation of $D^2(\mathbf{v})$

What follows relates to condition (\mathcal{L}_0) of Andresen and Spokoiny (2013). Our purpose is to bound $\|H_0^{-1}(D^2(\mathbf{v}) - D^2)H_0^{-1}\|$ for $\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})$. The result will depend on the smoothness of design basis Φ . It holds,

$$V(\Phi) = \sup_{\gamma \in \mathbb{R}^q, \|\gamma\|=1} \|\Phi\gamma\|_\infty, \quad \mathbf{v}_q \geq n^{1/2}V(\Phi),$$

and we only use v_q through inequality $\|\Phi\gamma\|_\infty \leq \frac{v_q}{\sqrt{n}}\|\gamma\|$ for each $\gamma \in \mathbb{R}^q$.

Lemma 4.6. *Let \mathbf{r} satisfies $\mathbf{r} \leq \sqrt{n}/w_{p,1}$. Then $\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})$ yields*

$$\begin{aligned} \|H_0^{-1}(D^2(\mathbf{v}) - D^2)H_0\| &\leq \delta(\mathbf{r}) \\ &\stackrel{\text{def}}{=} \mathbf{C}^{**} \sigma \frac{w_{p,1} + v_q + \frac{w_{p,2}v_q\mathbf{r}}{\sqrt{n}}}{\sqrt{n}} \{\mathbf{r} + \square\}, \end{aligned}$$

where

$$\begin{aligned} \square &= \|\mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\|, \\ \mathbf{C}^* &= (\sigma \vee 1) \times (\mathbf{C}_{f,1} \vee \mathbf{C}_{f,2} + 1), \\ \mathbf{C}^{**} &= 8\mathbf{C}^* \vee 3(\mathbf{C}^*)^2. \end{aligned}$$

Proof. For $\boldsymbol{\alpha} = (\boldsymbol{\beta}, \boldsymbol{\gamma}) \in \mathbb{R}^{p+q}$ with $\|\boldsymbol{\alpha}\| \leq 1$ we have

$$\begin{aligned} |\boldsymbol{\alpha}^\top H_0^{-1}(D^2(\mathbf{v}) - D^2)H_0^{-1}\boldsymbol{\alpha}| &\leq |\boldsymbol{\beta}^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}(D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2(\mathbf{v}) - D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2)D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\boldsymbol{\beta}| \\ &\quad + 2\sigma |\boldsymbol{\beta}^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}(A_{\boldsymbol{\theta}\boldsymbol{\eta}}(\mathbf{v}) - A_{\boldsymbol{\theta}\boldsymbol{\eta}})\boldsymbol{\gamma}| + \sigma^2 |\boldsymbol{\gamma}^\top (D_{\boldsymbol{\eta}\boldsymbol{\eta}}^2(\mathbf{v}) - D_{\boldsymbol{\eta}\boldsymbol{\eta}}^2)\boldsymbol{\gamma}|. \end{aligned}$$

We will proceed with each term separately. First,

$$\begin{aligned} &\boldsymbol{\beta}^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\{D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2(\mathbf{v}) - D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2\}D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\boldsymbol{\beta} \\ &= \boldsymbol{\beta}^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\{\Psi(\Phi\boldsymbol{\eta})\Psi(\Phi\boldsymbol{\eta})^\top - \Psi(\Phi\boldsymbol{\eta}^*)\Psi(\Phi\boldsymbol{\eta}^*)^\top\}D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\boldsymbol{\beta} \\ &= 2\boldsymbol{\gamma}^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\{\Psi(\Phi\boldsymbol{\eta}) - \Psi(\Phi\boldsymbol{\eta}^*)\}\Psi(\Phi\boldsymbol{\eta}^*)^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\boldsymbol{\gamma} + \|\{\Psi(\Phi\boldsymbol{\eta}) - \Psi(\Phi\boldsymbol{\eta}^*)\}^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\boldsymbol{\gamma}\|^2 \\ &\leq \frac{2\sigma w_{p,1}\mathbf{r}}{\sqrt{n}} + \left(\frac{\sigma w_{p,1}\mathbf{r}}{\sqrt{n}}\right)^2 \leq \frac{3\sigma w_{p,1}\mathbf{r}}{\sqrt{n}}. \end{aligned}$$

Second, we use decomposition

$$\begin{aligned} A_{\boldsymbol{\theta}\boldsymbol{\eta}}(\mathbf{v}) - A_{\boldsymbol{\theta}\boldsymbol{\eta}} &= \Psi(\Phi\boldsymbol{\eta}^*) \text{diag}\{(\Psi'(\Phi\boldsymbol{\eta}) - \Psi'(\Phi\boldsymbol{\eta}^*))^\top \boldsymbol{\theta}^*\} \Phi \\ &\quad + \Psi(\Phi\boldsymbol{\eta}^*) \text{diag}\{\Psi'(\Phi\boldsymbol{\eta}^*)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\} \Phi \\ &\quad + (\Psi(\Phi\boldsymbol{\eta}) - \Psi(\Phi\boldsymbol{\eta}^*)) \text{diag}\{\Psi'(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}\} \Phi \\ &\quad + \Psi'(\Phi\boldsymbol{\eta}) \text{diag}\{\mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}\} \Phi. \end{aligned}$$

Define the following vectors,

$$\begin{aligned}
\mathbf{a} &= (\Psi'(\Phi\boldsymbol{\eta}) - \Psi'(\Phi\boldsymbol{\eta}^*))^\top \boldsymbol{\theta}^*, & \|\mathbf{a}\| &\leq \mathbf{C}_{f,2}\mathbf{r}, \\
\mathbf{b} &= \Psi(\Phi\boldsymbol{\eta}^*)^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\boldsymbol{\beta}, & \|\mathbf{b}\|_\infty &\leq \mathbf{w}_{p,0}/\sqrt{n}, & \|\mathbf{b}\| &= 1, \\
\mathbf{c} &= \Psi'(\Phi\boldsymbol{\eta}^*)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*), & \|\mathbf{c}\|_\infty &\leq \mathbf{w}_{p,1}\mathbf{r}/\sqrt{n}, \\
\mathbf{d} &= (\Psi(\Phi\boldsymbol{\eta}) - \Psi(\Phi\boldsymbol{\eta}^*))^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\boldsymbol{\beta}, & \|\mathbf{d}\| &\leq \sigma\mathbf{w}_{p,1}\mathbf{r}/\sqrt{n}, \\
\mathbf{e} &= \Psi'(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}, & \|\mathbf{e}\|_\infty &\leq \mathbf{C}_{f,1}, \\
\mathbf{g} &= \Phi\boldsymbol{\gamma} & \|\mathbf{g}\| &\leq 1, \\
\mathbf{h} &= \Psi'(\Phi\boldsymbol{\eta})^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\boldsymbol{\beta}, & \|\mathbf{h}\|_\infty &\leq \mathbf{w}_{p,1}/\sqrt{n}, \\
\mathbf{k} &= \mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}, & \|\mathbf{k}\| &\leq \mathbf{C}^*\mathbf{r} + \square,
\end{aligned}$$

where the inequalities follow from Lemmas 4.1, 4.2 and 4.3. Then,

$$\begin{aligned}
|\boldsymbol{\beta}^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}(A_{\boldsymbol{\theta}\boldsymbol{\eta}}(\mathbf{v}) - A_{\boldsymbol{\theta}\boldsymbol{\eta}})\boldsymbol{\gamma}| &= |\mathbf{b}^\top \text{diag}(\mathbf{a})\mathbf{g} + \mathbf{b}^\top \text{diag}(\mathbf{c})\mathbf{g} + \mathbf{d}^\top \text{diag}(\mathbf{e})\mathbf{g} + \mathbf{h}^\top \text{diag}(\mathbf{k})\mathbf{g}| \\
&= \left| \sum_{i=1}^n b_i a_i g_i + b_i c_i g_i + d_i e_i g_i + h_i k_i g_i \right| \\
&\leq \|\mathbf{b}\|_\infty \|\mathbf{a}\| \|\mathbf{g}\| + \|\mathbf{b}\| \|\mathbf{c}\|_\infty \|\mathbf{g}\| + \|\mathbf{d}\| \|\mathbf{e}\|_\infty \|\mathbf{g}\| + \|\mathbf{h}\|_\infty \|\mathbf{k}\| \|\mathbf{g}\| \\
&\leq \{\mathbf{C}^* \mathbf{w}_{p,0} \mathbf{r} + \mathbf{w}_{p,1} \mathbf{r} + \mathbf{C}^* \mathbf{w}_{p,1} \mathbf{r} + \mathbf{w}_{p,1} (\mathbf{C}^* \mathbf{r} + \square)\} / \sqrt{n} \\
&\leq \frac{3\mathbf{C}^* \mathbf{w}_{p,1}}{\sqrt{n}} \{\mathbf{r} + \square\},
\end{aligned}$$

taking into account that $\mathbf{w}_{p,0} \leq \mathbf{w}_{p,1}$.

Third, we write,

$$\begin{aligned}
D_{\boldsymbol{\eta}\boldsymbol{\eta}}^2(\mathbf{v}) - D_{\boldsymbol{\eta}\boldsymbol{\eta}}^2 &= \Phi^\top \text{diag}\{(\Psi'(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta})^2 - (\Psi'(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*)^2\} \Phi \\
&\quad + \Phi^\top \text{diag}\{\Psi''(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}\} \text{diag}\{\mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}\} \Phi.
\end{aligned}$$

Define,

$$\begin{aligned}
\mathbf{a} &= \Psi'(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta} + \Psi'(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*, & \|\mathbf{a}\|_\infty &\leq 2\mathbf{C}^*, \\
\mathbf{b} &= (\Psi'(\Phi\boldsymbol{\eta}) - \Psi'(\Phi\boldsymbol{\eta}^*))^\top \boldsymbol{\theta}^*, & \|\mathbf{b}\| &\leq \mathbf{C}^*\mathbf{r}, \\
\mathbf{c} &= \Psi'(\Phi\boldsymbol{\eta})^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*), & \|\mathbf{c}\|_\infty &\leq \mathbf{w}_{p,1}\mathbf{r}/\sqrt{n}, \\
\mathbf{d} &= \Psi''(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}, & \|\mathbf{d}\|_\infty &\leq \mathbf{C}^* + \mathbf{w}_{p,2}\mathbf{r}/\sqrt{n}, \\
\mathbf{e} &= \mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}, & \|\mathbf{e}\| &\leq \mathbf{C}^*\mathbf{r} + \square, \\
\mathbf{g} &= \Phi\boldsymbol{\gamma}, & \|\mathbf{g}\|_\infty &\leq \mathbf{v}_q/\sqrt{n}, & \|\mathbf{g}\| &= 1.
\end{aligned}$$

Then,

$$\begin{aligned}
|\boldsymbol{\gamma}^\top (D_{\eta\eta}^2(\mathbf{v}) - D_{\eta\eta}^2)\boldsymbol{\gamma}| &= \left| \mathbf{g}^\top \text{diag}\{\mathbf{a}\} \text{diag}\{\mathbf{b} + \mathbf{c}\} \mathbf{g} + \mathbf{g}^\top \text{diag}\{\mathbf{d}\} \text{diag}\{\mathbf{e}\} \mathbf{g} \right| \\
&= \left| \sum_{i=1}^n (g_i^2 a_i b_i + g_i^2 a_i c_i + g_i^2 d_i e_i) \right| \\
&\leq \|\mathbf{a}\|_\infty \|\mathbf{g}\|_\infty \|\mathbf{g}\| \|\mathbf{b}\| + \|\mathbf{a}\|_\infty \|\mathbf{c}\|_\infty \|\mathbf{g}\|^2 + \|\mathbf{d}\|_\infty \|\mathbf{g}\|_\infty \|\mathbf{g}\| \|\mathbf{e}\| \\
&\leq \{4(\mathbf{C}^*)^2 \nu_q + 2\mathbf{C}^* \mathbf{w}_{p,1} + \mathbf{C}^* \mathbf{w}_{p,2} \nu_q / \sqrt{n}\} \frac{\mathbf{r} + \square}{n^{1/2}}
\end{aligned}$$

Putting all together gives the required bound. \square

4.4 Check of exponential moments of derivatives

The following two lemmas check $(\mathcal{E}\mathcal{D}_1)$ and $(\mathcal{E}\mathcal{D}_2)$ of Andresen and Spokoiny (2013), respectively.

Lemma 4.7. *For $\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})$ it holds for each $\boldsymbol{\gamma} \in \mathbb{R}^{p+q}$*

$$\log \mathbb{E} \exp \left(\boldsymbol{\gamma}^\top H_0^{-1} \nabla \zeta_\Phi(\mathbf{v}) \right) \leq \frac{\nu_1(\mathbf{r})^2 \|\boldsymbol{\gamma}\|^2}{2}$$

with $\nu_1(\mathbf{r}) = \sqrt{2} \nu_\varepsilon [1 + \sigma \mathbf{C}_{f,1} + \sigma \mathbf{w}_{p,1} \mathbf{r} / \sqrt{n}]$.

Proof. We have for $\boldsymbol{\gamma} = (\boldsymbol{\alpha}, \boldsymbol{\beta})$, $\boldsymbol{\alpha} \in \mathbb{R}^p$ and $\boldsymbol{\beta} \in \mathbb{R}^q$

$$\begin{aligned}
\log \mathbb{E} \exp \left(\boldsymbol{\gamma}^\top H_0^{-1} \nabla \zeta_\Phi(\mathbf{v}) \right) &= \log \mathbb{E} \exp \left(\boldsymbol{\alpha}^\top D_{\theta\theta}^{-1} \Psi(\Phi\boldsymbol{\eta}) \boldsymbol{\varepsilon} \right. \\
&\quad \left. + \boldsymbol{\beta}^\top \Phi^\top \boldsymbol{\nu} + \sigma \boldsymbol{\beta}^\top \Phi^\top \text{diag}\{\Psi'(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}\} \boldsymbol{\varepsilon} \right) \\
&\leq \frac{\nu_\varepsilon^2 \|\mathbf{u}\|^2}{2} + \frac{\nu_\varepsilon^2 \|\boldsymbol{\beta}\|^2}{2},
\end{aligned}$$

with $\mathbf{u} = \boldsymbol{\alpha}^\top D_{\theta\theta}^{-1} \Psi(\Phi\boldsymbol{\eta}) + \sigma \boldsymbol{\beta}^\top \Phi^\top \text{diag}\{\Psi'(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}\}$. By Lemma 4.3 we get

$$\begin{aligned}
\|\boldsymbol{\alpha}^\top D_{\theta\theta}^{-1} \Psi(\Phi\boldsymbol{\eta})\| &\leq \|\boldsymbol{\alpha}^\top D_{\theta\theta}^{-1} \Psi(\Phi\boldsymbol{\eta}^*)\| + \|\boldsymbol{\alpha}^\top D_{\theta\theta}^{-1} \{\Psi(\Phi\boldsymbol{\eta}) - \Psi(\Phi\boldsymbol{\eta}^*)\}\| \\
&\leq (1 + \sigma \mathbf{w}_{p,1} \mathbf{r} / \sqrt{n}) \|\boldsymbol{\alpha}\|
\end{aligned}$$

Then, we use Lemma 4.2 to get $\|\boldsymbol{\beta}^\top \Phi^\top \text{diag}\{\Psi'(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}\}\| \leq (\mathbf{C}_{f,1} + \mathbf{w}_{p,1} \mathbf{r} / \sqrt{n}) \|\boldsymbol{\beta}\|$.

Summing up we get the required statement. \square

Lemma 4.8. *Suppose $\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})$. Then it holds for each $\boldsymbol{\gamma}_{1,2} \in \mathbb{R}^{p+q}$*

$$\log \mathbb{E} \exp \left(\boldsymbol{\gamma}_1^\top H_0^{-1} \nabla^2 \zeta_\Phi(\mathbf{v}) H_0^{-1} \boldsymbol{\gamma}_2 \right) \leq \frac{\nu_2(\mathbf{r})^2 \|\boldsymbol{\gamma}_1\|^2 \|\boldsymbol{\gamma}_2\|^2}{2}$$

with $\nu_2(\mathbf{r})^2 = \sqrt{3} \nu_\varepsilon^2 \sigma [\mathbf{w}_{p,1} + \sigma (\mathbf{C}_{f,2} + \mathbf{w}_{p,2} \mathbf{r} / \sqrt{n}) \nu_q] / \sqrt{n}$.

Proof. Let $\gamma_i = (\alpha_i, \beta_i)$, $\alpha_i \in \mathbb{R}^p$ and $\beta_i \in \mathbb{R}^q$, $i = 1, 2$. Define then $\mathbf{b}_i = \Phi\beta_i$. By orthonormality of columns of Φ we have $\|\mathbf{b}_i\| = \|\beta_i\|$ and by condition (v) we have $\|\mathbf{b}_i\|_\infty \leq \frac{\nu_q}{\sqrt{n}}\|\beta_i\|$.

So, by generalized Hölder's inequality we have

$$\begin{aligned} & \log \mathbb{E} \exp \left(\gamma_1^\top H_0^{-1} \nabla^2 \zeta_\Phi(\mathbf{v}) H_0^{-1} \gamma_2 \right) \\ & \leq \log \mathbb{E} \exp \left(\sigma \alpha_1^\top D_{\theta\theta}^{-1} \Psi'(\Phi\eta) \text{diag}(\varepsilon) \mathbf{b}_2 + \sigma \alpha_2^\top D_{\theta\theta}^{-1} \Psi'(\Phi\eta) \text{diag}(\varepsilon) \mathbf{b}_1 \right. \\ & \quad \left. + \sigma^2 \mathbf{b}_1^\top \text{diag} \{ \Psi''(\Phi\eta)^\top \theta \} \text{diag}(\varepsilon) \mathbf{b}_2 \right) \\ & \leq \frac{1}{3} \log \mathbb{E} \exp \left(3 \sigma \alpha_1^\top D_{\theta\theta}^{-1} \Psi'(\Phi\eta) \text{diag}(\varepsilon) \mathbf{b}_2 \right) \\ & \quad + \frac{1}{3} \log \mathbb{E} \exp \left(3 \sigma \alpha_2^\top D_{\theta\theta}^{-1} \Psi'(\Phi\eta) \text{diag}(\varepsilon) \mathbf{b}_1 \right) \\ & \quad + \frac{1}{3} \log \mathbb{E} \exp \left(3 \sigma^2 \mathbf{b}_1^\top \text{diag} \{ \Psi''(\Phi\eta)^\top \theta \} \text{diag}(\varepsilon) \mathbf{b}_2 \right). \end{aligned}$$

Let's deal with each term separately. For $\alpha \in \mathbb{R}^p$, $\mathbf{b} = \Phi\beta \in \mathbb{R}^n$ we have

$$\alpha^\top D_{\theta\theta}^{-1} \Psi'(\mathbf{X}) \text{diag}(\varepsilon) \mathbf{b} = \mathbf{u}^\top \text{diag}(\varepsilon) \mathbf{b} = \sum_{i=1}^n u_i b_i \varepsilon_i,$$

where $\mathbf{u}^\top = \alpha^\top D_{\theta\theta}^{-1} \Psi'(\Phi\eta)$ by Lemma 4.1 satisfy $\|\mathbf{u}\|_\infty \leq w_{p,1} \|\alpha\| / \sqrt{n}$. Thus,

$$\begin{aligned} \log \mathbb{E} \exp \left(3 \sigma \alpha^\top D_{\theta\theta}^{-1} \Psi'(\Phi\eta) \text{diag}(\varepsilon) \mathbf{b} \right) &= \sum_{i=1}^n \log \mathbb{E} \exp(3 \sigma u_i b_i \varepsilon_i) \\ &\leq \frac{9\nu_\varepsilon^2 \sigma^2 \|\mathbf{u}\|_\infty^2 \|\mathbf{b}\|^2}{2} \\ &\leq \frac{9\nu_\varepsilon^2 \sigma^2 w_{p,1}^2 \|\alpha\|^2 \|\beta\|^2 / n}{2}. \end{aligned}$$

Further,

$$\mathbf{b}_1^\top \text{diag}(\Psi''(\Phi\eta)^\top \theta) \text{diag}(\varepsilon) \mathbf{b}_2 = \sum_{i=1}^n b_{1,i} k_i \varepsilon_i$$

with $\mathbf{k} = \text{diag} \{ \Psi''(\Phi\eta)^\top \theta \} \mathbf{b}_2$, which by Lemma 4.2 and condition (v) satisfy $\|\mathbf{k}\|_\infty \leq (\mathbb{C}_{f,2} + w_{p,2} \mathbf{r} / \sqrt{n}) \frac{\nu_q}{\sqrt{n}} \|\beta_2\|$. This brings us to

$$\log \mathbb{E} \exp \left(3 \sigma^2 \mathbf{b}_1^\top \text{diag}(\Psi''(\Phi\eta)^\top \theta) \text{diag}(\varepsilon) \mathbf{b}_2 \right) \leq \frac{9\sigma^4 (\mathbb{C}_{f,2} + w_{p,2} \mathbf{r} / \sqrt{n})^2 \frac{\mathbb{C}_\Phi^2}{n} \|\beta_1\|^2 \|\beta_2\|^2}{2}.$$

Bringing those inequalities together we get

$$\begin{aligned} \log \mathbb{E} \exp \left(\boldsymbol{\gamma}_1^\top H_0^{-1} \nabla^2 \zeta_\Phi(\boldsymbol{v}) H_0^{-1} \boldsymbol{\gamma}_2 \right) &\leq \frac{3\nu_\varepsilon^2 \sigma^2}{2} \left[\frac{\mathbf{w}_{p,1}^2}{n} (\|\boldsymbol{\alpha}_1\|^2 \|\boldsymbol{\beta}_2\|^2 + \|\boldsymbol{\alpha}_1\|^2 \|\boldsymbol{\beta}_2\|^2) \right. \\ &\quad \left. + \sigma^2 (\mathbf{C}_{f,2} + \mathbf{w}_{p,2} \mathbf{r} / \sqrt{n})^2 \frac{\mathbf{C}_\Phi^2}{n} \|\boldsymbol{\beta}_1\|^2 \|\boldsymbol{\beta}_2\|^2 \right] \\ &\leq \frac{3\nu_\varepsilon^2 \sigma^2 \left[\mathbf{w}_{p,1}^2 + \sigma^2 (\mathbf{C}_{f,2} + \mathbf{w}_{p,2} \mathbf{r} / \sqrt{n})^2 \mathbf{v}_q^2 \right] \|\boldsymbol{\gamma}_1\|^2 \|\boldsymbol{\gamma}_2\|^2}{2n}. \end{aligned}$$

□

4.5 Large deviation bound

In this section we state large deviation inequality for the full estimator $(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}})$. First, we check concentration of the plug-in estimator.

Lemma 4.9. *It holds with probability at least $1 - 4e^{-\mathbf{x}}$,*

$$\tilde{\boldsymbol{v}}^{(\text{pl})} = \left(\tilde{\boldsymbol{\theta}}^{(\text{pl})}, \tilde{\boldsymbol{\eta}}^{(\text{pl})} \right) \in \mathcal{Y}_o(\mathbf{r}_1),$$

where $\mathbf{r}_1 = 4(1 + \sigma \mathbf{C}_{f,1}) \nu_\varepsilon \sqrt{p + q + \mathbf{x}}$.

Proof. First of all, we have with probability $1 - 2e^{-\mathbf{x}}$,

$$\sigma^{-1} \|\tilde{\boldsymbol{\eta}}^{(\text{pl})} - \boldsymbol{\eta}^*\| \leq \nu_\varepsilon (\sqrt{q} + \sqrt{2\mathbf{x}}).$$

Then, using Lemma 4.3 with $\frac{\sigma \nu_\varepsilon \mathbf{w}_{p,1} (\sqrt{q} + \sqrt{2\mathbf{x}})}{\sqrt{n}} \leq \frac{1}{6}$ by condition (vii),

$$\tilde{\boldsymbol{\Psi}} \tilde{\boldsymbol{\Psi}}^\top \geq D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2 / 2, \quad \left\| (\tilde{\boldsymbol{\Psi}} - \boldsymbol{\Psi}(\Phi \boldsymbol{\eta}^*))^\top \boldsymbol{\theta}^* \right\| \leq \sigma \mathbf{C}_{f,1} \nu_\varepsilon (\sqrt{q} + \sqrt{2\mathbf{x}}).$$

So,

$$\tilde{\boldsymbol{\Psi}}^\top (\tilde{\boldsymbol{\theta}}^{(\text{pl})} - \boldsymbol{\theta}^*) = \tilde{\boldsymbol{\Psi}} (\tilde{\boldsymbol{\Psi}} \tilde{\boldsymbol{\Psi}}^\top)^{-1} \tilde{\boldsymbol{\Psi}} (\mathbf{Y} - \boldsymbol{\Psi}(\Phi \boldsymbol{\eta}^*))^\top \boldsymbol{\theta}^* + (\boldsymbol{\Psi}(\Phi \boldsymbol{\eta}^*) - \tilde{\boldsymbol{\Psi}})^\top \boldsymbol{\theta}^*,$$

and, since $\|(\tilde{\boldsymbol{\Psi}} \tilde{\boldsymbol{\Psi}}^\top)^{-1/2} \tilde{\boldsymbol{\Psi}}\| = 1$, we have

$$\left\| \tilde{\boldsymbol{\Psi}}^\top (\tilde{\boldsymbol{\theta}}^{(\text{pl})} - \boldsymbol{\theta}^*) \right\| \leq \left\| (\tilde{\boldsymbol{\Psi}} \tilde{\boldsymbol{\Psi}}^\top)^{-1/2} \tilde{\boldsymbol{\Psi}} \boldsymbol{\varepsilon} \right\| + \sigma \mathbf{C}_{f,1} \nu_\varepsilon (\sqrt{q} + \sqrt{2\mathbf{x}}),$$

where conditioned on \mathbf{Z} , we have $\left\| (\tilde{\boldsymbol{\Psi}} \tilde{\boldsymbol{\Psi}}^\top)^{-1/2} \tilde{\boldsymbol{\Psi}} \boldsymbol{\varepsilon} \right\| \leq \nu_\varepsilon (\sqrt{p} + \sqrt{2\mathbf{x}})$ with probability at least $1 - 2e^{-\mathbf{x}}$. □

Further we find a stationary point of log-likelihood as a maxima over the set

$$\mathcal{Y}_o(\mathbf{R}_0), \quad \mathbf{R}_0 \stackrel{\text{def}}{=} \frac{\sqrt{n}}{2\sigma \mathbf{w}_{p,1}}.$$

Lemma 4.10. *It holds with probability at least $1 - e^{-x}$,*

$$\sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{R}_0) \setminus \mathcal{Y}_\circ(\mathbf{r}_2)} L_\Phi(\mathbf{v}) - L_\Phi(\mathbf{v}^*) \leq 0, \quad (4.1)$$

where

$$\begin{aligned} \mathbf{r}_2 &= \mathbf{C}_0 \{ \sqrt{p+q+x} \vee \square \}, \\ \mathbf{C}_0 &= 96\nu_\varepsilon [2(\sigma \mathbf{C}_{f,1} \vee 1) + 1] (\sigma \mathbf{C}_{f,1} + 3/2). \end{aligned}$$

Proposition 4.11. *Let*

$$2\sigma \mathbf{C}_0 \mathbf{w}_{p,1} \{ \sqrt{p+q+x} \vee \square \} \leq \sqrt{n}.$$

Then, we have with probability at least $1 - 5e^{-x}$

$$(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}), (\boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*}) \in \mathcal{Y}_\circ(\mathbf{r}_0),$$

where $\mathbf{r}_0 = \mathbf{f}^{-1}\mathbf{F}(\mathbf{r}_1 + \mathbf{r}_2) + \mathbf{r}_1$.

Proof. The given inequality means $\mathbf{R}_0 > \mathbf{r}_2$, thus the maximum of $L_\Phi(\mathbf{v})$ over $\mathcal{Y}_\circ(\mathbf{R}_0)$ lies inside $\mathcal{Y}_\circ(\mathbf{r}_2)$ with probability at least $1 - e^{-x}$ and, therefore, is a stationary point. So, we have such a point at most $\mathbf{r}_1 + \mathbf{r}_2$ far from the plug-in estimator. Thus, the closest stationary point to plug-in satisfies w.h.p.

$$\|H_0(\tilde{\mathbf{v}} - \mathbf{v}^*)\| \leq \mathbf{f}^{-1}\mathbf{F}(\mathbf{r}_1 + \mathbf{r}_2) + \mathbf{r}_1.$$

Obviously, the same arguments work for $(\boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*})$. □

Remark 4.1. The inclusion $\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}_\circ(\mathbf{r}_0)$ implies

$$\hat{\mathbf{v}}_{\boldsymbol{\theta}^*} = \left(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*} - D_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-2} A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \right) \in \mathcal{Y}_\circ(3\mathbf{r}_0),$$

since

$$\begin{aligned} \|H_0(\hat{\mathbf{v}}_{\boldsymbol{\theta}^*} - \mathbf{v}^*)\| &\leq \mathbf{r}_0 + \sigma^{-1} \|\tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*} - \boldsymbol{\eta}^*\| + \|H_\sigma D_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-2} A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \\ &\leq 2\mathbf{r}_0 + \|H_\sigma D_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1}\| \times \|D_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1} A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\| \times \|D_{\boldsymbol{\theta}\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \\ &\leq 3\mathbf{r}_0, \end{aligned}$$

where $\|D_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-1} A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top D_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}\| \leq \rho < 1$ from Lemma (I). This strange object we need to show Wilks expansion.

To proof Lemma 4.10 we are going to straightforwardly use Theorem 2.1 of Andresen and Spokoiny (2013), which requires two conditions: the first one controls the stochastic part of $L_\Phi(\mathbf{v}) - L_\Phi(\mathbf{v}^*)$ and the second one controls it's expectation. Bounding the stochastic part requires check of the exponential moments of the score $\nabla\zeta_\Phi(\mathbf{v})$, which is done by Lemma 4.7. Now we need to check that $-\mathbb{E}[L_\Phi(\mathbf{v}) - L_\Phi(\mathbf{v}^*)]$ grows quadratically.

Lemma 4.12. *Let $A, \mathbf{b} > 0$ satisfy*

$$[4(\sigma\mathbf{C}_{f,1} \vee 1) + 2] \sqrt{\mathbf{b} + 2A^{-2}} \leq 1. \quad (4.2)$$

Then, $\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{R}_0)$ with $\mathbf{r} = \|H_0(\mathbf{v} - \mathbf{v}^)\| \geq A\Box$ implies*

$$-2\mathbb{E}[L_\Phi(\mathbf{v}) - L_\Phi(\mathbf{v}^*)] \geq \mathbf{b}\mathbf{r}^2.$$

Proof of Lemma 4.10. First,

$$\log \mathbb{E} \exp\{\boldsymbol{\gamma}^\top H_0^{-1} \nabla\zeta_\Phi(\mathbf{v})\} \leq \frac{\nu_1^2(\mathbf{R}_0) \|\boldsymbol{\gamma}\|^2}{2}, \quad \mathbf{v} \in \mathcal{Y}_\circ(\mathbf{R}_0),$$

where $\nu_1^2(\mathbf{R}_0)$ is given in Lemma 4.7. Since $\mathbf{R}_0 \leq \mathbf{w}_{p,1}^{-1} \sqrt{n}/2$, we can take

$$\nu_1(\mathbf{R}_0) \leq \nu_1 \stackrel{\text{def}}{=} \sqrt{2\nu_\varepsilon(3/2 + \sigma\mathbf{C}_{f,1})}.$$

Further, setting $A^2 = 2\mathbf{b}^{-1} = 16[2(\sigma\mathbf{C}_{f,1} \vee 1) + 1]^2$ in Lemma 4.12, we have by the theorem mentioned above, that (4.1) holds, when

$$\mathbf{r}_2 \geq \frac{6\nu_1}{\mathbf{b}} \sqrt{2(p+q) + \mathbf{x}},$$

which finishes the proof. \square

Proof of Lemma 4.12. We need to show that for $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta})$ with $\|H_0(\mathbf{v} - \mathbf{v}^*)\| = \mathbf{r}$ and $A\Box \leq \mathbf{r} \leq \mathbf{R}_0$ it holds,

$$\|\mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta}\|^2 + \sigma^{-2} \|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|^2 - \|\mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\|^2 \geq \mathbf{b}\mathbf{r}^2. \quad (4.3)$$

We have,

$$\begin{aligned} \mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta})^\top \boldsymbol{\theta} &= \Psi(\Phi\boldsymbol{\eta}^*)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + (\Psi(\Phi\boldsymbol{\eta}) - \Psi(\Phi\boldsymbol{\eta}^*))^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \\ &\quad + (\Psi(\Phi\boldsymbol{\eta}) - \Psi(\Phi\boldsymbol{\eta}^*))^\top \boldsymbol{\theta}^*. \end{aligned}$$

Denote,

$$\begin{aligned}
\mathbf{x}_\theta &= \Psi(\Phi\boldsymbol{\eta}^*)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*), & \mathbf{r}_\theta &\stackrel{\text{def}}{=} \|\mathbf{x}_\theta\|, \\
\mathbf{x}_\eta &= \sigma^{-2}(\boldsymbol{\eta} - \boldsymbol{\eta}^*), & \mathbf{r}_\eta &\stackrel{\text{def}}{=} \|\mathbf{x}_\eta\|, \\
\mathbf{y}_1 &= (\Psi(\Phi\boldsymbol{\eta}) - \Psi(\Phi\boldsymbol{\eta}^*))^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*), & \|\mathbf{y}_1\| &\leq \sigma n^{-1/2} \mathbf{w}_{p,1} \mathbf{r}_\theta \mathbf{r}_\eta, \\
\mathbf{y}_2 &= (\Psi(\Phi\boldsymbol{\eta}) - \Psi(\Phi\boldsymbol{\eta}^*))^\top \boldsymbol{\theta}^*, & \|\mathbf{y}_2\| &\leq \sigma \mathbf{C}_{f,1} \mathbf{r}_\eta, \\
\mathbf{y}_3 &= \mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*, & \|\mathbf{y}_3\| &= \square.
\end{aligned}$$

Then (4.3) rewrites as

$$\|\mathbf{x}_\theta + \mathbf{y}_1 + \mathbf{y}_2 + \mathbf{y}_3\|^2 + \mathbf{r}_\eta^2 - \square^2 \geq \mathbf{b} (\mathbf{r}_\theta^2 + \mathbf{r}_\eta^2).$$

If $\mathbf{r}_\eta^2 - \square^2 \geq \mathbf{b} (\mathbf{r}_\theta^2 + \mathbf{r}_\eta^2)$ then (4.3) is obviously satisfied. Otherwise, defining $\widehat{\mathbf{b}} \stackrel{\text{def}}{=} \mathbf{b} + A^{-2}$, we have

$$\mathbf{r}_\eta \leq \sqrt{2\widehat{\mathbf{b}}} \mathbf{r}_\theta, \quad \square \leq \sqrt{2} A^{-1} \mathbf{r}_\theta.$$

Further, $\sigma^{-1} \|\boldsymbol{\eta} - \boldsymbol{\eta}^*\| \leq (2\sigma \mathbf{w}_{p,1})^{-1} \sqrt{n}$ and therefore, $\|\mathbf{y}_1\| \leq \mathbf{r}_\theta/2$. So,

$$\mathbf{r}_\theta - \|\mathbf{y}_1\| - \|\mathbf{y}_2\| - \|\mathbf{y}_3\| \geq \left(1 - \frac{1}{2} - \sigma \mathbf{C}_{f,1} \sqrt{2\widehat{\mathbf{b}}} - \sqrt{2} A^{-1}\right) \mathbf{r}_\theta.$$

We need to show that the last expression to square is at least $\mathbf{b} \mathbf{r}_\theta^2 + \square^2 \leq \{\mathbf{b} + 2A^{-2}\} \mathbf{r}_\theta^2$ and the proof is finished by checking that (4.2) ensures

$$\frac{1}{2} - \sigma \mathbf{C}_{f,1} \sqrt{2\widehat{\mathbf{b}}} - \sqrt{2} A^{-1} \geq \sqrt{\widehat{\mathbf{b}} + A^{-2}}.$$

□

4.6 Local linear approximation of the score

We start with the linear approximation of gradient of the likelihood. Define for $\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{x})$ the following process

$$\check{\chi}(\mathbf{v}) \stackrel{\text{def}}{=} \check{D}_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1} \check{\nabla}_\theta L_\Phi(\mathbf{v}) - \check{\boldsymbol{\xi}} + \check{D}_{\boldsymbol{\theta}\boldsymbol{\theta}}(\mathbf{v} - \mathbf{v}^*),$$

where operator $\check{\nabla}_\theta = \nabla_\theta - A_{\boldsymbol{\theta}\boldsymbol{\eta}} D_{\boldsymbol{\eta}\boldsymbol{\eta}}^{-2} \nabla_\eta$, $\check{\boldsymbol{\xi}} = \check{D}_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1} \check{\nabla}_\theta \zeta_\Phi(\mathbf{v}^*)$ and $\check{D}_{\boldsymbol{\theta}\boldsymbol{\theta}}$ defined in (2.4).

In what follows we will also exploit function

$$z_{p+q}(\mathbf{x}) = \sqrt{2(p+q) + \mathbf{x}},$$

which is involved in the concentration inequalities for suprema of empirical processes.

Lemma 4.13. *On a set $C(\mathbf{x}, \mathbf{r}) \subset \Omega$ of probability at least $1 - e^{-\mathbf{x}}$ it holds*

$$\sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \|\check{\chi}(\mathbf{v})\| \leq \diamond(\mathbf{x}, \mathbf{r}) + \square,$$

where

$$\diamond(\mathbf{x}, \mathbf{r}) = 2(1 - \rho)^{-1/2} \{6\nu_2(\mathbf{r})z_{p+q}(\mathbf{x}) + \delta(\mathbf{r})\} \mathbf{r}$$

and $\square = \|\mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\|$ from the condition (vi).

Proof. Define the process, corresponding to the linear approximation of $\nabla L_\Phi(\mathbf{v})$ among the biased point \mathbf{v}^* ,

$$\chi(\mathbf{v}) \stackrel{\text{def}}{=} D^{-1} \{ \nabla L_\Phi(\mathbf{v}) - \nabla L_\Phi(\mathbf{v}^*) + D^2(\mathbf{v} - \mathbf{v}^*) \}.$$

We can split this vector into two terms, the expectation and the stochastic part. First, using $\mathbb{E}\nabla L_\Phi = \nabla \mathbb{E}L_\Phi$ we get

$$\mathbb{E}\chi(\mathbf{v}) = D^{-1} \{ \nabla \mathbb{E}L_\Phi(\mathbf{v}) - \nabla \mathbb{E}L_\Phi(\mathbf{v}^*) + D^2(\mathbf{v} - \mathbf{v}^*) \},$$

which using (\mathcal{L}_0) and (\mathcal{I}) can be easily bounded, see details in Andresen and Spokoiny (2013),

$$\|\mathbb{E}\chi(\mathbf{v})\| \leq 2(1 - \rho)^{-1/2} \mathbf{r} \delta(\mathbf{r}).$$

Using Theorem 4.1 of Andresen and Spokoiny (2013) along with Lemma 4.8, the value $\|H_0^{-1} \{ \nabla \zeta_\Phi(\mathbf{v}) - \nabla \zeta_\Phi(\mathbf{v}^*) \}\|$ is bounded by $6\nu_2(\mathbf{r})\sqrt{2(p+q)+\mathbf{x}}$ uniformly over $\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})$ with probability at least $1 - e^{-\mathbf{x}}$. To sum up, we get on a set $C(\mathbf{x}, \mathbf{r}) \subset \Omega$

$$\begin{aligned} \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \|\chi(\mathbf{v})\| &\leq \diamond(\mathbf{x}, \mathbf{r}) \\ &\stackrel{\text{def}}{=} 2(1 - \rho)^{-1/2} \{3\nu_2(\mathbf{r})z_{p+q}(\mathbf{x}) + \delta(\mathbf{r})\} \mathbf{r}, \end{aligned}$$

where $\mathbb{P}(C(\mathbf{x}, \mathbf{r})) \geq 1 - e^{-\mathbf{x}}$. Define,

$$\mathbf{b} = D^{-1} \nabla L_\Phi(\mathbf{v}^*), \quad \|\mathbf{b}\| \leq \|\mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\|,$$

with the inequality following from Lemma 4.5. Then, for $\check{\Pi}_\theta \stackrel{\text{def}}{=} \left(I_p \quad -A_{\theta\eta} D_{\eta\eta}^{-2} \right)$ we have $\check{\nabla}_\theta = \check{\Pi}_\theta \nabla$ and the proof is finished by

$$\check{\chi}(\mathbf{v}) = \check{D}_{\theta\theta}^{-1} \check{\Pi}_\theta D \{ \chi(\mathbf{v}) + \mathbf{b} \},$$

where we have $\|\check{D}_{\theta\theta}^{-1} \check{\Pi}_\theta D\| = 1$, since the matrix is orthonormal. \square

4.7 Semiparametric Fisher and square-root Wilks

The goal of the Fisher expansion is approximation $\check{D}_{\theta\theta}(\tilde{\theta} - \theta^*) \approx \check{\xi}$, which can be done as follows. Once we have large deviation bound $(\tilde{\theta}, \tilde{\eta}) \in \mathcal{Y}_o(\mathbf{r}_0)$ with high probability, we can apply Lemma 4.13 with $\mathbf{v} = \tilde{\mathbf{v}}$, which using $\nabla L_\Phi(\tilde{\mathbf{v}}) = 0$ gives

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

on the set $\{\tilde{\mathbf{v}} \in \mathcal{Y}_o(\mathbf{r}_0)\} \cap C(\mathbf{x}, \mathbf{r}_0)$, which is of probability at least $1 - 4e^{-\mathbf{x}}$.

Remark 4.2. We don't require \mathbf{r}_0 to be exactly as in Proposition 4.11, the global concentration can be held by any other means and we only require $\mathcal{Y}_o(\mathbf{r}_0)$ to contain the estimator $(\tilde{\theta}, \tilde{\eta})$ with high probability.

Semiparametric square-root Wilks approximates the distribution of $\{2L_\Phi(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*})\}^{1/2}$, by the the distribution of $\|\check{\xi}\|$. This might be shown via quadratic approximation of the likelihood. The Lemma 4.13 helps to ensure it.

Lemma 4.14. *Define,*

$$\begin{aligned} \alpha(\theta, \eta) &= L_\Phi(\theta, \eta) - L_\Phi(\theta^*, \eta - D_{\eta\eta}^{-2} A_{\theta\eta}^\top (\theta - \theta^*)) \\ &\quad - \{\check{D}_{\theta\theta}(\theta - \theta^*)\}^\top \check{\xi} + \|\check{D}_{\theta\theta}(\theta - \theta^*)\|^2/2. \end{aligned}$$

Then if $(\theta, \eta) \in \mathcal{Y}_o(\mathbf{x})$, it holds (pointwise)

$$\frac{|\alpha(\theta, \eta)|}{\|\check{D}_{\theta\theta}(\theta - \theta^*)\|} \leq \sup_{(\theta, \eta) \in \mathcal{Y}_o(\mathbf{x})} \|\chi(\theta, \eta)\|.$$

Proof. For a given pair (θ, η) define $\check{\eta} = \eta - D_{\eta\eta}^{-2} A_{\theta\eta}^\top \theta$ and $\check{L}_\Phi(\theta, \eta) = L_\Phi(\theta, \check{\eta})$, such that $\nabla_\theta \check{L}_\Phi(\theta, \eta) = \check{\nabla}_\theta L_\Phi(\theta, \check{\eta})$. Define the approximating process

$$\check{\alpha}(\theta, \eta) \stackrel{\text{def}}{=} \check{L}_\Phi(\theta, \eta) - \check{L}_\Phi(\theta^*, \eta) - \{\check{D}_{\theta\theta}(\theta - \theta^*)\}^\top \check{\xi} + \|\check{D}_{\theta\theta}(\theta - \theta^*)\|^2/2.$$

Note, that $\check{\alpha}(\theta^*, \eta) \equiv 0$ for all η . Thus, by mean value theorem $|\check{\alpha}(\theta, \eta)| \leq \|\check{D}_{\theta\theta}(\theta - \theta^*)\| \times \sup_{\theta^\circ \in [\theta, \theta^*]} \|\check{D}_{\theta\theta}^{-1} \check{\nabla}_\theta \check{\alpha}(\theta^\circ, \eta)\|$. Calculating the derivative gives

$$\check{D}_{\theta\theta}^{-1} \check{\nabla}_\theta \check{\alpha}(\theta, \eta) = \check{D}_{\theta\theta}^{-1} \check{\nabla}_\theta L_\Phi(\theta, \check{\eta}) - \check{\xi} + \check{D}_{\theta\theta}^2(\theta - \theta^*) = \chi(\theta, \check{\eta}),$$

which gives the desired result after noticing, that $\alpha(\theta, \check{\eta}) = \check{\alpha}(\theta, \eta)$. \square

As we mentioned in Remark 4.1, it holds $(\tilde{\theta}, \tilde{\eta}), (\tilde{\theta}, \tilde{\eta}_{\theta^*} - D_{\eta\eta}^{-2} A_{\theta\eta}^\top (\tilde{\theta} - \theta^*)) \in \mathcal{Y}_o(3\mathbf{r}_0)$ with probability $\geq 1 - 3e^{-\mathbf{x}}$. So, define

$$\diamond(\mathbf{x}) \stackrel{\text{def}}{=} \diamond(\mathbf{x}, 3\mathbf{r}_0) + \square.$$

Then by Lemmas 4.13, 4.14 we have on the set $C(\mathbf{x}, 3\mathbf{r}_0) \cap \{\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_o(\mathbf{r}_0)\}$ the following chain of inequalities, with $\tilde{\mathbf{u}} = \check{D}_{\theta\theta}(\tilde{\theta} - \theta^*)$

$$\begin{aligned} L_\Phi(\tilde{\mathbf{v}}) - L_\Phi(\tilde{\mathbf{v}}_{\theta^*}) &\leq L_\Phi(\tilde{\theta}, \tilde{\eta}) - L_\Phi(\theta^*, \tilde{\eta} - D_{\eta\eta}^{-2} A_{\theta\eta}^\top (\tilde{\theta} - \theta^*)) \\ &\leq \tilde{\mathbf{u}}^\top \check{\xi} - \|\tilde{\mathbf{u}}\|^2/2 + \diamond^*(\mathbf{x}) \|\tilde{\mathbf{u}}\| \\ &\leq \|\tilde{\mathbf{u}}\|^2/2 + \|\tilde{\mathbf{u}}\| \{\|\tilde{\mathbf{u}} - \check{\xi}\| + \diamond^*(\mathbf{x})\} \\ &\leq \|\tilde{\mathbf{u}}\|^2/2 + 2\|\tilde{\mathbf{u}}\| \diamond(\mathbf{x}), \end{aligned}$$

and on the other hand,

$$\begin{aligned} L_\Phi(\tilde{\mathbf{v}}) - L_\Phi(\tilde{\mathbf{v}}_{\theta^*}) &\geq L_\Phi(\tilde{\theta}, \tilde{\eta}_{\theta^*} - D_{\eta\eta}^{-2} A_{\theta\eta}(\tilde{\theta} - \theta^*)) - L_\Phi(\theta^*, \tilde{\eta}_{\theta^*}) \\ &\geq \|\tilde{\mathbf{u}}\|^2/2 - 2\|\tilde{\mathbf{u}}\| \diamond(\mathbf{x}), \end{aligned}$$

which together gives

$$\left| \{2L_\Phi(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*})\}^{1/2} - \|\tilde{\mathbf{u}}\| \right| \leq \frac{|2L_\Phi(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*}) - \|\tilde{\mathbf{u}}\|^2|}{\|\tilde{\mathbf{u}}\|} \leq 4\diamond(\mathbf{x}),$$

where we have used that on the same subset of Ω it holds $\|\tilde{\mathbf{u}} - \check{\xi}\| \leq \diamond^*(\mathbf{x})$ from the Fisher expansion. Using this inequality one last time gives

$$\left| \{2L_\Phi(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*})\}^{1/2} - \|\check{\xi}\| \right| \leq 5\diamond(\mathbf{x}).$$

4.8 Proof of Theorem 2.1

From Proposition 4.11 we have $\mathbf{r}_0 \asymp \sqrt{p+q+\mathbf{x}} \vee \square$. From Lemma 4.6 and Lemma 4.8 we have,

$$\begin{aligned} \diamond(\mathbf{x}) &\asymp \frac{\sigma(\mathbf{w}_{p,1} + \mathbf{v}_q) \mathbf{r}_0^2}{\sqrt{n}} + \frac{\sigma \mathbf{w}_{p,2} \mathbf{v}_q \mathbf{r}_0^3}{n} + \square, \\ &\lesssim \frac{\sigma(\mathbf{w}_{p,1} + \mathbf{v}_q)(p+q+\mathbf{x})}{\sqrt{n}} + \frac{\sigma \mathbf{w}_{p,2} \mathbf{v}_q (p+q+\mathbf{x})^{3/2}}{n} \\ &\quad + \square + \frac{\sigma(\mathbf{w}_{p,1} + \mathbf{v}_q) \square^2}{\sqrt{n}} + \frac{\sigma \mathbf{w}_{p,2} \mathbf{v}_q \square^3}{n} \\ &\lesssim \sigma \left(\mathbf{w}_{p,1} + \mathbf{v}_q + \mathbf{w}_{p,2} \mathbf{v}_q \sqrt{n^{-1}(p+q+\mathbf{x})} \right) (p+q+\mathbf{x}) n^{-1/2} + \square, \end{aligned}$$

by use of the inequalities given by condition (vii). It remains to check, that

$$\log \mathbb{E} \exp(\gamma^\top \check{\xi}) \leq \frac{\nu_\varepsilon^2 \|\gamma\|^2}{2}, \quad \gamma \in \mathbb{R}^p.$$

Note, that D^2 corresponds to $\text{Var}(\nabla\zeta_{\Phi}(\mathbf{v}^*))$, when the errors ε_i , ν_i are standard normal. In general, one can easily check, that for $\boldsymbol{\xi} = D^{-1}\nabla\zeta_{\Phi}(\mathbf{v}^*)$ and for each $\boldsymbol{\gamma} \in \mathbb{R}^{p+q}$ it holds

$$\log \mathbb{E} \exp(\boldsymbol{\gamma}^\top \boldsymbol{\xi}) \leq \frac{\nu_\varepsilon^2 \|\boldsymbol{\gamma}\|^2}{2}.$$

Now, the moments for the $\check{\boldsymbol{\xi}}$ follow by applying this inequality to $\check{D}_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1} \check{I}_{\boldsymbol{\theta}}^\top D \boldsymbol{\gamma}$, since $\check{D}_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1} \check{I}_{\boldsymbol{\theta}}^\top D$ is an orthonormal matrix.

4.9 Proof of Lemma 2.2

Define $N = \Phi^\top \Phi$ and $\tilde{\Phi} = \Phi N^{-1/2}$. Since N is symmetric, $\tilde{\Phi}^\top \tilde{\Phi} = N^{-1/2} \Phi \Phi^\top N^{-1/2} = I$. Moreover, applying trapezoidal rule for integrating $\gamma(t) = h_m(t)h_k(t)$, we get

$$\begin{aligned} |N_{mk} - \delta_{mk}| &= |\Phi_m^\top \Phi_k - \delta_{mk}| = \left| \frac{1}{n} \sum_{i=1}^n h_m(i/n)h_k(i/n) - \int_0^1 h_m(t)h_k(t)dt \right| \\ &= \left| \frac{1}{n} \sum_{i=1}^n \gamma(i/n) - \int_0^1 \gamma(t)dt \right| \\ &= \frac{|\gamma(1) - \gamma(0)|}{2n} + \left| \frac{1}{n} \sum_{i=1}^n \frac{\gamma((i-1)/n) + \gamma(i/n)}{2} - \int_0^1 \gamma(t)dt \right| \\ &\leq \frac{\|\gamma\|_\infty}{n} + \frac{\|\gamma''\|_\infty}{12n^2}. \end{aligned}$$

Using the inequalities $\|\gamma\|_\infty \leq A_0^2$ and

$$\|\gamma''\|_\infty = \|h_m h_k'' + h_m'' h_k + 2h_m' h_k'\|_\infty \leq A_0 A_2 m^2 + 2A_1^2 m k + A_0 A_2 k^2 \leq 2(A_0 A_2 + A_1^2) q^2,$$

we get that the absolute value of each element of the matrix $N - I$ is less or equal to $\delta = \frac{A_0^2}{n} + \frac{(A_0 A_2 + A_1^2) q^2}{6n^2}$. Thus,

$$\|N - I\| \leq \|N - I\|_{\text{Frob}} = \|\text{vec}(N - I)\|_2 \leq \sqrt{q^2} \|\text{vec}(N - I)\|_\infty \leq q\delta \leq \frac{1}{2},$$

with the last inequality following from the conditions of the lemma. Thus, $\|N^{-1}\| \leq 2$ and

$$C(\tilde{\Phi}) = n \max_{i \leq n} \|N^{-1/2} \tilde{\Phi}_i\|^2 \leq n \|N^{-1}\| \max_{i \leq n} \|\tilde{\Phi}_i\|^2 \leq 2C(\Phi).$$

4.10 Proof of Lemma 2.4

For $h \in C^1[0, 1]$ and $N \in \mathbb{N}$ we have,

$$\int_0^1 h(x) dx \leq \frac{1}{N} \sum_{j=1}^N \min_{x \in [\frac{j-1}{N}, \frac{j}{N})} h(x) + \frac{\|h'\|_\infty}{N}.$$

Take $N < \delta^{-1}p^2$. By condition of the lemma, each interval $\left[\frac{j-1}{N}, \frac{j}{N}\right)$ contains at least $c_0n\left(\frac{1}{N} - \delta p^{-2}\right)$ points. So,

$$\frac{1}{N} \min_{x \in \left[\frac{j-1}{N}, \frac{j}{N}\right)} h(x) \leq \frac{(1 - N\delta p^{-2})^{-1}}{c_0n} \sum_{X_i^* \in \left[\frac{j-1}{N}, \frac{j}{N}\right)} h(X_i^*),$$

which, by summing over $j = 1, \dots, N$, gives

$$\int_0^1 h(x) dx \leq \frac{(c_0(1 - N\delta p^{-2}))^{-1}}{n} \sum_{i=1}^n h(X_i^*) + \frac{\|h'\|_\infty}{N}.$$

Fix arbitrary unit vector $\gamma \in \mathbb{R}^p$ and set $h(x) = (\Psi(x)^\top \gamma)^2$. By (a) we have

$$\int_0^1 h(x) dx = 1, \quad \|h'\|_\infty \leq 2A^2p^2,$$

and taking $N(\delta) = \lfloor (1 - \delta)\delta^{-1}p^2 \rfloor$

$$\frac{1}{n} \sum_{i=1}^n h(X_i^*) \geq c_0(1 - \alpha(\delta)), \quad \alpha(\delta) = \left(1 - \frac{2A^2p^2}{N(\delta)}\right) (1 - N(\delta)\delta p^{-2}).$$

Since $N(\delta) \geq \delta^{-1}p^2 - p^2 - 1$,

$$\alpha(\delta) \leq (\delta + \delta p^{-2}) \left(1 - \frac{2A^2\delta}{1 - \delta}\right).$$

It is left to notice, that $\frac{1}{n} \sum_{i=1}^n h(X_i^*) = \gamma^\top [n^{-1}\Psi(\mathbf{X}^*)\Psi(\mathbf{X}^*)^\top] \gamma$.

4.11 Proof of Lemma 2.5

For any $h \in C^1[0, 1]$ it holds

$$\|h\|_2^2 \geq \frac{1}{n} \sum_{j=1}^n \max_{x \in \left[\frac{j-1}{n}, \frac{j}{n}\right]} h(x)^2 - \frac{2\|h\|_\infty \|h'\|_\infty}{n}.$$

The amount of points X_i^* lying in $\left[\frac{j-1}{n}, \frac{j}{n}\right]$ is bounded by C_0 , so it holds

$$\max_{x \in \left[\frac{j-1}{n}, \frac{j}{n}\right]} h(x)^2 \geq C_0^{-1} \sum_{X_i^* \in \left[\frac{j-1}{n}, \frac{j}{n}\right]} h(X_i^*)^2,$$

which brings us to

$$\|h\|_2^2 \geq -\frac{2\|h\|_\infty \|h'\|_\infty}{n} + \frac{1}{C_0n} \sum_{i \leq n} h(X_i^*)^2.$$

Applying this to $h = f - f_p$, using $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$ and Lemma 2.3, gives

$$\begin{aligned} \|\mathbf{f}^* - \Psi(\mathbf{X}^*)^\top \boldsymbol{\theta}^*\| &\leq \mathbf{c}_0^{1/2} \sqrt{n} \|f - f_p\|_2 + (2\mathbf{C}_0)^{1/2} \sqrt{\|f - f_q\|_\infty \|f' - f'_q\|_\infty} \\ &\leq \mathbf{c}_0^{1/2} Q \sqrt{np}^{-\beta} + \frac{(2\mathbf{C}_0)^{1/2} A Q}{\sqrt{2\beta - 3}} p^{-\beta+1}. \end{aligned}$$

Further, since $\|\Phi \boldsymbol{\eta}^* - \mathbf{X}^*\|^2 = \sum_{i \leq n} (g_q(i/n) - g(i/n))^2$, we use arguments above and Lemma 2.3 to show

$$\begin{aligned} \|\Phi \boldsymbol{\eta}^* - \mathbf{X}^*\| &\leq \sqrt{n} \|g_q - g\|_2 + \sqrt{2 \|g_q - g\|_\infty \|(g_q - g)'\|_\infty} \\ &\leq Q \sqrt{n} q^{-\beta x} + \frac{\sqrt{2} A Q}{\sqrt{2\beta_X - 3}} q^{-\beta x + 1}. \end{aligned}$$

4.12 Proof of Proposition 2.6

First we check the conditions of Theorem 2.1.

(i). Due to (a) and Lemma 2.3, we have for $\beta > 2.5$ and $s = 0, 1, 2$

$$\|f_p^{(s)}\|_\infty \leq \mathbf{c}_{f,s} = \frac{A Q}{\sqrt{2(\beta - s) - 1}}.$$

(ii). We need to transfer condition (c) from the matrix $\mathbb{F} = \Psi(\mathbf{X}^*) \Psi(\mathbf{X}^*)^\top$ to the matrix $D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2 = \Psi(\Phi \boldsymbol{\eta}^*) \Psi(\Phi \boldsymbol{\eta}^*)^\top$. By Lemma 4.1 it holds,

$$\begin{aligned} \|\mathbb{F}^{-1/2} [\Psi(\mathbf{X}^*) - \Psi(\Phi \boldsymbol{\eta}^*)]\|_{\text{op}} &\leq n^{-1/2} \mathbf{w}_{p,1} \|\mathbf{X}^* - \Phi \boldsymbol{\eta}^*\| \\ &\leq \mathbf{c} p^{3/2} q^{-\beta x}, \end{aligned}$$

where the last inequality follows from $\mathbf{w}_{p,1} = \mathbf{f}^{-1/2} A p^{3/2}$. So, for a unit vector $\boldsymbol{\gamma} \in \mathbb{R}^p$, since $\|\boldsymbol{\gamma}^\top \mathbb{F}^{-1/2} \Psi(\mathbf{X}^*)\| = 1$, it holds

$$|\boldsymbol{\gamma}^\top \mathbb{F}^{-1/2} (D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2 - \mathbb{F}) \mathbb{F}^{-1/2} \boldsymbol{\gamma}| \leq \delta (2 + \delta),$$

where $\delta = \mathbf{c} p^{3/2} q^{-\beta x}$. We have then,

$$D_{\boldsymbol{\theta}\boldsymbol{\theta}}^2 \geq (1 - \delta(2 + \delta)) \mathbb{F} \geq \mathbf{f} (1 - \delta(2 + \delta)) (n I_p),$$

which ensures condition (ii), since $p^{3/2} q^{-\beta x} = n^{-(\beta-3/2)/(2\beta+1)}$ is small.

(iii). By (a) we can take $\mu_{m,s} = A m^s$.

(iv). Using Lemma 2.2, we get this condition satisfied with $\mathbf{C}_{\tilde{\Phi}} = 2A\sqrt{q}$. Indeed

$$\mathbf{C}_{\tilde{\Phi}} = n \|\Phi_i\|^2 = \sum_{m=1}^q \phi_m^2(i/m) \leq A^2 q.$$

(v). We have by (2.13)

$$\begin{aligned} \square &= \|\mathbf{f}^* - \Psi(\Phi\boldsymbol{\eta}^*)^\top \boldsymbol{\theta}^*\| \leq \mathbf{C}_{f,1} \|\Phi\boldsymbol{\eta}^* - \mathbf{X}^*\| + \|\mathbf{f}^* - \Psi(\mathbf{X}^*)^\top \boldsymbol{\theta}^*\| \\ &\leq \mathbf{C} \sqrt{n}(p^{-\beta} + q^{-\beta x}), \end{aligned}$$

and $\square \leq \sqrt{n}$ for large enough p and q .

(vii). Taking $\mathbf{w}_{p,s} = Ap^{s+\frac{1}{2}}$ the conditions are satisfied, since $p^{3/2}\sqrt{(p+q)/n}$ is small.

Using $n^{-1}\check{D}_{\boldsymbol{\theta}\boldsymbol{\theta}}^2 \geq (1-\rho)\mathbf{f}I_p$, which is given by Lemma 4.4 and condition (ii), Theorem 2.1 provides us with

$$\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\| \lesssim_{\mathcal{P}} \sqrt{\frac{p}{n}} + \Delta_{p,q},$$

where

$$\Delta_{p,q} = \frac{\sigma(p^{3/2} + q^{1/2})(p+q)}{n} + \sigma(p^{3/2} + q^{1/2})(p^{-\beta} + q^{-\beta x})^2 + q^{-\beta x} + p^{-\beta x}.$$

Since $\tilde{f}_{p,q}(x) = \Psi(x)^\top \tilde{\boldsymbol{\theta}}$ and $f_p(x) = \Psi(x)^\top \boldsymbol{\theta}^*$, we have $\|\tilde{f}_{p,q} - f_p\|_2 = \|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|$. Further, using $\|f - f_p\|_2 \lesssim p^{-\beta}$, given by Sobolev condition, we arrive at

$$\begin{aligned} \|f - \tilde{f}_{p,q}\|_2 &\lesssim_{\mathcal{P}} \sqrt{\frac{p}{n}} + p^{-\beta} + q^{-\beta x} \\ &\quad + \left\{ \frac{\sigma(p^{3/2} + q^{1/2})(p+q)}{n} + \sigma(p^{3/2} + q^{1/2})(p^{-\beta} + q^{-\beta x})^2 \right\}. \end{aligned}$$

The choice of p and q provides

$$\sqrt{\frac{p}{n}} = p^{-\beta} = q^{-\beta x} = n^{-\frac{\beta}{2\beta+1}}.$$

Next, we deal with the second term from the figure braces,

$$(p^{3/2} + q^{1/2})(p^{-\beta} + q^{-\beta x})^2 \asymp n^{\frac{\frac{3}{2}\sqrt{\frac{\beta}{2\beta_X}} - 2\beta}{2\beta+1}} \leq n^{\frac{-\beta}{2\beta+1}},$$

where the inequality holds since $\beta \geq 2 > 3/2$ and $2\beta_X \geq 3 > 1$. Finally, we let us deal with the term

$$\begin{aligned} \frac{(p^{3/2} + q^{1/2})(p+q)}{n} &\asymp \frac{p^{5/2} + q^{3/2} + p^{3/2}q}{n} \\ &= n^{\left(\frac{5}{2}\sqrt{\frac{3\beta}{2\beta_X}}\sqrt{\left(\frac{3}{2} + \frac{\beta}{\beta_X}\right)} - 2\beta - 1\right)/(2\beta+1)} \leq r_n(\beta) = n^{-\beta/(2\beta+1)}, \end{aligned}$$

with the inequality equivalent to the union of the following inequalities

$$\left\{ \begin{array}{l} \frac{5}{2} - 2\beta - 1 \leq -\beta, \\ \frac{3\beta}{2\beta_X} - 2\beta - 1 \leq -\beta, \\ \frac{3}{2} + \frac{\beta}{\beta_X} - 2\beta - 1 \leq -\beta \end{array} \right. \Leftrightarrow \left\{ \begin{array}{l} \frac{3}{2} \leq \beta, \\ \frac{3}{2}\beta \leq \beta\beta_X + \beta_X, \\ \frac{1}{2}\beta_X + \beta \leq \beta\beta_X \end{array} \right. \Leftrightarrow \left\{ \begin{array}{l} \beta \geq \frac{3}{2}, \\ (\beta + 1)(\beta_X - 3/2) \geq -\frac{3}{2}, \\ (\beta - 1/2)(\beta_X - 1) \geq \frac{1}{2}, \end{array} \right.$$

where all of the last one are true, since $\beta \geq 2$ and $\beta_X \geq 3/2$.

References

- Andresen, A. and Spokoiny, V. (2013). Critical dimension in profile semiparametric estimation. Manuscript. arXiv:1303.4640.
- Butucea, C. and Taupin, M.-L. (2008). New M -estimators in semi-parametric regression with errors in variables. *Ann. Inst. Henri Poincaré, Probab. Stat.*, 44(3):393–421.
- Carroll, R., Ruppert, D., and Stefanski, L. (1995). Measurement error in nonlinear models. number 63 in monographs on statistics and applied probability.
- Cheng, C.-L. and Van Ness, J. W. (1999). *Statistical regression with measurement error*. John Wiley & Sons.
- Fan, J., Truong, Y. K., and Wang, Y. (1991). *Nonparametric function estimation involving errors-in-variables*. Springer.
- Fuller, W. A. (2009). *Measurement error models*, volume 305. John Wiley & Sons.
- Juditsky, A. B., Lepski, O. V., and Tsybakov, A. B. (2009). Nonparametric estimation of composite functions. *The Annals of Statistics*, pages 1360–1404.
- Schneeweiß, H. and Mittag, H.-J. (1986). *Lineare Modelle mit fehlerbehafteten Daten*. Physica-Verlag Heidelberg-Wien.
- Spokoiny, V. (2012). Parametric estimation. Finite sample theory. *Ann. Statist.*, 40(6):2877–2909. arXiv:1111.3029.
- Tsybakov, A. (2009). *Introduction to Nonparametric estimation*. Springer New York.
- Wansbeek, T. J. and Meijer, E. (2000). *Measurement error and latent variables in econometrics*, volume 37. Elsevier Amsterdam.

