

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

Preprint

ISSN 2198-5855

**Two convergence results for an alternation maximization
procedure**

Andreas Andresen , Vladimir Spokoiny

submitted: January 9, 2015

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

No. 2061
Berlin 2015



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

Key words and phrases. alternating procedure, EM-algorithm, M-estimation, profile maximum likelihood, local linear approximation, spread, local concentration.

This work was partially supported by DFG Research Units 1735 "Structural Inference in Statistics: Adaptation and Efficiency".

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

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

Andresen and Spokoiny's (2013) "critical dimension in semiparametric estimation" provide a technique for the finite sample analysis of profile M-estimators. This paper uses very similar ideas to derive two convergence results for the alternating procedure to approximate the maximizer of random functionals such as the realized log likelihood in MLE estimation. We manage to show that the sequence attains the same deviation properties as shown for the profile M-estimator in Andresen and Spokoiny (2013), i.e. a finite sample Wilks and Fisher theorem. Further under slightly stronger smoothness constraints on the random functional we can show nearly linear convergence to the global maximizer if the starting point for the procedure is well chosen.

1 Introduction

This paper presents a convergence result for an alternating maximization procedure to approximate M-estimators. Let $\mathbb{Y} \in \mathcal{Y}$ denote some observed random data, and \mathbb{P} denote the data distribution. In the semiparametric profile M-estimation framework the target of analysis is

$$\boldsymbol{\theta}^* = \Pi_{\boldsymbol{\theta}} \mathbf{v}^* = \Pi_{\boldsymbol{\theta}} \underset{\mathbf{v}}{\operatorname{argmax}} \mathbb{E}_{\mathbb{P}} \mathcal{L}(\mathbf{v}, \mathbb{Y}), \quad (1.1)$$

where $\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, $\Pi_{\boldsymbol{\theta}} : \mathcal{Y} \rightarrow \mathbb{R}^p$ is a projection and where \mathcal{Y} is some high dimensional or even infinite dimensional parameter space. This paper focuses on finite dimensional parameter spaces $\mathcal{Y} \subseteq \mathbb{R}^{p^*}$ with $p^* = p + m \in \mathbb{N}$ being the full dimension, as infinite dimensional maximization problem are computationally anyways not feasible. A prominent way of estimating $\boldsymbol{\theta}^*$ is the profile M-estimator (pME)

$$\tilde{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \tilde{\mathbf{v}} \stackrel{\text{def}}{=} \underset{(\boldsymbol{\theta}, \boldsymbol{\eta})}{\operatorname{argmax}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}).$$

The alternating maximization procedure is used in situations where a direct computation of the full maximum estimator (ME) $\tilde{\mathbf{v}} \in \mathbb{R}^{p^*}$ is not feasible or simply very difficult to implement. Consider for example the task to calculate the pME where with scalar random observations $\mathbb{Y} = (y_i)_{i=1}^n \subset \mathbb{R}$, parameter $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^p \times \mathbb{R}^m$ and a function basis $(e_k) \subset L^2(\mathbb{R})$

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}) = -\frac{1}{2} \sum_{i=1}^n \left| y_i - \sum_{k=0}^m \boldsymbol{\eta}_k e_k(\mathbf{X}_i^\top \boldsymbol{\theta}) \right|^2.$$

In this case the maximization problem is high dimensional and non-convex (see Section 3 for more details). But for fixed $\boldsymbol{\theta} \in S_1 \subset \mathbb{R}^p$ maximization with respect to $\boldsymbol{\eta} \in \mathbb{R}^m$ is rather simple while for fixed $\boldsymbol{\eta} \in \mathbb{R}^m$ the maximization with respect to $\boldsymbol{\theta} \in \mathbb{R}^p$ can be feasible for low $p \in \mathbb{N}$. This motivates the following iterative procedure. Given some (data dependent) functional $\mathcal{L} : \mathbb{R}^p \times \mathbb{R}^m \rightarrow \mathbb{R}$ and an initial guess $\tilde{\mathbf{v}}_0 \in \mathbb{R}^{p+m}$ set for $k \in \mathbb{N}$

$$\begin{aligned} \tilde{\mathbf{v}}_{k,k+1} &\stackrel{\text{def}}{=} (\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_{k+1}) = \left(\tilde{\boldsymbol{\theta}}_k, \underset{\boldsymbol{\eta} \in \mathbb{R}^m}{\operatorname{argmax}} \mathcal{L}(\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\eta}) \right), \\ \tilde{\mathbf{v}}_{k,k} &\stackrel{\text{def}}{=} (\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k) = \left(\underset{\boldsymbol{\theta} \in \mathbb{R}^p}{\operatorname{argmax}} \mathcal{L}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k), \tilde{\boldsymbol{\eta}}_k \right). \end{aligned} \quad (1.2)$$

The so called "alternation maximization procedure" (or minimization) is a widely applied algorithm in many parameter estimation tasks (see [9], [13], [10] or [17]). Some natural questions arise: Does the sequence $(\tilde{\boldsymbol{\theta}}_k)$ converge to a limit that satisfies the same statistical properties as the profile estimator? And if the answer is yes, after how many steps does the sequence acquire these properties? Under what circumstances does the sequence actually converge to the global maximizer $\tilde{\mathbf{v}}$? This problem is hard because the behavior of each step of the sequence is determined by the actual finite sample realization of the functional $\mathcal{L}(\cdot, \mathbb{Y})$. To the authors' knowledge no general "convergence" result is available that answers the questions from above except for the treatment of specific models (see again [9], [13], [10] or [17]).

We address this difficulty via employing new finite sample techniques of [2] and [14] which allow to answer the above questions: with growing iteration number $k \in \mathbb{N}$ the estimators $\tilde{\theta}_k$ attain the same statistical properties as the profile M-estimator and Theorem 2.2 provides a choice of the necessary number of steps $K \in \mathbb{N}$. Under slightly stronger conditions on the structure of the model we can give a convergence result to the global maximizer that does not rely on unimodality. Further we can address the important question under which ratio of full dimension $p^* = p + m \in \mathbb{N}$ to sample size $n \in \mathbb{N}$ the sequence behaves as desired. For instance for smooth \mathcal{L} our results become sharp if p^*/\sqrt{n} is small and convergence to the full maximizer already occurs if p^*/n is small.

The alternation maximization procedure can be understood as a special case of the Expectation Maximization algorithm (EM algorithm) as we will illustrate below. The EM algorithm itself was derived by [5] who generalized particular versions of this approach and presented a variety of problems where its application can be fruitful; for a brief history of the EM algorithm see [11] (Sect. 1.8). We briefly explain the EM algorithm. Take observations $(\mathbb{X}) \sim IP_\theta$ for some parametric family $(IP_\theta, \theta \in \Theta)$. Assume that a parameter $\theta \in \Theta$ is to be estimated as maximizer of the functional $\mathcal{L}_c(\mathbb{X}, \theta) \in \mathbb{R}$, but that only $\mathbb{Y} \in \mathcal{Y}$ is observed, where $\mathbb{Y} = f_Y(\mathbb{X})$ is the image of the complete data set $\mathbb{X} \in \mathcal{X}$ under some map $f_Y : \mathcal{X} \rightarrow \mathcal{Y}$. Prominent examples for the map f_Y are projections onto some components of \mathbb{X} if both are vectors. The information lost under the map can be regarded as missing data or latent variables. As a direct maximization of the functional is impossible without knowledge of \mathbb{X} the EM algorithm serves as a workaround. It consists of the iteration of tow steps: starting with some initial guess $\tilde{\theta}_0$ the kth "Expectation step" derives the functional Q via

$$Q(\theta, \theta_k) = \mathbb{E}_{\theta_k}[\mathcal{L}_c(\mathbb{X}, \theta) | \mathbb{Y}],$$

which means that on the right hand side the conditional expectation is calculated under the distribution IP_{θ_k} . The kth "Maximation step" then simply locates the maximizer θ_{k+1} of Q .

Since the algorithm is very popular in applications a lot of research on its behaviour has been done. We are only dealing with a special case of this procedure so we restrict our selves to citing the well known convergence result by [16]. Wu presents regularity conditions that ensure that $\mathcal{L}(\theta_{k+1}) \geq \mathcal{L}(\theta_k)$ where

$$\mathcal{L}(\theta, \mathbb{Y}) \stackrel{\text{def}}{=} \log \int_{\{\mathbb{X} | \mathbb{Y} = f_Y(\mathbb{X})\}} \exp \mathcal{L}_c(\mathbb{X}, \theta) d\mathbb{X},$$

such that $\mathcal{L}(\theta_k) \rightarrow \mathcal{L}^*$ for some limit value $\mathcal{L}^* > 0$, that may depend on the starting point θ_0 . Additionally Wu gives conditions that guarantee that the sequence θ_k (possibly a sequence of sets) converges to $C(\mathcal{L}^*) \stackrel{\text{def}}{=} \{\theta | \mathcal{L}(\theta) = \mathcal{L}^*\}$. [5] show that the speed of convergence is linear in the case of point valued θ_k and of some differentiability criterion being met. A limitation of these results is that it is not clear whether $\mathcal{L}^* = \sup \mathcal{L}(\theta)$ and thus it is not guaranteed that $C(\mathcal{L}^*)$ is the desired MLE and not just some local maximum. Of course this problem disappears if $\mathcal{L}(\cdot)$ is unimodal and the regularity conditions are met but this assumption may be too restrictive.

In a recent work [3] present a new way of addressing the properties of the EM sequence in a very general i.i.d. setting, based on concavity of $\theta \mapsto \mathbb{E}_{\theta^*}[\mathcal{L}_c(\mathbb{X}, \theta)]$. They show that if

additional to concavity the functional \mathcal{L}_c is smooth enough (First order stability) and if for a sample (\mathbf{Y}_i) with high probability an uniform bound holds of the kind

$$\sup_{\theta \in B_r(\theta^*)} \left| \sum_{i=1}^n \operatorname{argmax}_{\theta^\circ} \mathbb{E}_\theta[\mathcal{L}_c(\mathbb{X}, \theta^\circ) | \mathbf{Y}_i] - \operatorname{argmax}_{\theta^\circ} \mathbb{E}_{\theta^*}[\mathbb{E}_\theta[\mathcal{L}_c(\mathbb{X}, \theta^\circ) | \mathbb{Y}]] \right| \leq \epsilon_n \quad (1.3)$$

that then with high probability and some $\rho < 1$

$$\|\tilde{\theta}_k - \theta^*\| \leq \rho^k \|\theta_0 - \theta^*\| + C\epsilon_n. \quad (1.4)$$

Unfortunately this does not answer our two questions to full satisfaction. First the bound (1.3) is rather high level and has to be checked for each model, while we seek (and find) properties of the functional - such as smoothness and bounds on the moments of its gradient - that lead to comparably desirable behavior. Further with (1.4) it remains unclear whether for large $k \in \mathbb{N}$ the alternating sequence satisfies a Fisher expansion or whether a Wilks type phenomenon occurs. In particular it remains open which ratio of dimension to sample size ensures good performance of the procedure. Also the actual convergence of $\tilde{\theta}_k \rightarrow \theta^*$ is not implied, as the right hand side in (1.4) is bounded from below by $C\epsilon_n > 0$.

Remark 1.1. In the context of the alternating procedure the bound (1.3) would read

$$\max_{\theta^\circ \in B_r(\theta^*)} \left| \operatorname{argmax}_{\theta} \mathcal{L}(\theta, \tilde{\eta}_{\theta^\circ}) - \operatorname{argmax}_{\theta} \mathbb{E} \mathcal{L}(\theta, \tilde{\eta}_{\theta^\circ}) \right| \leq \epsilon_n,$$

which is still difficult to check.

To see that the procedure (1.2) is a special case of the EM algorithm denote in the notation from above $\mathbb{X} = (\operatorname{argmax}_{\eta} \mathcal{L}\{(\theta, \eta), \mathbb{Y}\}, \mathbb{Y})$ - where θ is the parameter specifying the distribution \mathbb{P}_θ - and $f_Y(\mathbb{X}) = \mathbb{Y}$. Then with $\mathcal{L}_c(\theta, \mathbb{X}) = \mathcal{L}_c(\theta, \eta, \mathbb{Y}) \stackrel{\text{def}}{=} \mathcal{L}(\theta, \eta)$

$$Q(\theta, \tilde{\theta}_{k-1}) = \mathbb{E}_{\tilde{\theta}_{k-1}}[\mathcal{L}_c(\theta, \mathbb{X}) | \mathbb{Y}] = \mathcal{L}_c\left(\theta, \operatorname{argmax}_{\eta} \mathcal{L}\{(\tilde{\theta}_{k-1}, \eta), \mathbb{Y}\}, \mathbb{Y}\right) = \mathcal{L}(\theta, \tilde{\eta}_k),$$

and thus the resulting sequence is the same as in (1.2). Consequently the convergence results from above apply to our problem if the involved regularity criteria are met. But as noted these results do not tell us if the limit of the sequence $(\tilde{\theta}_k)$ actually is the profile and the statistical properties of limit points are not clear without too restrictive assumptions on \mathcal{L} and the data.

This work fills this gap for a wide range of settings. Our main result can be summarized as follows: Under a set of regularity conditions on the data and the functional \mathcal{L} points of the sequence $(\tilde{\theta}_k)$ behave for large iteration number $k \in \mathbb{N}$ like the pME. To be more precise we show in Theorem 2.2 that when the initial guess $\tilde{\theta}_0 \in \mathcal{Y}$ is good enough, then the step estimator sequence $(\tilde{\theta}_k)$ satisfies with high probability

$$\begin{aligned} \|\check{D}(\tilde{\theta}_k - \theta^*) - \check{\xi}\|^2 &\leq \epsilon(p^* + \rho^k R_0), \\ \left| \max_{\eta} \mathcal{L}(\tilde{\theta}_k, \eta) - \max_{\eta} \mathcal{L}(\theta^*, \eta) - \|\check{\xi}\|^2/2 \right| &\leq (p + \mathbf{x})^{1/2} \epsilon(p^* + \rho^k R_0), \end{aligned}$$

where $\rho < 1$ and $\epsilon > 0$ is some small number, for example $\epsilon = Cp^*/\sqrt{n}$ in the smooth i.i.d setting. Further $R_0 > 0$ is a bound related to the quality of the initial guess. The random variable $\xi \in \mathbb{R}^p$ and the matrix $\check{D} \in \mathbb{R}^{p \times p}$ are related to the efficient influence function in semiparametric models and its covariance. These are up to $\rho^k R_0$ the same properties as those proven for the pME in [2] under nearly the same set of conditions. Further in our second main result we manage to show under slightly stronger smoothness conditions that $(\tilde{\theta}_k, \tilde{\eta}_k)$ approaches the ME \tilde{v} with nearly linear convergence speed, i.e. $\|\mathcal{D}((\theta_k, \eta_k) - \tilde{v})\| \leq \tau^{k/\log(k)}$ with some $0 < \tau < 1$ and $\mathcal{D}^2 = \mathbb{E}\nabla^2 \mathcal{L}(\mathbf{v}^*)$ (see Theorem 2.4).

In the following we write $\tilde{v}_{k,k(+1)}$ in statements that are true for both $\tilde{v}_{k,k+1}$ and $\tilde{v}_{k,k}$. Also we do not specify whether the elements of the resulting sequence are sets or single points. All statements made about properties of $\tilde{v}_{k,k(+1)}$ are to be understood in the sense that they hold for “every point of $\tilde{v}_{k,k(+1)}$ ”.

1.1 Idea of the proof

To motivate the approach first consider the toy model

$$\mathbb{Y} = \mathbf{v}^* + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, \mathbb{F}_{\mathbf{v}^*}^{-2}), \quad \mathbb{F}_{\mathbf{v}^*}^2 =: \begin{pmatrix} \mathbb{F}_{\theta^*}^2 & A \\ A^\top & \mathbb{F}_{\eta^*}^2 \end{pmatrix}.$$

In this case we set \mathcal{L} to be the true log likelihood of the observations

$$\mathcal{L}(\mathbf{v}, \mathbb{Y}) = -\|\mathbb{F}(\mathbf{v}^* - \mathbb{Y})\|^2/2.$$

With any starting initial guess $\tilde{v}_0 \in \mathbb{R}^{p+m}$ we obtain from (1.2) for $k \in \mathbb{N}$ and the usual first order criterion of maximality the following two equations

$$\begin{aligned} \mathbb{F}_{\theta^*}(\tilde{\theta}_k - \theta^*) &= I_{\theta^*} \epsilon_\theta + \mathbb{F}_{\theta^*}^{-1} A (\tilde{\eta}_k - \eta^*), \\ \mathbb{F}_{\eta^*}(\tilde{\eta}_{k+1} - \eta^*) &= I_{\eta^*} \epsilon_\eta + \mathbb{F}_{\eta^*}^{-1} A^\top (\tilde{\theta}_k - \theta^*). \end{aligned}$$

Combining these two equations we derive, assuming $\|\mathbb{F}_{\theta^*}^{-1} A \mathbb{F}_{\eta^*}^{-2} A^\top I_{\theta^*}^{-1}\| =: \|\mathbf{M}_0\| = \nu < 1$

$$\begin{aligned} \mathbb{F}_{\theta^*}(\tilde{\theta}_k - \theta^*) &= \mathbb{F}_{\theta^*}^{-1}(\mathbb{F}_{\theta^*}^2 \epsilon_\theta - A \epsilon_\eta) + \mathbb{F}_{\theta^*}^{-1} A \mathbb{F}_{\eta^*}^{-1} A^\top \mathbb{F}_{\theta^*}^{-1} \mathbb{F}_{\theta^*}(\tilde{\theta}_{k-1} - \theta^*) \\ &= \sum_{l=1}^k \mathbf{M}_0^{k-l} \mathbb{F}_{\theta^*}^{-1}(\mathbb{F}_{\theta^*}^2 \epsilon_\theta - A \epsilon_\eta) \\ &\quad + \mathbf{M}_0^k \mathbb{F}_{\theta^*}(\tilde{\theta}_0 - \theta^*) \rightarrow \mathbb{F}_{\theta^*}(\hat{\theta} - \theta^*). \end{aligned}$$

Because the limit $\hat{\theta}$ is independent of the initial point \tilde{v}_0 and because the profile $\tilde{\theta}$ is a fix point of the procedure the unique limit satisfies $\hat{\theta} = \tilde{\theta}$. This argument is based on the fact that in this setting the functional is quadratic such that the gradient satisfies

$$\nabla \mathcal{L}(\mathbf{v}) = \mathbb{F}_{\mathbf{v}^*}^2(\mathbf{v} - \mathbf{v}^*) + \mathbb{F}_{\mathbf{v}^*}^2 \epsilon.$$

Any smooth function is quadratic around its maximizer which motivates a local linear approximation of the gradient of the functional \mathcal{L} to derive our results with similar arguments. This is done in the proof of Theorem 2.2.

First it is ensured that the whole sequence $(\tilde{\mathbf{v}}_{k,k+1})_{k \in \mathbb{N}_0}$ satisfies for some $R_0 > 0$

$$\{\tilde{\mathbf{v}}_{k,k+1}, k \in \mathbb{N}_0\} \subset \{\|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \leq R_0\}, \quad (1.5)$$

where $\mathcal{D}^2 \stackrel{\text{def}}{=} \nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}^*)$ (see Theorem 4.3). In the second step we approximate with $\zeta = \mathcal{L} - \mathbb{E} \mathcal{L}$

$$\mathcal{L}(\mathbf{v}, \mathbf{v}^*) = \nabla \zeta(\mathbf{v}^*)(\mathbf{v} - \mathbf{v}^*) - \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|^2/2 + \alpha(\mathbf{v}, \mathbf{v}^*), \quad (1.6)$$

where $\alpha(\mathbf{v}, \mathbf{v}^*)$ is defined by (1.6). Similar to the toy case above this allows using the first order criterion of maximality and (1.5) to obtain a bound of the kind

$$\begin{aligned} \|\mathcal{D}(\mathbf{v}_{k,k} - \mathbf{v}^*)\| &\leq \mathfrak{c} \sum_{l=0}^k \rho^l (\|\mathcal{D}^{-1} \nabla \zeta(\mathbf{v}^*)\| + |\alpha(\mathbf{v}_{l,l}, \mathbf{v}^*)|) \\ &\leq \mathfrak{c}_1 (\|\mathcal{D}^{-1} \nabla \zeta(\mathbf{v}^*)\| + \epsilon(R_0)) + \rho^k R_0 \stackrel{\text{def}}{=} \mathfrak{r}_k. \end{aligned}$$

This is done in Lemma 4.5 using results from [2] to show that $\epsilon(R_0)$ is small. Finally the same arguments as in [2] allow to obtain our main result using that with high probability for all $k \in \mathbb{N}_0$ $\tilde{\mathbf{v}}_{k,k} \in \{\|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \leq \mathfrak{r}_k\}$. For the convergence result similar arguments are used. The only difference is that instead of (1.6) we use the approximation

$$\mathcal{L}(\mathbf{v}, \tilde{\mathbf{v}}) = -\|\mathcal{D}(\mathbf{v} - \tilde{\mathbf{v}})\|^2/2 + \alpha'(\mathbf{v}, \tilde{\mathbf{v}}),$$

exploiting that $\nabla \mathcal{L}(\tilde{\mathbf{v}}) \equiv 0$, which allows to obtain actual convergence to the ME.

It is worthy to point out two technical challenges of the analysis. First the sketched approach relies on (1.5). As all estimators $(\tilde{\mathbf{v}}_{k,k+1})$ are random this means that we need with some small $\beta > 0$

$$\mathbb{P} \left(\bigcap_{k \in \mathbb{N}_0} \left\{ \tilde{\mathbf{v}}_{k,k}, \tilde{\mathbf{v}}_{k,k+1} \in \{\|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \leq R_0\} \right\} \right) \geq 1 - \beta.$$

This is not trivial but the result of Theorem 4.3 serves the result thanks to $\mathcal{L}(\tilde{\mathbf{v}}_{k,k+1}) \geq \mathcal{L}(\tilde{\mathbf{v}}_0)$. Second the main result 2.2 is formulated to hold for all $k \in \mathbb{N}_0$. This implies the need of a bound of the kind

$$\mathbb{P} \left(\bigcap_{k \in \mathbb{N}_0} \left\{ \left\| \check{D}^{-1} \{ \check{\nabla} \zeta(\tilde{\mathbf{v}}_{k,k}) - \check{\nabla} \zeta(\mathbf{v}^*) \} \right\| \leq \epsilon(\mathfrak{r}_k) \right\} \right) \geq 1 - \beta,$$

with some small $\epsilon(\mathfrak{r}) > 0$ that is decreasing if $\mathfrak{r} > 0$ shrinks. Again this is not trivial and not a direct implication of the results of [2] or [14]. We manage to derive this result in the desired way in Theorem 8.2, which is an adapted version of Theorem D.1 of [2] based on Corollary 2.5 of [14].

2 Main results

2.1 Conditions

This section collects the conditions imposed on the model. We use the same set of assumptions as in [2] and this section closely follows Section 2.1 of that paper.

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 (1.1). The matrix \mathcal{D}^2 can be defined as follows:

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

Here and in what follows we implicitly assume that the log-functional 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 we can interchange $\nabla \mathbb{E}\mathcal{L} = \mathbb{E}\nabla \mathcal{L}$ on $\mathcal{Y}_\circ(\mathbf{r}_0)$, where $\mathcal{Y}_\circ(\mathbf{r})$ is defined in (2.1) and $\mathbf{r}_0 > 0$ in (2.4). It is worth mentioning that $\mathcal{D}^2 = \mathcal{V}^2 \stackrel{\text{def}}{=} \text{Cov}(\nabla \mathcal{L}(\mathbf{v}^*))$ if the model $\mathbf{Y} \sim \mathbb{P}_{\mathbf{v}^*} \in (\mathbb{P}_{\mathbf{v}})$ is correctly specified and sufficiently regular; see e.g. [7].

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 $\rho < 1$

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

Remark 2.1. The condition (\mathcal{I}) allows to introduce 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

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

and (\mathcal{I}) ensures that the matrix \check{D}^2 is well posed.

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

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

The following two conditions quantify the smoothness properties on $\mathcal{Y}_\circ(\mathbf{r})$ of the expected log-functional $\mathbb{E}\mathcal{L}(\mathbf{v})$ and of the stochastic component $\zeta(\mathbf{v}) = \mathcal{L}(\mathbf{v}) - \mathbb{E}\mathcal{L}(\mathbf{v})$.

($\check{\mathcal{L}}$) For each $r \leq r_0$, there is a constant $\delta(r)$ such that it holds on the set $\mathcal{Y}_o(r)$:

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

Remark 2.2. 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 projected gradient $\check{\nabla}_\theta \mathbb{E}\mathcal{L}(\mathbf{v})$ which is defined as

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

Under condition ($\check{\mathcal{L}}_0$) it follows from the second order Taylor expansion for any $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_o(r)$ (see Lemma B.1 of [2])

$$\|\check{D}^{-1}(\check{\nabla} \mathbb{E}\mathcal{L}(\mathbf{v}) - \check{\nabla} \mathbb{E}\mathcal{L}(\mathbf{v}^*)) - \check{D}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \delta(r)r. \quad (2.2)$$

In the proofs we actually only need the condition (2.2) which in some cases can be weaker than ($\check{\mathcal{L}}_0$).

The next condition concerns the regularity of the stochastic component $\zeta(\mathbf{v}) \stackrel{\text{def}}{=} \mathcal{L}(\mathbf{v}) - \mathbb{E}\mathcal{L}(\mathbf{v})$. Similarly to [14], we implicitly assume that the stochastic component $\zeta(\mathbf{v})$ is a separable stochastic process.

($\check{\mathcal{D}}_1$) For all $0 < r < r_0$, there exists a constant $\omega \leq 1/2$ such that for all $|\mu| \leq \check{g}$ and $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_o(r)$

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

The above conditions allow to derive the main result once the accuracy of the sequence is established. We include another condition that allows to control the deviation behavior of $\|\check{D}^{-1}\check{\nabla}\zeta(\mathbf{v}^*)\|$. To present this condition define the covariance matrix $\mathcal{V}^2 \in \mathbb{R}^{p^* \times p^*}$ and $\check{V}^2 \in \mathbb{R}^{p \times p}$

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

$\mathcal{V}^2 \in \mathbb{R}^{p^* \times p^*}$ describes the variability of the process $\mathcal{L}(\mathbf{v})$ around the central point \mathbf{v}^o .

($\check{\mathcal{D}}_0$) There exist constants $\nu_0 > 0$ and $\check{g} > 0$ such that for all $|\mu| \leq \check{g}$

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

So far we only presented conditions that allow to treat the properties of $\tilde{\boldsymbol{\theta}}_k$ on local sets $\mathcal{Y}_o(r_k)$. To show that r_k is not too large the following, stronger conditions are employed:

(\mathcal{L}_0) For each $r \leq r_0$, there is a constant $\delta(r)$ such that it holds on the set $\mathcal{Y}_o(r)$:

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

($\mathcal{E}\mathcal{D}_1$) There exists a constant $\omega \leq 1/2$, such that for all $|\mu| \leq g$ and all $0 < r < r_0$

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

($\mathcal{E}\mathcal{D}_0$) There exist constants $\nu_0 > 0$ and $g > 0$ such that for all $|\mu| \leq g$

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

It is important to note, that the constants $\check{\omega}, \check{\delta}(r), \check{\nu}$ and $\omega, \delta(r), \nu$ in the respective weak and strong version can differ substantially and may depend on the full dimension $p^* \in \mathbb{N}$ in less or more severe ways ($AH^{-2}\nabla_\eta \mathcal{L}$ might be quite smooth while $\nabla_\eta \mathcal{L}$ could be less regular). This is why we use both sets of conditions where they suit best, although the list of assumptions becomes rather long. If a short list is preferred the following lemma shows, that the stronger conditions imply the weaker ones from above:

Lemma 2.1. *[[2], Lemma 2.1] Assume (\mathcal{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*

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

Finally we present two conditions that allow to ensure that with a high probability the sequence $(\mathbf{v}_{k, k(+1)})$ stays close to \mathbf{v}^* if the initial guess $\tilde{\mathbf{v}}_0$ lands close to \mathbf{v}^* . These conditions have to be satisfied on the whole set $\mathcal{Y} \subseteq \mathbb{R}^{p^*}$.

($\mathcal{L}r$) For any $r > r_0$ there exists a value $b(r) > 0$, such that

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

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

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

We impose one further merely technical condition:

(\mathbf{B}_1) We assume for all $r \geq \frac{6\nu_0}{b}\sqrt{\mathbf{x} + 4p^*}$

$$1 + \sqrt{\mathbf{x} + 4p^*} \leq \frac{3\nu_r^2}{b}g(r).$$

Remark 2.3. Without this the calculation of $R_0(\mathbf{x})$ in Section 4.1 would become technically more involved, without that further insight would be gained.

Remark 2.4. For a discussion on how restrictive these conditions are we refer the reader to Remark 2.8 and 2.9 of [2].

2.2 Introduction of important objects

In this section we introduce all objects and bounds that are relevant for Theorem 2.2. This section is quite technical but necessary to understand the results.

First consider the $p^* \times p^*$ matrices \mathcal{D}^2 and \mathcal{V}^2 from Section 2.1, which could be defined similarly to the Fisher information matrix:

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

We represent the information and covariance matrix in block form:

$$\mathcal{D}^2 = \begin{pmatrix} D^2 & A \\ A^\top & H^2 \end{pmatrix}, \quad \mathcal{V}^2 = \begin{pmatrix} V^2 & E \\ E^\top & Q^2 \end{pmatrix}.$$

A crucial object is the constant $0 \leq \rho$ defined by

$$\|D^{-1}AH^{-1}\|^2 \stackrel{\text{def}}{=} \rho,$$

which we assume to be smaller 1 ($\|\cdot\|$ here and everywhere denotes the spectral norm when its argument is a matrix). It determines the speed of convergence of the alternating procedure (see Theorem 2.2). Define also the local sets

$$\begin{aligned} \mathcal{Y}_o(\mathbf{r}) &\stackrel{\text{def}}{=} \{\mathbf{v} : (\mathbf{v} - \mathbf{v}^*)^\top \mathcal{D}^2 (\mathbf{v} - \mathbf{v}^*) \leq \mathbf{r}^2\}, \\ \tilde{\mathcal{Y}}_o(\mathbf{r}) &\stackrel{\text{def}}{=} \{\mathbf{v} : (\mathbf{v} - \tilde{\mathbf{v}})^\top \mathcal{D}^2 (\mathbf{v} - \tilde{\mathbf{v}}) \leq \mathbf{r}^2\}, \end{aligned}$$

and the radius $r_0 > 0$ via

$$r_0(\mathbf{x}) \stackrel{\text{def}}{=} \inf_{r \geq 0} \left\{ \mathbb{P} \left(\underset{\substack{\mathbf{v} \in \mathcal{Y} \\ \Pi_\theta \mathbf{v} = \theta^*}}{\text{argmax}} \mathcal{L}(\mathbf{v}), \tilde{\mathbf{v}} \in \mathcal{Y}_o(\mathbf{r}) \right) \geq 1 - e^{-\mathbf{x}} \right\}. \quad (2.3)$$

Remark 2.5. This radius can be determined using conditions (\mathcal{L}_r) and (\mathcal{E}_r) of Section 2.1 and Theorem 4.3 which would yield $r_0(\mathbf{x}) = C\sqrt{\mathbf{x} + p^*}$.

Further introduce the $p \times p$ matrix \check{D} and the p -vectors $\check{\nabla}_\theta$ and $\check{\xi}$ as

$$\check{D}^2 = D^2 - AH^{-2}A^\top, \quad \check{\nabla}_\theta = \nabla_\theta - AH^{-2}\nabla_\eta, \quad \check{\xi} = \check{D}^{-1}\check{\nabla}_\theta,$$

and the matrices

$$\mathbb{B}^2 \stackrel{\text{def}}{=} \mathcal{D}^{-1}\mathcal{V}^2\mathcal{D}^{-1}, \quad \mathbb{B}_\theta \stackrel{\text{def}}{=} D^{-1}V^2D^{-1}, \quad \mathbb{B}_\eta \stackrel{\text{def}}{=} H^{-1}Q^2H^{-1}.$$

Remark 2.6. The random variable $\check{\xi} \in \mathbb{R}^p$ is related to the efficient influence function in semiparametric models. If the model is regular and correctly specified \check{D}^2 is the covariance of the efficient influence function and its inverse the semiparametric Cramer-Rao lower bound for regular estimators. The matrices $\check{B}, \check{B}_\theta, \check{B}_\eta$ describe the miss specification of the model and are related to the White-statistic.

For our estimations we need the constant

$$\mathfrak{z}(\mathbf{x}) \stackrel{\text{def}}{=} \mathfrak{z}(\mathbf{x}, \check{B}) \vee \mathfrak{z}_Q(\mathbf{x}, 4p^*) \approx \sqrt{p^* + \mathbf{x}},$$

where $\mathfrak{z}(\mathbf{x}, \cdot)$ is explained in Section 7 and $\mathfrak{z}_Q(\mathbf{x}, \cdot)$ is defined in Equation (8.2).

Remark 2.7. The constant $\mathfrak{z}(\mathbf{x})$ is only introduced for ease of notation. This makes some bounds less sharp but allows to address all terms that are of order $\sqrt{p^* + \mathbf{x}}$ with one symbol. The constant $\mathfrak{z}(\mathbf{x}, \check{B})$ is comparable to the " $1 - e^{-\mathbf{x}}$ " quantile of the norm of $\mathcal{D}^{-1}\mathcal{V}\mathbb{X}$, where $\mathbb{X} \sim \mathcal{N}(0, Id_{p^*})$, i.e. it is of order of the trace of \check{B} . The constant $\mathfrak{z}_Q(\mathbf{x}, \mathbb{Q})$ arises as an exponential deviation bound for the supremum of a smooth process over a set with complexity described by \mathbb{Q} .

To bound the deviations of the points of the sequence $(\tilde{\mathbf{v}}_{k,k(+1)})$ we need the following radius:

$$R_0(\mathbf{x}, K_0) \stackrel{\text{def}}{=} \mathfrak{z}(\mathbf{x}) \vee \frac{6\nu_0}{b(1-\rho)} \sqrt{\mathbf{x} + 2.4p^* + \frac{b^2}{9\nu_0^2} K_0(\mathbf{x})}, \quad (2.4)$$

which will ensure $\{\tilde{\mathbf{v}}_0, \tilde{\mathbf{v}}_{0,1}, \dots\} \subset \mathcal{Y}_\circ(R_0)$, where $K_0(\mathbf{x}) > 0$ is defined as

$$K_0(\mathbf{x}) \stackrel{\text{def}}{=} \inf_{K>0} \{ \mathbb{P}(\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -K) \geq \beta(\mathbf{x}) \},$$

for some $\beta(\mathbf{x}) \rightarrow 0$ as $\mathbf{x} \rightarrow \infty$, see condition (A_1) in 2.3. Finally define the *parametric uniform spread* and the *semiparametric uniform spread*

$$\begin{aligned} \diamond_Q(\mathbf{r}, \mathbf{x}) &\stackrel{\text{def}}{=} \{ \delta(\mathbf{r})\mathbf{r} + 6\nu_1\omega(\mathfrak{z}_Q(\mathbf{x}, 4p^*) + 2\mathbf{r}^2) \}, \\ \check{\diamond}_Q(\mathbf{r}, \mathbf{x}) &\stackrel{\text{def}}{=} \frac{8}{(1-\rho^2)^2} \check{\delta}(\mathbf{r})\mathbf{r} + 6\nu_1\check{\omega}(\mathfrak{z}_Q(\mathbf{x}, 2p^* + 2p)^2 + 2\mathbf{r}^2). \end{aligned} \quad (2.5)$$

Remark 2.8. This object is central to our analysis as it describes the accuracy of our main result of Theorem 2.2. It is small for not too large \mathbf{r} , if $\check{\omega}, \check{\delta}$ from conditions $(\check{\mathcal{E}}\mathcal{D}_1)$, $(\check{\mathcal{L}}_0)$ from Section 2.1 are small (with Lemma 2.1 it suffices that ω, δ from $(\mathcal{E}\mathcal{D}_1)$, (\mathcal{L}_0) are small). $\check{\diamond}_Q(\mathbf{r}, \mathbf{x})$ is structurally slightly different from $\check{\diamond}(\mathbf{r}, \mathbf{x})$ in [2] as it is based on Theorem 8.2 and allows a uniform in k formulation of our main result Theorem 2.2, but for moderate $\mathbf{x} \in \mathbb{R}_+$ they are of similar size.

2.3 Dependence on initial guess

Our main theorem is only valid under the conditions from Section 2.1 and under some constraints on the quality of the initial guess $\tilde{\mathbf{v}}_0 \in \mathbb{R}^{p^*}$ which we denote by (A_1) , (A_2) and (A_3) :

- (A₁) With probability greater $1 - \beta_{(\mathbf{A})}(\mathbf{x})$ the initial guess satisfies $\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -K_0(\mathbf{x})$ for some $K_0(\mathbf{x}) \geq 0$.
- (A₂) The conditions $(\check{\mathcal{E}}\mathcal{D}_1)$, $(\check{\mathcal{L}}_0)$, $(\mathcal{E}\mathcal{D}_1)$ and (\mathcal{L}_0) from Section 2.1 hold for all $\mathbf{r} \leq R_0(\mathbf{x}, K_0)$ where R_0 is defined in (2.4) with $\beta(\mathbf{x}) = \beta_{(\mathbf{A})}(\mathbf{x})$.
- (A₃) There is some $\epsilon > 0$ such that $\delta(\mathbf{r})/\mathbf{r} \vee 12\nu_1\omega \leq \epsilon$ for all $\mathbf{r} \leq R_0$. Further $K_0(\mathbf{x}) \in \mathbb{R}$ and $\epsilon > 0$ are small enough to ensure

$$c(\epsilon, \mathfrak{z}(\mathbf{x})) \stackrel{\text{def}}{=} \epsilon 7\mathcal{C}(\rho) \frac{1}{1-\rho} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2) < 1, \quad (2.6)$$

$$c(\epsilon, R_0) \stackrel{\text{def}}{=} \epsilon 7\mathcal{C}(\rho) \frac{1}{1-\rho} R_0 < 1, \quad (2.7)$$

with

$$\mathcal{C}(\rho) \stackrel{\text{def}}{=} 2\sqrt{2}(1 + \sqrt{\rho})(1 - \sqrt{\rho})^{-1}. \quad (2.8)$$

Remark 2.9. One way of obtaining condition (A₁) is to show that $\tilde{\mathbf{v}} \in \mathcal{Y}_\circ(R_K)$ with probability greater $1 - \beta_{(\mathbf{A})}(\mathbf{x})$ for some finite $R_K(\mathbf{x}) \in \mathbb{R}$ and $0 \leq \beta_{(\mathbf{A})}(\mathbf{x}) < 1$. Then (see Section 4.1)

$$K_0(\mathbf{x}) \stackrel{\text{def}}{=} (1/2 + 12\nu_0\omega)R_K^2 + (\delta(R_K) + \mathfrak{z}(\mathbf{x}))R_K + 6\nu_0\omega\mathfrak{z}(\mathbf{x})^2.$$

Condition (A₁) is specified by conditions (A₂) and (A₃) and is fundamental, as it allows with dominating probability to concentrate the analysis on a local set $\mathcal{Y}_\circ(R_0(\mathbf{x}))$ (see Theorem 4.3). Conditions (A₂) and (A₃) impose a bound on $R_0(\mathbf{x})$ and thus on K_0 from (A₁). These conditions boil down to $\delta(R_0) + \omega R_0$ being significantly smaller than 1. Condition (A₃) ensures that the quality of the main result from [2] can be attained, i.e. that $\check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}) \approx \check{\diamond}(\mathbf{r}_0, \mathbf{x})$ under rather mild conditions on the size R_0 , as we only need ϵR_0 to be small. A violation of (A₂) would make it impossible to apply Theorem 8.1 the backbone of our proofs.

Remark 2.10. In the case of iid observations with sample size n one often has $\delta(R_0) + \omega R_0 \leq CR_0(\mathbf{x})/\sqrt{n}$ which suggests at first glance that (A₂) and (A₃) are only a question of the sample size. But note that in case of iid observations the functional satisfies $n \approx -\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*)$ such that the conditions (A₂) and (A₃) are not satisfied automatically with sufficiently large sample size. They are true conditions on the quality of the first guess.

2.4 Statistical properties of the alternating sequence

In this Section we present our main theorem in full rigor, i.e. that the limit of the alternating sequence satisfies a finite sample Wilks Theorem and Fisher expansion.

Theorem 2.2. *Assume that the conditions $(\mathcal{E}\mathcal{D}_0)$, $(\mathcal{E}\mathcal{D}_1)$, (\mathcal{L}_0) , (\mathcal{L}_r) and $(\mathcal{E}r)$ of Section 2.1 are met with a constant $\mathbf{b}(\mathbf{r}) \equiv \mathbf{b}$ and where $\mathcal{V}^2 = \text{Cov}(\nabla\mathcal{L}(\mathbf{v}^*))$, $\mathcal{D}^2 =$*

$-\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}^*)$ and where $\mathbf{v}^\circ = \mathbf{v}^*$. Assume that $(\check{\mathcal{D}}_1)$ and $(\check{\mathcal{L}}_0)$ are met. Further assume (B_1) and that the initial guess satisfies (A_1) and (A_2) of Section 2.3. Then it holds with probability greater $1 - 8e^{-x} - \beta_{(\mathbf{A})}$ for all $k \in \mathbb{N}$

$$\|\check{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\| \leq \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}), \quad (2.9)$$

$$\begin{aligned} |2\check{L}(\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}^*) - \|\check{\boldsymbol{\xi}}\|^2| &\leq 8 \left(\|\check{\boldsymbol{\xi}}\| + \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}) \right) \check{\diamond}_Q(2(1 + \rho)\mathbf{r}_k, \mathbf{x}) \\ &\quad + \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x})^2, \end{aligned} \quad (2.10)$$

where

$$\mathbf{r}_k \leq 2\sqrt{2}(1 - \sqrt{\rho})^{-1} \{(\mathfrak{z}(\mathbf{x}) + \diamond_Q(\mathbf{R}_0, \mathbf{x})) + (1 + \sqrt{\rho})\rho^k \mathbf{R}_0(\mathbf{x})\}.$$

If further condition (A_3) is satisfied then (2.9) and (2.10) are met with

$$\begin{aligned} \mathbf{r}_k &\leq \mathbf{C}(\rho) (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2) + \epsilon \frac{7^2 \mathbf{C}(\rho)^4}{1 - c(\epsilon, \mathfrak{z}(\mathbf{x}))} \left(\frac{1}{1 - \rho} \right) (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^2 \\ &\quad + \rho^k \left(\mathbf{C}(\rho) \mathbf{R}_0 + \epsilon \frac{7^2 \mathbf{C}(\rho)^4}{1 - c(\epsilon, \mathbf{R}_0)} \left(\frac{1}{\rho^{-1} - 1} \right) \mathbf{R}_0^2 \right). \end{aligned}$$

In particular this means that if

$$k \geq \frac{2 \log(\mathfrak{z}(\mathbf{x})) - \log\{2\mathbf{R}_0(\mathbf{x}, \mathbf{K}_0)\}}{\log(\rho)},$$

we have with $\mathfrak{z}(\mathbf{x})^2 \leq \mathbf{C}_3(p^* + \mathbf{x})$

$$\check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}) \approx \check{\diamond}_Q(\mathbf{C}\sqrt{p^* + \mathbf{x}}, \mathbf{x}).$$

Remark 2.11. Note that the results are very similar to those in [2] for the profile M estimator $\tilde{\boldsymbol{\theta}}$. This is evident after noting that (ignoring terms of the order $\epsilon \mathfrak{z}(\mathbf{x})$)

$$\mathbf{r}_k \lesssim \mathbf{C}(\rho) (\mathfrak{z}(\mathbf{x}) + \rho^k (\mathbf{R}_0 + \mathbf{C}\epsilon \mathbf{R}_0^2)),$$

which for large $k \in \mathbb{N}$ means $\mathbf{r}_k \lesssim \mathbf{C}(\rho) \mathfrak{z}(\mathbf{x})$.

Remark 2.12. Concerning the properties of $\check{\boldsymbol{\xi}} \in \mathbb{R}^p$ we repeat remark 2.1 of [2]. In the case of the correct model specification the deviation properties of the quadratic form $\|\check{\boldsymbol{\xi}}\|^2 = \|\check{D}^{-1} \check{\nabla}_\theta\|^2$ are essentially the same as of a chi-square random variable with p degrees of freedom; see Theorem 7.1 in the appendix. In the case of a possible model misspecification with, the behavior of the quadratic form $\|\check{\boldsymbol{\xi}}\|^2$ will depend on the characteristics of the matrix $\check{B} \stackrel{\text{def}}{=} \check{D}^{-1} \text{Cov}(\check{\nabla} \mathcal{L}(\mathbf{v}^*)) \check{D}^{-1}$; see again Theorem 7.1. Moreover, in the asymptotic setup the vector $\check{\boldsymbol{\xi}}$ is asymptotically standard normal; see Section 2.2. of [2] for the i.i.d. case.

Remark 2.13. These results allow to derive some important corollaries like concentration and confidence sets (see [14], Section 3.2).

Remark 2.14. In general an exact numerical computation of

$$\theta(\boldsymbol{\eta}) \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \mathbb{R}^p} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}), \text{ or } \eta(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\eta} \in \mathbb{R}^m} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}),$$

is not possible. Define $\widehat{\theta}(\boldsymbol{\eta})$ and $\widehat{\eta}(\boldsymbol{\theta})$ as the numerical approximations to $\theta(\boldsymbol{\eta})$ and $\eta(\boldsymbol{\theta})$ and assume that

$$\begin{aligned} \|D(\widehat{\theta}(\boldsymbol{\eta}) - \theta(\boldsymbol{\eta}))\| &\leq \tau, \text{ for all } \boldsymbol{\eta} \in \mathcal{Y}_{\circ, \boldsymbol{\eta}}(\mathbf{R}_0) \stackrel{\text{def}}{=} \{\boldsymbol{v} \in \mathcal{Y}_{\circ}(\mathbf{R}_0), \Pi_{\boldsymbol{\eta}} \boldsymbol{v} = \boldsymbol{\eta}\}, \\ \|H(\widehat{\eta}(\boldsymbol{\theta}) - \eta(\boldsymbol{\theta}))\| &\leq \tau, \text{ for all } \boldsymbol{\theta} \in \mathcal{Y}_{\circ, \boldsymbol{\theta}}(\mathbf{R}_0) \stackrel{\text{def}}{=} \{\boldsymbol{v} \in \mathcal{Y}_{\circ}(\mathbf{R}_0), \Pi_{\boldsymbol{\theta}} \boldsymbol{v} = \boldsymbol{\theta}\}. \end{aligned}$$

Then we can easily modify the proof of Theorem 2.2 via adding $C(\rho)\tau$ to the error terms and the radii \mathbf{r}_k , where $C(\rho)$ is some rational function of ρ .

Remark 2.15. Note that under condition (A_3) the size of \mathbf{r}_k for $k \rightarrow \infty$ does not depend on $\mathbf{R}_0 > 0$. So as long as $\epsilon \mathbf{R}_0$ is small enough the quality of the initial guess no longer affects the statistical properties of the sequence $(\boldsymbol{\theta}_k)$ for large $k \in \mathbb{N}$.

2.5 Convergence to the ME

Even though Theorem 2.2 tells us, that the statistical properties of the alternating sequence resemble those of its target, the profile ME, it is an interesting question if the underlying approach allows to qualify conditions under which the sequence actually attains the maximizer $\widetilde{\boldsymbol{v}}$. Without further assumptions Theorem 2.2 yields the following Corollary:

Corollary 2.3. *Under the assumptions of Theorem 2.2 it holds with probability greater $1 - 8e^{-x} - \beta_{(\mathbf{A})}$*

$$\|\check{D}(\widetilde{\boldsymbol{\theta}} - \widetilde{\boldsymbol{\theta}}_k)\| \leq \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}) + \check{\diamond}(\mathbf{r}_0, \mathbf{x}),$$

where $\mathbf{r}_0 > 0$ is defined in (2.3) and

$$\check{\diamond}(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \frac{8}{(1 - \rho^2)^2} \check{\delta}(\mathbf{r}) \mathbf{r} + 6\nu_1 \check{\omega} \mathfrak{z}_1(\mathbf{x}, 2p^* + 2p) \mathbf{r}.$$

Remark 2.16. The value $\mathfrak{z}_1(\mathbf{x}, \cdot)$ is defined in (2.11).

Corollary 2.3 is a first step in the direction of an actual convergence result but the gap $\check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}) + \check{\diamond}(\mathbf{r}_0, \mathbf{x})$ is not a zero sequence in $k \in \mathbb{N}$. It turns out that it is possible to prove convergence to the ME with the cost of assuming more smoothness of the functional \mathcal{L} and using the right bound for the maximal eigenvalue of the hessian $\nabla^2 \mathcal{L}(\boldsymbol{v}^*)$.

Consider the following condition, that basically quantifies how "well behaved" the second derivative $\nabla^2(\mathcal{L} - \mathbb{E}\mathcal{L})$ is:

($\mathcal{E}\mathcal{D}_2$) There exists a constant $\omega \leq 1/2$, such that for all $|\mu| \leq g$ and all $0 < r < r_0$

$$\sup_{\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_o(\mathbf{r})} \sup_{\|\gamma_1\|=1} \sup_{\|\gamma_2\|=1} \log \mathbb{E} \exp \left\{ \frac{\mu \gamma_1^\top \mathcal{D}^{-1} \{ \nabla^2 \zeta(\mathbf{v}) - \nabla^2 \zeta(\mathbf{v}') \} \gamma_2}{\omega_2 \|\mathcal{D}(\mathbf{v} - \mathbf{v}')\|} \right\} \leq \frac{\nu_2^2 \mu^2}{2}.$$

Define $\mathfrak{z}(\mathbf{x}, \nabla^2 \mathcal{L}(\mathbf{v}^*))$ via

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

and $\varkappa(\mathbf{x}, R_0)$

$$\varkappa(\mathbf{x}, R_0) \stackrel{\text{def}}{=} \frac{2\sqrt{2}(1 + \sqrt{\rho})}{\sqrt{1 - \rho}} [\delta(R_0) + 9\omega_2 \nu_2 \|\mathcal{D}^{-1}\| \mathfrak{z}_1(\mathbf{x}, 6p^*) R_0 + \|\mathcal{D}^{-1}\| \mathfrak{z}(\mathbf{x}, \nabla^2 \mathcal{L}(\mathbf{v}^*))],$$

where $\mathfrak{z}_1(\mathbf{x}, \cdot)$ satisfies (see Theorem 9.2)

$$\mathfrak{z}_1(\mathbf{x}, \mathbb{Q}) = \begin{cases} \sqrt{2(\mathbf{x} + \mathbb{Q})} & \text{if } \sqrt{2(\mathbf{x} + \mathbb{Q})} \leq g_0, \\ g_0^{-1}(\mathbf{x} + \mathbb{Q}) + g_0/2 & \text{otherwise.} \end{cases} \quad (2.11)$$

Remark 2.17. For the case that $\mathcal{L}(\mathbf{v}) = \sum_{i=1}^n \ell_i(\mathbf{v})$ with a sum of independent marginal functionals $\ell_i : \mathcal{Y} \rightarrow \mathbb{R}$ we can use Corollary 3.7 of [15] to obtain

$$\mathfrak{z}(\mathbf{x}, \nabla^2 \mathcal{L}(\mathbf{v}^*)) = \sqrt{2\tau\nu_3} \sqrt{\mathbf{x} + p^*},$$

if with a sequence of matrices $(\mathbf{A}_i) \subset \mathbb{R}^{p^* \times p^*}$

$$\log \mathbb{E} \exp \lambda \nabla^2 \ell_i(\mathbf{v}^*) \leq \nu_3^2 \lambda^2 / 2 \mathbf{A}_i, \quad \left\| \sum_{i=1}^n \mathbf{A}_i \right\| \leq \tau.$$

Remark 2.18. In the case of smooth i.i.d models this means that $\varkappa(\mathbf{x}, R_0) \leq \mathbb{C}(R_0 + \mathbf{x} + \log(p^*)) / \sqrt{n} + \mathbb{C}R_0 \sqrt{\mathbf{x} + p^*} / n$. This means that $\varkappa(\mathbf{x}, R_0) = O((\mathbf{x} + R_0 + \log(p^*)) / \sqrt{n})$ if $p^* + \mathbf{x} = o(n)$.

With these definitions we can prove the following Theorem:

Theorem 2.4. Let the conditions $(\mathcal{E}\mathcal{D}_2)$, (\mathcal{L}_0) , (\mathcal{L}_r) and $(\mathcal{E}\mathbf{r})$ be met with a constant $\mathfrak{b}(\mathbf{r}) \equiv \mathfrak{b}$ and where $\mathcal{D}^2 = -\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}^*)$ and $\mathbf{v}^* = \mathbf{v}^\circ$. Further suppose (B_1) and that the initial guess satisfies (A_1) and (A_2) . Assume that $\varkappa(\mathbf{x}, R_0) < (1 - \rho)$. Then

$$\mathbb{P} \left(\bigcap_{k \in \mathbb{N}} \{ \mathbf{v}_{k, k(+1)} \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^*) \} \right) \geq 1 - 3e^{-\mathbf{x}} - \beta(\mathbf{A}),$$

where

$$\mathbf{r}_k^* \leq \begin{cases} \rho^k 2\sqrt{2} \frac{1}{1 - \varkappa(\mathbf{x}, R_0)k} \tilde{R}_0, & \varkappa(\mathbf{x}, R_0)k \leq 1, \\ 2 \frac{1 - \rho}{\varkappa(\mathbf{x}, R_0)} \tau(\mathbf{x})^{k/\log(k)} \tilde{R}_0, & \text{otherwise,} \end{cases} \quad (2.12)$$

with $\widetilde{R}_0 \stackrel{\text{def}}{=} R_0 + r_0$ and

$$\tau(\mathbf{x}) \stackrel{\text{def}}{=} \left(\frac{\varkappa(\mathbf{x}, R_0)}{1 - \rho} \right)^{L(k)} < 1$$

$$L(k) \stackrel{\text{def}}{=} \left\lfloor \frac{\log(1/\rho) - \frac{1}{k} (\log(2\sqrt{2}) - \log(\varkappa(\mathbf{x}, R_0)k - 1))}{\left(1 + \frac{1}{\log(k)} \log(1 - \rho)\right)} \right\rfloor \in \mathbb{N},$$

where $\lfloor x \rfloor \in \mathbb{N}$ denotes the largest natural number smaller than $x > 0$.

Remark 2.19. This means that we obtain nearly linear convergence to the global maximizer $\widetilde{\mathbf{v}}$.

Remark 2.20. As in Remark 2.14 if no exact numerical computation of the stepwise maximizers is possible we can easily modify the proof of Theorem 2.4 via adding $C(\rho)\tau$ to $\varkappa(\mathbf{x}, R_0)$, to address that case.

2.6 Critical dimension

In parallel to [2] we want to address 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 are accurate if the spread function $\check{\diamond}_Q(\mathbf{r}_k, \mathbf{x})$ from (2.5) is small. The critical size of p^* then depends on the exact bounds on $\check{\delta}(\cdot)$ and $\check{\omega}$. In the i.i.d setting $\check{\delta}(\mathbf{r})/\mathbf{r} \asymp \check{\omega} \asymp 1/\sqrt{n}$ such that $\check{\diamond}(\mathbf{r}_k, \mathbf{x}) \asymp p^*/\sqrt{n}$ for large $k \in \mathbb{N}$. In other words, one needs that “ p^{*2}/n is small” to obtain an accurate non asymptotic version of the Wilks phenomenon and the Fisher Theorem for the limit of the alternating sequence. This is not surprising because good performance of the ME itself can only be guaranteed if “ p^{*2}/n is small”, as is shown in [2]. There are examples where the pME only satisfies a Wilks- or Fisher result if “ p^{*2}/n is small”, such that in any of those settings the alternating sequence started in the global maximizer does not admit an accurate Wilks- or Fisher expansion.

Interesting enough the constrain $\varkappa(\mathbf{x}, R_0) < (1 - \rho)$ of Theorem 2.4 for the convergence of the sequence to the global maximizer means that one needs $p^*/n \ll 1$ in the smooth i.i.d. setting if $R_0 \leq C_{R_0} \sqrt{p^* + \bar{x}}$. Further Theorem 2.4 states a lower bound for the speed of convergence that in the smooth i.i.d. setting decreases if p^*/n grows. Unfortunately we were unable to find an example that meets the conditions of Section 2.1 and where no convergence occurs if p^*/n tends to infinity. So whether this dimension effect on the convergence is an artifact of our proofs or indeed a property of the alternating procedure remains an open question.

3 Application to single index model

We illustrate how the results of Theorem 2.2 and Theorem 2.4 can be applied in Single Index modeling. Consider the following model

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

for some $f : \mathbb{R} \rightarrow \mathbb{R}$ and $\boldsymbol{\theta}^* \in S_1^{p,+} \subset \mathbb{R}^p$ and with i.i.d errors $\varepsilon_i \in \mathbb{R}$, $\text{Var}(\varepsilon_i) = \sigma^2$ 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, [4] show the asymptotic efficiency of the general semiparametric maximum-functional estimator for particular examples and in [6] 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 $e_0 \stackrel{\text{def}}{=} \psi : [-s_{\mathbf{X}}, s_{\mathbf{X}}] \mapsto \mathbb{R}$ we set for $k = (2^{j_k} - 1)13 + r_k$ with $j_k \in \mathbb{N}_0$ and $r_k \in \{0, \dots, (2^{j_k} - 1)13 - 1\}$

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

A candidate to estimate $\boldsymbol{\theta}^*$ is the profile ME

$$\tilde{\boldsymbol{\theta}}_m \stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \operatorname{argmax}_{(\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathcal{Y}_m} \mathcal{L}_m(\boldsymbol{\theta}, \boldsymbol{\eta}),$$

where

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

and where $\mathcal{Y}_m \subset S_1^{p,+} \times B_{\mathbf{r}^\circ}^m \subset \mathbb{R}^p \times \mathbb{R}^m$ where $B_{\mathbf{r}^\circ}^m \subset \mathbb{R}^m$ denotes the centered ball of radius $\mathbf{r}^\circ > 0$ for some $\mathbf{r}^\circ > 0$. [8] analyzed a very similar estimator in a more general setting based on a kernel estimation of $\mathbb{E}[y | f(\boldsymbol{\theta}^\top \mathbf{X})]$ instead of using a parametric sieve approximation $\sum_{k=0}^m \boldsymbol{\eta}_k e_k$. He showed \sqrt{n} -consistency and asymptotic normality of the proposed estimator.

In this setting a direct computation of $\tilde{\boldsymbol{v}}$ becomes involved, as the maximization problem is high dimensional and not convex. But as noted in the introduction the maximization with respect to $\boldsymbol{\eta}$ for given $\boldsymbol{\theta}$ is high dimensional but convex and consequently feasible. Further for moderate $p \in \mathbb{N}$ the maximization with respect to $\boldsymbol{\theta}$ for fixed $\boldsymbol{\eta}$ is computationally realistic. So an alternating maximization procedure is applicable. To show that it behaves in a desired way we apply the technique presented above.

For the initial guess $\tilde{\boldsymbol{v}}_0 \in \mathcal{Y}$ one can use a simple grid search. For this generate a uniform grid $G_N \stackrel{\text{def}}{=} (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \subset S_1^+$ and define

$$\tilde{\boldsymbol{v}}_0 \stackrel{\text{def}}{=} \operatorname{argmax}_{\substack{(\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathcal{Y} \\ \boldsymbol{\theta} \in G_N}} \mathcal{L}(\boldsymbol{v}). \quad (3.1)$$

Note that given the grid the above maximizer is easily obtained. Simply calculate

$$\tilde{\boldsymbol{\eta}}_{0,k} \stackrel{\text{def}}{=} \operatorname{argmax} \mathcal{L}(\boldsymbol{\theta}_k, \boldsymbol{\eta}) = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{e} \mathbf{e}^\top (\mathbf{X}_i^\top \boldsymbol{\theta}_k) \right)^{-1} \frac{1}{n} \sum_{i=1}^n y_i \mathbf{e}^\top (\mathbf{X}_i^\top \boldsymbol{\theta}_k) \in \mathbb{R}^m \quad (3.2)$$

where by abuse of notation $\mathbf{e} = (\mathbf{e}_1, \dots, \mathbf{e}_m) \in \mathbb{R}^m$. Now observe that

$$\tilde{\mathbf{v}}_0 = \operatorname{argmax}_{k=1, \dots, N} \mathcal{L}(\boldsymbol{\theta}_k, \tilde{\boldsymbol{\eta}}_{0,k}).$$

Define $\tau \stackrel{\text{def}}{=} \sup_{\boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in G_N} \|\boldsymbol{\theta} - \boldsymbol{\theta}^\circ\|$.

To apply the result presented in Theorem 2.2 and Theorem 2.4 we need a list of assumptions denoted by (\mathcal{A}) . We start with conditions on the regressors $\mathbf{X} \in \mathbb{R}^p$:

(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 $B_{s_{\mathbf{X}}}(0) \subset \mathbb{R}^p$ and Lipschitz continuous on $B_{s_{\mathbf{X}}}(0) \subset \mathbb{R}^p$ with Lipschitz constant $L_{d_{\mathbf{X}}} > 0$. Further we assume that for any $\boldsymbol{\theta} \perp \boldsymbol{\theta}^*$ with $\|\boldsymbol{\theta}\| = 1$ we have $\operatorname{Var}(\mathbf{X}^\top \boldsymbol{\theta} \mid \mathbf{X}^\top \boldsymbol{\theta}^*) > \sigma_{\mathbf{X} \mid \boldsymbol{\theta}^*}^2$ for some constant $\sigma_{\mathbf{X} \mid \boldsymbol{\theta}^*}^2 > 0$ that does not depend on $\mathbb{X}^\top \boldsymbol{\theta}^* \in \mathbb{R}$. Also the density $d_{\mathbb{X}} : \mathbb{R}^p \rightarrow \mathbb{R}$ 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 $0 < c_{d_{\mathbf{X}}} \leq C_{d_{\mathbf{X}}} < \infty$.

(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 $\mathbb{E}[\varepsilon_i] = 0$, $\operatorname{Cov}(\varepsilon_i) = \sigma^2$ and satisfy for all $|\mu| \leq \tilde{g}$ for some $\tilde{g} > 0$ and some $\tilde{\nu}_{\mathbf{r}} > 0$

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

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

Proposition 3.1. *Let $\tau = o(p^{*-3/2})$ and $p^{*5}/n \rightarrow 0$. With initial guess given by Equation (3.1) and for $\mathbf{x} \leq 2\tilde{\nu}^2 \tilde{g}^2 n$ the alternating sequence satisfies (2.9) and (2.10) with probability greater $1 - 9 \exp\{-\mathbf{x}\}$ and where with some constant $C_{\diamond} \in \mathbb{R}$*

$$\check{\diamond}_Q(\mathbf{r}, \mathbf{x}) \leq \frac{C_{\diamond} (p^* + \mathbf{x})^{3/2}}{\sqrt{n}} (\mathbf{r}^2 + p^* + \mathbf{x}).$$

Remark 3.1. The constraint $\tau = o(p^{*-3/2})$ implies that for the calculation of the initial guess the vector $\tilde{\boldsymbol{\eta}}_{0,l}$ of (3.2) and the functional $\mathcal{L}(\cdot)$ have to be evaluated $N = p^{*3(p-1)/2}$ times.

Proposition 3.2. Take the initial guess given by Equation (3.1). Assume (\mathcal{A}) but use a three times continuously differentiable wavelet basis. Further assume that $p^{*4}/n \rightarrow 0$ and $\tau = o(p^{*-3/2})$. Let $\mathbf{x} > 0$ be chosen such that

$$\mathbf{x} \leq \frac{1}{2} (\tilde{\nu}^2 n \tilde{\mathbf{g}}^2 - \log(p^*)) \wedge p^*.$$

Then we get the claim of Theorem 2.4 with $\beta_{(\mathbf{A})} = e^{-\mathbf{x}}$ and

$$\kappa(\mathbf{x}, R_0) = O(\tau m^{3/2} + \sqrt{\tau \mathbf{x}} m^{3/2} / n^{1/4}) + O(p^{*2} / \sqrt{n}) \rightarrow 0,$$

for moderate choice of $\mathbf{x} > 0$.

For details see [1].

4 Proof of Theorem 2.2

In this section we will proof Theorem 2.2. Before we start with the actual proof we want to explain the agenda. The first step of the proof is to find a desirable set $\Omega(\mathbf{x}) \subset \Omega$ of high probability, on which a linear approximation of the gradient of the functional $\mathcal{L}(\mathbf{v})$ can be carried out with sufficient accuracy. Once this set is found all subsequent analysis concerns events in $\Omega(\mathbf{x}) \subset \Omega$.

For this purpose define for some $K \in \mathbb{N}$ the set

$$\Omega(\mathbf{x}) = \bigcap_{k=0}^K (C_{k,k} \cap C_{k,k+1}) \cap C(\nabla) \cap \{\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -K_0(\mathbf{x})\}, \text{ where} \quad (4.1)$$

$$C_{k,k(+1)} = \left\{ \|\mathcal{D}(\tilde{\mathbf{v}}_{k,k(+1)} - \mathbf{v}^*)\| \leq R_0(\mathbf{x}), \|D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*)\| \leq R_0(\mathbf{x}), \right. \\ \left. \|H(\tilde{\boldsymbol{\eta}}_{k(+1)} - \boldsymbol{\eta}^*)\| \leq R_0(\mathbf{x}) \right\},$$

$$C(\nabla) = \bigcap_{\mathbf{r} \leq R_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{I}_o(\mathbf{r})} \left\{ \frac{1}{6\omega\nu_1} \|\mathcal{Y}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 4p^*)^2 \right\} \\ \bigcap_{\mathbf{r} \leq 4R_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{I}_o(\mathbf{r})} \left\{ \frac{1}{6\check{\omega}\check{\nu}_1} \|\check{\mathcal{Y}}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 2p^* + 2p)^2 \right\}$$

$$\cap \left\{ \max\{\|\mathcal{D}^{-1}\nabla\mathcal{L}\|, \|D^{-1}\nabla_{\boldsymbol{\theta}}\mathcal{L}\|, \|H^{-1}\nabla_{\boldsymbol{\eta}}\mathcal{L}\|\} \leq \mathfrak{z}(\mathbf{x}) \right\}$$

$$\cap \{\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \in \mathcal{I}_o(\mathbf{r}_0(\mathbf{x}))\}.$$

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

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

the parametric normalized stochastic gradient gap $\mathcal{Y}(\mathbf{v})$ is defined as

$$\mathcal{Y}(\mathbf{v}) = \mathcal{D}_0^{-1}\left(\nabla\zeta(\mathbf{v}) - \nabla\zeta(\mathbf{v}^*)\right),$$

and $\mathbf{r}_0(\mathbf{x}) > 0$ is chosen such that $\mathbb{P}(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_o(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}$, where

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

Remark 4.1. We intersect the set with the event $\{\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_o(\mathbf{r}_0)\}$ where we a priory demand $\mathbf{r}_0(\mathbf{x}) > 0$ to be chosen such that $\mathbb{P}(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_o(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}$. Note that condition $(\mathcal{E}\mathbf{r})$ together with $(\mathcal{L}\mathbf{r})$ allow to set $\sqrt{p^* + \mathbf{x}} \approx \mathbf{r}_0 \leq \mathbf{R}_0$ (see Theorem 4.3).

In Section 4.1 we show that this set is of probability greater $1 - 8e^{-\mathbf{x}} - \beta_{(\mathbf{A})}$. We want to explain the purpose of this set along the architecture of the proof of our main theorem.

$\{\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -\mathbf{K}_0(\mathbf{x})\}$: This set ensures, that the first guess satisfies $\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -\mathbf{K}_0(\mathbf{x})$, which means that it is close enough to the target $\mathbf{v}^* \in \mathbb{R}^{p^*}$. This fact allows us to obtain an a priori bound for the deviation of the sequence $(\tilde{\mathbf{v}}_{k,k+1}) \subset \mathcal{Y}$ from $\mathbf{v}^* \in \mathcal{Y}_o(\mathbf{R}_0)$ with Theorem 4.3.

$\{\mathcal{D}(\tilde{\mathbf{v}}_{k,k+1} - \mathbf{v}^*) \leq \mathbf{R}_0(\mathbf{x})\}$: As just mentioned this event is of high probability due to $\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -\mathbf{K}_0(\mathbf{x})$ and Theorem 4.3. This allows to concentrate the analysis on the set $\mathcal{Y}_o(\mathbf{R}_0)$ on which Taylor expansions of the functional $\mathcal{L} : \mathbb{R}^{p^*} \rightarrow \mathbb{R}$ become accurate.

$C(\nabla)$: This set ensures that on $\Omega(\mathbf{x}) \subset \Omega$ all occurring random quadratic forms and stochastic errors are controlled by $\mathfrak{z}(\mathbf{x}) \in \mathbb{R}$. Consequently we can derive in the proof of Lemma 4.5 an a priori bound of the form $\|\mathcal{D}(\tilde{\mathbf{v}}_{k,k+1} - \mathbf{v}^*)\| \leq \mathbf{r}_k$ for a decreasing sequence of radii $(\mathbf{r}_k) \subset \mathbb{R}_+$ satisfying $\limsup_{k \rightarrow \infty} \mathbf{r}_k = \mathbf{C}\mathfrak{z}(\mathbf{x})$. Further this set allows to obtain in Lemma 4.7 the bounds for all $k \in \mathbb{N}$.

On $\Omega(\mathbf{x}) \subset \Omega$ we find $\tilde{\mathbf{v}}_{k,k+1} \in \mathcal{Y}_o(\mathbf{r}_k)$ such that we can follow the arguments of Theorem 2.2 of [2] to obtain the desired result with accuracy measured by $\check{\diamond}_Q(\mathbf{r}_k, \mathbf{x})$.

4.1 Probability of desirable set

Here we show that the set $\Omega(\mathbf{x})$ actually is of probability greater $1 - 8e^{-\mathbf{x}} - \beta_{(\mathbf{A})}$. We prove the following two Lemmas, which together yield the claim.

Lemma 4.1. *The set $C(\nabla)$ satisfies*

$$\mathbb{P}(C(\nabla)) \geq 1 - 7e^{-\mathbf{x}}.$$

Proof. The proof is similar to the proof of Theorem 3.1 in [14]. Denote

$$\begin{aligned}\mathcal{A} &\stackrel{\text{def}}{=} \bigcap_{r \leq R_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_\circ(r)} \left\{ \frac{1}{6\omega\nu_1} \|\mathfrak{y}(\mathbf{v})\| - 2r^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 4p^*)^2 \right\} \\ \mathcal{B} &\stackrel{\text{def}}{=} \bigcap_{r \leq 4R_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_\circ(r)} \left\{ \frac{1}{6\check{\omega}\check{\nu}_1} \|\check{\mathfrak{y}}(\mathbf{v})\| - 2r^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 2p^* + 2p)^2 \right\} \\ \mathcal{C} &\stackrel{\text{def}}{=} \left\{ \max\{\|\mathcal{D}^{-1}\nabla\mathcal{L}\|, \|D^{-1}\nabla_\theta\mathcal{L}\|, \|H^{-1}\nabla_\eta\mathcal{L}\|\} \leq \mathfrak{z}(\mathbf{x}) \right\}.\end{aligned}$$

We estimate

$$\begin{aligned}\mathbb{P}(\mathcal{C}(\nabla)) &\geq 1 - \mathbb{P}(\mathcal{A}^c) - \mathbb{P}(\mathcal{B}^c) - \mathbb{P}(\mathcal{C}^c) \\ &\quad - \mathbb{P}(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \notin \mathcal{Y}_\circ(r_0)) - \mathbb{P}\left(\|\check{D}^{-1}\check{\nabla}_\theta\|^2 > \mathfrak{z}(\mathbf{x}, \check{B}_\theta)\right).\end{aligned}$$

We bound using for both terms Theorem 8.2 which is applicable due to $(\mathcal{E}\mathcal{D}_1)$ and $(\check{\mathcal{E}}\mathcal{D}_1)$:

$$\mathbb{P}(\mathcal{A}^c) \leq e^{-\mathbf{x}}, \quad \mathbb{P}(\mathcal{B}^c) \leq e^{-\mathbf{x}}.$$

For the set $\mathcal{C} \subset \Omega$ observe that we can use (\mathcal{I}) and Lemma 4.2 to find

$$\|H^{-1}\nabla_\eta\| \vee \|D^{-1}\nabla_\theta\| \leq \|\mathcal{D}^{-1}\nabla\|.$$

This implies that

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

Using the deviation properties of quadratic forms as sketched in Section 7 we find

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

By the choice of $\mathfrak{z}(\mathbf{x}) > 0$ and $r_0 > 0$ this gives the claim. \square

We cite Lemma B.2 of [2]:

Lemma 4.2. *Let*

$$\begin{aligned}\mathcal{D}^2 &= \begin{pmatrix} D^2 & A \\ A^\top & H^2 \end{pmatrix} \in \mathbb{R}^{(p+p) \times (p+p)}, \quad D \in \mathbb{R}^{p \times p}, \quad H \in \mathbb{R}^{m \times m} \text{ invertible,} \\ &\quad \|D^{-1}AH^{-1}\| < 1.\end{aligned}$$

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

The next step is to show that the set $\bigcap_{k=1}^K (C_{k,k} \cap C_{k,k+1})$ has high probability, that is independent of the number of necessary steps. A close look at the proof of Theorem 4.1 of [14] shows that it actually yields the following modified version:

Theorem 4.3 ([14], Theorem 4.1). *Suppose (\mathcal{E}_r) and (\mathcal{L}_r) with $b(r) \equiv b$. 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 r_0 and any $r \geq r_0$, the following conditions are fulfilled:

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

then

$$\mathbb{P}(\mathcal{Y}(K) \subseteq \mathcal{Y}_o(r_0)) \geq 1 - e^{-x}.$$

Note that with (\mathcal{I})

$$\|D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*)\| \vee \|H(\tilde{\boldsymbol{\eta}}_{k(+1)} - \boldsymbol{\eta}^*)\| \leq \frac{1}{1 - \rho} \|\mathcal{D}(\tilde{\mathbf{v}}_{k,k(+1)} - \mathbf{v}^*)\|.$$

With assumption (B_1) and

$$\mathbf{R}_0(\mathbf{x}) = \frac{6\nu_0}{b(1 - \rho)} \sqrt{x + \mathbb{Q} + \frac{b}{9\nu_0^2} \mathbf{K}_0(\mathbf{x})},$$

this implies the desired result as $\mathcal{L}(\mathbf{v}_{k,k(+1)}, \mathbf{v}^*) \geq \mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*)$ such that with Theorem 4.3

$$\begin{aligned} \mathbb{P} \left(\bigcap_{k=0}^K (C_{k,k} \cap C_{k,k+1}) \right) &\geq \mathbb{P} \left(\bigcap_{k=0}^K (C_{k,k} \cap C_{k,k+1}) \cap \{\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -\mathbf{K}_0\} \right) \\ &\quad - \mathbb{P}(\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \leq -\mathbf{K}_0) \\ &\geq \mathbb{P} \left\{ \mathcal{Y}(\mathbf{K}_0(\mathbf{x})) \subset \mathcal{Y}_o \left((1 - \rho) \mathbf{R}_0(\mathbf{x}) \right) \right\} - \beta_{(\mathbf{A})} \\ &\geq 1 - e^{-x} - \beta_{(\mathbf{A})}. \end{aligned}$$

Remark 4.2. This also shows that the sets of maximizers $(\tilde{\mathbf{v}}_{k,k(+1)})$ are nonempty and well defined since the maximization always takes place on compact sets of the form $\{\boldsymbol{\theta} \in \mathbb{R}^p, (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathcal{Y}_o(\mathbf{R}_0)\}$ or $\{\boldsymbol{\eta} \in \mathbb{R}^m, (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathcal{Y}_o(\mathbf{R}_0)\}$.

To address the claim of remark 2.9 we present the following Lemma:

Lemma 4.4. *On the set $C(\nabla) \cap \{\tilde{\mathbf{v}}_0 \in \mathcal{Y}_o(R_K)\}$ it holds*

$$\mathcal{L}(\mathbf{v}_0, \mathbf{v}^*) \geq -(1/2 + 12\nu_0\omega)R_K^2 - (\delta(R_K) + \mathfrak{z}(\mathbf{x}))R_K - 6\nu_0\omega\mathfrak{z}(\mathbf{x})^2.$$

Proof. With similar arguments as in the proof of Lemma 4.5 we have on $C(\nabla) \subset \Omega$ that

$$\begin{aligned}
\mathcal{L}(\mathbf{v}_0, \mathbf{v}^*) &\geq \mathbb{E}[\mathcal{L}(\mathbf{v}_0, \mathbf{v}^*)] - \|\mathcal{D}^{-1}\nabla\zeta(\mathbf{v}^*)\|R_K - |\{\nabla\zeta(\widehat{\mathbf{v}}) - \nabla\zeta(\mathbf{v}^*)\}(\mathbf{v}_0 - \mathbf{v}^*)| \\
&\geq -\|\mathcal{D}(\mathbf{v}_0 - \mathbf{v}^*)\|^2/2 - \|\mathcal{D}^{-1}\nabla\zeta(\mathbf{v}^*)\|R_K \\
&\quad - \|\mathcal{D}^{-1}\{\nabla\mathcal{L}(\widehat{\mathbf{v}}) - \nabla\mathcal{L}(\mathbf{v}^*)\}\|R_K - R_K\delta(R_K) \\
&\geq -(1/2 + 12\nu_0\omega)R_K^2 - (\delta(R_K) + \mathfrak{z}(\mathbf{x}))R_K - 6\nu_0\omega\mathfrak{z}(\mathbf{x})^2.
\end{aligned}$$

□

4.2 Proof convergence

We derive the a priori bound $\widetilde{\mathbf{v}}_{k,k(+)1} \in \mathcal{Y}_\circ(\mathbf{r}_k)$ with an adequately decreasing sequence $(\mathbf{r}_k) \subset \mathbb{R}_+$ using the argument of Section 1.1, where $\limsup \mathbf{r}_k \approx \mathfrak{z}(\mathbf{x})$.

Lemma 4.5. *Assume that*

$$\Omega(\mathbf{x}) \subseteq \bigcap_{k \in \mathbb{N}} \left\{ \mathbf{v}_{k,k(+)1} \in \mathcal{Y}_\circ(\mathbf{r}_k^{(l)}) \right\}.$$

Then under the assumptions of Theorem 2.2 we get on $\Omega(\mathbf{x})$ for all $k \in \mathbb{N}_0$

$$\begin{aligned}
\|\mathcal{D}(\widetilde{\mathbf{v}}_{k,k(+)1} - \mathbf{v}^*)\| &\leq 2\sqrt{2}(1 - \sqrt{\rho})^{-1} (\mathfrak{z}(\mathbf{x}) + (1 + \sqrt{\rho})\rho^k R_0(\mathbf{x})) \\
&\quad + 2\sqrt{2}(1 + \sqrt{\rho}) \sum_{r=0}^{k-1} \rho^r \diamond_Q(\mathbf{r}_r^{(l)}) \\
&=: \mathbf{r}_k^{(l+1)}.
\end{aligned}$$

Proof. 1. We first show that on $\Omega(\mathbf{x})$

$$\begin{aligned}
D(\widetilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) &= D^{-1}\nabla_{\boldsymbol{\theta}}\mathcal{L}(\mathbf{v}^*) - D^{-1}A(\widetilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*) + \boldsymbol{\tau}(\mathbf{r}_k^{(l)}), \\
H(\widetilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*) &= H^{-1}\nabla_{\boldsymbol{\eta}}\mathcal{L}(\mathbf{v}^*) - H^{-1}A^\top(\widetilde{\boldsymbol{\theta}}_{k-1} - \boldsymbol{\theta}^*) + \boldsymbol{\tau}(\mathbf{r}_k^{(l)}),
\end{aligned} \tag{4.2}$$

where

$$\|\boldsymbol{\tau}(\mathbf{r})\| \leq \diamond_Q(\mathbf{r}, \mathbf{x}) = \{\delta(\mathbf{r})\mathbf{r} + 6\nu_1\omega(\mathfrak{z}_Q(\mathbf{x}, 4\rho^*) + 2\mathbf{r}^2)\}.$$

The proof is the same in each step for both statements such that we only prove the first one. The arguments presented here are similar to those of Theorem D.1 in [2]. By assumption on $\Omega(\mathbf{x})$ we have $\widetilde{\mathbf{v}}_{k,k(+)1} \in \mathcal{Y}_\circ(\mathbf{r}_k^{(l)})$. Define with $\zeta = \mathcal{L} - \mathbb{E}\mathcal{L}$

$$\alpha(\mathbf{v}, \mathbf{v}^*) := \mathcal{L}(\mathbf{v}, \mathbf{v}^*) - (\nabla\zeta(\mathbf{v}^*)(\mathbf{v} - \mathbf{v}^*) - \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|^2/2).$$

Note that

$$\begin{aligned}\mathcal{L}(\mathbf{v}, \mathbf{v}^*) &= \nabla\zeta(\mathbf{v}^*)(\mathbf{v} - \mathbf{v}^*) - \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|^2/2 + \alpha(\mathbf{v}, \mathbf{v}^*) \\ &= \nabla_{\boldsymbol{\theta}}\zeta(\mathbf{v}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \|\mathcal{D}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2 + (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top A(\boldsymbol{\eta} - \boldsymbol{\eta}^*) \\ &\quad + \nabla_{\boldsymbol{\eta}}\zeta(\mathbf{v}^*)(\boldsymbol{\eta} - \boldsymbol{\eta}^*) - \|\mathcal{H}(\boldsymbol{\eta} - \boldsymbol{\eta}^*)\|^2/2 + \alpha(\mathbf{v}, \mathbf{v}^*).\end{aligned}$$

Setting $\nabla_{\boldsymbol{\theta}}\mathcal{L}(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k) = 0$ we find

$$\mathcal{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \mathcal{D}^{-1}(\nabla_{\boldsymbol{\theta}}\zeta(\mathbf{v}^*) - A(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*)) = \mathcal{D}^{-1}\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \mathbf{v}^*).$$

As we assume that $\tilde{\mathbf{v}}_{k,k} \in \mathcal{Y}_o(\mathbf{R}_0)$ it suffices to show that with dominating probability

$$\sup_{(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \in \mathcal{Y}_o(\mathbf{R}_0)} \|\mathcal{U}_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k)\| \leq \diamond(\mathbf{r}_k^{(l)}),$$

where

$$\mathcal{U}_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \stackrel{\text{def}}{=} \mathcal{D}^{-1}\{\nabla_{\boldsymbol{\theta}}\mathcal{L}(\tilde{\mathbf{v}}_{k,k}) - \nabla_{\boldsymbol{\theta}}\mathcal{L}(\mathbf{v}^*) - \mathcal{D}^2(\boldsymbol{\theta} - \boldsymbol{\theta}^*) - A(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*)\}.$$

To see this note first that with Lemma 4.2 $\|\mathcal{D}^{-1}\Pi_{\boldsymbol{\theta}}\mathcal{D}\mathbf{v}\| \leq \|\mathcal{D}^{-1}\mathcal{D}\mathbf{v}\|$. This gives by condition (\mathcal{L}_0) , Lemma 4.2 and Taylor expansion

$$\begin{aligned}\sup_{(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{E}\mathcal{U}_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k)\| &\leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{D}^{-1}\Pi_{\boldsymbol{\theta}}(\nabla\mathcal{E}\mathcal{L}(\mathbf{v}) - \nabla\mathcal{E}\mathcal{L}(\mathbf{v}^*) - \mathcal{D}(\mathbf{v} - \mathbf{v}^*))\| \\ &\leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{D}^{-1}\Pi_{\boldsymbol{\theta}}\mathcal{D}\| \|\mathcal{D}^{-1}\nabla^2\mathcal{E}\mathcal{L}(\mathbf{v})^2\mathcal{D}^{-1} - I_{p^*}\|^{1/2}\mathbf{r} \\ &\leq \delta(\mathbf{r})\mathbf{r}.\end{aligned}$$

For the remainder note that again with Lemma 4.2

$$\|\mathcal{D}^{-1}(\nabla_{\boldsymbol{\theta}}\zeta(\mathbf{v}) - \nabla_{\boldsymbol{\theta}}\zeta(\mathbf{v}^*))\| \leq \|\mathcal{D}^{-1}(\nabla\zeta(\mathbf{v}) - \nabla\zeta(\mathbf{v}^*))\|.$$

This yields that on $\Omega(\mathbf{x})$

$$\begin{aligned}\sup_{(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{U}_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) - \mathcal{E}\mathcal{U}_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k)\| &\leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{D}^{-1}(\nabla_{\boldsymbol{\theta}}\zeta(\mathbf{v}) - \nabla_{\boldsymbol{\theta}}\zeta(\mathbf{v}^*))\| \\ &\leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \left\{ \frac{1}{6\nu_1\omega} \|\mathcal{Y}(\mathbf{v})\| \right\} 6\nu_1\omega \leq 6\nu_1\omega \{3Q(\mathbf{x}, 4p^*) + 2\mathbf{r}^2\}.\end{aligned}$$

Using the same argument for $\tilde{\boldsymbol{\eta}}_k$ gives the claim.

2. We prove the apriori bound for the distance of the k. estimator to the oracle

$$\|\mathcal{D}(\tilde{\mathbf{v}}_{k,k(+1)} - \mathbf{v}^*)\| \leq \mathbf{r}_k^{(l+1)}.$$

To see this we first use the inequality

$$\|\mathcal{D}(\tilde{\mathbf{v}}_{k,k(+1)} - \mathbf{v}^*)\| \leq \sqrt{2}\|D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*)\| + \sqrt{2}\|H(\tilde{\boldsymbol{\eta}}_{k(+1)} - \boldsymbol{\eta}^*)\|.$$

Now we find with (4.2)

$$\begin{aligned} \|D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*)\| &\leq \|D^{-1}\nabla_{\boldsymbol{\theta}}\mathcal{L}(\mathbf{v}^*)\| + \|D^{-1}A(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*)\| + \|\boldsymbol{\tau}(\mathbf{r}_k^{(l)})\| \\ &\leq \|D^{-1}\nabla_{\boldsymbol{\theta}}\mathcal{L}(\mathbf{v}^*)\| + \|D^{-1}AH^{-1}\| \|H(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*)\| + \|\boldsymbol{\tau}(\mathbf{r}_k^{(l)})\|. \end{aligned}$$

Next we use that on $\Omega(\mathbf{x})$

$$\|D^{-1}AH^{-1}\| \leq \sqrt{\rho}, \quad \|D^{-1}\nabla_{\boldsymbol{\theta}}\mathcal{L}(\mathbf{v}^*)\| \leq \mathfrak{z}(\mathbf{x}), \quad \|H^{-1}\nabla_{\boldsymbol{\eta}}\mathcal{L}(\mathbf{v}^*)\| \leq \mathfrak{z}(\mathbf{x}),$$

and

$$\|H(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*)\| \leq \|H^{-1}\nabla_{\boldsymbol{\eta}}\mathcal{L}(\mathbf{v}^*)\| + \|H^{-1}A^\top(\tilde{\boldsymbol{\theta}}_{k-1} - \boldsymbol{\theta}^*)\| + \|\boldsymbol{\tau}(\mathbf{r}_k^{(l)})\|,$$

to derive the recursive formula

$$\|D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*)\| \leq (1 + \sqrt{\rho}) \left(\mathfrak{z}(\mathbf{x}) + \|\boldsymbol{\tau}(\mathbf{r}_k^{(l)})\| \right) + \rho \|D(\tilde{\boldsymbol{\theta}}_{k-1} - \boldsymbol{\theta}^*)\|.$$

Deriving the analogous formula for $\|H(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*)\|$ and solving the recursion gives the claim. \square

Lemma 4.6. *Assume the same as in Theorem 2.2. Then we get*

$$\Omega(\mathbf{x}) \subseteq \bigcap_{k \in \mathbb{N}} \left\{ \mathbf{v}_{k,k(+1)} \in \Upsilon_{\circ}(\mathbf{r}_k^{(1)}) \right\},$$

where

$$\mathbf{r}_k^{(1)} \leq 2\sqrt{2}(1 - \sqrt{\rho})^{-1} \left\{ (\mathfrak{z}(\mathbf{x}) + \diamond_Q(\mathbf{R}_0, \mathbf{x})) + (1 + \sqrt{\rho})\rho^k \mathbf{R}_0(\mathbf{x}) \right\}. \quad (4.3)$$

Further assume that $\delta(\mathbf{r})/\mathbf{r} \vee 12\nu_1\omega \leq \epsilon$ and that (2.6) and (2.7) are met with $\mathbf{C}(\rho)$ defined in (2.8). Then

$$\Omega(\mathbf{x}) \subseteq \bigcap_{k \in \mathbb{N}} \left\{ \mathbf{v}_{k,k(+1)} \in \Upsilon_{\circ}(\mathbf{r}_k^*) \right\},$$

where

$$\begin{aligned} \mathbf{r}_k^* &\leq \mathbf{C}(\rho) \left(\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2 \right) + \epsilon \frac{7^2 \mathbf{C}(\rho)^4}{1 - c(\epsilon, \mathfrak{z}(\mathbf{x}))} \left(\frac{1}{1 - \rho} \right) \left(\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2 \right)^2 \\ &\quad + \rho^k \left(\mathbf{C}(\rho) \mathbf{R}_0 + \epsilon \frac{7^2 \mathbf{C}(\rho)^4}{1 - c(\epsilon, \mathbf{R}_0)} \left(\frac{1}{\rho^{-1} - 1} \right) \mathbf{R}_0^2 \right). \end{aligned} \quad (4.4)$$

Proof. We proof this claim via induction. On $\Omega(\mathbf{x})$ we have

$$\mathbf{v}_{k,k(+1)} \in \mathcal{Y}_\circ(\mathbf{R}_0), \quad \text{set } \mathbf{r}_k^{(0)} \stackrel{\text{def}}{=} \mathbf{R}_0.$$

Now with Lemma 4.5 we find that

$$\Omega(\mathbf{x}) \subseteq \bigcap_{k \in \mathbb{N}} \left\{ \mathbf{v}_{k,k(+1)} \in \mathcal{Y}_\circ(\mathbf{r}_k^{(l)}) \right\} \text{ implies } \Omega(\mathbf{x}) \subseteq \bigcap_{k \in \mathbb{N}} \left\{ \mathbf{v}_{k,k(+1)} \in \mathcal{Y}_\circ(\mathbf{r}_k^{(l+1)}) \right\},$$

where

$$\begin{aligned} \mathbf{r}_k^{(l)} &\leq 2\sqrt{2}(1 - \sqrt{\rho})^{-1} (\mathfrak{z}(\mathbf{x}) + (1 + \sqrt{\rho})\rho^k \mathbf{R}_0(\mathbf{x})) \\ &\quad + 2\sqrt{2}(1 + \sqrt{\rho}) \sum_{r=0}^{k-1} \rho^r \diamond_Q(\mathbf{r}_r^{(l-1)}, \mathbf{x}). \end{aligned}$$

Setting $l = 1$ this gives

$$\mathbf{r}_k^{(1)} \leq 2\sqrt{2}(1 - \sqrt{\rho})^{-1} \left\{ (\mathfrak{z}(\mathbf{x}) + \diamond_Q(\mathbf{R}_0, \mathbf{x})) + (1 + \sqrt{\rho})\rho^k \mathbf{R}_0(\mathbf{x}) \right\},$$

which gives (4.3). For the second claim we show that

$$\Omega(\mathbf{x}) \subseteq \bigcap_{k \in \mathbb{N}} \left\{ \mathbf{v}_{k,k(+1)} \in \mathcal{Y}_\circ \left(\limsup_{l \rightarrow \infty} \mathbf{r}_k^{(l)} \right) \right\} \subseteq \bigcap_{k \in \mathbb{N}} \left\{ \mathbf{v}_{k,k(+1)} \in \mathcal{Y}_\circ(\mathbf{r}_k^*) \right\}.$$

So we have to show that $\limsup_{l \rightarrow \infty} \mathbf{r}_k^{(l)} \leq \mathbf{r}_k^*$ from (4.4). For this we use $\delta(r)/r \vee 12\nu_1\omega \leq \epsilon$ to estimate further

$$\begin{aligned} \mathbf{r}_k^{(l)} &\leq 2\sqrt{2}(1 - \sqrt{\rho})^{-1} (\mathfrak{z}(\mathbf{x}) + (1 + \sqrt{\rho})\rho^k \mathbf{R}_0(\mathbf{x})) \\ &\quad + 2\sqrt{2}(1 + \sqrt{\rho})\epsilon \sum_{r=0}^{k-1} \rho^r \left((\mathbf{r}_{k-r}^{(l-1)})^2 + \mathfrak{z}(\mathbf{x})^2 \right) \\ &\leq 2\sqrt{2}(1 - \sqrt{\rho})^{-1} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2 + (1 + \sqrt{\rho})\rho^k \mathbf{R}_0(\mathbf{x})) \\ &\quad + 2\sqrt{2}(1 + \sqrt{\rho})\epsilon \sum_{r=0}^{k-1} \rho^r (\mathbf{r}_{k-r}^{(l-1)})^2 \\ &\leq C(\rho) \left\{ (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2) + \rho^k \mathbf{R}_0 + \epsilon \sum_{r=0}^{k-1} \rho^r (\mathbf{r}_{k-r}^{(l-1)})^2 \right\}, \end{aligned}$$

where $C(\rho) > 0$ is defined in (2.8). We set

$$A_{s,k}^{(l)} \stackrel{\text{def}}{=} \sum_{r_1=0}^{k-1} \rho^{r_1} \left(\sum_{r_2=0}^{k-r_1-1} \rho^{r_2} \left(\dots \sum_{r_s=0}^{k-r_1-\dots-r_{s-1}-1} \rho^{r_s} (\mathbf{r}_{k-r_1-\dots-r_s}^{(l-1)})^2 \dots \right) \right)^2.$$

Claim

$$\begin{aligned}
A_{s,k}^{(l)} &\leq 7^{\sum_{t=0}^{s-1} 2^t} \mathbf{C}(\rho)^{2^s} \left\{ \left(\frac{1}{1-\rho} \right)^{\sum_{t=0}^{s-1} 2^t} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^{2^s} \right. \\
&\quad \left. + \rho^k \left(\frac{1}{\rho^{-1}-1} \right)^{\sum_{t=0}^{s-1} 2^t} \mathbf{R}_0^{2^s} \right\} \\
&\quad + 7^{\sum_{t=0}^{s-1} 2^t} (\mathbf{C}(\rho)\epsilon)^{2^s} A_{s+1,k}^{(l-1)}.
\end{aligned} \tag{4.5}$$

We proof this claim via induction. Clearly

$$\begin{aligned}
A_{1,k}^{(l)} &= \sum_{r_1=0}^{k-1} \rho^{r_1} (\mathbf{r}_{k-r_1}^{(l-1)})^2 \leq 7\mathbf{C}(\rho)^2 \sum_{r_1=0}^{k-1} \rho^{r_1} \left\{ (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^2 + \rho^{2(k-r_1)} \mathbf{R}_0^2 \right\} \\
&\quad + 7\mathbf{C}(\rho)^2 \epsilon^2 \sum_{r_1=0}^{k-1} \rho^{r_1} \left(\sum_{r_2=0}^{k-r_1-r_2-1} \rho^{r_2} (\mathbf{r}_{k-r_1-r_2}^{(l-2)})^2 \right)^2 \\
&\leq 7\mathbf{C}(\rho)^2 \left\{ \frac{1}{1-\rho} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^2 + \frac{\rho^k}{\rho^{-1}-1} \mathbf{R}_0^2 \right\} \\
&\quad + 7\mathbf{C}(\rho)^2 \epsilon^2 A_{2,k}^{(l-1)}.
\end{aligned}$$

Further

$$\begin{aligned}
A_{s,k}^{(l)} &\stackrel{\text{def}}{=} \sum_{r_1=0}^{k-1} \rho^{r_1} \left(\sum_{r_2=0}^{k-r_1-1} \rho^{r_2} \left(\dots \sum_{r_s=0}^{k-r_1-\dots-r_{s-1}-1} \rho^{r_s} (\mathbf{r}_{k-r_1-\dots-r_s}^{(l-1)})^2 \dots \right) \right)^2 \\
&= \sum_{r_1=0}^{k-1} \rho^{r_1} \left(A_{s-1,k-r_1}^{(l)} \right)^2.
\end{aligned} \tag{4.6}$$

Plugging in (4.5) we get for $s \geq 2$

$$\begin{aligned}
A_{s,k}^{(l)} &\leq \sum_{r_1=0}^{k-1} \rho^{r_1} \left(7^{\sum_{t=0}^{s-2} 2^t} \mathbf{C}(\rho)^{2^{s-1}} \left\{ \left(\frac{1}{1-\rho} \right)^{\sum_{t=0}^{s-2} 2^t} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^{2^{s-1}} \right. \right. \\
&\quad \left. \left. + \rho^k \left(\frac{1}{\rho^{-1}-1} \right)^{\sum_{t=0}^{s-2} 2^t} \mathbf{R}_0^{2^{s-1}} \right\} \right. \\
&\quad \left. + 7^{\sum_{t=0}^{s-2} 2^t} (\mathbf{C}(\rho)\epsilon)^{2^{s-1}} A_{s,k-r_1}^{(l-1)} \right)^2.
\end{aligned}$$

Shifting the index this gives

$$A_{s,k}^{(l)} \leq 7 \sum_{r_1=0}^{k-1} \rho^{r_1} \left(7^{\sum_{t=1}^{s-1} 2^t} \mathbf{C}(\rho)^{2^s} \left\{ \left(\frac{1}{1-\rho} \right)^{\sum_{t=1}^{s-1} 2^{t-1}} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^{2^s} \right. \right. \\ \left. \left. + \rho^k \left(\frac{1}{\rho^{-1}-1} \right)^{\sum_{t=1}^{s-1} 2^t} \mathbf{R}_0^{2^s} \right\} \right. \\ \left. + 7^{\sum_{t=1}^{s-1} 2^t} (\mathbf{C}(\rho)\epsilon)^{2^s} (A_{s,k-r_1}^{(l-1)})^2 \right).$$

Direct calculation then leads to

$$A_{s,k}^{(l)} \leq 7^{\sum_{t=0}^{s-1} 2^t} \mathbf{C}(\rho)^{2^s} \left\{ \left(\frac{1}{1-\rho} \right)^{\sum_{t=0}^{s-1} 2^t} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^{2^s} \right. \\ \left. + \rho^k \left(\frac{1}{\rho^{-1}-1} \right)^{\sum_{t=0}^{s-1} 2^t} \mathbf{R}_0^{2^s} \right\} \\ + 7^{\sum_{t=0}^{s-1} 2^t} (\mathbf{C}(\rho)\epsilon)^{2^s} \sum_{r_1=0}^{k-1} \rho^{r_1} (A_{s,k-r_1}^{(l-1)})^2,$$

which gives (4.5) with (4.6). Similarly we can prove

$$A_{s,k}^{(1)} = \left(\frac{1}{1-\rho} \right)^{2^{s-1}} \mathbf{R}_0^{2^s}.$$

Abbreviate

$$\lambda_s \stackrel{\text{def}}{=} 7^{2^s-1} \mathbf{C}(\rho)^{2^s}, \quad \beta_s \stackrel{\text{def}}{=} 7^{2^s-1} (\mathbf{C}(\rho)\epsilon)^{2^s}, \\ \mathfrak{z}_s(\mathbf{x}) \stackrel{\text{def}}{=} \left(\frac{1}{1-\rho} \right)^{2^{s-1}} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^{2^s}, \quad \mathbf{R}_s \stackrel{\text{def}}{=} \left(\frac{1}{\rho^{-1}-1} \right)^{2^{s-1}} \mathbf{R}_0^{2^s}.$$

Then

$$\mathbf{r}_k^{(l)} \leq \mathbf{C}(\rho) \left\{ (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2) + \rho^k \mathbf{R}_0 + \epsilon A_{1,k}^{(l)} \right\} \\ \leq \sum_{s=0}^{l-1} \lambda_s \prod_{r=0}^{s-1} \beta_r \mathfrak{z}_s(\mathbf{x}) + \rho^k \sum_{s=0}^{l-1} \lambda_s \prod_{r=0}^{s-1} \beta_r \mathbf{R}_s + \prod_{r=0}^{l-1} \beta_r \mathbf{R}_l. \quad (4.7)$$

We estimate further

$$\begin{aligned}
& \sum_{s=0}^{l-1} \lambda_s \prod_{r=0}^{s-1} \beta_r \mathfrak{z}_s(\mathbf{x}) - \mathbf{C}(\rho) (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2) = \sum_{s=1}^{l-1} \lambda_s \prod_{r=0}^{s-1} \beta_r \mathfrak{z}_s(\mathbf{x}) \\
& \leq \sum_{s=1}^{l-1} 7^{2^s} \mathbf{C}(\rho)^{2^{s+1}} \epsilon^{2^s-1} \left(\frac{1}{1-\rho} \right)^{2^s-1} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^{2^s} \\
& = \epsilon 7^2 \mathbf{C}(\rho)^4 \left(\frac{1}{1-\rho} \right) (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^2 \sum_{s=1}^{l-1} \left(\epsilon 7 \mathbf{C}(\rho) \frac{1}{1-\rho} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2) \right)^{2^s-1}.
\end{aligned}$$

Assuming (2.6) this gives

$$\begin{aligned}
\sum_{s=0}^{l-1} \lambda_s \prod_{r=0}^{s-1} \beta_r \mathfrak{z}_s(\mathbf{x}) & \leq \mathbf{C}(\rho) (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2) \\
& + \epsilon \frac{7^2 \mathbf{C}(\rho)^4}{1 - c(\epsilon, \mathfrak{z}(\mathbf{x}))} \left(\frac{1}{1-\rho} \right) (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^2.
\end{aligned}$$

With the same argument we find under (2.7) that

$$\rho^k \sum_{s=0}^{l-1} \lambda_s \prod_{r=0}^{s-1} \beta_r \mathbf{R}_s \leq \rho^k \left(\mathbf{C}(\rho) \mathbf{R}_0 + \epsilon \frac{7^2 \mathbf{C}(\rho)^4}{1 - c(\epsilon, \mathbf{R}_0)} \left(\frac{1}{\rho^{-1} - 1} \right) \mathbf{R}_0^2 \right).$$

Additionally (2.7) implies

$$\prod_{r=0}^{l-1} \beta_r \mathbf{R}_l \leq \left(\epsilon 7 \mathbf{C}(\rho) \frac{1}{\rho^{-1} - 1} \right)^{2^{l-1}} \mathbf{R}_0^{2^l} \rightarrow 0.$$

Plugging these bounds into (4.7) and letting $l \rightarrow \infty$ gives the claim. \square

4.3 Result after convergence

In the previous section we showed that

$$\begin{aligned}
\Omega(\mathbf{x}) & \subset \bigcap_{\mathbf{r} \leq 4\mathbf{R}_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \left\{ \frac{1}{6\tilde{\omega}\tilde{\nu}_1} \|\check{\mathfrak{y}}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 2p^* + 2p)^2 \right\} \\
& \cap \bigcap_{k \in \mathbb{N}} \left\{ \mathbf{v}_{k,k} \in \mathcal{Y}_\circ(\mathbf{r}_k^{(\cdot)}), \mathbf{v}_{k,k+1} \in \mathcal{Y}_\circ(\mathbf{r}_k^{(\cdot)}) \right\} \cap \{ \tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_\circ(\mathbf{r}_0) \},
\end{aligned}$$

where $\mathbf{r}_k^{(\cdot)}$ is defined in (4.4) or (4.3). The claim of Theorem 2.2 follows with the following lemma:

Lemma 4.7. Assume $(\check{\mathcal{D}}_1)$, $(\check{\mathcal{L}}_0)$, and (\mathcal{I}) with a central point $\mathbf{v}^\circ = \mathbf{v}^*$ and $\mathcal{D}^2 = \nabla^2 \mathcal{E}\mathcal{L}(\mathbf{v}^*)$. Then it holds on $\Omega(\mathbf{x}) \subseteq \Omega$ that for all $k \in \mathbb{N}$

$$\|\check{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\| \leq \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}), \quad (4.8)$$

$$\begin{aligned} |2\check{L}(\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}^*) - \|\check{\boldsymbol{\xi}}\|^2| &\leq 8 \left(\|\check{D}^{-1}\check{\nabla}\| + \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}) \right) \check{\diamond}_Q(2(1+\rho)\mathbf{r}_k, \mathbf{x}) \\ &\quad + \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x})^2, \end{aligned} \quad (4.9)$$

where the spread $\check{\diamond}(\mathbf{r}, \mathbf{x})$ is defined in (2.5) and where

$$\mathbf{r}_k \stackrel{\text{def}}{=} \mathbf{r}_k^{(\cdot)} \vee \mathbf{r}_0.$$

Proof. The proof is nearly the same as that of Theorem 2.2 of [2] which is inspired by the proof of Theorem 1 of [12]. So we only sketch it and refer the reader to [2] for the skipped arguments. We define

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

Note that

$$\nabla_{\boldsymbol{\theta}_1} l(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\eta}) = \check{\nabla}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_1, \boldsymbol{\eta} + H^{-2}A^\top(\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1)), \quad \tilde{\boldsymbol{\theta}}_k = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} l(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k),$$

such that $\check{\nabla}_{\boldsymbol{\theta}} \mathcal{L}(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k) = 0$. This gives

$$\|\check{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\| = \|\check{D}^{-1}\check{\nabla}\mathcal{L}(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k) - \check{D}^{-1}\check{\nabla}\mathcal{L}(\mathbf{v}^*) + \check{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*)\|.$$

Now the right hand side can be bounded just as in the proof of Theorem 2.2 of [2]. This gives (4.8).

For (4.9) we can represent:

$$\check{L}(\tilde{\boldsymbol{\theta}}_k) - \check{L}(\boldsymbol{\theta}^*) = l(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_{k+1}) - l(\boldsymbol{\theta}^*, \boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*}),$$

where

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

Due to the definition of $\tilde{\boldsymbol{\theta}}_k$ and $\tilde{\boldsymbol{\eta}}_{k+1}$

$$l(\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*}) - l(\boldsymbol{\theta}^*, \boldsymbol{\theta}^*, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^*}) \leq \check{L}(\tilde{\boldsymbol{\theta}}_k) - \check{L}(\boldsymbol{\theta}^*) \leq l(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_{k+1}) - l(\boldsymbol{\theta}^*, \tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_{k+1}).$$

Again the remaining steps are exactly the same as in the proof of Theorem 2.2 of [2].

□

5 Proof of Corollary 2.3

Proof. Note that with the argument of Section 4.1 $\mathbb{P}(\Omega'(\mathbf{x})) \geq 1 - 8e^{-x} - \beta_{(\mathbf{A})}$ where with $\Omega(\mathbf{x})$ from (4.1)

$$\Omega'(\mathbf{x}) = \Omega(\mathbf{x}) \cap \{\tilde{\mathbf{v}} \in \mathcal{Y}_\circ(\mathbf{r}_0)\}.$$

On $\Omega'(\mathbf{x})$ it holds due to Theorem 2.2 and due to Theorem 2.1 of [2]

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

Now the claim follows with the triangular inequality. \square

6 Proof of Theorem 2.4

We prove this Theorem in a similar manner to the convergence result in Lemma 4.5. Redefine the set $\Omega(\mathbf{x})$

$$\Omega(\mathbf{x}) \stackrel{\text{def}}{=} \bigcap_{k=0}^K (C_{k,k} \cap C_{k,k+1}) \cap C(\nabla) \cap \{\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -K_0(\mathbf{x})\}, \text{ where}$$

$$C_{k,k(+1)} = \left\{ \|\mathcal{D}(\tilde{\mathbf{v}}_{k,k(+1)} - \mathbf{v}^*)\| \leq R_0(\mathbf{x}), \|\mathcal{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*)\| \leq R_0(\mathbf{x}), \right. \\ \left. \|H(\tilde{\boldsymbol{\eta}}_{k(+1)} - \boldsymbol{\eta}^*)\| \leq R_0(\mathbf{x}) \right\},$$

$$C(\nabla) = \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_\circ(R_0(\mathbf{x}))} \|\mathcal{Y}(\nabla^2)(\mathbf{v})\| \leq 9\nu_2\omega_2\mathfrak{z}_1(\mathbf{x}, 6p^*)R_0(\mathbf{x}) \right\} \\ \cap \{\|\mathcal{D}^{-1}\nabla^2\zeta(\mathbf{v}^*)\| \leq \mathfrak{z}(\mathbf{x}, \nabla^2\zeta(\mathbf{v}^*))\}.$$

where

$$\mathcal{Y}(\nabla^2)(\mathbf{v}) \stackrel{\text{def}}{=} \mathcal{D}^{-1}(\nabla^2\zeta(\mathbf{v}) - \nabla^2\zeta(\mathbf{v}^*)) \in \mathbb{R}^{p^*2}.$$

We see that on $\Omega(\mathbf{x})$

$$\mathbf{v}_{k,k(+1)} \in \tilde{\mathcal{Y}}_\circ(R_0) \stackrel{\text{def}}{=} \{\|\mathcal{D}(\mathbf{v} - \tilde{\mathbf{v}})\| \leq R_0 + \mathbf{r}_0\} \cap \mathcal{Y}_\circ(R_0).$$

Lemma 6.1. *Under the conditions of Theorem 2.4*

$$\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 3e^{-x} - \beta_{(\mathbf{A})}.$$

Proof. The proof is very similar to the one presented in Section 4.1, so we only give a sketch. By assumption

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

and due to $(\mathcal{E}\mathcal{D}_2)$ with Theorem 9.2

$$\mathbb{P} \left(\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{R}_0(\mathbf{x}))} \|\mathcal{Y}(\nabla^2)(\mathbf{v})\| \leq 9\nu_2\omega_2\mathfrak{z}_1(\mathbf{x}, 6p^*)\mathbf{R}_0(\mathbf{x}) \right) \geq 1 - e^{-\mathbf{x}}.$$

□

Lemma 6.2. Assume for some sequence $(\mathbf{r}_k^{(l)})$ that

$$\bigcap_{k \in \mathbb{N}} \left\{ \|\mathcal{D}(\tilde{\mathbf{v}}_{k,k+1}) - \tilde{\mathbf{v}}\| \leq \mathbf{r}_k^{(l)} \right\} \subseteq \Omega(\mathbf{x}).$$

Then we get on $\Omega(\mathbf{x})$

$$\begin{aligned} \|\mathcal{D}(\tilde{\mathbf{v}}_{k,k+1}) - \tilde{\mathbf{v}}\| &\leq 2\sqrt{2}(1 + \sqrt{\rho}) \sum_{r=0}^{k-1} \rho^r \|\boldsymbol{\tau}(\mathbf{r}_{k-r}^{(l)})\| + 2\sqrt{2}\rho^k(\mathbf{R}_0 + \mathbf{r}_0), \\ &=: \mathbf{r}_k^{(l+1)}. \end{aligned} \tag{6.1}$$

where

$$\|\boldsymbol{\tau}(\mathbf{r})\| \leq [\delta(\mathbf{R}_0) + 9\nu_2\omega_2\|\mathcal{D}^{-1}\|\mathfrak{z}_1(\mathbf{x}, 6p^*)\mathbf{R}_0 + \|\mathcal{D}^{-1}\|\mathfrak{z}(\mathbf{x}, \nabla^2\zeta(\mathbf{v}^*))] \mathbf{r}.$$

Proof. 1. We first show that on $\Omega(\mathbf{x})$

$$\begin{aligned} D(\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}) &= -D^{-1}A(\tilde{\boldsymbol{\eta}}_k - \tilde{\boldsymbol{\eta}}) + \boldsymbol{\tau}(\mathbf{r}_k^{(l)}), \\ H(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*) &= -H^{-1}A^\top(\tilde{\boldsymbol{\theta}}_{k-1} - \tilde{\boldsymbol{\theta}}) + \boldsymbol{\tau}(\mathbf{r}_k^{(l)}), \end{aligned}$$

The proof is very similar to that of Lemma 4.5. Define

$$\alpha(\mathbf{v}, \tilde{\mathbf{v}}) := \mathcal{L}(\mathbf{v}, \tilde{\mathbf{v}}) + \|\mathcal{D}(\mathbf{v} - \tilde{\mathbf{v}})\|^2/2.$$

Note that

$$\begin{aligned} \mathcal{L}(\mathbf{v}, \tilde{\mathbf{v}}) &= \nabla\mathcal{L}(\mathbf{v}) - \|\mathcal{D}(\mathbf{v} - \tilde{\mathbf{v}})\|^2/2 + \alpha(\mathbf{v}, \mathbf{v}^*) \\ &= -\|D(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})\|^2/2 + (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top A(\boldsymbol{\eta} - \tilde{\boldsymbol{\eta}}) \\ &\quad - \|H(\boldsymbol{\eta} - \tilde{\boldsymbol{\eta}})\|^2/2 + \alpha(\mathbf{v}, \tilde{\mathbf{v}}). \end{aligned}$$

Setting $\nabla_{\boldsymbol{\theta}}\mathcal{L}(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k) = 0$ we find

$$D(\tilde{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}) = D^{-1}A(\tilde{\boldsymbol{\eta}}_k - \tilde{\boldsymbol{\eta}}) + D^{-1}\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \tilde{\mathbf{v}}).$$

We want to show

$$\sup_{(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^{(l)}) \cap \mathcal{Y}_o(\mathbf{R}_0)} D^{-1}\nabla_{\boldsymbol{\theta}}\alpha((\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k), \tilde{\mathbf{v}}) \leq \|\boldsymbol{\tau}(\mathbf{r}_k^{(l)})\|,$$

where

$$D^{-1}\nabla_{\theta}\alpha(\mathbf{v}, \tilde{\mathbf{v}}) \stackrel{\text{def}}{=} D^{-1}\{\nabla_{\theta}\mathcal{L}(\mathbf{v}) - D^2(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) - A(\tilde{\boldsymbol{\eta}}_k - \tilde{\boldsymbol{\eta}})\}.$$

To see this note that by assumption we have $\Omega(\mathbf{x}) \subseteq \{\tilde{\mathbf{v}} \in \mathcal{Y}_o(\mathbf{r}_o)\} \subseteq \{\tilde{\mathbf{v}} \in \mathcal{Y}_o(\mathbf{R}_o)\}$. By condition (\mathcal{L}_o) , Lemma 4.2 and Taylor expansion we have

$$\begin{aligned} & \sup_{(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^{(l)}) \cap \mathcal{Y}_o(\mathbf{R}_o)} \|\mathbb{E}\mathcal{U}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k)\| \\ & \leq \sup_{\mathbf{v} \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^{(l)}) \cap \mathcal{Y}_o(\mathbf{R}_o)} \|D^{-1}\Pi_{\boldsymbol{\theta}}\left(\nabla\mathbb{E}\mathcal{L}(\mathbf{v}) - \nabla\mathbb{E}\mathcal{L}(\tilde{\mathbf{v}}) - \mathcal{D}(\mathbf{v} - \mathbf{v}^*)\right)\| \\ & \leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{R}_o)} \|D^{-1}\Pi_{\boldsymbol{\theta}}\mathcal{D}\| \|\mathcal{D}^{-1}\nabla^2\mathbb{E}\mathcal{L}(\mathbf{v})\mathcal{D}^{-1} - I_{p^*}\| \mathbf{r}_k^{(l)} \\ & \leq \delta(\mathbf{R}_o)\mathbf{r}_k^{(l)}. \end{aligned}$$

For the remainder note that with $\zeta = \mathcal{L} - \mathbb{E}\mathcal{L}$ on $\Omega(\mathbf{x})$ using Lemma 4.2 we can bound

$$\begin{aligned} & \sup_{(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^{(l)}) \cap \mathcal{Y}_o(\mathbf{R}_o)} \left\| \mathcal{U}_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) - \mathbb{E}\mathcal{U}_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \right\| \\ & \leq \sup_{\mathbf{v} \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^{(l)}) \cap \mathcal{Y}_o(\mathbf{R}_o)} \left\| D^{-1}\left(\nabla_{\boldsymbol{\theta}}\zeta(\mathbf{v}) - \nabla_{\boldsymbol{\theta}}\zeta(\tilde{\mathbf{v}})\right) \right\| \\ & \leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{x})} \left\| \mathcal{D}^{-1}\nabla^2\zeta(\mathbf{v})\mathcal{D}^{-1} \right\| \mathbf{r}_k^{(l)} \\ & \leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{R}_o)} \left\{ \frac{1}{9\nu_2\omega_2} \left\| \mathcal{D}^{-1}\left(\nabla^2\zeta(\mathbf{v}) - \nabla^2\zeta(\mathbf{v}^*)\right)\mathcal{D}^{-1} \right\| \right\} 6\nu_1\omega\mathbf{r}_k^{(l)} \\ & \quad + \left\{ \left\| \mathcal{D}^{-1}\nabla^2\zeta(\mathbf{v}^*)\mathcal{D}^{-1} \right\| \right\} \mathbf{r}_k^{(l)} \\ & \leq [9\nu_2\omega_2\|\mathcal{D}^{-1}\|\mathfrak{z}_1(\mathbf{x}, 6p^*)\mathbf{R}_o + \|\mathcal{D}^{-1}\|\mathfrak{z}(\mathbf{x}, \nabla^2\zeta(\mathbf{v}^*))] \mathbf{r}_k^{(l)}. \end{aligned}$$

Using the same argument for $\tilde{\boldsymbol{\eta}}_k$ gives the claim.

Now the claim follows as in the proof of Lemma 4.5. \square

Lemma 6.3. Assume that $\delta(\mathbf{r})/\mathbf{r} \vee 9\nu_2\omega_2 \vee \|\mathcal{D}^{-1}\| \leq \epsilon_2$. Further assume that $\varkappa(\mathbf{x}, \mathbf{R}_o) < 1 - \rho$ where

$$\begin{aligned} \varkappa(\mathbf{x}, \mathbf{R}_o) \stackrel{\text{def}}{=} & \frac{2\sqrt{2}(1 + \sqrt{\rho})}{\sqrt{1 - \rho}} \left(\delta(\mathbf{R}_o) + 9\omega_2\nu_2\|\mathcal{D}^{-1}\|\mathfrak{z}_1(\mathbf{x}, 6p^*)\mathbf{R}_o \right. \\ & \left. + \|\mathcal{D}^{-1}\|\mathfrak{z}(\mathbf{x}, \nabla^2\mathcal{L}(\mathbf{v}^*)) \right). \end{aligned}$$

Then

$$\Omega(\mathbf{x}) \subseteq \bigcap_{k \in \mathbb{N}} \left\{ \mathbf{v}_{k,k(+1)} \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k) \right\},$$

where $(\mathbf{r}_k)_{k \in \mathbb{N}}$ satisfy the bound (2.12).

Proof. Define for all $k \in \mathbb{N}_0$ the sequence $\mathbf{r}_k^{(0)} = \mathbf{R}_0$. We estimate

$$\|\boldsymbol{\tau}(\mathbf{r}_k^{(l)})\| \leq \frac{1}{\sqrt{1-\rho}} (\delta(\mathbf{R}_0) + 6\nu_1\omega_2\|\mathcal{D}^{-1}\|\mathfrak{J}_1(\mathbf{x}, 6p^*)\mathbf{R}_0 + \|\mathcal{D}^{-1}\|\mathfrak{J}(\mathbf{x}, \mathcal{B}(\nabla^2))\mathbf{r}_k^{(l)}),$$

such that by definition

$$2\sqrt{2}(1 + \sqrt{\rho}) \sum_{r=0}^{k-1} \rho^r \|\boldsymbol{\tau}(\mathbf{r}_{k-r}^{(l)})\| \leq \varkappa(\mathbf{x}, \mathbf{R}_0) \sum_{r=0}^{k-1} \rho^r \mathbf{r}_{k-r}^{(l)}.$$

Plugging in the recursive formula for $\mathbf{r}_k^{(l)}$ from (6.1) and denoting $\widetilde{\mathbf{R}}_0 \stackrel{\text{def}}{=} \mathbf{R}_0 + \mathbf{r}_0$ we find

$$\begin{aligned} \mathbf{r}_k^{(l)} &\leq \varkappa(\mathbf{x}, \mathbf{R}_0) \sum_{r=0}^{k-1} \rho^r \mathbf{r}_{k-r}^{(l-1)} + 2\sqrt{2}\rho^k \widetilde{\mathbf{R}}_0 \\ &\leq \varkappa(\mathbf{x}, \mathbf{R}_0) \sum_{r=0}^{k-1} \rho^r \left(\varkappa(\mathbf{x}, \mathbf{R}_0) \sum_{s=0}^{k-r-1} \rho^s \mathbf{r}_{k-r-s}^{(l-2)} + 2\rho^{k-r} \widetilde{\mathbf{R}}_0 \right) + 2\sqrt{2}\rho^k \widetilde{\mathbf{R}}_0 \\ &\leq \varkappa(\mathbf{x}, \mathbf{R}_0)^2 \sum_{r=0}^{k-1} \rho^r \sum_{s=0}^{k-r-1} \rho^s \mathbf{r}_{k-r-s}^{(l-2)} + 2\sqrt{2}\rho^k \widetilde{\mathbf{R}}_0 (\varkappa(\mathbf{x}, \mathbf{R}_0)k + 1) \\ &\leq \varkappa(\mathbf{x}, \mathbf{R}_0)^2 \sum_{r=0}^{k-1} \rho^r \sum_{s=0}^{k-r-1} \rho^s \left(\varkappa(\mathbf{x}, \mathbf{R}_0) \sum_{t=0}^{k-r-s-1} \rho^t \mathbf{r}_{k-r-s-t}^{(l-3)} + 2\rho^{k-r-s} \widetilde{\mathbf{R}}_0 \right) \\ &\quad + 2\sqrt{2}\rho^k \widetilde{\mathbf{R}}_0 (\varkappa(\mathbf{x}, \mathbf{R}_0)k + 1) \\ &\leq \varkappa(\mathbf{x}, \mathbf{R}_0)^3 \sum_{r=0}^{k-1} \rho^r \sum_{s=0}^{k-r-1} \rho^s \mathbf{r}_{k-r-s}^{(l-3)} + 2\sqrt{2}\rho^k \widetilde{\mathbf{R}}_0 (\varkappa(\mathbf{x}, \mathbf{R}_0)^2 k^2 + \varkappa(\mathbf{x}, \mathbf{R}_0)k + 1). \end{aligned}$$

By induction this gives for $l \in \mathbb{N}$

$$\begin{aligned}
\mathbf{r}_k^{(l)} &\leq \varkappa(\mathbf{x}, \mathbf{R}_0)^l \sum_{r_1=0}^{k-1} \rho^{r_1} \sum_{r_2=0}^{k-r_1-1} \rho^{r_2} \dots \sum_{r_l=0}^{k-\sum_{s=1}^{l-1} r_s-1} \rho^{r_l} \widetilde{\mathbf{R}}_0 \\
&\quad + 2\sqrt{2}\rho^k \widetilde{\mathbf{R}}_0 \sum_{s=0}^{l-1} \varkappa(\mathbf{x}, \mathbf{R}_0)^s k^s \\
&\leq \left(\left(\frac{\varkappa(\mathbf{x}, \mathbf{R}_0)}{1-\rho} \right)^l + 2\sqrt{2}\rho^k \sum_{s=0}^{l-1} (\varkappa(\mathbf{x}, \mathbf{R}_0)k)^s \right) \widetilde{\mathbf{R}}_0 \\
&\leq \begin{cases} \left(\left(\frac{\varkappa(\mathbf{x}, \mathbf{R}_0)}{1-\rho} \right)^l + 2\sqrt{2}\rho^k \frac{1}{1-\varkappa(\mathbf{x}, \mathbf{R}_0)k} \right) \widetilde{\mathbf{R}}_0, & \varkappa(\mathbf{x}, \mathbf{R}_0)k \leq 1, \\ \varkappa(\mathbf{x}, \mathbf{R}_0)^l \left(\left(\frac{1}{1-\rho} \right)^l + 2\sqrt{2}\rho^k \frac{k^l}{\varkappa(\mathbf{x}, \mathbf{R}_0)k-1} \right) \widetilde{\mathbf{R}}_0, & \text{otherwise.} \end{cases}
\end{aligned}$$

By Lemma 6.2

$$\Omega(\mathbf{x}) \subset \bigcap_{k \in \mathbb{N}_0} \bigcap_{l \in \mathbb{N}} \left\{ \tilde{\mathbf{v}}_{k, k+(+1)} \in \tilde{\mathcal{Y}}_0(\mathbf{r}_k^{(l)}) \right\}.$$

Set if $\varkappa(\mathbf{x}, \mathbf{R}_0)/(1-\rho) < 1$

$$l(k) \stackrel{\text{def}}{=} \begin{cases} \infty, & \varkappa(\mathbf{x}, \mathbf{R}_0)k \leq 1, \\ \frac{k \log(\rho) + \log(2\sqrt{2}) - \log(\varkappa(\mathbf{x}, \mathbf{R}_0)k - 1)}{-\log(1-\rho) - \log(k)}, & \text{otherwise.} \end{cases}$$

Then with $\mathbf{r}_k^* \stackrel{\text{def}}{=} \mathbf{r}_k^{(l(k))}$ we get

$$\Omega(\mathbf{x}) \subset \bigcap_{k \in \mathbb{N}_0} \left\{ \tilde{\mathbf{v}}_{k, k+(+1)} \in \tilde{\mathcal{Y}}_0(\mathbf{r}_k^*) \right\}, \quad \mathbf{r}_k^* \leq \begin{cases} \frac{\rho^k 2\sqrt{2}}{1-\varkappa(\mathbf{x}, \mathbf{R}_0)k} \widetilde{\mathbf{R}}_0, & \varkappa(\mathbf{x}, \mathbf{R}_0)k \leq 1, \\ 2 \left(\frac{\varkappa(\mathbf{x}, \mathbf{R}_0)}{1-\rho} \right)^{\frac{k}{\log(k)} L(k)-1} \widetilde{\mathbf{R}}_0, & \text{otherwise,} \end{cases}$$

as claimed. \square

7 Deviation bounds for quadratic forms

This section is the same as Section A of [2]. The following general result from [14] helps to control the deviation for quadratic forms of type $\|\mathbf{B}\boldsymbol{\xi}\|^2$ for a given positive matrix \mathbf{B} and a random vector $\boldsymbol{\xi}$. It will be used several times in our proofs. Suppose that

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

For a symmetric matrix \mathbf{B} , define

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

For ease of presentation, suppose that $g^2 \geq 2p_B$. The other case only changes the constants in the inequalities. Note that $\|\xi\|^2 = \eta^\top \mathbb{B} \eta$. Define $\mu_c = 2/3$ and

$$\begin{aligned} g_c &\stackrel{\text{def}}{=} \sqrt{g^2 - \mu_c p_B}, \\ 2(x_c + 2) &\stackrel{\text{def}}{=} (g^2/\mu_c - p_B)/\lambda^* + \log \det(I_p - \mu_c \mathbb{B}/\lambda^*). \end{aligned}$$

Proposition 7.1. *Let (ED_0) hold with $\nu_0 = 1$ and $g^2 \geq 2p_B$. Then for each $x > 0$*

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

where $\mathfrak{z}(x, \mathbb{B})$ is defined by

$$\mathfrak{z}^2(\mathbb{B}, x) \stackrel{\text{def}}{=} \begin{cases} p_B + 2v_B(x+1)^{1/2}, & x+1 \leq v_B/(18\lambda^*), \\ p_B + 6\lambda^*(x+1), & v_B/(18\lambda^*) < x+1 \leq x_c + 2, \\ |y_c + 2\lambda^*(x - x_c + 1)/g_c|^2, & x > x_c + 1, \end{cases}$$

with $y_c^2 \leq p_B + 6\lambda^*(x_c + 2)$.

8 A uniform bound for the norm of a random process

We want to derive for a random process $\check{y}(\mathbf{v}) \in \mathbb{R}^p$ a bound of the kind

$$\mathbb{P} \left(\sup_{\mathbf{r} \leq \mathbf{r}^*} \sup_{\mathbf{v} \in \mathcal{V}_c(\mathbf{r})} \left\{ \frac{1}{\omega} \|\check{y}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \geq \mathfrak{C}_{\mathfrak{z}Q}(x, p^*) \right) \leq e^{-x}.$$

This is a slightly stronger result than the one derived in Section D of [2] but the ideas employed here are very similar.

We want to apply Corollary 2.5 of the supplement of [14] which we cite here as a Theorem. Note that we slightly generalized the formulation of the theorem, to make it applicable in our setting. The proof remains the same.

Theorem 8.1. *Let $(U(\mathbf{r}))_{0 \leq \mathbf{r} \leq \mathbf{r}^*} \subset \mathbb{R}^p$ be a sequence of balls around \mathbf{v}^* induced by the metric $d(\cdot, \cdot)$. Let a random real valued process $\mathcal{U}(\mathbf{r}, \mathbf{v})$ fulfill for any $0 \leq \mathbf{r} \leq \mathbf{r}^*$ that $\mathcal{U}(\mathbf{r}, \mathbf{v}^*) = 0$ and*

(Ed) *For any $\mathbf{v}, \mathbf{v}^\circ \in U(\mathbf{r})$*

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{r}, \mathbf{v}) - \mathcal{U}(\mathbf{r}, \mathbf{v}^\circ)}{d(\mathbf{v}, \mathbf{v}^\circ)} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq g. \quad (8.1)$$

Finally assume that $\sup_{\mathbf{v} \in U(\mathbf{r})} (\mathcal{U}(\mathbf{r}, \mathbf{v}))$ increases in \mathbf{r} . Then with probability greater $1 - e^{-x}$

$$\sup_{\mathbf{v} \in U(\mathbf{r})} \left\{ \frac{1}{3\nu_1} \mathcal{U}(\mathbf{r}, \mathbf{v}) - d(\mathbf{v}, \mathbf{v}^*)^2 \right\} \leq \mathfrak{z}Q(x, p^*)^2,$$

where $\mathfrak{z}_Q(\mathbf{x}, p^*) \stackrel{\text{def}}{=} \mathbb{Q}(U(\mathbf{r}^*))$ denotes the entropy of the set $U(\mathbf{r}^*) \subset \mathbb{R}^p$ and where with $g_0 = \nu_0 g$ and for some $Q > 0$

$$\mathfrak{z}_Q(\mathbf{x}, Q)^2 \stackrel{\text{def}}{=} \begin{cases} (1 + \sqrt{\mathbf{x} + Q})^2 & \text{if } 1 + \sqrt{\mathbf{x} + Q} \leq g_0, \\ 1 + \{2g_0^{-1}(\mathbf{x} + Q) + g_0\}^2 & \text{otherwise.} \end{cases} \quad (8.2)$$

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

$$\sup_{\|\mathbf{u}\| \leq 1} \log \mathbb{E} \exp \left\{ \lambda \frac{\mathbf{u}^\top (\check{\mathfrak{Y}}(\mathbf{v}) - \check{\mathfrak{Y}}(\mathbf{v}^\circ))}{\omega \|\mathcal{D}(\mathbf{v} - \mathbf{v}^\circ)\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}. \quad (8.3)$$

Remark 8.1. In the setting of Theorem 2.2 we have

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

and condition (8.3) becomes $(\mathcal{E}\mathcal{D}_1)$ from 2.1.

Theorem 8.2. Let a random p -vector process $\check{\mathfrak{Y}}(\mathbf{v})$ fulfill $\check{\mathfrak{Y}}(\mathbf{v}^*) = 0$ and let condition (8.3) be satisfied. Then for each $0 \leq \mathbf{r} \leq \mathbf{r}^*$, on a set of probability greater $1 - e^{-\mathbf{x}}$

$$\sup_{\mathbf{r} \leq \mathbf{r}^*} \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \left\{ \frac{1}{6\omega\nu_1} \|\check{\mathfrak{Y}}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 2p^* + 2p)^2,$$

with $g_0 = \nu_0 g$.

Remark 8.2. Note that the entropy of the original set $\mathcal{Y}_o(\mathbf{r}) \subset \mathbb{R}^{p^*}$ is equal to $2p^*$. So in order to control the norm $\|\check{\mathfrak{Y}}(\mathbf{v})\|$ one only pays with the additional sumand $2p$.

Proof. In what follows, we use the representation

$$\|\check{\mathfrak{Y}}(\mathbf{v})\| = \omega \sup_{\|\mathbf{u}\| \leq \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \frac{1}{\omega \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \mathbf{u}^\top \check{\mathfrak{Y}}(\mathbf{v}).$$

This implies

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\check{\mathfrak{Y}}(\mathbf{v})\| = \omega \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \sup_{\|\mathbf{u}\| \leq \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \frac{1}{\omega \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \mathbf{u}^\top \check{\mathfrak{Y}}(\mathbf{v}).$$

Due to Lemma 8.3 the process $\mathcal{U}(\mathbf{r}, \mathbf{v}, \mathbf{u}) \stackrel{\text{def}}{=} \frac{1}{\omega \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \mathbf{u}^\top \check{\mathfrak{Y}}(\mathbf{v})$ satisfies condition $(\mathcal{E}d)$ (see (8.1)) as process on $U(\mathbf{r}^*)$ where

$$U(\mathbf{r}) \stackrel{\text{def}}{=} \mathcal{Y}_o(\mathbf{r}) \times B_{\mathbf{r}}(0). \quad (8.4)$$

Further $\sup_{(\mathbf{v}, \mathbf{u}) \in U(\mathbf{r})} \mathcal{U}(\mathbf{r}, \mathbf{v}, \mathbf{u})$ is increasing in \mathbf{r} . This allows to apply Theorem 8.2 to obtain the desired result. Set $d((\mathbf{v}, \mathbf{u}), (\mathbf{v}^\circ, \mathbf{u}^\circ))^2 = \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|^2 + \|\mathbf{u} - \mathbf{u}^\circ\|^2$. We get on a set of probability greater $1 - e^{-x}$

$$\begin{aligned} & \sup_{(\mathbf{v}, \mathbf{u}) \in U(\mathbf{r}^*)} \left\{ \frac{1}{6\omega\nu_1 \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \mathbf{u}^\top \check{\mathcal{Y}}(\mathbf{v}) - \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|^2 - \|\mathbf{u}\|^2 \right\} \\ & \leq \mathfrak{J}_Q(\mathbf{x}, \mathbb{Q}(U(\mathbf{r}^*))). \end{aligned}$$

The constant $\mathbb{Q}(U(\mathbf{r}^*)) > 0$ quantifies the complexity of the set $U(\mathbf{r}^*) \subset \mathbb{R}^{p^*} \times \mathbb{R}^p$. We point out that for compact $M \subset \mathbb{R}^{p^*}$ we have $\mathbb{Q}(M) = 2p^*$ (see Supplement of [14], Lemma 2.10). This gives $\mathbb{Q}(U) = 2p^* + 2p$. Finally observe that

$$\begin{aligned} & \sup_{\mathbf{r} \leq \mathbf{r}^*} \sup_{\mathbf{v} \in \mathcal{V}_\circ(\mathbf{r})} \left\{ \frac{1}{6\omega\nu_1} \|\check{\mathcal{Y}}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \\ & \leq \sup_{\mathbf{r} \leq \mathbf{r}^*} \sup_{(\mathbf{v}, \mathbf{u}) \in U(\mathbf{r})} \left\{ \frac{1}{6\omega\nu_1 \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \mathbf{u}^\top \check{\mathcal{Y}}(\mathbf{v}) - \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|^2 - \|\mathbf{u}\|^2 \right\} \\ & = \sup_{(\mathbf{v}, \mathbf{u}) \in U(\mathbf{r}^*)} \left\{ \frac{1}{6\omega\nu_1 \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \mathbf{u}^\top \check{\mathcal{Y}}(\mathbf{v}) - \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|^2 - \|\mathbf{u}\|^2 \right\}. \end{aligned}$$

□

Lemma 8.3. Suppose that $\check{\mathcal{Y}}(\mathbf{v})$ satisfies for each $\|\mathbf{u}\| \leq 1$ and $|\lambda| \leq g$ the inequality (8.3). Then the process $\mathcal{U}(\mathbf{v}, \mathbf{u}) = \frac{1}{2\omega \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \check{\mathcal{Y}}(\mathbf{v})^\top \mathbf{u}_1$ satisfies $(\mathcal{E}d)$ from (8.1) with $|\lambda| \leq g/2$, $d((\mathbf{v}, \mathbf{u}), (\mathbf{v}^\circ, \mathbf{u}^\circ))^2 = \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|^2 + \|\mathbf{u} - \mathbf{u}^\circ\|^2$, $\nu = 2\nu_0$ and $U \subset \mathbb{R}^{p^*+p}$ defined in (8.4), i.e. for any $(\mathbf{v}, \mathbf{u}_1), (\mathbf{v}^\circ, \mathbf{u}_2) \in U$

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}, \mathbf{u}_1) - \mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_2)}{d((\mathbf{v}, \mathbf{u}_1), (\mathbf{v}^\circ, \mathbf{u}_2))} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq g/2.$$

Proof. Let $(\mathbf{v}, \mathbf{u}_1), (\mathbf{v}^\circ, \mathbf{u}_2) \in U$ and w.l.o.g. $\mathbf{u}_1 \leq \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \leq \|\mathcal{D}(\mathbf{v}^\circ - \mathbf{v}^*)\|$. By

the Hölder inequality and (8.3), we find

$$\begin{aligned}
& \log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}, \mathbf{u}_1) - \mathcal{U}(\mathbf{v}, \mathbf{u}_2)}{d((\mathbf{v}, \mathbf{u}_1), (\mathbf{v}^\circ, \mathbf{u}_2))} \right\} \\
&= \log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}, \mathbf{u}_1) - \mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_1) + \mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_1) - \mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_2)}{d((\mathbf{v}, \mathbf{u}_1), (\mathbf{v}^\circ, \mathbf{u}_2))} \right\} \\
&\leq \frac{1}{2} \log \mathbb{E} \exp \left\{ 2\lambda \frac{\mathbf{u}_1^\top \left(\frac{1}{\|\mathcal{D}(\mathbf{v}-\mathbf{v}^*)\|} \check{\mathcal{Y}}(\mathbf{v}) - \frac{1}{\|\mathcal{D}(\mathbf{v}^\circ-\mathbf{v}^*)\|} \check{\mathcal{Y}}(\mathbf{v}^\circ) \right)}{\omega \|\mathcal{D}(\mathbf{v}-\mathbf{v}^\circ)\|} \right\} \\
&\quad + \frac{1}{2} \log \mathbb{E} \exp \left\{ 2\lambda \frac{(\mathbf{u}_1^\top - \mathbf{u}_2^\top) \check{\mathcal{Y}}(\mathbf{v}^\circ)}{\omega \|\mathbf{u}_1 - \mathbf{u}_2\| \|\mathcal{D}(\mathbf{v}-\mathbf{v}^*)\|} \right\} \\
&\leq \sup_{\|\mathbf{u}\| \leq 1} \frac{1}{2} \log \mathbb{E} \exp \left\{ 2\lambda \frac{\mathbf{u}^\top (\check{\mathcal{Y}}(\mathbf{v}) - \check{\mathcal{Y}}(\mathbf{v}^\circ))}{\omega \|\mathcal{D}(\mathbf{v}-\mathbf{v}^\circ)\|} \right\} \\
&\quad + \sup_{\|\mathbf{u}\| \leq 1} \frac{1}{2} \log \mathbb{E} \exp \left\{ 2\lambda \frac{\mathbf{u}^\top (\check{\mathcal{Y}}(\mathbf{v}^\circ) - \check{\mathcal{Y}}(\mathbf{v}^*))}{\omega \|\mathcal{D}(\mathbf{v}-\mathbf{v}^*)\|} \right\} \\
&\leq \frac{4\nu_0^2 \lambda^2}{2}, \quad \lambda \leq g/2.
\end{aligned}$$

□

9 A bound for the spectral norm of a random matrix process

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

$$\mathbb{P} \left(\sup_{\mathbf{v} \in \mathcal{U}_\circ(\mathbf{r})} \left\{ \|\check{\mathcal{Y}}(\mathbf{v})\| \right\} \geq C\omega_2 \mathfrak{z}_1(\mathbf{x}, p^*) \mathbf{r} \right) \leq e^{-\mathbf{x}}.$$

We derive such a bound in a very similar manner to Theorem E.1 of [2].

We want to apply Corollary 2.2 of the supplement of [14]. Again we slightly generalized the formulation but the proof remains the same.

Corollary 9.1. *Let $(U(\mathbf{r}))_{0 \leq \mathbf{r} \leq \mathbf{r}^*} \subset \mathbb{R}^p$ be a sequence of balls around \mathbf{v}^* induced by the metric $d(\cdot, \cdot)$. Let a random real valued process $\mathcal{U}(\mathbf{v})$ fulfill that $\mathcal{U}(\mathbf{v}^*) = 0$ and*

(Ed) *For any $\mathbf{v}, \mathbf{v}^\circ \in U(\mathbf{r})$*

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)}{d(\mathbf{v}, \mathbf{v}^\circ)} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq g. \quad (9.1)$$

Then for each $0 \leq \mathbf{r} \leq \mathbf{r}^$, on a set of probability greater $1 - e^{-\mathbf{x}}$*

$$\sup_{\mathbf{v} \in U(\mathbf{r})} \mathcal{U}(\mathbf{v}) \leq 3\nu_1 \mathfrak{z}_1(\mathbf{x}, p^*)^2 d(\mathbf{v}, \mathbf{v}^*),$$

where $\mathfrak{z}_1(\mathbf{x}, p^*) \stackrel{\text{def}}{=} \mathbb{Q}(U(\mathbf{r}^*))$ denotes the entropy of the set $U(\mathbf{r}^*) \subset \mathbb{R}^p$ and where with $\mathfrak{g}_0 = \nu_0 \mathfrak{g}$ and for some $\mathbb{Q} > 0$

$$\mathfrak{z}_1(\mathbf{x}, \mathbb{Q}) \stackrel{\text{def}}{=} \begin{cases} \sqrt{2(\mathbf{x} + \mathbb{Q})} & \text{if } \sqrt{2(\mathbf{x} + \mathbb{Q})} \leq \mathfrak{g}_0, \\ \mathfrak{g}_0^{-1}(\mathbf{x} + \mathbb{Q}) + \mathfrak{g}_0/2 & \text{otherwise.} \end{cases}$$

To use this result let $\mathcal{Y}(\mathbf{v})$ be a smooth centered random process with values in $\mathbb{R}^{p^* \times p^*}$ and let $\mathcal{D} : \mathbb{R}^{p^*} \rightarrow \mathbb{R}^{p^*}$ be some linear operator. We aim at bounding the maximum of the spectral 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}^* . Suppose that $\mathcal{Y}(\mathbf{v})$ satisfies $\mathcal{Y}(\mathbf{v}^*) = 0$ and for each $0 < \mathbf{r} < \mathbf{r}^*$ and for all pairs $\mathbf{v}, \mathbf{v}^o \in \mathcal{Y}_o(\mathbf{r}) = \{\mathbf{v} \in \mathcal{Y} : \|\mathbf{v} - \mathbf{v}^*\|_{\mathcal{Y}} \leq \mathbf{r}\} \subset \mathbb{R}^{p^*}$

$$\sup_{\|\mathbf{u}_1\| \leq 1} \sup_{\|\mathbf{u}_2\| \leq 1} \log \mathbb{E} \exp \left\{ \lambda \frac{\mathbf{u}_1^\top (\mathcal{Y}(\mathbf{v}) - \mathcal{Y}(\mathbf{v}^o)) \mathbf{u}_2}{\omega_2 \|\mathcal{D}(\mathbf{v} - \mathbf{v}^o)\|} \right\} \leq \frac{\nu_2^2 \lambda^2}{2}. \quad (9.2)$$

Remark 9.1. In the setting of Theorem 2.4 we have $\|\mathbf{v} - \mathbf{v}^o\|_{\mathcal{Y}} = \|\mathcal{D}(\mathbf{v} - \mathbf{v}^o)\|$ and

$$\mathcal{Y}(\mathbf{v}) = \mathcal{D}^{-1} \nabla^2 \zeta(\mathbf{v}) - \mathcal{D}^{-1} \nabla^2 \zeta(\mathbf{v}^*),$$

and condition (9.2) becomes $(\mathcal{E}\mathcal{D}_2)$ from 2.1.

Theorem 9.2. Let a random process $\mathcal{Y}(\mathbf{v}) \in \mathbb{R}^{p^* \times p^*}$ fulfill $\mathcal{Y}(\mathbf{v}^*) = 0$ and let condition (9.2) be satisfied. Then for each $0 \leq \mathbf{r} \leq \mathbf{r}^*$, on a set of probability greater $1 - e^{-\mathbf{x}}$

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| \leq 9\omega_2 \nu_2 \mathfrak{z}_1(\mathbf{x}, 6p^*) \mathbf{r},$$

with $\mathfrak{g}_0 = \nu_0 \mathfrak{g}$.

Remark 9.2. Note that the entropy of the original set $\mathcal{Y}_o(\mathbf{r}) \subset \mathbb{R}^{p^*}$ is multiplied by 3. So in order to control the spectral norm $\|\mathcal{Y}(\mathbf{v})\|$ one only pays with this factor.

Proof. In what follows, we use the representation

$$\|\mathcal{Y}(\mathbf{v})\| = \omega_2 \sup_{\|\mathbf{u}_2\| \leq \mathbf{r}} \sup_{\|\mathbf{u}_1\| \leq \mathbf{r}} \frac{1}{\omega_2 \mathbf{r}^2} \mathbf{u}_1^\top \check{\mathcal{Y}}(\mathbf{v}) \mathbf{u}_2.$$

This implies

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| = \omega \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \sup_{\|\mathbf{u}_2\| \leq \mathbf{r}} \sup_{\|\mathbf{u}_1\| \leq \mathbf{r}} \frac{1}{\omega \mathbf{r}^2} \mathbf{u}_1^\top \check{\mathcal{Y}}(\mathbf{v}) \mathbf{u}_2.$$

Due to Lemma 9.3 the process $\mathcal{U}(\mathbf{v}) \stackrel{\text{def}}{=} \frac{1}{\omega \mathbf{r}^2} \mathbf{u}_1^\top \mathcal{Y}(\mathbf{v}) \mathbf{u}_2$ satisfies condition $(\mathcal{E}d)$ (see (9.1)) as process on

$$U(\mathbf{r}) \stackrel{\text{def}}{=} \mathcal{Y}_o(\mathbf{r}) \times B_{\mathbf{r}}(0) \times B_{\mathbf{r}}(0) \subset \mathbb{R}^{3p^*}. \quad (9.3)$$

This allows to apply Corollary 9.1 to obtain the desired result. We get on a set of probability greater $1 - e^{-x}$

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| \leq \sup_{(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2) \in U(\mathbf{r})} \left\{ \frac{1}{\mathbf{r}^2} \mathbf{u}_1^\top \mathcal{Y}(\mathbf{v}) \mathbf{u}_2 \right\} \leq 9\omega_2 \nu_2 \mathfrak{z}_1 \left(\mathbf{x}, \mathbb{Q}(U(\mathbf{r}^*)) \right) \mathbf{r}.$$

The constant $\mathbb{Q}(U(\mathbf{r})) > 0$ quantifies the complexity of the set $U(\mathbf{r}) \subset \mathbb{R}^{3p^*}$. We point out that for compact $M \subset \mathbb{R}^{3p^*}$ we have $\mathbb{Q}(M) = 6p^*$ (see Supplement of [14], Lemma 2.10). This gives the claim. \square

Lemma 9.3. *Suppose that $\mathcal{Y}(\mathbf{v}) \in \mathbb{R}^{p^* \times p^*}$ satisfies $\mathcal{Y}(\mathbf{v}^*) = 0$ and for each $\|\mathbf{u}_1\| \leq 1$, $\|\mathbf{u}_2\| \leq 1$ and $|\lambda| \leq \mathfrak{g}$ the inequality (9.2). Then the process*

$$\mathcal{U}(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2) = \frac{1}{2\omega_2 \mathbf{r}^2} \mathbf{u}_1^\top \mathcal{Y}(\mathbf{v}) \mathbf{u}_2$$

satisfies $(\mathcal{E}d)$ from (9.1) with $U \subset \mathbb{R}^{3p^*}$ defined in (9.3), with $|\lambda| \leq \mathfrak{g}/3$ and with

$$d((\mathbf{v}, \mathbf{u}), (\mathbf{v}^\circ, \mathbf{u}^\circ))^2 = \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|^2 + \|\mathbf{u}_1 - \mathbf{u}_1^\circ\|^2 + \|\mathbf{u}_2 - \mathbf{u}_2^\circ\|^2,$$

i.e. for any $(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2), (\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ) \in U$

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2) - \mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ)}{d((\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2), (\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ))} \right\} \leq \frac{9\nu_2^2 \lambda^2}{2}, \quad |\lambda| \leq \mathfrak{g}/3.$$

Proof. Let $(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2), (\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ) \in U$. By the Hölder inequality and (9.2), we find

$$\begin{aligned}
& \log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2) - \mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ)}{d((\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2), (\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ))} \right\} \\
&= \log \mathbb{E} \exp \left\{ \lambda \left(\frac{\mathcal{U}(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2) - \mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_1, \mathbf{u}_2)}{d((\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2), (\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ))} + \frac{\mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_1, \mathbf{u}_2) - \mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2)}{d((\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2), (\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ))} \right. \right. \\
&\quad \left. \left. + \frac{\mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2) - \mathcal{U}(\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ)}{d((\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2), (\mathbf{v}^\circ, \mathbf{u}_1^\circ, \mathbf{u}_2^\circ))} \right) \right\} \\
&\leq \frac{1}{3} \log \mathbb{E} \exp \left\{ 3\lambda \frac{\mathbf{u}_1^\top (\frac{1}{r^2} \check{\mathcal{Y}}(\mathbf{v}) - \frac{1}{r^2} \check{\mathcal{Y}}(\mathbf{v}^\circ)) \mathbf{u}_2}{\omega_2 \|\mathcal{D}(\mathbf{v} - \mathbf{v}^\circ)\|} \right\} \\
&\quad + \frac{1}{3} \log \mathbb{E} \exp \left\{ 3\lambda \frac{(\mathbf{u}_1 - \mathbf{u}_1^\circ)^\top \mathcal{Y}(\mathbf{v}^\circ) \mathbf{u}_2}{\omega_2 \|\mathbf{u}_1 - \mathbf{u}_2\| r^2} \right\} \\
&\quad + \frac{1}{3} \log \mathbb{E} \exp \left\{ 3\lambda \frac{(\mathbf{u}_1^\circ)^\top \mathcal{Y}(\mathbf{v}^\circ) (\mathbf{u}_2 - \mathbf{u}_2^\circ)}{\omega_2 \|\mathbf{u}_1 - \mathbf{u}_2\| r^2} \right\} \\
&\leq \frac{1}{3} \sup_{\|\mathbf{u}_1\| \leq 1} \sup_{\|\mathbf{u}_2\| \leq 1} \log \mathbb{E} \exp \left\{ 3\lambda \frac{\mathbf{u}_1^\top (\mathcal{Y}(\mathbf{v}) - \mathcal{Y}(\mathbf{v}^\circ)) \mathbf{u}_2}{\omega_2 \|\mathcal{D}(\mathbf{v} - \mathbf{v}^\circ)\|} \right\} \\
&\quad + \frac{2}{3} \sup_{\|\mathbf{u}_1\| \leq 1} \sup_{\|\mathbf{u}_2\| \leq 1} \log \mathbb{E} \exp \left\{ 3\lambda \frac{\mathbf{u}_1^\top (\mathcal{Y}(\mathbf{v}^\circ) - \mathcal{Y}(\mathbf{v}^*)) \mathbf{u}_2}{\omega_2 \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \right\} \\
&\leq \frac{9\nu_2^2 \lambda^2}{2}, \quad \lambda \leq g/3.
\end{aligned}$$

□

References

- [1] Andresen, A. (2014). Finite sample analysis of profile m-estimation in the single index model. *arXiv:1406.4052*.
- [2] Andresen, A. and Spokoiny, V. (2014). Critical dimension in profile semiparametric estimation. *arXiv:1303.4640*.
- [3] Balakrishnan, S., Wainwright, M. J., and Yu, B. (2014). Statistical guarantees for the em algorithm: From population to sample-based analysis. *arXiv: 1408.2156*.
- [4] Delecroix., M., Haerdle, W., and Hristache, M. (1997). Efficient estimation in single-index regression. Technical report, SFB 373, Humboldt Univ. Berlin.
- [5] Dempster, A., Laird, N., and Rubin, D. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society, Series B*, 39:1–38.
- [6] Haerdle, W., Hall, P., and Ichimura, H. (1993). Optimal smoothing in single-index models. *Ann. Statist.*, 21:157–178.
- [7] Ibragimov, I. and Khas'minskij, R. (1981). *Statistical estimation. Asymptotic theory. Transl. from the Russian by Samuel Kotz*. New York - Heidelberg -Berlin: Springer-Verlag .
- [8] Ichimura, H. (1993). Semiparametric least squares (sls) and weighted sls estimation of single- index models. *J Econometrics*, 58:71–120.
- [9] Jain, P., Netrapalli, P., and Sanghavi, S. (2013). Low-rank matrix completion using alternating minimization. *STOC*, pages 665–674.
- [10] Keshavan, R. H., Montanari, A., and Oh, S. (2010). Matrix completion from few entries. *IEEE Transactions on Information Theory*, 56(6):2980–2998.
- [11] McLachlan, G. and Krishnan, T. (1997). *The EM Algorithm and Extensions*. Wiley, New York.
- [12] Murphy, S. A. and Van der Vaart, A. W. (1999). Observed information in semi-parametric models. *Bernoulli*, 5(3):381–412.
- [13] Netrapalli, P., Jain, P., and Sanghavi, S. (2013). Phase retrieval using alternating minimization. *NIPS*, pages 2796–2804.
- [14] Spokoiny, V. (2012). Parametric estimation. Finite sample theory. *Ann. Statist.*, 40(6):2877–2909. *arXiv:1111.3029*.
- [15] Tropp, J. A. (2012). User-friendly tail bounds for sums of random matrices. *Foundations of Computation Mathematics*, 12:389–434.
- [16] Wu, C. (1983). On the convergence properties of the em algorithm. *Annals of Statistics*, 11:95–103.
- [17] Yi, X., Caramanis, C., and Sanghavi, S. (2013). Alternating minimization for mixed linear regression. *arXiv: 1310.3745*.