

Convergence of an Alternating Maximization Procedure

Andreas Andresen

ANDREAS.ANDRESEN@POSTEO.DE

*Weierstrass-Institute,
Mohrenstr. 39,
10117 Berlin, Germany*

Vladimir Spokoiny

SPOKOINY@WIAS-BERLIN.DE

*Weierstrass-Institute and Humboldt University Berlin,
Higher School of Economics, IITP RAN, MIPT Moscow,
Mohrenstr. 39, 10117 Berlin, Germany*

Editor: Jie Peng

Abstract

We derive two convergence results for a sequential alternating maximization 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 by Andresen and Spokoiny (2013), that means 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.

Keywords: alternating maximization, alternating minimization, profile maximum likelihood, EM-algorithm, M-estimation, local linear approximation, local concentration, semi-parametric

Contents

1	Introduction	3
1.1	Relation to the EM Algorithm	6
1.2	Linear Series Estimators	8
1.3	Finite sample Wilks and Fisher Theorems	9
2	Main Results	10
2.1	Conditions	11
2.1.1	Smoothness	11
2.1.2	Complexity	12
2.1.3	Moments	13
2.1.4	Conditions for the Full Model	14
2.1.5	Quadratic Drift Beats Linear Fluctuation	14
2.2	Dependence on Initial Guess	15
2.3	Introduction of Important Objects	16
2.4	Statistical Properties of the Alternating Sequence	17
2.5	Convergence to the ME	18
2.6	Critical Dimension	20
3	Application to Single Index Model	20
4	Proof of Theorem 7	23
4.1	Idea of the Proof	23
4.2	A Desirable Set	24
4.3	Probability of Desirable Set	26
4.4	Proof Convergence	30
4.5	Result after Convergence	37
5	Proof of Corollary 13	38
6	Proof of Theorem 14	38
A	Deviation Bounds for Quadratic Forms	46
B	A Uniform Bound for the Norm of a Random Process	46
C	A Bound for the Spectral Norm of a Random Matrix Process	49

1. Introduction

This paper presents two convergence results 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

$$\boldsymbol{\theta}^* = \Pi_{\boldsymbol{\theta}} \boldsymbol{v}^* = \Pi_{\boldsymbol{\theta}} \underset{\boldsymbol{v}}{\operatorname{argmax}} \mathbb{E}_{\mathbb{P}} \mathcal{L}(\boldsymbol{v}, \mathbb{Y}), \quad (1)$$

where $\mathcal{L} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ is an appropriate functional, $\Pi_{\boldsymbol{\theta}} : \mathcal{Y} \rightarrow \mathbb{R}^p$ is a projection and where \mathcal{Y} is some high dimensional or even infinite dimensional parameter space. A prominent way of estimating $\boldsymbol{\theta}^*$ is the profile M-estimator (pME)

$$\tilde{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \tilde{\boldsymbol{v}} \stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \underset{(\boldsymbol{\theta}, \boldsymbol{\eta})}{\operatorname{argmax}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}). \quad (2)$$

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 problems are computationally not feasible. This is motivated by the sieve M-estimation technique, which projects the estimation problem to a finite dimensional submodel - see Section 1.2 for details.

The alternating maximization procedure is used in situations where a direct computation of the full maximum estimator (ME) $\tilde{\boldsymbol{v}} \in \mathbb{R}^{p^*}$ 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 $\boldsymbol{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^p \times \mathbb{R}^m$ and a function basis $(\boldsymbol{e}_k) \subset L^2(\mathbb{R})$

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

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 $\tilde{\boldsymbol{v}}_0 \in \mathbb{R}^{p+m}$ set for $k \in \mathbb{N}$

$$\begin{aligned} \tilde{\boldsymbol{v}}_{k,k+1} &\stackrel{\text{def}}{=} (\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_{k+1}) = \left(\tilde{\boldsymbol{\theta}}_k, \underset{\boldsymbol{\eta} \in \mathbb{R}^m}{\operatorname{argmax}} \mathcal{L}(\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\eta}) \right), \\ \tilde{\boldsymbol{v}}_{k,k} &\stackrel{\text{def}}{=} (\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k) = \left(\underset{\boldsymbol{\theta} \in \mathbb{R}^p}{\operatorname{argmax}} \mathcal{L}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k), \tilde{\boldsymbol{\eta}}_k \right). \end{aligned} \quad (3)$$

The so called "alternating maximization procedure" (or minimization) is a widely applied algorithm in many parameter estimation tasks (see Jain, Netrapalli, and Sanghavi, 2013; Netrapalli, Jain, and Sanghavi, 2013; Keshavan, Montanari, and Oh, 2010; Yi, Caramanis, and Sanghavi, 2013). 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 $\tilde{\boldsymbol{\nu}}$? 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 again (see Jain, Netrapalli, and Sanghavi, 2013; Netrapalli, Jain, and Sanghavi, 2013; Keshavan, Montanari, and Oh, 2010; Yi, Caramanis, and Sanghavi, 2013) or variants of the procedure (Cheng, 2013).

We address this difficulty via employing new finite sample techniques (Andresen and Spokoiny, 2014; Spokoiny, 2012), which allow to answer the above questions: with growing iteration number $k \in \mathbb{N}$ the estimators $\tilde{\boldsymbol{\theta}}_k$ attain the same statistical properties as the profile M-estimator and Theorem 7 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 alternating maximization procedure can be understood as a special case of the Expectation Maximization algorithm (EM algorithm) as we illustrate in Section 1.1. The EM algorithm itself was derived in a work of Dempster, Laird, and Rubin (1977) where particular versions of this approach are generalized. This paper (Dempster, Laird, and Rubin, 1977) also contains a variety of problems where an application of the EM algorithm can be fruitful; for a brief history of the EM algorithm (see McLachlan and Krishnan, 1997, Section 1.8). We briefly explain the EM algorithm in Section 1.1.

Since the EM algorithm is very popular in applications a lot of research on its behavior 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 Wu Wu (1983), which is still state of the art in most settings. Unfortunately Wu's result - as most convergence results on these iterative procedures - only ensures convergence to some set of local maximizers or fixpoints of the procedure. Only in very special cases like unimodality can actual convergence to the maximizer be ensured.

In a recent work (Balakrishnan, Wainwright, and Yu, 2014) a new way is presented of addressing the properties of the EM sequence in a very general i.i.d. setting, based on concavity of \mathcal{L} . They assume that the functional \mathcal{L} is concave and smooth enough (First order stability) and that for a sample $(\mathbf{Y}_i)_{i=1, \dots, n}$ with high probability an uniform bound is satisfied of the kind

$$\max_{\boldsymbol{\theta}^\circ \in B_r(\boldsymbol{\theta}^*)} \left| \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^\circ}) - \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E} \mathcal{L}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_{\boldsymbol{\theta}^\circ}) \right| \leq \epsilon_n. \quad (4)$$

Under these assumptions, with high probability and some $\nu < 1$ they show

$$\|\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*\| \leq \nu^k \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\| + C\epsilon_n. \quad (5)$$

Unfortunately this does not answer our two questions to full satisfaction. First the bound (4) 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 (5) 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 $\tilde{\boldsymbol{\theta}}_k \rightarrow \tilde{\boldsymbol{\theta}}$ is not addressed.

These results 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 $(\tilde{\boldsymbol{\theta}}_k)$ actually is the profile and the statistical properties of limit points are not clear without too restrictive assumptions on \mathcal{L} and the data.

Another new work (Cheng, 2013) contains the analysis of a slightly altered algorithm in a very general semiparametric asymptotic framework. Instead of alternatingly maximizing the functional \mathcal{L} a kind of gradient decent procedure for the profile likelihood $pl(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta})$ is analyzed, i.e. they define

$$\boldsymbol{\theta}_k \stackrel{\text{def}}{=} \boldsymbol{\theta}_{k-1} + \widehat{\mathbf{D}}(\boldsymbol{\theta}_{k-1})^{-2} \ell(\boldsymbol{\theta}_{k-1}), \tag{6}$$

where $\widehat{\boldsymbol{\eta}}(\cdot)$ is an estimator of $\operatorname{argmax}_{\boldsymbol{\eta}} \mathbb{E} \mathcal{L}(\cdot, \boldsymbol{\eta})$, $\widehat{\mathbf{D}}(\cdot)$ is an estimator of $\nabla^2 \mathbb{E} \max_{\boldsymbol{\eta}} \mathcal{L}(\cdot, \boldsymbol{\eta})$ and $\ell(\cdot)$ is an estimator of $\nabla \max_{\boldsymbol{\eta}} \mathcal{L}(\cdot, \boldsymbol{\eta})$. Under common regularity conditions it is shown, that $\|\boldsymbol{\theta}_k - \tilde{\boldsymbol{\theta}}\| = o_{\mathbb{P}}(1/\sqrt{n})$ if $k(n) \in \mathbb{N}$ is chosen such that the rate of the initial guess $\boldsymbol{\theta}_0$ - obtained via a stochastic grid search - and the rate of the estimator of the nuisance parameter are addressed. These results resemble very much what is aimed for in this work but it is important to note a series of differences between the results of that work and the present paper. First and most importantly, the treated algorithm in that paper (Cheng, 2013) is similar in virtue to the alternating procedure, but in fact is a different procedure. It is a gradient descend scheme and involves a very careful data driven choice of step sizes when carrying out the estimations necessary in (6) and in that sense differs substantially from the simple and direct alternating maximization. Also the estimating step of the nuisance component is not object of the analysis but assumed to be sufficiently good for the arguments to go through. Finally the results of (Cheng, 2013) are purely asymptotic.

In this work we carry out a finite sample analysis for the alternating maximization procedure in (3). Instead of a general semiparametric framework we address sieve profile estimators also called *finite dimensional linear series estimation* (see Chen, 2007; Andersen and Spokoiny, 2014), see Section 1.2 for more details. In this setting the bias of estimation - induced by projection the full model to a finite dimensional sub model - can be treated separately and the model becomes finite dimensional as far as the algorithm is concerned. This allows a very careful and explicit analysis of the behavior of the procedure. In particular the speed of convergence can be linked to characteristics of the information matrix $-\nabla^2 \mathbb{E} \mathcal{L}(\boldsymbol{v}^*)$ - namely to the constant $\nu < 1$ in (20) - and to the the full dimension of the projected parameter space. The resulting number of iterations necessary for efficient estimation can be given in a rather simple and closed form. Finally our results are nonasymptotic which in this context is crucial as a clear comparison of the computational and the estimation error for finite samples is needed for reasonable inference.

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 $(\tilde{\boldsymbol{\theta}}_k)$ behave for large iteration number $k \in \mathbb{N}$ like the pME. To be more precise we show in Theorem 7 that if the initial guess $\tilde{\boldsymbol{v}}_0 \in \Upsilon$ is good enough the step estimator sequence $(\tilde{\boldsymbol{\theta}}_k)$ satisfies with high probability

$$\|\check{\mathbf{D}}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\|^2 \leq \epsilon(p^* + \nu^k R_0), \quad (7)$$

$$\left| \max_{\boldsymbol{\eta}} \mathcal{L}(\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\eta}) - \max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}^*, \boldsymbol{\eta}) - \|\check{\boldsymbol{\xi}}\|^2/2 \right| \leq (p + \mathbf{x})^{1/2} \epsilon(p^* + \nu^k R_0), \quad (8)$$

where $\nu < 1$ is introduced in (20) and $\epsilon > 0$ is some small number, for example $\epsilon = \mathbf{C}/\sqrt{n}$ in the smooth i.i.d setting. Further $R_0 > 0$ is a bound related to the quality of the initial guess. Generally it is proportional to the full dimension and in this way the rate with which the full nuisance can be estimated affects the speed of the convergence of the procedure. The random variable $\check{\boldsymbol{\xi}} \in \mathbb{R}^p$ and the matrix $\check{\mathbf{D}} \in \mathbb{R}^{p \times p}$ are related to the efficient influence function in semiparametric models and its covariance. These are up to $\nu^k R_0$ the same properties as those proven for the pME by Andersen and Spokoiny (2014) under nearly the same set of conditions. Up to the finite sample bounds on the right hand sides this means that the estimating points of the procedure admit a Fisher expansion - in other words are asymptotical normal - and a Wilks expansion. Consequently the usual inference procedures based on confidence and concentration sets can be applied to these estimators. Further in our second main result we manage to show under slightly stronger smoothness conditions that $(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k)$ approaches the ME $(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}) = \operatorname{argmax} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}^*)$ with nearly linear convergence speed, i.e. $\|\mathcal{D}((\boldsymbol{\theta}_k, \boldsymbol{\eta}_k) - \tilde{\boldsymbol{v}})\| \leq \tau^{k/\log(k)}$ with some $0 < \tau < 1$ and $\mathcal{D}^2 = -\mathbb{E}\nabla^2 \mathcal{L}(\boldsymbol{v}^*)$ (see Theorem 14).

To clarify we want to mention that the term convergence refers to the behavior of the sequence $(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k)$ when the number of iterations $k \in \mathbb{N}$ tends to infinity. We show $(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k) \rightarrow (\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}})$. This has to be distinguished from the usual stochastic convergence results of the M-estimator $(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}})$ towards the target $(\boldsymbol{\theta}^*, \boldsymbol{\eta}^*)$ or the weak convergence to a normal distribution as the sample size increases. Our setup is assuming finite sample size such that even with $k \rightarrow \infty$ there remains a gap between $(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k)$ and $(\boldsymbol{\theta}^*, \boldsymbol{\eta}^*)$ and between $\tilde{\boldsymbol{\theta}}_k$ and $\boldsymbol{\theta}^*$ that is related to the parametric and semiparametric Cramer-Rao lower bounds respectively. This is why we can in the finite sample setting only hope to obtain convergence of the alternating procedure to $(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}})$ but not to $(\boldsymbol{\theta}^*, \boldsymbol{\eta}^*)$. But for a growing sample size (7) implies the weak convergence results also for the estimator $\tilde{\boldsymbol{\theta}}_k$ when $k(n) \in \mathbb{N}$ is large enough and ϵp^* vanishes (see Section 1.3).

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

1.1 Relation to the EM Algorithm

In the introduction we claimed that the alternating procedure analyzed in this work is related to the EM algorithm. In this section we want to elaborate on that.

First we explain the EM algorithm. Consider data $(\mathbb{X}) \sim \mathbb{P}_\theta$ for some parametric family $(\mathbb{P}_\theta, \theta \in \Theta)$. Assume that a parameter $\theta \in \Theta$ is to be estimated as maximizer of the functional $\mathcal{L}_c(\theta, \mathbb{X}) \in \mathbb{R}$, 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 subspaces of \mathcal{X} if both \mathcal{Y}, \mathcal{X} are vector spaces. 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 two steps: starting with some initial guess $\tilde{\theta}_0$ the k th "Expectation step" derives the functional Q via

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

which means that on the right hand side the conditional expectation is calculated under the distribution \mathbb{P}_{θ_k} . The k th "Maximization step" then simply locates the maximizer θ_{k+1} of Q .

Now we can present the convergence result of Wu (1983) in more detail. Wu presents regularity conditions that ensure that $\mathcal{L}(\theta_{k+1}|\mathbb{Y}) \geq \mathcal{L}(\theta_k|\mathbb{Y})$ where

$$\mathcal{L}(\theta|\mathbb{Y}) \stackrel{\text{def}}{=} \mathbb{E}[\mathcal{L}_c(\theta, \mathbb{X})|\mathbb{Y} = f_Y(\mathbb{X})],$$

such that $\mathcal{L}(\theta_k|\mathbb{Y}) \rightarrow \mathcal{L}^*$ for some limit value $\mathcal{L}^*(\mathbb{Y}) > 0$, 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 $C(\mathcal{L}^*) \stackrel{\text{def}}{=} \{\theta | \mathcal{L}(\theta|\mathbb{Y}) = \mathcal{L}^*(\mathbb{Y})\}$. Dempster, Laird, and Rubin (1977) 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 $\mathcal{L}^*(\mathbb{Y}) = \sup \mathcal{L}(\theta|\mathbb{Y})$ 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 $\mathcal{L}(\cdot|\mathbb{Y})$ is unimodal and the regularity conditions are met but this assumption may be too restrictive.

To see that the procedure (3) is a special case of the EM algorithm we have to find the right triplet $(\mathbb{X}, f_Y, \mathcal{L}_c)$. For this we take $\mathbb{X} = (\mathbf{Z}, \mathbb{Y})$ with $\mathbf{Z} \sim \text{argmax}_\eta \mathcal{L}\{(\theta, \eta), \mathbb{Y}\}$ under \mathbb{P}_θ . Further we set $f_Y(\mathbb{X}) = \mathbb{Y}$ and $\mathcal{L}_c(\theta, \mathbf{X}) \stackrel{\text{def}}{=} \mathcal{L}(\theta, \eta, \mathbb{Y})$, where $\mathbf{X} = (\eta, \mathbb{Y})$. Then we find

$$\begin{aligned} Q(\theta, \tilde{\theta}^{(k-1)}) &= \mathbb{E}_{\tilde{\theta}^{(k-1)}}[\mathcal{L}_c(\theta, \mathbb{X})|\mathbb{Y}] \\ &= \mathbb{E}_{\tilde{\theta}^{(k-1)}}\left[\mathcal{L}_c\left(\theta, \text{argmax}_\eta \mathcal{L}\{(\tilde{\theta}^{(k-1)}, \eta), \mathbb{Y}\}, \mathbb{Y}\right)\middle|\mathbb{Y}\right] \\ &= \mathcal{L}_c\left(\theta, \text{argmax}_\eta \mathcal{L}\{(\tilde{\theta}^{(k-1)}, \eta), \mathbb{Y}\}, \mathbb{Y}\right) \\ &= \mathcal{L}(\theta, \tilde{\eta}^{(k)}, \mathbb{Y}), \end{aligned}$$

and thus the resulting sequence is the same as in (3).

1.2 Linear Series Estimators

In semiparametric models the profile M estimator $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^p$ from Equation (2) cannot be calculated in practice if the full model is infinite dimensional. There are various ways to circumvent this problem. Next to non parametric estimation and plugin of the nuisance $\boldsymbol{\eta} \in \mathcal{X}$ a prominent approach is the so called sieve technique that motivates the setting in this work.

The sieve approach was systematically introduced by Grenander (see Grenander, 1981, Chapter 8) and consists in choosing a suitable sequence of subsets $(\mathcal{Y}_m)_{m=1}^\infty \subset \mathcal{Y}$ such that for each $\mathbf{v} \in \mathcal{Y}$ there exists a sequence $\Pi_m(\mathbf{v}) \subset \mathcal{Y}_m$ with $\|\mathbf{v} - \Pi_m(\mathbf{v})\| \rightarrow 0$. Furthermore the sets $\mathcal{Y}_m \subset \mathcal{Y}$ have to be such that $\sup_{\mathbf{v} \in \mathcal{Y}_m} \mathcal{L}(\mathbf{v})$ can be calculated in practice. In the setting of semiparametric M-Estimation we assume $\mathcal{Y} = \mathcal{Y}_\boldsymbol{\theta} \times \mathcal{Y}_\boldsymbol{\eta} \subseteq \mathbb{R}^p \times \mathcal{X}$ with some infinite dimensional separable Hilbert space \mathcal{X} and countable basis $\{\mathbf{e}_1, \mathbf{e}_2, \dots\} \subset \mathcal{X}$. We set $\mathcal{Y}_m = \mathcal{Y}_\boldsymbol{\theta} \times \Pi_m \mathcal{Y}_\boldsymbol{\eta}$, where $\Pi_m : \mathcal{X} \rightarrow \mathcal{X}_m$ denotes the orthogonal projection onto $\mathcal{X}_m \stackrel{\text{def}}{=} \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_m)$. For each $m \in \mathbb{N}$ the sieve profile M-estimator is defined as

$$\tilde{\boldsymbol{\theta}}_m \stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \tilde{\mathbf{v}}_m \stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \underset{\substack{\boldsymbol{\theta} \in \mathbb{R}^p \\ \boldsymbol{\eta} \in \mathbb{R}^m}}{\text{argmax}} \mathcal{L} \left(\boldsymbol{\theta}, \sum_{k=1}^m \eta_k \mathbf{e}_k \right).$$

This means that for the calculation of the estimator $\tilde{\boldsymbol{\theta}}_m$ only a finite dimensional setting has to be considered. In our analysis we will focus on the behavior of the alternating maximization procedure in that case.

But of course the projection onto a finite dimensional submodel induces an approximation bias ” $\mathbf{v}^* - \mathbf{v}_m^*$ ” where

$$(\boldsymbol{\theta}_m^*, \boldsymbol{\eta}_m^*) \stackrel{\text{def}}{=} \mathbf{v}_m^* \stackrel{\text{def}}{=} \underset{\substack{\boldsymbol{\theta} \in \mathbb{R}^p \\ \boldsymbol{\eta} \in \mathbb{R}^m}}{\text{argmax}} \mathbb{E} \mathcal{L} \left(\boldsymbol{\theta}, \sum_{k=1}^m \eta_k \mathbf{e}_k \right).$$

In (Andresen and Spokoiny, 2014) it is explained in detail how this bias can be treated. Once the bias is controlled this leads for each $m \in \mathbb{N}$ to bounds of the kind

$$\begin{aligned} \left\| \check{\mathbf{D}}(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}_m \right\| &\leq \check{\diamond}(\mathbf{x}) + \alpha(m), \\ \left| \max_{\boldsymbol{\eta} \in \Pi_m \mathcal{Y}_\boldsymbol{\eta}} \mathcal{L}(\tilde{\boldsymbol{\theta}}_m, \boldsymbol{\eta}) - \max_{\boldsymbol{\eta} \in \Pi_m \mathcal{Y}_\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}^*, \boldsymbol{\eta}) - \|\check{\boldsymbol{\xi}}_m\|^2 \right| &\leq p \check{\diamond}(\mathbf{x}) + \alpha(m), \end{aligned}$$

where $\alpha(m) \geq 0$ quantifies the impact of the bias ” $\mathbf{v}^* - \mathbf{v}_m^*$ ”. The choice of $m \in \mathbb{N}$ then has to balance the two terms $\check{\diamond}(\mathbf{x})$ and $\alpha(m)$, which leads to common optimal choices for the dimension based on the ”smoothness” of the nuisance component $\boldsymbol{\eta}^* \in \mathcal{X}$. To ease notation we drop the \cdot_m in the following, as the treatment of the bias can be done separately, (see Andresen and Spokoiny, 2014).

1.3 Finite sample Wilks and Fisher Theorems

Before we present our main results we want to explain what type of results we aim at and how they can be interpreted. Hopefully this will ease the understanding and will make some of the apparently cumbersome notation more intelligible.

Usually in asymptotic treatments of semiparametric M-estimators like $\tilde{\boldsymbol{\theta}}$ in (2) the aim is to derive statements of the kind

$$\sqrt{n}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \check{d}^{-1}\check{\boldsymbol{\xi}} = o_{\mathbb{P}}(1), \quad (9)$$

$$\max_{\boldsymbol{\eta}} \mathcal{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\eta}) - \max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}^*, \boldsymbol{\eta}) - \|\check{\boldsymbol{\xi}}\|^2/2 = o_{\mathbb{P}}(1), \quad (10)$$

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

where $n \in \mathbb{N}$ denotes the sample size. The random variable $\check{\boldsymbol{\xi}} \in \mathbb{R}^p$ is called semiparametric score. Below we will briefly explain its derivation along with the explanation of the matrices $\check{v}^2, \check{d}^2 \in \mathbb{R}^{p \times p}$. But before, we sketch how (9) and (10) can be used for the construction of asymptotic confidence sets that yield statistical tests. Given the matrices \check{v}^2, \check{d}^2 the construction works as follows. Let $q_{\alpha}^2 > 0$ be an α -level quantile of a $\chi_p^2(\check{d}^{-2}\check{v}^2\check{d}^{-2})$ -distribution. Set

$$\mathcal{E}(q_{\alpha}) = \{\boldsymbol{\theta} : \sqrt{n}\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}\| \leq q_{\alpha}\}; \quad (11)$$

then one can use (9) to show

$$\mathbb{P}\{\boldsymbol{\theta}^* \notin \mathcal{E}(q_{\alpha})\} = \mathbb{P}\{\sqrt{n}\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\| \geq q_{\alpha}\} \rightarrow 1 - \alpha.$$

Similarly one can exploit (10).

Now we explain the definition of \check{v}^2 and \check{d}^2 . Introduce

$$\begin{aligned} n\check{v}^{-2}(\mathbf{v}) &\stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \text{Cov}(\nabla \mathcal{L}(\mathbf{v}))^{-1} \Pi_{\boldsymbol{\theta}}^{\top}, \\ n\check{d}^{-2}(\mathbf{v}) &\stackrel{\text{def}}{=} -\Pi_{\boldsymbol{\theta}} (\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}))^{-1} \Pi_{\boldsymbol{\theta}}^{\top}, \end{aligned}$$

where $\Pi_{\boldsymbol{\theta}}$ is the orthogonal projection onto the $\boldsymbol{\theta}$ -component in \mathbb{R}^p and $\Pi_{\boldsymbol{\theta}}^{\top}$ is its dual operator. Then $\check{v}^2 = \check{v}^2(\mathbf{v}^*)$ and $\check{d}^2 = \check{d}^2(\mathbf{v}^*)$.

Remark 1 *Note that these two matrices coincide if the functional \mathcal{L} was the complete loglikelihood of the observations and that then \check{d}^2 would equal the covariance of the efficient influence function (see Kosorok, 2005, for more details).*

For the definition of the semiparametric score $\check{\boldsymbol{\xi}} \in \mathbb{R}^p$ consider

$$nd^2(\mathbf{v}) = \begin{pmatrix} d^2(\mathbf{v}) & a(\mathbf{v}) \\ a^{\top}(\mathbf{v}) & h^2(\mathbf{v}) \end{pmatrix} \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}) \in \mathbb{R}^{p^* \times p^*}.$$

Then

$$\begin{aligned}\check{\xi} &\stackrel{\text{def}}{=} \frac{1}{\sqrt{n}} \check{d}(\mathbf{v}^*) (1 - \mathbb{E}_\epsilon) \Pi_{\boldsymbol{\theta}} d^{-2}(\mathbf{v}^*) \nabla \mathcal{L}(\mathbf{v}^*) \\ &= \frac{1}{\sqrt{n}} (1 - \mathbb{E}_\epsilon) \check{d}^{-1}(\mathbf{v}^*) \{ \nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{v}^*) - ah^{-2}(\mathbf{v}^*) \nabla_{\boldsymbol{\eta}} \mathcal{L}(\mathbf{v}^*) \},\end{aligned}$$

This random variable is related to the efficient influence function in semiparametric estimation and it plays the role that the usual score $\nabla \mathcal{L}(\mathbf{v}^*)$ plays in the setting of parametric M-estimation.

In this work we derive (9) and (10) for $\tilde{\boldsymbol{\theta}}_k$ instead of $\tilde{\boldsymbol{\theta}}$ but with finite sample bounds for the terms on the right-hand sides of (9) and (10). To be more precise we derive statements of the following kind. With probability greater than $1 - \mathbb{C}e^{-x}$

$$\left\| \sqrt{n}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{d}^{-1} \check{\xi} \right\| \leq \epsilon(p^* + \nu^k R_0), \quad (12)$$

$$\left| \max_{\boldsymbol{\eta}} \mathcal{L}(\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\eta}) - \max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}^*, \boldsymbol{\eta}) - \|\check{\xi}\|^2/2 \right| = (p + \mathbf{x})^{1/2} \epsilon(p^* + \nu^k R_0), \quad (13)$$

with some small value $\epsilon > 0$ as in (7) and (8). Note that for vanishing right hand sides these equations imply (9) and (10) when n tends to infity. Using the scheme in (11) the bounds (12) and (13) allow the construction of (conservative) finite sample "confidence sets". Assume that (approximate) quantiles q_α for $\|\check{\xi}\|$ are available, i.e. that with some small $\epsilon > 0$ and any $\alpha \in [0, 1]$

$$\mathbb{P}(\|\check{\xi}\| \leq q_\alpha) \in (\alpha - \delta, \alpha + \delta),$$

then with some generic constant $\mathbb{C} > 0$ (see Andresen and Spokoiny, 2014, Remark 2.13)

$$\begin{aligned}\alpha + \delta + \mathbb{C}e^{-x} &\leq \mathbb{P} \left\{ \boldsymbol{\theta}^* \in \mathcal{E}(q_\alpha + \epsilon(p^* + \nu^k R_0)) \right\}, \\ \mathbb{P} \left\{ \boldsymbol{\theta}^* \in \mathcal{E}(q_\alpha - \epsilon(p^* + \nu^k R_0)) \right\} &\leq \alpha - \delta - \mathbb{C}e^{-x}.\end{aligned}$$

The important achievement is that one can make approximate confidence statements for the estimators of the alternating procedure and this even in the finite sample case, without ignoring "hopefully small enough" terms. As remarked above such approximate quantiles could be attained via an plug-in-estimation of $\check{d}^{-2} \check{v}^2 \check{d}^{-2}$ combined with a Gaussian approximation or a bootstrap.

2. Main Results

This Section contains the thorough presentation of our main convergence results. It involves the introduction of various technicalities and may appear a bit cumbersome on first read. We recommend to carefully read Section 1.3 first to ease understanding.

2.1 Conditions

This section collects the conditions imposed on the model. We use the same set of assumptions as Andresen and Spokoiny (2014) and this section closely follows Section 2.1 of that paper.

Let the full dimension of the parameter space be finite, i.e. $p^* < \infty$. Our conditions involve the symmetric positive definite *information matrix* $\mathcal{D}_0^2 \in \mathbb{R}^{p^* \times p^*}$ and a central point $\mathbf{v}^* \in \mathbb{R}^{p^*}$. To ease presentation in this paper we identify \mathbf{v}^* with the “true point” from (1) and define

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

where we assume, that the second derivative exists. In the context of semiparametric estimation, it is convenient to represent the *information matrix* in block form:

$$\mathcal{D}_0^2 = \begin{pmatrix} \mathcal{D}_0^2 & \mathbf{A}_0 \\ \mathbf{A}_0^\top & \mathbf{H}_0^2 \end{pmatrix}.$$

First we state an *identifiability condition*, which basically imposes that \mathcal{D}^2 is positive definite. Note that in this work $\|\cdot\|$ always denotes the spectral norm when its argument is a matrix.

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

$$\|\mathbf{H}_0^{-1} \mathbf{A}_0^\top \mathcal{D}_0^{-1}\| \leq \sqrt{\nu}. \quad (14)$$

The condition (\mathcal{I}) allows to introduce the important $p \times p$ efficient information matrix $\check{\mathcal{D}}_0^2$ which is defined as the inverse of the $\boldsymbol{\theta}$ -block of the inverse of the full dimensional matrix \mathcal{D}_0^2 . The exact formula is given by

$$\check{\mathcal{D}}_0^2 \stackrel{\text{def}}{=} \mathcal{D}_0^2 - \mathbf{A}_0 \mathbf{H}^{-2} \mathbf{A}_0^\top,$$

and (\mathcal{I}) ensures that the matrix $\check{\mathcal{D}}_0^2$ is positive definite such that $\check{\mathcal{D}}$ is well defined. At the same time $\nu < 1$ ensures that the alternating sequence actually converges. As can be seen in Theorem 7 the speed of convergence is linear in ν .

Using the matrix \mathcal{D}_0^2 and the central point $\mathbf{v}^* \in \mathbb{R}^{p^*}$, we define the local set $\mathcal{Y}_o(\mathbf{r}) \subset \mathcal{Y} \subseteq \mathbb{R}^{p^*}$ with some $\mathbf{r} \geq 0$:

$$\mathcal{Y}_o(\mathbf{r}) \stackrel{\text{def}}{=} \{\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathcal{Y} : \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}\}. \quad (15)$$

2.1.1 SMOOTHNESS

Usually in the context of regular M-estimation one assumes local quadraticity, i.e. that

$$pl(\boldsymbol{\theta}) = \nabla pl(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta}^* - \boldsymbol{\theta}^*)^\perp \nabla^2 pl(\boldsymbol{\theta}^*)(\boldsymbol{\theta}^* - \boldsymbol{\theta}^*) + o_{\mathbb{P}}(1),$$

for $\boldsymbol{\theta} \in \mathbb{R}^p$ close enough to $\boldsymbol{\theta}^*$. The functional $pl(\cdot)$ in this context denotes the profile-functional and is defined as

$$pl(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\eta}),$$

see for instance (Cheng, 2013).

In our setting we need a precise bound for the accuracy of a quadratic approximation of the expected profile functional $\mathbb{E}pl$. We want to bound the error of a local linear approximation of the projected gradient $\check{\nabla}_{\boldsymbol{\theta}} \mathbb{E}\mathcal{L}(\boldsymbol{v})$ which is defined as

$$\check{\nabla}_{\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}} - A_0 H_0^{-2} \nabla_{\boldsymbol{\eta}}.$$

Instead of local quadraticity of the profile functional we impose that for $(\boldsymbol{\theta}, \boldsymbol{\eta}) = \boldsymbol{v} \in \mathcal{Y}_o(\mathbf{r})$ - with the local set $\mathcal{Y}_o(\mathbf{r})$ defined in (15) - with some small $\check{\epsilon} > 0$

$$\left\| \check{\mathcal{D}}^{-1} \left(\check{\nabla} \mathbb{E}\mathcal{L}(\boldsymbol{v}) - \check{\nabla} \mathbb{E}\mathcal{L}(\boldsymbol{v}^*) \right) - \check{\mathcal{D}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right\| \leq \check{\epsilon} \mathbf{r}^2.$$

The following condition serves a less high level way of checking that such an approximation holds (see Andresen and Spokoiny, 2014, Lemma B.1)

($\check{\mathcal{L}}_0$) For each $\mathbf{r} \leq \mathbf{r}_0$, there is a constant $\check{\epsilon} > 0$ such that it holds on the set $\mathcal{Y}_o(\mathbf{r})$:

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

If $\mathbb{E}\mathcal{L}$ is three times continuously differentiable one obtains $\check{\epsilon} \leq \mathbf{C} \|\mathcal{D}^{-1}\|$. In i.i.d. models one usually has $\|\mathcal{D}^{-1}\| = O(1/\sqrt{n})$.

Remark 2 Here and in what follows we implicitly assume that the function $\mathcal{L}(\boldsymbol{v}): \mathbb{R}^{p^*} \rightarrow \mathbb{R}$ is sufficiently smooth in $\boldsymbol{v} \in \mathbb{R}^{p^*}$, $\nabla \mathcal{L}(\boldsymbol{v}) \in \mathbb{R}^{p^*}$ stands for the gradient and $\nabla^2 \mathbb{E}\mathcal{L}(\boldsymbol{v}) \in \mathbb{R}^{p^* \times p^*}$ for the Hessian of the expectation $\mathbb{E}\mathcal{L}: \mathbb{R}^{p^*} \rightarrow \mathbb{R}$ at $\boldsymbol{v} \in \mathbb{R}^{p^*}$. By smooth enough we mean that we can interchange $\nabla \mathbb{E}\mathcal{L} = \mathbb{E} \nabla \mathcal{L}$ on $\mathcal{Y}_o(R_0)$, where $\mathcal{Y}_o(\mathbf{r})$ is defined in (15) and $R_0 > 0$ in (19). It is worth mentioning that $\mathcal{D}_0^2 = \mathcal{V}_0^2 \stackrel{\text{def}}{=} \text{Cov}(\nabla \mathcal{L}(\boldsymbol{v}^*))$ if the model $\mathbf{Y} \sim \mathbb{P}_{\boldsymbol{v}^*} \in (\mathbb{P}_{\boldsymbol{v}})$ is correctly specified and sufficiently regular; see e.g. (Ibragimov and Khas'minskij, 1981).

2.1.2 COMPLEXITY

The usual approach to gain asymptotic control on profile M-estimators would be to assume that the class

$$\{\check{\nabla} \mathcal{L}(\boldsymbol{v}), \boldsymbol{v} \in \mathcal{Y}_o(\mathbf{r})\},$$

is \mathbb{P} -Donsker (see Cheng, 2013).

To pin down the estimator sequence $(\tilde{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\eta}}_k)_{k \in \mathbb{N}}$ and to obtain finite sample results we use a more specific approach, which is based on a new finite sample approach (Spokoiny,

2012; Andresen and Spokoiny, 2014). First note that we assume that - as far as the alternating procedure is concerned - the model is finite dimensional, which means that $\mathcal{Y}_\circ(\mathbf{r}) \subset \mathcal{Y} \subset \mathbb{R}^{p^*}$ in (15) is automatically compact for finite radius $\mathbf{r} > 0$. If we can ensure the right smoothness and moment conditions on $\check{\nabla}\mathcal{L}(\mathbf{v})$, we automatically obtain that the above class is \mathbb{P} -Donsker. But using the new techniques (Spokoiny, 2012; Andresen and Spokoiny, 2014) we manage to obtain slightly stronger bounds that are useful in a finite sample setting.

To understand the next condition consider first the definition of a subgaussian random vector. A random vector $\mathbf{X} \in \mathbb{R}^p$ is called subgaussian, if for any $\mu \in \mathbb{R}$ and some $\nu > 0$

$$\sup_{\|\gamma\| \leq 1} \log \mathbb{E} \exp \left\{ \mu \gamma^\top \mathbf{X} \right\} \leq \nu^2 \mu^2 / 2. \quad (16)$$

Obviously this is a strong condition. As it turns out in many situations it is sufficient to assume subexponentiality, which simply relaxes subgaussianity to demanding that (16) is met for any $|\mu| \leq \mathbf{g}$ with some $\mathbf{g} > 0$. In this way many distributions that would be excluded by assuming subgaussianity can still be treated.

Our next condition combines subexponentiality with a smoothness constraint on the stochastic component of $\check{\nabla}\mathcal{L}(\mathbf{v})$. It assumes that - with some $\check{\epsilon} > 0$ - the random vector

$$\frac{1}{\check{\epsilon} \|\mathcal{D}(\mathbf{v} - \mathbf{v}')\|} \check{\mathbb{D}}^{-1} \{ \check{\nabla}_\theta \zeta(\mathbf{v}) - \check{\nabla}_\theta \zeta(\mathbf{v}') \} \in \mathbb{R}^p,$$

is uniformly subexponential for $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_\circ(\mathbf{r})$, with the stochastic component defined as $\zeta(\mathbf{v}) \stackrel{\text{def}}{=} \mathcal{L}(\mathbf{v}) - \mathbb{E}\mathcal{L}(\mathbf{v})$. It reads:

($\check{\mathcal{E}}\mathcal{D}_1$) For all $0 < \mathbf{r} < \mathbf{r}_0$, there exists a constant $\check{\epsilon} \leq 1/2$ such that for all $|\mu| \leq \check{\mathbf{g}}$ and $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_\circ(\mathbf{r})$

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

To convey more intuition consider for some pair $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_\circ(\mathbf{r})$ the function

$$\psi_\gamma(t) \stackrel{\text{def}}{=} \gamma^\top \check{\mathbb{D}}^{-1} \{ \check{\nabla}_\theta \zeta(\mathbf{v}) - \check{\nabla}_\theta \zeta(\mathbf{v} + t(\mathbf{v}' - \mathbf{v})) \}.$$

Then ($\check{\mathcal{E}}\mathcal{D}_1$) is met with $\check{\epsilon} \leq \|\mathcal{D}^{-1}\|$, if - uniformly in γ - for any pair $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_\circ(\mathbf{r})$ the corresponding function $\psi_\gamma : [0, 1] \rightarrow \mathbb{R}$ is Lipschitz continuous and the Lipschitz constant a subexponential random variable.

2.1.3 MOMENTS

We need another condition that allows to control the deviation behavior of $\|\check{\mathbb{D}}^{-1} \check{\nabla} \zeta(\mathbf{v}^*)\|$. To present this condition define the covariance matrix $\mathcal{V}_0^2 \in \mathbb{R}^{p^* \times p^*}$ and $\check{V}^2 \in \mathbb{R}^{p \times p}$

$$\mathcal{V}_0^2 \stackrel{\text{def}}{=} \text{Cov} \{ \nabla \mathcal{L}(\mathbf{v}^*) \}, \quad \check{V}^2 = \text{Cov}(\check{\nabla}_\theta \zeta(\mathbf{v}^*)).$$

We impose subexponential moments on $\check{V}^{-1} \check{\nabla}_\theta \zeta(\mathbf{v}^*)$:

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

$$\sup_{\|\gamma\| \leq 1} \log \mathbb{E} \exp \left\{ \mu \gamma^\top \check{V}^{-1} \check{\nabla}_{\theta} \zeta(\mathbf{v}^*) \right\} \leq \frac{\check{\nu}_0^2 \mu^2}{2}.$$

2.1.4 CONDITIONS FOR THE FULL MODEL

So far we only presented conditions that allow to treat the properties of $\tilde{\theta}_k$ on local sets $\mathcal{Y}_o(\mathbf{r}_k)$, for some sequence $(\mathbf{r}_k)_{k \in \mathbb{N}}$. To show that the sequence of estimators $(\mathbf{v}_k)_{k \in \mathbb{N}}$ satisfies $\mathbf{v}_k \in \mathcal{Y}_o(\mathbf{r}_k)$ for an appropriately decreasing sequence $(\mathbf{r}_k)_{k \in \mathbb{N}}$ the following, stronger conditions are employed, which can be interpreted just as the previous ones.

(\mathcal{L}_0) For each $\mathbf{r} \leq \mathbf{r}_0$, there is a constant $\epsilon > 0$ such that it holds on the set $\mathcal{Y}_o(\mathbf{r})$:

$$\|\mathcal{D}_0^{-1} \{ \nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}) \} \mathcal{D}_0^{-1} - I_{p^*}\| \leq \epsilon \mathbf{r}.$$

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

$$\sup_{\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_o(\mathbf{r})} \sup_{\|\gamma\|=1} \log \mathbb{E} \exp \left\{ \frac{\mu \gamma^\top \mathcal{D}_0^{-1} \{ \nabla \zeta(\mathbf{v}) - \nabla \zeta(\mathbf{v}') \}}{\epsilon \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}')\|} \right\} \leq \frac{\nu_1^2 \mu^2}{2}.$$

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

$$\sup_{\|\gamma\| \leq 1} \log \mathbb{E} \exp \left\{ \mu \gamma^\top \mathcal{V}_0^{-1} \nabla \zeta(\mathbf{v}^*) \right\} \leq \frac{\nu_0^2 \mu^2}{2}.$$

It is important to note, that the constants $\check{\epsilon}, \check{\nu}$ and ϵ, ν 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 ($\mathcal{A}H^{-2} \nabla_{\eta} \mathcal{L}$ might be quite smooth while $\nabla_{\eta} \mathcal{L}$ could be less regular).

For the convergence statement in Theorem 14 we additionally need the following condition, that controls the moments and the smoothness of the process $\nabla^2(\mathcal{L} - \mathbb{E}\mathcal{L})$:

($\mathcal{E}\mathcal{D}_2$) There exists a constant $\epsilon_2 \leq 1/2$, such that for all $|\mu| \leq g$ and all $0 < \mathbf{r} < \mathbf{r}_0$

$$\sup_{\mathbf{v}, \mathbf{v}' \in \mathcal{Y}_o(\mathbf{r})} \sup_{\|\gamma_1\|=1} \sup_{\|\gamma_2\|=1} \log \mathbb{E} \exp \left\{ \frac{\mu \gamma_1^\top \mathcal{D}^{-1} \{ \nabla^2 \zeta(\mathbf{v}) - \nabla^2 \zeta(\mathbf{v}') \} \gamma_2}{\epsilon_2 \|\mathcal{D}(\mathbf{v} - \mathbf{v}')\|} \right\} \leq \frac{\nu_2^2 \mu^2}{2}.$$

2.1.5 QUADRATIC DRIFT BEATS LINEAR FLUCTUATION

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 $\mathcal{Y} \subseteq \mathbb{R}^{p^*}$.

The first condition imposes that $\mathbb{E}\mathcal{L}(\mathbf{v})$ decreases nearly quadratically as the distance of \mathbf{v} to \mathbf{v}^* grows.

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

$$\mathbb{E}[\mathcal{L}(\mathbf{v}) - \mathcal{L}(\mathbf{v}^*)] \leq \mathbf{b} \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2.$$

The next condition bounds moments of the full gradient $\nabla \mathcal{L}(\mathbf{v})$ - again via subexponentiality.

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

$$\sup_{\mathbf{v} \in \mathcal{Y}_0(\mathbf{r})} \sup_{\mu \leq \mathbf{g}(\mathbf{r})} \sup_{\|\gamma\| \leq 1} \log \mathbb{E} \exp \left\{ \mu \gamma^\top \mathcal{D}^{-1} \nabla \zeta(\mathbf{v}) \right\} \leq \frac{\nu_{\mathbf{r}}^2 \mu^2}{2}.$$

We impose one further merely technical condition:

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

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

Remark 3 Without this the calculation of $R_0(\mathbf{x})$ in Section 4.3 would become technically more involved, without that further insight would be gained.

Remark 4 The condition ($\mathcal{E}\mathbf{r}$) can be substantially relaxed to $\mathbf{b} = \mathbf{b}(\mathbf{r}) > 0$ that decreases to 0 as $\mathbf{r} \rightarrow \infty$. We avoid the resulting technicalities and refer the reader to the original publication for the non constant case (see Spokoiny, 2012, Theorem 4.1).

2.2 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):

(\mathbf{A}_1) With probability greater $1 - \beta$ the initial guess satisfies $\mathcal{L}(\tilde{\mathbf{v}}_0) - \mathcal{L}(\mathbf{v}^*) \geq -K_0$ for some $K_0 \geq 0$.

(\mathbf{A}_2) 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 $\mathbf{r} \leq R_0(\mathbf{x})$ where $R_0(\mathbf{x})$ can be bounded with (see (19))

$$R_0(\mathbf{x}) \leq \mathbf{C} \sqrt{\mathbf{x} + p^* + K_0}.$$

(\mathbf{A}_3) $K_0 \in \mathbb{R}$ and $\epsilon > 0$ are small enough to ensure

$$\epsilon \mathbf{C}(\nu) R_0 < 1, \tag{17}$$

with

$$\mathbf{C}(\nu) \stackrel{\text{def}}{=} \frac{16\sqrt{2}(1 + \sqrt{\nu})}{(1 - \nu)(1 - \sqrt{\nu})}, \tag{18}$$

where $\nu > 0$ is defined in (14).

Condition (A_1) allows to concentrate the analysis on a local set $\{\|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \leq R_0(\mathbf{x})\} \subset \mathcal{T}$ with dominating probability (see Theorem 28). Conditions (A_2) and (A_3) ensure that this neighborhood is small enough to imply convergence of the procedure. They impose a bound on $R_0(\mathbf{x})$ and thus on K_0 from (A_1) . These conditions boil down to $\epsilon\sqrt{K_0}$ being significantly smaller than 1, which is a quantification of the quality of the first guess. There are numerous ways to initiate the procedure. In Section 3 we use a grid search and show that for a sufficiently fine grid these conditions can be met in the treated model.

Remark 5 *One way of obtaining condition (A_1) is to show that $\|\mathcal{D}(\tilde{\mathbf{v}} - \mathbf{v}^*)\| \leq R$ with probability greater $1 - \beta$ for some finite $R \in \mathbb{R}$ and $0 \leq \beta(R) < 1$. Then one can use - with some constant $\mathbf{C} > 0$ -*

$$K_0 \leq (1/2 + \epsilon(1 + 12\nu_0))(R + \mathbf{C}\sqrt{\mathbf{x} + p^*})^2,$$

as we show in Lemma 31.

Remark 6 *The precise definition of $R_0(\mathbf{x}) > 0$ reads*

$$R_0(\mathbf{x}) \stackrel{\text{def}}{=} \mathfrak{z}(\mathbf{x}) \vee \frac{6\nu_{\mathbf{r}}}{\mathbf{b}(1 - \nu)} \sqrt{\mathbf{x} + 2.4p^* + \frac{\mathbf{b}^2}{9\nu_{\mathbf{r}}^2} K_0}, \quad (19)$$

with the term

$$\mathfrak{z}(\mathbf{x}) \approx \mathbf{C}\sqrt{p^* + \mathbf{x}},$$

which is defined in (30).

2.3 Introduction of Important Objects

In this section we collect the most important objects and bounds that are relevant for Theorem 7. Remember the $p^* \times p^*$ information matrix \mathcal{D}^2 from Section 2.1, which is defined similarly to the Fisher information matrix:

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

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

$$\|\mathcal{D}^{-1} \mathcal{A} \mathcal{H}^{-1}\|^2 \stackrel{\text{def}}{=} \nu, \quad (20)$$

which we assume with condition (\mathcal{I}) to be smaller 1. It determines the speed of convergence of the alternating procedure (see Theorem 7).

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

$$\begin{aligned} \check{\mathcal{D}}^2 &= \mathcal{D}^2 - \mathcal{A} \mathcal{H}^{-2} \mathcal{A}^\top, \\ \check{\boldsymbol{\xi}} &= \check{\mathcal{D}}^{-1} \check{\nabla}_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{v}^*), \quad \check{\nabla}_{\boldsymbol{\theta}} \mathcal{L} = \nabla_{\boldsymbol{\theta}} \mathcal{L} - \mathcal{A} \mathcal{H}^{-2} \nabla_{\boldsymbol{\eta}} \mathcal{L}. \end{aligned}$$

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.

We define the *semiparametric uniform spread*

$$\check{\diamond}_Q(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \check{\epsilon} \left\{ \frac{16}{(1-\nu^2)^2} \mathbf{r}^2 + \mathfrak{z}_Q(\mathbf{x}, 2p^* + 2p)^2 \right\} (1 + 6\nu_1^2). \quad (21)$$

This object is central for our analysis as it describes the accuracy of our main results. It is small if $\check{\epsilon}(\mathbf{r}^2 + \mathbf{x} + p^*)$ is small, since $\mathfrak{z}_Q(\mathbf{x}, p^*) \approx \sqrt{p^* + \mathbf{x}}$ (see its definition in Equation (42)).

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 7 *Assume that the conditions $(\mathcal{L}_{\mathbf{r}})$, $(\mathcal{E}_{\mathbf{r}})$ and (B_1) of Section 2.1 are met. Further assume (A_1) , (A_2) and (A_3) of Section 2.2. Then it holds with probability greater $1 - 8e^{-\mathbf{x}} - \beta$ for all $k \in \mathbb{N}$*

$$\|\check{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\xi}\| \leq \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}), \quad (22)$$

$$\left| \max_{\boldsymbol{\eta}} \mathcal{L}(\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\eta}) - \max_{\boldsymbol{\eta}} \mathcal{L}(\boldsymbol{\theta}^*, \boldsymbol{\eta}) - \|\check{\xi}\|^2/2 \right| \leq 5 \left(\|\check{\xi}\| + \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}) \right) \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}), \quad (23)$$

where

$$\mathbf{r}_k \leq \mathbf{C} \left(\sqrt{p^* + \mathbf{x}} + \nu^k R_0 \right),$$

with a constant \mathbf{C} that depends on $\nu < 1$ and $1 - \mathbf{C}(\nu)\epsilon R_0 > 0$. In particular this means that if

$$k \geq \frac{\log(p^* + \mathbf{x}) - \log\{R_0\}}{\log(\nu)},$$

we have

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

Remark 8 *Note that with linear convergence speed this leads to statements about $\tilde{\boldsymbol{\theta}}_k$ that are very similar to those in (Andresen and Spokoiny, 2014) for the profile M estimator $\tilde{\boldsymbol{\theta}}$.*

Remark 9 *Concerning the properties of $\check{\xi} \in \mathbb{R}^p$ we repeat remark 2.1 of (Andresen and Spokoiny, 2014). In case of correct model specification the deviation properties of the quadratic form $\|\check{\xi}\|^2 = \|\check{D}^{-1}\check{\nabla}_{\boldsymbol{\theta}}\|^2$ are essentially the same as those of a chi-square random variable with p degrees of freedom; see Theorem 39 in the appendix. In the case of a possible*

model misspecification the behavior of the quadratic form $\|\check{\xi}\|^2$ will depend on the characteristics of the matrix $\check{B} \stackrel{\text{def}}{=} \text{Cov}(\check{\xi})$; see again Theorem 39. Moreover, in the asymptotic setup the vector $\check{\xi}$ is asymptotically standard normal; see Section 2.2. of (Andresen and Spokoiny, 2014) for the i.i.d. case.

Remark 10 These results allow to derive some important corollaries like concentration and confidence sets (see Section 1.3).

Remark 11 In general an exact numerical computation of

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

is not possible. Define $\hat{\theta}(\boldsymbol{\eta})$ and $\hat{\eta}(\boldsymbol{\theta})$ as the numerical approximations to $\theta(\boldsymbol{\eta})$ and $\eta(\boldsymbol{\theta})$ and assume that - with the local set $\mathcal{Y}_\circ(\mathbf{x})$ defined in (15) -

$$\|D(\hat{\theta}(\boldsymbol{\eta}) - \theta(\boldsymbol{\eta}))\| \leq \tau, \quad \text{for all } \boldsymbol{\eta} \in \mathcal{Y}_{\circ, \boldsymbol{\eta}}(R_0) \stackrel{\text{def}}{=} \{\boldsymbol{v} \in \mathcal{Y}_\circ(R_0), \Pi_{\boldsymbol{\eta}} \boldsymbol{v} = \boldsymbol{\eta}\},$$

$$\|H(\hat{\eta}(\boldsymbol{\theta}) - \eta(\boldsymbol{\theta}))\| \leq \tau, \quad \text{for all } \boldsymbol{\theta} \in \mathcal{Y}_{\circ, \boldsymbol{\theta}}(R_0) \stackrel{\text{def}}{=} \{\boldsymbol{v} \in \mathcal{Y}_\circ(R_0), \Pi_{\boldsymbol{\theta}} \boldsymbol{v} = \boldsymbol{\theta}\}.$$

Then we can easily modify the proof of Theorem 7 via adding $\mathcal{C}(\nu)\tau$ to the error terms and the radii \mathbf{r}_k , where $\mathcal{C}(\nu)$ is some rational function of ν .

2.5 Convergence to the ME

Even though Theorem 7 tells us that the statistical properties of the alternating sequence resemble those of its target - the profile ME $\tilde{\boldsymbol{\theta}}$ - it is an interesting question if the underlying approach allows to qualify conditions under which the sequence actually attains the maximizer $\tilde{\boldsymbol{v}}$.

Define the radius $\mathbf{r}_0(\mathbf{x}) > 0$ to be the smallest radius $\mathbf{r} > 0$ such that $\mathbb{P}(\tilde{\boldsymbol{v}}, \tilde{\boldsymbol{v}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}_\circ(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}$, where

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

Remark 12 This radius can be determined using conditions $(\mathcal{L}_{\mathbf{r}})$ and $(\mathcal{E}_{\mathbf{r}})$ of Section 2.1 and Theorem 28 which would yield $\mathbf{r}_0(\mathbf{x}) = O(\sqrt{\mathbf{x} + p^*})$.

Without further assumptions Theorem 7 yields the following Corollary:

Corollary 13 Under the assumptions of Theorem 7 it holds with probability greater $1 - 8e^{-\mathbf{x}} - \beta$

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

Corollary 13 is a first step in the direction of an actual convergence result but the gap $\check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}) + \check{\diamond}_Q(\mathbf{r}_0, \mathbf{x})$ is not a zero sequence in $k \in \mathbb{N}$. It turns out that it is possible to

prove convergence to the ME at the cost of assuming more smoothness of the functional \mathcal{L} and using the right bound for the maximal eigenvalue of the Hessian $\nabla^2\mathcal{L}(\mathbf{v}^*)$.

Define $\mathfrak{z}_2(\mathbf{x}, \nabla^2\mathcal{L}(\mathbf{v}^*))$ via

$$\mathbb{P}\left(\|\mathcal{D}^{-1}\nabla^2\mathcal{L}(\mathbf{v}^*)\| \geq \mathfrak{z}_2(\mathbf{x}, \nabla^2\mathcal{L}(\mathbf{v}^*))\right) \leq e^{-\mathbf{x}},$$

and $\kappa(\mathbf{x}, R_0)$ as

$$\kappa(\mathbf{x}, R_0) \stackrel{\text{def}}{=} \frac{2\sqrt{2}(1+\sqrt{\nu})}{\sqrt{1-\nu}} \left\{ \epsilon R_0 + 9\epsilon_2\nu_2\|\mathcal{D}^{-1}\|\mathfrak{z}_1(\mathbf{x}, 6p^*)R_0 + \|\mathcal{D}^{-1}\|\mathfrak{z}_2(\mathbf{x}, \nabla^2\mathcal{L}(\mathbf{v}^*)) \right\},$$

where $\mathfrak{z}_1(\mathbf{x}, \cdot) \approx \sqrt{\mathbf{x} + p^*}$, see (46). With these definitions we can prove the following Theorem:

Theorem 14 *Let the conditions (\mathcal{L}_r) , $(\mathcal{E}r)$ and (B_1) be met. Further suppose (A_1) and (A_2) , with $(\mathcal{E}\mathcal{D}_2)$ instead of $(\mathcal{E}\mathcal{D}_1)$. Assume that $\kappa(\mathbf{x}, R_0) < (1 - \nu)$. Then*

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

where

$$\mathbf{r}_k^* \leq \begin{cases} \nu^k \frac{4\sqrt{2}}{1-\kappa(\mathbf{x}, R_0)k} R_0, & \kappa(\mathbf{x}, R_0)k \leq 1, \\ \nu^{\frac{k}{\log(k)}} \log\left(\frac{1-\nu}{\kappa(\mathbf{x}, R_0)}\right) c_k R_0, & \text{otherwise,} \end{cases} \quad (24)$$

with some sequence $(c_k) \in \mathbb{N}$, where $0 < c_k \rightarrow 2$.

Remark 15 *This means that we obtain nearly linear convergence to the global maximizer $\tilde{\mathbf{v}}$.*

Remark 16 *As in Remark 11 if no exact numerical computation of the stepwise maximizers is possible we can easily modify the proof of Theorem 14 via adding $\mathcal{C}(\nu)\tau$ to $\kappa(\mathbf{x}, R_0)$ to address that case.*

Remark 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 (Tropp, 2012) to obtain*

$$\mathfrak{z}_2(\mathbf{x}, \nabla^2\mathcal{L}(\mathbf{v}^*)) = \sqrt{2\tau\nu}\sqrt{\mathbf{x} + \log(p^*)},$$

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

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

In the case of smooth i.i.d models this means that

$$\kappa(\mathbf{x}, R_0) \leq \frac{\mathcal{C}}{\sqrt{n}}(\mathbf{x} + R_0 + \log(p^*)),$$

if $p^* + \mathbf{x} = o(n)$.

Remark 18 *It may happen that $\kappa(\mathbf{x}, R_0)/(1 - \nu)$ is very close to or even larger than 1. But a close look at the proof of Theorem 14 reveals that this can be improved using Lemma 33. For this purpose bound $\mathbf{r}_k^* \leq \mathbf{C}^*(\mathfrak{z}(\mathbf{x}) + \nu^k R_0)$ with \mathbf{r}_k^* defined in (33) and with some constant $\mathbf{C}^* > 0$. Then the result of Theorem 14 is true with $\kappa(\mathbf{x}, \mathbf{C}^* \sqrt{p^* + \bar{\mathbf{x}}})$ instead of $\kappa(\mathbf{x}, R_0)$ and with probability greater $1 - 10e^{-x}$. See Remark 38 for more details.*

2.6 Critical Dimension

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 7 are accurate if the spread function $\check{\diamond}_Q(\mathbf{r}_k, \mathbf{x})$ from (21) is small. The critical size of p_n then depends on the exact bounds on $\epsilon, \check{\epsilon}$. In the i.i.d setting we have $\epsilon \asymp \check{\epsilon} \asymp 1/\sqrt{n}$ such that $\check{\diamond}(\mathbf{r}_k, \mathbf{x}) \asymp p_n/\sqrt{n}$ for large $k \in \mathbb{N}$. In other words, one needs that “ p_n^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_n^2/n is small”, as is shown by Andresen and Spokoiny (2014). There are examples where the pME only satisfies a Wilks- or Fisher result if “ p_n^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 $\kappa(\mathbf{x}, R_0) < (1 - \nu)$ of Theorem 14 for the convergence of the sequence to the global maximizer means that one needs $p_n/n \ll 1$ in the smooth i.i.d. setting if $R_0 \leq \mathbf{C}_{R_0} \sqrt{p_n + \bar{\mathbf{x}}}$. Further Theorem 14 states a lower bound for the speed of convergence that in the smooth i.i.d. setting decreases if p_n/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/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 7 and Theorem 14 can be applied in Single Index modeling. This section is based on (Andresen, 2015). See that paper for a more detailed presentation.

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 non-parametric 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, (Delecroix. et al., 1997) show the asymptotic efficiency of the general semiparametric maximum-functional estimator for particular examples and in (Haerdle et al., 1993) the right choice of band-

width for the nonparametric estimation of the link function is analyzed. We want to use this model to illustrate our theoretical results.

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 $s > 0$. This allows to approximate $f \in \{f : [-s, s] \mapsto \mathbb{R}\}$ by an orthonormal C^3 -Daubechies-wavelet basis $(\mathbf{e}_k)_{k \in \mathbb{N}}$ on the interval.

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

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

where

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

and where $\Upsilon_m = S_1^{p,+} \times \mathbb{R}^m$. Again we will suppress the sub index \cdot_m in the following.

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

$$\tilde{\boldsymbol{v}}_0 \stackrel{\text{def}}{=} \operatorname{argmax}_{\substack{(\boldsymbol{\theta}, \boldsymbol{\eta}) \in \Upsilon \\ \boldsymbol{\theta} \in G_N}} \mathcal{L}(\boldsymbol{v}). \quad (25)$$

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

$$\tilde{\boldsymbol{\eta}}_{0,k} \stackrel{\text{def}}{=} \operatorname{argmax} \mathcal{L}(\boldsymbol{\theta}_k, \boldsymbol{\eta}) = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{e} \mathbf{e}^\top (\mathbf{X}_i^\top \boldsymbol{\theta}_k) \right)^{-1} \frac{1}{n} \sum_{i=1}^n y_i \mathbf{e}^\top (\mathbf{X}_i^\top \boldsymbol{\theta}_k) \in \mathbb{R}^m, \quad (26)$$

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

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

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

To apply the result presented in Theorem 7 and Theorem 14 we need a list of assumptions denoted by (\mathcal{A}) .

(Cond \mathbf{X}) The random variables $(\mathbf{X}_i)_{i=1, \dots, n} \subset \mathbb{R}^p$ are i.i.d and bounded with distribution denoted by $\mathbb{P}^{\mathbf{X}}$ and independent of $(\varepsilon_i)_{i=1, \dots, n} \subset \mathbb{R}$.

The measure $\mathbb{P}^{\mathbf{X}}$ is absolutely continuous with respect to the Lebesgue measure. The

Lebesgue density $p_{\mathbf{X}}$ of $\mathbb{P}^{\mathbf{X}}$ is Lipschitz continuous and positive on $B_s(0) \subset \mathbb{R}^p$. For any pair $\boldsymbol{\theta} \in S_1^{+,p}$ with $\boldsymbol{\theta} \perp \boldsymbol{\theta}^*$ we have almost surely

$$\text{Var}(\mathbf{X}^\top \boldsymbol{\theta} | \mathbf{X}^\top \boldsymbol{\theta}^*) > 0.$$

(**Cond_f**) For some $\boldsymbol{\eta}^* \in B_{\mathbf{r}^\circ}(0) \subset l^2 \stackrel{\text{def}}{=} \{(u_k)_{k \in \mathbb{N}} : \sum_{k=1}^{\infty} u_k^2 < \infty\}$

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

where $\|f'\|_\infty < \infty$ and $\|f''\|_\infty < \infty$ and where with some $\alpha > 2$

$$\sum_{k=0}^{\infty} k^{2\alpha} \eta_k^{*2} < \infty.$$

On some interval $[t_0 - h, t_0 + h] \subseteq [-s + s]$ with $h > 0$ it holds true that

$$|f'(t)| > 0.$$

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

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

Remark 19 Note that our assumptions in terms of moments and smoothness are quite common in this model. For instance (Haerdle et al., 1993) assume that the density $p_{\mathbf{X}}$ of the regressors (\mathbf{X}_i) is twice continuously differentiable, that f has two bounded derivatives and that the errors (ε_i) are centered with bounded polynomial moments of arbitrary degree.

Remark 20 $\text{Var}(\mathbf{X}^\top \boldsymbol{\theta}^\circ | \mathbf{X}^\top \boldsymbol{\theta}^*) = 0$ would mean that $\mathbf{X}^\top \boldsymbol{\theta}^\circ = a(\mathbf{X}^\top \boldsymbol{\theta}^*)$ for some measurable function $a : \mathbb{R} \rightarrow \mathbb{R}$. But then we would have for any $(\alpha, \beta) \in \mathbb{R}^2$ with $\alpha^2 + \beta^2 = 1$ that

$$f(\mathbf{X}^\top (\alpha \boldsymbol{\theta}^* + \beta \boldsymbol{\theta}^\circ)) = f(\alpha \mathbf{X}^\top \boldsymbol{\theta}^* + \beta a(\mathbf{X}^\top \boldsymbol{\theta}^*)) \stackrel{\text{def}}{=} f_{\alpha, \beta}^\circ(\mathbf{X}^\top \boldsymbol{\theta}^*),$$

such that the problem would no longer be identifiable. Also $|f'(t)| > 0$ on some interval is necessary for identifiability of $\boldsymbol{\theta}^*$.

Proposition 21 Let $\tau = o(p^{*-3/2})$ and $p^{*5}/n \rightarrow 0$. With initial guess given by Equation (25) and for not too large $\mathbf{x} > 0$ the alternating sequence satisfies (22) and (23) with probability greater $1 - 9 \exp\{-\mathbf{x}\}$ and where with some constant $\mathbf{C} \in \mathbb{R}$

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

Proposition 22 *Take the initial guess given by Equation (25). Assume (A). Further assume that $p^{*4}/n \rightarrow 0$ and $\tau = o(p^{*-3/2})$. Then we get the claim of Theorem 14 with $\beta = e^{-\mathbf{x}}$ and*

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

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

Remark 23 *The constraint $\tau = o(p^{*-3/2})$ implies that for the calculation of the initial guess the vector $\tilde{\boldsymbol{\eta}}_{0,l}$ of (26) and the functional $\mathcal{L}(\cdot)$ have to be evaluated $N = p^{*3(p-1)/2}$ times.*

For details and proofs see (Andresen, 2015).

4. Proof of Theorem 7

In this section we present the proof of Theorem 7. As the proof is quite technical and complex we want to first explain the basic ideas of the proof. In a second section we will outline more clearly the steps of the proof. Finally we carry out each of these steps which combine to yield the proof.

4.1 Idea of the Proof

To ease the understanding of what follows in the subsequent sections we want to illustrate the central ideas with a simple model. Consider for some positive definite matrix $\mathcal{D} \in \mathbb{R}^{p^* \times p^*}$ and some vector $\mathbf{v}^* = (\boldsymbol{\theta}^*, \boldsymbol{\eta}^*) \in \mathbb{R}^{p+m} = \mathbb{R}^{p^*}$ the model

$$\mathbb{Y} = \mathbf{v}^* + \boldsymbol{\varepsilon} \in \mathbb{R}^{p^*}, \text{ where } \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathcal{D}^{-2}), \quad \mathcal{D}^2 = \begin{pmatrix} \mathcal{D}^2 & A \\ A^\top & \mathbb{H}^2 \end{pmatrix} \in \mathbb{R}^{p^* \times p^*}.$$

Set \mathcal{L} to be the true log likelihood of the observations, i.e.

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

With any starting initial guess $\tilde{\mathbf{v}}_0 \in \mathbb{R}^{p+m}$ we obtain from (3) for $k \in \mathbb{N}$ and the usual first order criterion of maximality the following two equations

$$\begin{aligned} \mathcal{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) &= \mathcal{D}\boldsymbol{\varepsilon}_\theta + \mathcal{D}^{-1}A(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*), \\ \mathbb{H}(\tilde{\boldsymbol{\eta}}_{k+1} - \boldsymbol{\eta}^*) &= \mathbb{H}\boldsymbol{\varepsilon}_\eta + \mathbb{H}^{-1}A^\top(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*). \end{aligned}$$

Combining these two equations we derive, assuming $\|\mathcal{D}^{-1}A\mathbb{H}^{-2}A^\top\mathcal{D}^{-1}\| \stackrel{\text{def}}{=} \|\mathbf{M}_0\| = \nu < 1$

$$\begin{aligned} \mathcal{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) &= \mathcal{D}^{-1}(\mathcal{D}^2\boldsymbol{\varepsilon}_\theta - A\boldsymbol{\varepsilon}_\eta) + \mathcal{D}^{-1}A\mathbb{H}^{-1}A^\top\mathcal{D}^{-1}\mathcal{D}(\tilde{\boldsymbol{\theta}}_{k-1} - \boldsymbol{\theta}^*) \\ &= \sum_{l=1}^k \mathbf{M}_0^{k-l}\mathcal{D}^{-1}(\mathcal{D}^2\boldsymbol{\varepsilon}_\theta - A\boldsymbol{\varepsilon}_\eta) + \mathbf{M}_0^k\mathcal{D}(\tilde{\boldsymbol{\theta}}_0 - \boldsymbol{\theta}^*) \\ &\rightarrow \mathbf{X} \stackrel{\text{def}}{=} \mathcal{D}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*). \end{aligned}$$

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

$$\nabla \mathcal{L}(\boldsymbol{v}) = \mathcal{D}_{\boldsymbol{v}^*}^2(\boldsymbol{v} - \boldsymbol{v}^*) + \mathcal{D}_{\boldsymbol{v}^*}^2 \boldsymbol{\varepsilon}.$$

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 7 .

First it is ensured that the whole sequence $(\widetilde{\boldsymbol{v}}_{k,k(+1)})_{k \in \mathbb{N}_0}$ satisfies for some $R_0(\mathbf{x}) > 0$ and with probability greater than $1 - e^{-x}$

$$\{\widetilde{\boldsymbol{v}}_{k,k(+1)}, k \in \mathbb{N}_0\} \subset \{\|\mathcal{D}(\boldsymbol{v} - \boldsymbol{v}^*)\| \leq R_0(\mathbf{x})\}, \quad (27)$$

where $\mathcal{D}^2 \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E} \mathcal{L}(\boldsymbol{v}^*)$ (see Theorem 28). In the second step we approximate with $\zeta = \mathcal{L} - \mathbb{E} \mathcal{L}$

$$\mathcal{L}(\boldsymbol{v}) - \mathcal{L}(\boldsymbol{v}^*) = \nabla \zeta(\boldsymbol{v}^*)(\boldsymbol{v} - \boldsymbol{v}^*) - \|\mathcal{D}(\boldsymbol{v} - \boldsymbol{v}^*)\|^2/2 + \alpha(\boldsymbol{v}, \boldsymbol{v}^*), \quad (28)$$

where $\alpha(\boldsymbol{v}, \boldsymbol{v}^*)$ is defined by (28). Similar to the toy case above this allows using the first order criterion of maximality and (27) to obtain a bound of the kind

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

This is done in Lemma 32 using results from (Andresen and Spokoiny, 2014) to show that $\epsilon(R_0)$ is small. Finally the same arguments as in (Andresen and Spokoiny, 2014) allow to obtain our main result using that with high probability for all $k \in \mathbb{N}_0$ $\widetilde{\boldsymbol{v}}_{k,k} \in \{\|\mathcal{D}(\boldsymbol{v} - \boldsymbol{v}^*)\| \leq \mathbf{r}_k\}$. For the convergence result similar arguments are used. The only difference is that instead of (28) we use the approximation

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

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

4.2 A Desirable Set

In this section we will explain the agenda of the proof. 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 $\mathcal{L}(\boldsymbol{v})$ 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 - with the local set $\mathcal{Y}_o(\mathbf{r})$ defined in (15) - for some $K \in \mathbb{N}$ the set

$$\begin{aligned} \Omega(\mathbf{x}) &= \bigcap_{k=0}^K (C_{k,k} \cap C_{k,k+1}) \cap C(\nabla) \cap \{\mathcal{L}(\tilde{\mathbf{v}}_0) - \mathcal{L}(\mathbf{v}^*) \geq -K_0\}, \text{ where} \quad (29) \\ C_{k,k+1} &= \left\{ \|\mathcal{D}(\tilde{\mathbf{v}}_{k,k+1}) - \mathbf{v}^*\| \leq R_0(\mathbf{x}), \|\mathcal{D}(\tilde{\boldsymbol{\theta}}_k) - \boldsymbol{\theta}^*\| \leq R_0(\mathbf{x}), \right. \\ &\quad \left. \|\mathcal{H}(\tilde{\boldsymbol{\eta}}_{k+1}) - \boldsymbol{\eta}^*\| \leq R_0(\mathbf{x}) \right\}, \\ C(\nabla) &= \bigcap_{\mathbf{r} \leq R_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \left\{ \frac{1}{6\epsilon\nu_1} \|\mathcal{Y}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 4p^*)^2 \right\} \\ &\quad \bigcap_{\mathbf{r} \leq 4R_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \left\{ \frac{1}{6\check{\epsilon}\check{\nu}_1} \|\check{\mathcal{Y}}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 2p^* + 2p)^2 \right\} \\ &\quad \cap \left\{ \max\{\|\mathcal{D}^{-1}\nabla\mathcal{L}\|, \|\mathcal{D}^{-1}\nabla_{\boldsymbol{\theta}}\mathcal{L}\|, \|\mathcal{H}^{-1}\nabla_{\boldsymbol{\eta}}\mathcal{L}\|\} \leq \mathfrak{z}(\mathbf{x}) \right\} \\ &\quad \cap \{\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \in \mathcal{Y}_o(\mathbf{r}_0(\mathbf{x}))\}. \end{aligned}$$

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

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

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

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

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

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

The constant $\mathfrak{z}(\mathbf{x})$ in the definition of $C(\nabla)$ 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. It is defined as

$$\mathfrak{z}(\mathbf{x}) \stackrel{\text{def}}{=} \mathfrak{z}(\mathbf{x}, \mathcal{D}^{-1}\mathcal{V}^2\mathcal{D}^{-1}) \vee \mathfrak{z}_Q(\mathbf{x}, 4p^*) \approx \sqrt{p^* + \mathbf{x}}, \quad (30)$$

where $\mathcal{V}^2 \in \mathbb{R}^{p^* \times p^*}$ denotes the *covariance matrix* from Section 2.1

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

Remark 24 $\mathfrak{z}(\mathbf{x}, \cdot)$ is explained in more detail in Section A and $\mathfrak{z}_Q(\mathbf{x}, \cdot)$ is defined in Equation (42). The constant $\mathfrak{z}(\mathbf{x}, \mathcal{B})$ is comparable to the " $1 - e^{-\mathbf{x}}$ "-quantile of the norm

of \mathbf{X} , where $\mathbf{X} \sim \mathcal{N}(0, \mathbb{B})$, i.e. it is of order of the trace of \mathbb{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} .

Remark 25 We intersect the set with the event $\{\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_\circ(\mathbf{r}_0)\}$ where we a priori demand $\mathbf{r}_0(\mathbf{x}) > 0$ to be chosen such that $\mathbb{P}(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_\circ(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}$. Note that condition $(\mathcal{E}\mathbf{r})$ together with $(\mathcal{L}\mathbf{r})$ allow to set $\sqrt{p^* + \mathbf{x}} \approx \mathbf{r}_0 \leq R_0$ (see Theorem 28).

In Section 4.3 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.

$\{\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -\mathbf{K}_0\}$: This set ensures, that the first guess satisfies $\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -\mathbf{K}_0$, 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 $\mathbf{v}^* \in \mathcal{Y}$ with Theorem 28.

$\{\mathcal{D}(\tilde{\mathbf{v}}_{k,k(+1)} - \mathbf{v}^*) \leq \mathbf{R}_0(\mathbf{x})\}$: As just mentioned this event is of high probability due to $\mathcal{L}(\tilde{\mathbf{v}}_0, \mathbf{v}^*) \geq -\mathbf{K}_0$ and Theorem 28. This allows to concentrate the analysis on the set $\mathcal{Y}_\circ(\mathbf{R}_0)$ 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 32 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 34 the bounds for all $k \in \mathbb{N}$.

On $\Omega(\mathbf{x}) \subset \Omega$ we find $\tilde{\mathbf{v}}_{k,k(+1)} \in \mathcal{Y}_\circ(\mathbf{r}_k)$ such that we can follow the arguments of Theorem 2.2 of (Andresen and Spokoiny, 2014) to obtain the desired result with accuracy measured by $\check{\diamond}_Q(\mathbf{r}_k, \mathbf{x})$.

The sketch in figure 4.2 illustrates the behavior of the first steps of the procedure. The axes correspond to the θ - or η -subspaces respectively. The two ellipsoids with center \mathbf{v}^* and solid frame represent the local sets $\mathcal{Y}_\circ(R_K) \subset \mathcal{Y}_\circ(\mathbf{R}_0)$, with $R_K > 0$ from Remark 5. We see that the initial guess $\tilde{\mathbf{v}}_0$ lies in $\mathcal{Y}_\circ(R)$. The elements $(\mathbf{v}_{k,k(+1)})$ of the alternating sequence all land inside of the respective $\mathcal{Y}_\circ(\mathbf{r}_k)$, which are represented by shrinking ellipsoids centered in \mathbf{v}^* with dotted frames. Note that not the set $\mathcal{Y}_\circ(R_K)$ but $\mathcal{Y}_\circ(\mathbf{R}_0)$ contains all points of the sequence.

4.3 Probability of Desirable Set

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

Lemma 26 *The set $C(\nabla)$ satisfies*

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

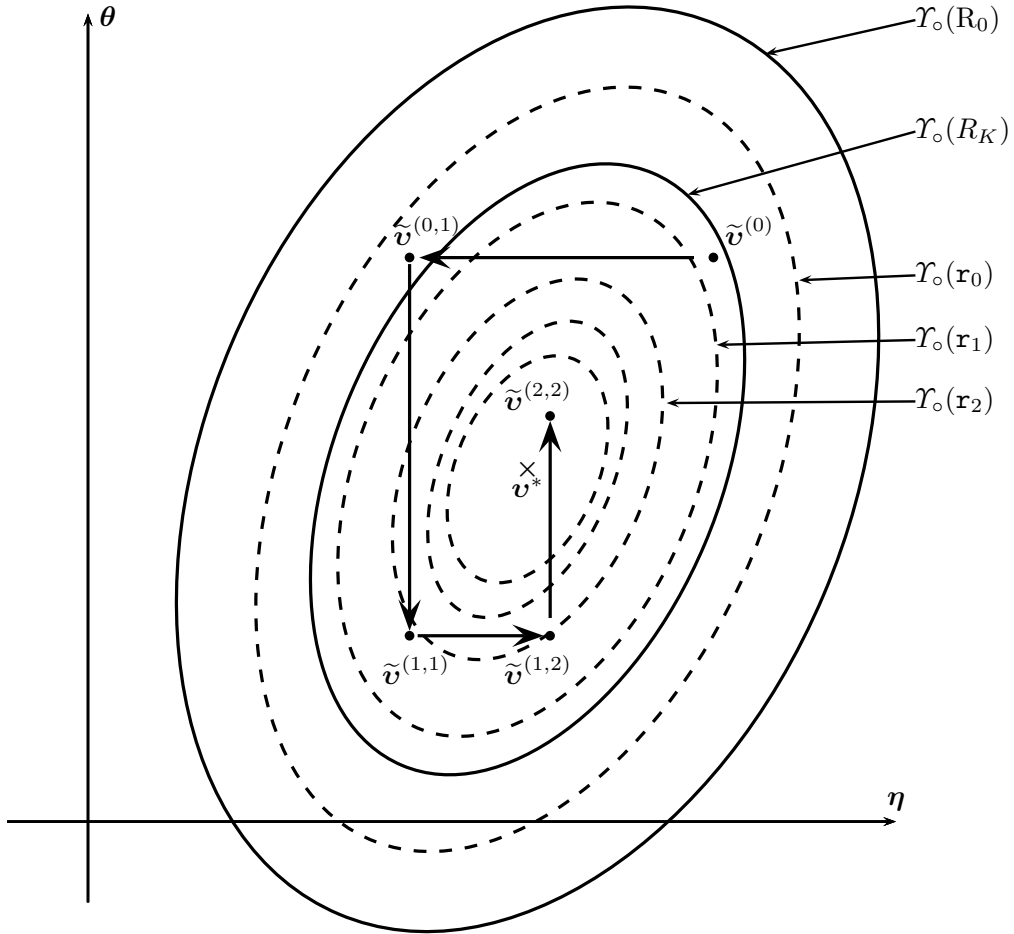


Figure 1: The behavior of the procedure for the first 4 steps of the alternating algorithm.

Proof The proof is similar to the proof of Theorem 3.1 in (Spokoiny, 2012). Denote

$$\begin{aligned}\mathcal{A} &\stackrel{\text{def}}{=} \bigcap_{\mathbf{r} \leq \mathbf{R}_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \left\{ \frac{1}{6\epsilon\nu_1} \|\mathfrak{y}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 4p^*)^2 \right\} \\ \mathcal{B} &\stackrel{\text{def}}{=} \bigcap_{\mathbf{r} \leq 4\mathbf{R}_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \left\{ \frac{1}{6\check{\epsilon}\check{\nu}_1} \|\check{\mathfrak{y}}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \leq \mathfrak{z}_Q(\mathbf{x}, 2p^* + 2p)^2 \right\} \\ \mathcal{C} &\stackrel{\text{def}}{=} \left\{ \max\{\|\mathcal{D}^{-1}\nabla\mathcal{L}\|, \|\mathcal{D}^{-1}\nabla_{\boldsymbol{\theta}}\mathcal{L}\|, \|\mathbb{H}^{-1}\nabla_{\boldsymbol{\eta}}\mathcal{L}\|\} \leq \mathfrak{z}(\mathbf{x}) \right\}.\end{aligned}$$

We estimate

$$\begin{aligned}\mathbb{P}(\mathcal{C}(\nabla)) &\geq 1 - \mathbb{P}(\mathcal{A}^c) - \mathbb{P}(\mathcal{B}^c) - \mathbb{P}(\mathcal{C}^c) \\ &\quad - \mathbb{P}(\tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\boldsymbol{\theta}^*} \notin \mathcal{Y}_\circ(\mathbf{r}_0)) - \mathbb{P}(\|\check{\boldsymbol{\xi}}\| > \mathfrak{z}(\mathbf{x}, \text{Cov}(\check{\boldsymbol{\xi}}))).\end{aligned}$$

We bound using for both terms Theorem 42 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 27 to find

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

This implies that

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

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

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

By the choice of $\mathfrak{z}(\mathbf{x}) > 0$ and $\mathbf{r}_0 > 0$ this gives the claim. ■

We cite Lemma B.2 of (Andresen and Spokoiny, 2014):

Lemma 27 *Let*

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

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

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 (Spokoiny, 2012) shows that it actually yields the following modified version:

Theorem 28 ((Spokoiny, 2012), Theorem 4.1) *Suppose (\mathcal{E}_r) and (\mathcal{L}_r) with $\mathbf{b}(r) \equiv \mathbf{b}$. Further define the following random set*

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

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

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

then

$$\mathbb{P}(\Upsilon(K) \subseteq \Upsilon_{\circ}(\mathbf{r}_0)) \geq 1 - e^{-\mathbf{x}}.$$

Note that with (\mathcal{I})

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

With assumption (B_1) and

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

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 28

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

Remark 29 *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 $\{\boldsymbol{\theta} \in \mathbb{R}^p, (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \Upsilon_{\circ}(R_0)\}$ or $\{\boldsymbol{\eta} \in \mathbb{R}^m, (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \Upsilon_{\circ}(R_0)\}$.*

Remark 30 *To address the claim of Remark 5 we present the following lemma:*

Lemma 31 *On the set $C(\nabla) \cap \{\tilde{\mathbf{v}}_0 \in \Upsilon_\circ(R_K)\}$ it holds*

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

Proof *With similar arguments as in the proof of Lemma 32 we have on $C(\nabla) \subset \Omega$ that*

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

■

4.4 Proof Convergence

We derive the a priori bound $\tilde{\mathbf{v}}_{k,k(+1)} \in \Upsilon_\circ(\mathbf{r}_k)$ with an adequately decreasing sequence $(\mathbf{r}_k) \subset \mathbb{R}_+$ using the argument of Section 4.1, where $\limsup \mathbf{r}_k \approx \mathfrak{z}(\mathbf{x})$. For this purpose we define the *parametric uniform spread*

$$\diamond_Q(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \epsilon \{2\mathbf{r}^2 + \mathfrak{z}_Q(\mathbf{x}, 4p^*)^2\} (1 + 6\nu_1^2). \quad (31)$$

Lemma 32 *Assume that for some sequence $(\mathbf{r}_k^{(l)})_{k \in \mathbb{N}}$*

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

Then under the assumptions of Theorem 7 we get on $\Omega(\mathbf{x})$ for all $k \in \mathbb{N}_0$

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

Proof

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

$$\mathcal{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) = \mathcal{D}^{-1}\nabla_{\boldsymbol{\theta}}\mathcal{L}(\mathbf{v}^*) - \mathcal{D}^{-1}\mathcal{A}(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*) + \boldsymbol{\tau}(\mathbf{r}_k^{(l)}), \quad (32)$$

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

where - with $\diamond_Q(\mathbf{r}, \mathbf{x})$ defined in (31) -

$$\|\boldsymbol{\tau}(\mathbf{r})\| \leq \diamond_Q(\mathbf{r}, \mathbf{x}).$$

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 (Andresen and Spokoiny, 2014). By assumption on $\Omega(\mathbf{x})$ we have $\tilde{\mathbf{v}}_{k,k+1} \in \mathcal{Y}_\circ(\mathbf{r}_k^{(l)})$. Define with $\zeta = \mathcal{L} - \mathbb{E}\mathcal{L}$

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

Note that

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

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

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

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

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

where

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

To see this note first that with Lemma 27 $\|\mathcal{D}^{-1}\Pi_{\boldsymbol{\theta}}\mathcal{D}\mathbf{v}\| \leq \|\mathcal{D}^{-1}\mathcal{D}\mathbf{v}\|$. This gives by condition (\mathcal{L}_0) , Lemma 27 and Taylor expansion

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

For the remainder note that again with Lemma 27

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

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

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

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

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

To see this we first use the inequality

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

Now we find with (32)

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

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

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

and

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

to derive the recursive formula

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

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

Lemma 33 *Assume the same as in Theorem 7 . Further assume that (17) is met with $\mathfrak{C}(\nu)$ defined in (18). Then*

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

where - with $\mathfrak{C}(\nu) \geq 0$ defined in (18) -

$$\begin{aligned} \mathbf{r}_k^* &\leq \left(\mathfrak{C}(\nu) + \frac{4\epsilon\mathfrak{C}(\nu)^4\mathfrak{z}(\mathbf{x})}{1 - \epsilon\mathfrak{C}(\nu)\mathfrak{z}(\mathbf{x})} \right) 2\mathfrak{z}(\mathbf{x}) \\ &\quad + \nu^k \left(\mathfrak{C}(\nu) + \frac{4\epsilon\mathfrak{C}(\nu)^4 R_0}{1 - \epsilon\mathfrak{C}(\nu)R_0} \right) R_0. \end{aligned} \tag{33}$$

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 32 we find that if

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

that then

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

where

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

Setting $l = 1$ this gives

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

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 $\limsup_{l \rightarrow \infty} \mathbf{r}_k^{(l)} \leq \mathbf{r}_k^*$ in (33). For this we estimate further

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

where $\mathbf{C}(\nu) > 0$ is defined in (18). We set

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

Claim

$$\begin{aligned}
 A_{s,k}^{(l)} &\leq 4^{\sum_{t=0}^{s-1} 2^t} \mathbf{C}(\nu)^{2^s} \left\{ \left(\frac{1}{1-\nu} \right)^{\sum_{t=0}^{s-1} 2^t} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^{2^s} \right. \\
 &\quad \left. + \nu^k \left(\frac{1}{\nu^{-1}-1} \right)^{\sum_{t=0}^{s-1} 2^t} \mathbf{R}_0^{2^s} \right\} \\
 &\quad + 4^{\sum_{t=0}^{s-1} 2^t} (\mathbf{C}(\nu)\epsilon)^{2^s} A_{s+1,k}^{(l-1)}.
 \end{aligned} \tag{34}$$

We proof this claim via induction. Clearly

$$\begin{aligned}
 A_{1,k}^{(l)} &= \sum_{r_1=0}^{k-1} \nu^{r_1} (\mathbf{r}_{k-r_1}^{(l-1)})^2 \leq 4\mathbf{C}(\nu)^2 \sum_{r_1=0}^{k-1} \nu^{r_1} \left\{ (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^2 + \nu^{2(k-r_1)} \mathbf{R}_0^2 \right\} \\
 &\quad + 4\mathbf{C}(\nu)^2 \epsilon^2 \sum_{r_1=0}^{k-1} \nu^{r_1} \left(\sum_{r_2=0}^{k-r_1-r_2-1} \nu^{r_2} (\mathbf{r}_{k-r_1-r_2}^{(l-2)})^2 \right)^2 \\
 &\leq 4\mathbf{C}(\nu)^2 \left\{ \frac{1}{1-\nu} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^2 + \frac{\nu^k}{\nu^{-1}-1} \mathbf{R}_0^2 \right\} \\
 &\quad + 4\mathbf{C}(\nu)^2 \epsilon^2 A_{2,k}^{(l-1)}.
 \end{aligned}$$

Furthermore

$$\begin{aligned}
 A_{s,k}^{(l)} &\stackrel{\text{def}}{=} \sum_{r_1=0}^{k-1} \nu^{r_1} \left(\sum_{r_2=0}^{k-r_1-1} \nu^{r_2} \left(\dots \sum_{r_s=0}^{k-r_1-\dots-r_{s-1}-1} \nu^{r_s} (\mathbf{r}_{k-r_1-\dots-r_s}^{(l-1)})^2 \dots \right)^2 \right)^2 \\
 &= \sum_{r_1=0}^{k-1} \nu^{r_1} \left(A_{s-1,k-r_1}^{(l)} \right)^2. \tag{35}
 \end{aligned}$$

Plugging in (34) we get for $s \geq 2$

$$\begin{aligned}
 A_{s,k}^{(l)} &\leq \sum_{r_1=0}^{k-1} \nu^{r_1} \left(4^{\sum_{t=0}^{s-2} 2^t} \mathbf{C}(\nu)^{2^{s-1}} \left\{ \left(\frac{1}{1-\nu} \right)^{\sum_{t=0}^{s-2} 2^t} (\mathfrak{J}(\mathbf{x}) + \epsilon \mathfrak{J}(\mathbf{x})^2)^{2^{s-1}} \right. \right. \\
 &\quad \left. \left. + \nu^k \left(\frac{1}{\nu^{-1}-1} \right)^{\sum_{t=0}^{s-2} 2^t} \mathbf{R}_0^{2^{s-1}} \right\} \right. \\
 &\quad \left. + 4^{\sum_{t=0}^{s-2} 2^t} (\mathbf{C}(\nu)\epsilon)^{2^{s-1}} A_{s,k-r_1}^{(l-1)} \right)^2.
 \end{aligned}$$

Shifting the index this gives

$$\begin{aligned}
 A_{s,k}^{(l)} &\leq 4 \sum_{r_1=0}^{k-1} \nu^{r_1} \left(4^{\sum_{t=1}^{s-1} 2^t} \mathbf{C}(\nu)^{2^s} \left\{ \left(\frac{1}{1-\nu} \right)^{\sum_{t=1}^{s-1} 2^{t-1}} (\mathfrak{J}(\mathbf{x}) + \epsilon \mathfrak{J}(\mathbf{x})^2)^{2^s} \right. \right. \\
 &\quad \left. \left. + \nu^k \left(\frac{1}{\nu^{-1}-1} \right)^{\sum_{t=1}^{s-1} 2^t} \mathbf{R}_0^{2^s} \right\} \right. \\
 &\quad \left. + 4^{\sum_{t=1}^{s-1} 2^t} (\mathbf{C}(\nu)\epsilon)^{2^s} (A_{s,k-r_1}^{(l-1)})^2 \right).
 \end{aligned}$$

Direct calculation then leads to

$$\begin{aligned}
 A_{s,k}^{(l)} &\leq 4^{\sum_{t=0}^{s-1} 2^t} \mathbf{C}(\nu)^{2^s} \left\{ \left(\frac{1}{1-\nu} \right)^{\sum_{t=0}^{s-1} 2^t} (\mathfrak{J}(\mathbf{x}) + \epsilon \mathfrak{J}(\mathbf{x})^2)^{2^s} \right. \\
 &\quad \left. + \nu^k \left(\frac{1}{\nu^{-1}-1} \right)^{\sum_{t=0}^{s-1} 2^t} \mathbf{R}_0^{2^s} \right\} \\
 &\quad + 4^{\sum_{t=0}^{s-1} 2^t} (\mathbf{C}(\nu)\epsilon)^{2^s} \sum_{r_1=0}^{k-1} \nu^{r_1} (A_{s,k-r_1}^{(l-1)})^2,
 \end{aligned}$$

which gives (34) with (35). Similarly we can prove

$$A_{s,k}^{(1)} = \left(\frac{1}{1-\nu} \right)^{2^{s-1}} R_0^{2^s}.$$

Abbreviate

$$\begin{aligned} \lambda_s &\stackrel{\text{def}}{=} 4^{2^s-1} \mathbf{C}(\nu)^{2^s}, \quad \beta_s \stackrel{\text{def}}{=} 4^{2^s-1} (\mathbf{C}(\nu)\epsilon)^{2^s}, \\ \mathfrak{z}_s(\mathbf{x}) &\stackrel{\text{def}}{=} \left(\frac{1}{1-\nu} \right)^{2^{s-1}} (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2)^{2^s}, \quad R_s \stackrel{\text{def}}{=} \left(\frac{1}{\nu^{-1}-1} \right)^{2^{s-1}} R_0^{2^s}. \end{aligned}$$

Then

$$\begin{aligned} \mathbf{r}_k^{(l)} &\leq \mathbf{C}(\nu) \left\{ (\mathfrak{z}(\mathbf{x}) + \epsilon \mathfrak{z}(\mathbf{x})^2) + \nu^k R_0 + \epsilon A_{1,k}^{(l)} \right\} \\ &\leq \sum_{s=0}^{l-1} \lambda_s \prod_{r=0}^{s-1} \beta_r \mathfrak{z}_s(\mathbf{x}) + \nu^k \sum_{s=0}^{l-1} \lambda_s \prod_{r=0}^{s-1} \beta_r R_s + \prod_{r=0}^{l-1} \beta_r R_l. \end{aligned} \quad (36)$$

We estimate further

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

Assuming (17) and the definition of $R_0 > \mathfrak{z}(\mathbf{x})$ this gives

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

With the same argument we find under (17) that

$$\nu^k \sum_{s=0}^{l-1} \lambda_s \prod_{r=0}^{s-1} \beta_r R_s \leq \nu^k \left(\mathbf{C}(\nu) + \frac{4\epsilon \mathbf{C}(\nu)^4 R_0}{1 - \epsilon \mathbf{C}(\nu) R_0} \right) R_0.$$

Additionally (17) implies

$$\prod_{r=0}^{l-1} \beta_r R_r \leq \left(\epsilon 4 \mathbf{C}(\nu) \frac{1}{\nu^{-1}-1} \right)^{2^{l-1}} R_0^{2^l} \rightarrow 0.$$

Plugging these bounds into (36) and letting $l \rightarrow \infty$ gives the claim. \blacksquare

4.5 Result after Convergence

In the previous section we showed that

$$\begin{aligned} \Omega(\mathbf{x}) \subset \bigcap_{\mathbf{r} \leq 4R_0(\mathbf{x})} \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \left\{ \frac{1}{6\check{\epsilon}\check{\nu}_1} \|\check{\mathcal{Y}}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \leq 3Q(\mathbf{x}, 2p^* + 2p)^2 \right\} \\ \cap \bigcap_{k \in \mathbb{N}} \{ \mathbf{v}_{k,k} \in \mathcal{Y}_\circ(\mathbf{r}_k^*), \mathbf{v}_{k,k+1} \in \mathcal{Y}_\circ(\mathbf{r}_k^*) \} \cap \{ \tilde{\mathbf{v}}, \tilde{\mathbf{v}}_{\theta^*} \in \mathcal{Y}_\circ(\mathbf{r}_0) \}, \end{aligned}$$

where \mathbf{r}_k^* is defined in (33). The claim of Theorem 7 follows with the following lemma:

Lemma 34 *Assume $(\check{\mathcal{D}}_1)$, $(\check{\mathcal{L}}_0)$, and (\mathcal{I}) . Then it holds on $\Omega(\mathbf{x}) \subseteq \epsilon$ that for all $k \in \mathbb{N}$*

$$\|\check{D}(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\boldsymbol{\xi}}\| \leq \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}), \quad (37)$$

$$|2\check{L}(\tilde{\boldsymbol{\theta}}_k, \boldsymbol{\theta}^*) - \|\check{\boldsymbol{\xi}}\|^2| \leq 5 \left(\|\check{D}^{-1}\check{\nabla}\| + \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}) \right) \check{\diamond}_Q(\mathbf{r}_k, \mathbf{x}), \quad (38)$$

where the spread $\check{\diamond}(\mathbf{r}, \mathbf{x})$ is defined in (21) and where

$$\mathbf{r}_k \stackrel{\text{def}}{=} \mathbf{r}_k^* \vee \mathbf{r}_0.$$

Proof The proof is nearly the same as that of Theorem 2.2 of (Andresen and Spokoiny, 2014) which is inspired by the proof of Theorem 1 of (Murphy and van der Vaart, 2000). So we only sketch it and refer the reader to (Andresen and Spokoiny, 2014) for the skipped arguments. We define

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

Note that

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

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

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

Now the right hand side can be bounded just as in the proof of Theorem 2.2 of (Andresen and Spokoiny, 2014). This gives (37).

For (38) we can represent:

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

where

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

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 (Andresen and Spokoiny, 2014). ■

5. Proof of Corollary 13

Proof Note that with the argument of Section 4.3 $\mathbb{P}(\epsilon'(\mathbf{x})) \geq 1 - 8e^{-\mathbf{x}} - \beta$ where with $\Omega(\mathbf{x})$ from (29)

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

On $\epsilon'(\mathbf{x})$ it holds due to Theorem 7 and due to Theorem 2.1 of (Andresen and Spokoiny, 2014)

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

Now the claim follows with the triangular inequality and noting that $\check{\diamond}(\mathbf{r}_0, \mathbf{x}) \leq \check{\diamond}_Q(\mathbf{r}_0, \mathbf{x})$. ■

6. Proof of Theorem 14

We prove this Theorem in a similar manner to the convergence result in Lemma 32. Redefine the set $\Omega(\mathbf{x})$

$$\begin{aligned} \Omega(\mathbf{x}) &\stackrel{\text{def}}{=} \bigcap_{k=0}^K (C_{k,k} \cap C_{k,k+1}) \cap C(\nabla) \cap \{\mathcal{L}(\tilde{\boldsymbol{v}}_0) - \mathcal{L}(\boldsymbol{v}^*) \geq -K_0\}, \quad \text{where} \quad (39) \\ C_{k,k+1} &= \left\{ \|\mathcal{D}(\tilde{\boldsymbol{v}}_{k,k+1}) - \boldsymbol{v}^*\| \leq R_0(\mathbf{x}), \|\mathcal{D}(\tilde{\boldsymbol{\theta}}_k) - \boldsymbol{\theta}^*\| \leq R_0(\mathbf{x}), \right. \\ &\quad \left. \|\mathcal{H}(\tilde{\boldsymbol{\eta}}_{k+1}) - \boldsymbol{\eta}^*\| \leq R_0(\mathbf{x}) \right\}, \\ C(\nabla) &= \left\{ \sup_{\boldsymbol{v} \in \mathcal{Y}_o(R_0(\mathbf{x}))} \|\mathcal{Y}(\nabla^2)(\boldsymbol{v})\| \leq 9\nu_2\epsilon_2\mathfrak{z}_1(\mathbf{x}, 6p^*)R_0(\mathbf{x}) \right\} \\ &\quad \cap \{\|\mathcal{D}^{-1}\nabla^2\zeta(\boldsymbol{v}^*)\| \leq \mathfrak{z}(\mathbf{x}, \nabla^2\zeta(\boldsymbol{v}^*))\}. \end{aligned}$$

where

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

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

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

Lemma 35 *Under the conditions of Theorem 14*

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

Proof The proof is very similar to the one presented in Section 4.3, so we only give a sketch. By assumption

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

and due to $(\mathcal{E}\mathcal{D}_2)$ with Theorem 47

$$\mathbb{P}\left(\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{R}_0(\mathbf{x}))} \|\mathfrak{y}(\nabla^2)(\mathbf{v})\| \leq 9\nu_2 \epsilon_2 \mathfrak{z}_1(\mathbf{x}, 6p^*) \mathbf{R}_0(\mathbf{x})\right) \geq 1 - e^{-\mathbf{x}}.$$

■

Lemma 36 *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 $\Omega(\mathbf{x})$

$$\begin{aligned} \|\mathcal{D}(\tilde{\mathbf{v}}_{k,k+1}) - \tilde{\mathbf{v}}\| &\leq 2\sqrt{2}(1 + \sqrt{\nu}) \sum_{r=0}^{k-1} \nu^r \|\boldsymbol{\tau}(\mathbf{r}_{k-r}^{(l)})\| + 2\sqrt{2}\nu^k (\mathbf{R}_0 + \mathbf{r}_0), \\ &=: \mathbf{r}_k^{(l+1)}. \end{aligned} \tag{40}$$

where

$$\|\boldsymbol{\tau}(\mathbf{r})\| \leq [\epsilon \mathbf{R}_0 + 9\nu_2 \epsilon_2 \|\mathcal{D}^{-1}\| \mathfrak{z}_1(\mathbf{x}, 6p^*) \mathbf{R}_0 + \|\mathcal{D}^{-1}\| \mathfrak{z}(\mathbf{x}, \nabla^2 \zeta(\mathbf{v}^*))] \mathbf{r}.$$

Proof

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

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

The proof is very similar to that of Lemma 32. Define

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

Note that

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

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

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

We want to show

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

where

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

To see this note that by assumption we have $\Omega(\mathbf{x}) \subseteq \{\tilde{\mathbf{v}} \in \mathcal{Y}_o(\mathbf{r}_0)\} \subseteq \{\tilde{\mathbf{v}} \in \mathcal{Y}_o(\mathbf{R}_0)\}$. By condition (\mathcal{L}_0) , Lemma 27 and Taylor expansion we have

$$\begin{aligned} &\sup_{(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^{(l)}) \cap \mathcal{Y}_o(\mathbf{R}_0)} \|\mathbb{E} \mathcal{U}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k)\| \\ &\leq \sup_{\mathbf{v} \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^{(l)}) \cap \mathcal{Y}_o(\mathbf{R}_0)} \|\mathbf{D}^{-1} \Pi_{\boldsymbol{\theta}} \left(\nabla \mathbb{E} \mathcal{L}(\mathbf{v}) - \nabla \mathbb{E} \mathcal{L}(\tilde{\mathbf{v}}) - \mathcal{D}(\mathbf{v} - \mathbf{v}^*) \right)\| \\ &\leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{R}_0)} \|\mathbf{D}^{-1} \Pi_{\boldsymbol{\theta}} \mathcal{D}\| \|\mathcal{D}^{-1} \nabla^2 \mathbb{E} \mathcal{L}(\mathbf{v}) \mathcal{D}^{-1} - I_{p^*}\| \|\mathbf{r}_k^{(l)}\| \\ &\leq \epsilon \mathbf{r}_k^{(l)2}. \end{aligned}$$

For the remainder note that with $\zeta = \mathcal{L} - \mathbb{E}\mathcal{L}$ on $\Omega(\mathbf{x})$ using Lemma 27 we can bound

$$\begin{aligned}
 & \sup_{(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^{(l)}) \cap \mathcal{Y}_o(\mathbf{R}_0)} \left\| \mathcal{U}_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) - \mathbb{E}\mathcal{U}_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}_k) \right\| \\
 & \leq \sup_{\mathbf{v} \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^{(l)}) \cap \mathcal{Y}_o(\mathbf{R}_0)} \left\| \mathcal{D}^{-1} \left(\nabla_{\boldsymbol{\theta}} \zeta(\mathbf{v}) - \nabla_{\boldsymbol{\theta}} \zeta(\tilde{\mathbf{v}}) \right) \right\| \\
 & \leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \left\| \mathcal{D}^{-1} \nabla^2 \zeta(\mathbf{v}) \mathcal{D}^{-1} \right\| \mathbf{r}_k^{(l)} \\
 & \leq \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{R}_0)} \left\{ \frac{1}{9\nu_2\epsilon_2} \left\| \mathcal{D}^{-1} \left(\nabla^2 \zeta(\mathbf{v}) - \nabla^2 \zeta(\mathbf{v}^*) \right) \mathcal{D}^{-1} \right\| \right\} 6\nu_1\epsilon \mathbf{r}_k^{(l)} \\
 & \quad + \left\{ \left\| \mathcal{D}^{-1} \nabla^2 \zeta(\mathbf{v}^*) \mathcal{D}^{-1} \right\| \right\} \mathbf{r}_k^{(l)} \\
 & \leq \left[9\nu_2\epsilon_2 \left\| \mathcal{D}^{-1} \right\| \mathfrak{z}_1(\mathbf{x}, 6p^*) \mathbf{R}_0 + \left\| \mathcal{D}^{-1} \right\| \mathfrak{z}(\mathbf{x}, \nabla^2 \zeta(\mathbf{v}^*)) \right] \mathbf{r}_k^{(l)}.
 \end{aligned}$$

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

Now the claim follows as in the proof of Lemma 32. ■

Lemma 37 *Assume that $\kappa(\mathbf{x}, \mathbf{R}_0) < 1 - \nu$ where*

$$\begin{aligned}
 \kappa(\mathbf{x}, \mathbf{R}_0) \stackrel{\text{def}}{=} & \frac{2\sqrt{2}(1 + \sqrt{\nu})}{\sqrt{1 - \nu}} \left(\epsilon \mathbf{R}_0 + 9\epsilon_2\nu_2 \left\| \mathcal{D}^{-1} \right\| \mathfrak{z}_1(\mathbf{x}, 6p^*) \mathbf{R}_0 \right. \\
 & \left. + \left\| \mathcal{D}^{-1} \right\| \mathfrak{z}_2(\mathbf{x}, \nabla^2 \mathcal{L}(\mathbf{v}^*)) \right).
 \end{aligned}$$

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 (24).

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

$$\left\| \boldsymbol{\tau}(\mathbf{r}_k^{(l)}) \right\| \leq \frac{1}{\sqrt{1 - \nu}} \left(\epsilon \mathbf{R}_0 + 6\nu_1\epsilon_2 \left\| \mathcal{D}^{-1} \right\| \mathfrak{z}_1(\mathbf{x}, 6p^*) \mathbf{R}_0 + \left\| \mathcal{D}^{-1} \right\| \mathfrak{z}(\mathbf{x}, \mathcal{B}(\nabla^2)) \right) \mathbf{r}_k^{(l)},$$

such that by definition

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

Plugging in the recursive formula for $\mathbf{r}_k^{(l)}$ from (40) and denoting $\tilde{\mathbf{R}}_0 \stackrel{\text{def}}{=} \mathbf{R}_0 + \mathbf{r}_0$ we find

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

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

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

By Lemma 36

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

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

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

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

$$\Omega(\mathbf{x}) \subset \bigcap_{k \in \mathbb{N}_0} \left\{ \tilde{\mathbf{v}}_{k,k(+1)} \in \tilde{\Upsilon}_o(\mathbf{r}_k^*) \right\},$$

$$\mathbf{r}_k^* \leq \begin{cases} \frac{\nu^k 2\sqrt{2}}{1 - \kappa(\mathbf{x}, \mathbf{R}_0)k} \tilde{\mathbf{R}}_0, & \kappa(\mathbf{x}, \mathbf{R}_0)k \leq 1, \\ 2 \left(\frac{\kappa(\mathbf{x}, \mathbf{R}_0)}{1 - \nu} \right)^{\frac{k}{\log(k)} L(k) - 1} \tilde{\mathbf{R}}_0, & \text{otherwise.} \end{cases}$$

The sequence $L(k) > 0$ is defined as

$$L(k) \stackrel{\text{def}}{=} \left\lfloor \frac{\log(1/\nu) - \frac{1}{k} (\log(2\sqrt{2}) - \log(\kappa(\mathbf{x}, \mathbf{R}_0)k - 1))}{1 + \frac{1}{\log(k)} \log(1 - \nu)} \right\rfloor \in \mathbb{N},$$

where $\lfloor x \rfloor \in \mathbb{N}_0$ denotes the largest natural number smaller than $x > 0$. To ensure that $L(k) > 0$ we assume that $k \log(1/\nu) - \log(2\sqrt{2}) > k$. Further as $\kappa(\mathbf{x}, \mathbf{R}_0) < (1 - \nu)$ and $L(k)$ is only relevant once $\kappa(\mathbf{x}, \mathbf{R}_0)k > 1$ it follows that

$$0 < 1 + \frac{1}{\log(k)} \log(1 - \nu) < 1.$$

Then

$$L(k) \geq \log(1/\nu) - \frac{1}{k} (\log(2\sqrt{2}) - \log(\kappa(\mathbf{x}, \mathbf{R}_0)k - 1)) > 1.$$

Consequently

$$\begin{aligned} \left(\frac{\kappa(\mathbf{x}, \mathbf{R}_0)}{1 - \nu} \right)^{\frac{k}{\log(k)} L(k)} &\leq \nu^{\frac{k}{\log(k)} \log\left(\frac{1 - \nu}{\kappa(\mathbf{x}, \mathbf{R}_0)}\right)} \left(\frac{\kappa(\mathbf{x}, \mathbf{R}_0)}{1 - \nu} \right)^{-\frac{1}{\log(k)} (\log(2\sqrt{2}) - \log(\kappa(\mathbf{x}, \mathbf{R}_0)k - 1))} \\ &\stackrel{\text{def}}{=} \nu^{\frac{k}{\log(k)} \log\left(\frac{1 - \nu}{\kappa(\mathbf{x}, \mathbf{R}_0)}\right)} c_k, \end{aligned}$$

where $c_k \rightarrow \frac{\kappa(\mathbf{x}, \mathbf{R}_0)}{1 - \nu}$. Finally note that $\tilde{\mathbf{R}}_0 \leq 2\mathbf{R}_0$ and the proof is complete. \blacksquare

Remark 38 *As pointed out in Remark 18 the above result can be improved. Redefine $\Omega(\mathbf{x})$ as the intersection of the two sets in (29) and (39). Then $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 10e^{-x}$. Also redefine*

$$\begin{aligned} \kappa(\mathbf{x}, \mathbf{r}) \stackrel{\text{def}}{=} \frac{2\sqrt{2}(1 + \sqrt{\nu})}{\sqrt{1 - \nu}} &\left(\epsilon \mathbf{r} + 3\epsilon_2 \|\mathcal{D}^{-1}\| \mathfrak{z}_1(\mathbf{x}, 6p^*) \mathbf{r} \right. \\ &\left. + \|\mathcal{D}^{-1}\| \mathfrak{z}_2(\mathbf{x}, \nabla^2 \mathcal{L}(\mathbf{v}^*)) \right). \end{aligned}$$

By the arguments of the proof of Theorem 7 we find with \mathbf{r}_k^* defined in (33)

$$\Omega(\mathbf{x}) \subset \bigcap_{k \in \mathbb{N}} \{ \mathbf{v}_{k,k(+1)} \in \Upsilon_o(\mathbf{r}_k^*) \}.$$

Using this in Lemma 36 instead of $\cap_{k \in \mathbb{N}} \{\mathbf{v}_{k,k(+1)} \in \mathcal{Y}_o(R_0)\}$ we can bound

$$\begin{aligned} \|\boldsymbol{\tau}(\mathbf{r}_k^{(l)})\| &\leq \frac{1}{\sqrt{1-\nu}} \left(\epsilon \mathbf{r}_k^* + 6\nu_1 \epsilon_2 \|\mathcal{D}^{-1}\| \mathfrak{z}_1(\mathbf{x}, 6p^*) \mathbf{r}_k^* \right. \\ &\quad \left. + \|\mathcal{D}^{-1}\| \mathfrak{z}(\mathbf{x}, \mathcal{B}(\nabla^2)) \right) \mathbf{r}_k^{(l)}. \end{aligned}$$

Consequently, representing $\mathbf{r}_k^* = \mathbf{C}(\mathfrak{z}(\mathbf{x}) + \nu^k R_0)$ we find

$$\begin{aligned} \mathbf{r}_k^{(l)} &\leq \kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x})) \sum_{r=0}^{k-1} \nu^r \mathbf{r}_{k-r}^{(l-1)} + 2\sqrt{2}\nu^k (R_0 + \mathbf{r}_0) \\ &\quad + \mathbf{C}\epsilon(1 + \|\mathcal{D}^{-1}\| \mathfrak{z}_1(\mathbf{x}, 6p^*)) \sum_{r=0}^{k-1} \nu^k R_0 \mathbf{r}_{k-r}^{(l-1)} \\ &\leq \kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x})) \sum_{r=0}^{k-1} \nu^r \mathbf{r}_{k-r}^{(l-1)} + \mathbf{C}_1 \epsilon R_0 k \nu^k (R_0 + \mathbf{r}_0), \end{aligned}$$

where $\mathbf{C}_1 \leq 2\sqrt{2} + \mathbf{C}(1 + \|\mathcal{D}^{-1}\| \mathfrak{z}_1(\mathbf{x}, 6p^*))$. With the same arguments as in the proof of Lemma 37 we infer

$$\begin{aligned} &\mathbf{r}_k^{(l)} / (R_0 + \mathbf{r}_0) \\ &\leq \begin{cases} \left(\left(\frac{\kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))}{1-\nu} \right)^l + k \nu^k \frac{\mathbf{C}_1 \epsilon R_0}{1-\kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))k} \right), & \kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))k \leq 1, \\ \kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))^l \left(\left(\frac{1}{1-\nu} \right)^l + \nu^k \frac{k^{l+1} \mathbf{C}_1 \epsilon R_0}{\kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))k-1} \right), & \text{otherwise.} \end{cases} \end{aligned}$$

Set

$$l(k) \stackrel{\text{def}}{=} \begin{cases} \infty, & \kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))k \leq 1, \\ \frac{k \log(\nu) + \log(\mathbf{C}_1 \epsilon R_0) - \log(\kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))k-1) - \log(k)}{-\log(1-\nu) - \log(k)}, & \text{otherwise.} \end{cases}$$

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

$$\begin{aligned} \Omega(\mathbf{x}) &\subset \bigcap_{k \in \mathbb{N}_0} \left\{ \tilde{\mathbf{v}}_{k,k(+1)} \in \tilde{\mathcal{Y}}_o(\mathbf{r}_k^*) \right\}, \\ \mathbf{r}_k^* &\leq \begin{cases} \frac{\nu^k 2\sqrt{2}}{1-\kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))k} (R_0 + \mathbf{r}_0), & \kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))k \leq 1, \\ 2 \left(\frac{\kappa(\mathbf{x}, \mathbf{C}\mathfrak{z}(\mathbf{x}))}{1-\nu} \right)^{\frac{k}{\log(k)} L(k)-1} (R_0 + \mathbf{r}_0), & \text{otherwise.} \end{cases} \end{aligned}$$

This gives the claim.

Acknowledgments

The first author was supported by Research Units 1735 "Structural Inference in Statistics: Adaptation and Efficiency". The research was partly supported by the Russian Science Foundation grant (project 14-50-00150). Financial support by the German Research Foundation (DFG) through the Collaborative Research Center 649 "Economic Risk" is gratefully acknowledged.

Appendix A. Deviation Bounds for Quadratic Forms

This section is the same as Section A of Andresen and Spokoiny (2014). The following general result from Spokoiny (2012) helps to control the deviation for quadratic forms of type $\|\mathcal{B}\boldsymbol{\xi}\|^2$ for a given positive matrix \mathcal{B} and a random vector $\boldsymbol{\xi}$. It will be used several times in our proofs. Suppose that

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

For a symmetric matrix \mathcal{B} , define

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

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

$$\begin{aligned} \mathfrak{g}_c &\stackrel{\text{def}}{=} \sqrt{\mathfrak{g}^2 - \mu_c \mathfrak{p}_B}, \\ 2(\mathfrak{x}_c + 2) &\stackrel{\text{def}}{=} (\mathfrak{g}^2/\mu_c - \mathfrak{p}_B)/\lambda^* + \log \det(\mathcal{I}_p - \mu_c \mathcal{B}/\lambda^*). \end{aligned}$$

Proposition 39 *Let (ED_0) hold with $\nu_0 = 1$ and $\mathfrak{g}^2 \geq 2\mathfrak{p}_B$. Then for each $\mathfrak{x} > 0$*

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

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

$$\mathfrak{z}^2(\mathcal{B}, \mathfrak{x}) \stackrel{\text{def}}{=} \begin{cases} \mathfrak{p}_B + 2\mathfrak{v}_B(\mathfrak{x} + 1)^{1/2}, & \mathfrak{x} + 1 \leq \mathfrak{v}_B/(18\lambda^*), \\ \mathfrak{p}_B + 6\lambda^*(\mathfrak{x} + 1), & \mathfrak{v}_B/(18\lambda^*) < \mathfrak{x} + 1 \leq \mathfrak{x}_c + 2, \\ |y_c + 2\lambda^*(\mathfrak{x} - \mathfrak{x}_c + 1)/\mathfrak{g}_c|^2, & \mathfrak{x} > \mathfrak{x}_c + 1, \end{cases}$$

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

Appendix B. A Uniform Bound for the Norm of a Random Process

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

$$\mathbb{P} \left(\sup_{\mathbf{r} \leq \mathbf{r}^*} \sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \left\{ \frac{1}{\epsilon} \|\check{\mathcal{Y}}(\mathbf{v})\| - 2\mathbf{r}^2 \right\} \geq \mathfrak{C}\mathfrak{z}_Q(\mathbf{x}, p^*) \right) \leq e^{-\mathfrak{x}}.$$

This is a slightly stronger result than the one derived in Section D of (Andresen and Spokoiny, 2014) but the ideas employed here are very similar.

We want to apply Corollary 2.5 of the supplement of Spokoiny (2012) 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 40 *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 $\mathbf{v}, \mathbf{v}^\circ \in U(\mathbf{r})$

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{r}, \mathbf{v}) - \mathcal{U}(\mathbf{r}, \mathbf{v}^\circ)}{d(\mathbf{v}, \mathbf{v}^\circ)} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathbf{g}. \quad (41)$$

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

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

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 $\mathbf{g}_0 = \nu_0 \mathbf{g}$ and for some $\mathbb{Q} > 0$

$$\mathfrak{z}_Q(\mathbf{x}, \mathbb{Q})^2 \stackrel{\text{def}}{=} \begin{cases} (1 + \sqrt{\mathbf{x} + \mathbb{Q}})^2 & \text{if } 1 + \sqrt{\mathbf{x} + \mathbb{Q}} \leq \mathbf{g}_0, \\ 1 + \{2\mathbf{g}_0^{-1}(\mathbf{x} + \mathbb{Q}) + \mathbf{g}_0\}^2 & \text{otherwise.} \end{cases} \quad (42)$$

To use this result let $\check{\mathcal{Y}}(\mathbf{v})$ 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 $\|\check{\mathcal{Y}}(\mathbf{v})\|$ over a vicinity $\mathcal{Y}_\circ(\mathbf{r}) \stackrel{\text{def}}{=} \{\|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}\}$ of \mathbf{v}^* . Suppose that $\check{\mathcal{Y}}(\mathbf{v})$ satisfies for each $0 < \mathbf{r} < \mathbf{r}^*$ and for all pairs $\mathbf{v}, \mathbf{v}^\circ \in \mathcal{Y}_\circ(\mathbf{r}) = \{\mathbf{v} \in \mathcal{Y} : \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}\} \subset \mathbb{R}^{p^*}$

$$\sup_{\|\mathbf{u}\| \leq 1} \log \mathbb{E} \exp \left\{ \lambda \frac{\mathbf{u}^\top (\check{\mathcal{Y}}(\mathbf{v}) - \check{\mathcal{Y}}(\mathbf{v}^\circ))}{\epsilon \|\mathcal{D}(\mathbf{v} - \mathbf{v}^\circ)\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}. \quad (43)$$

Remark 41 In the setting of Theorem 7 we have

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

and condition (43) becomes $(\mathcal{E}\mathcal{D}_1)$ from 2.1.

Theorem 42 Let a random p -vector process $\check{\mathcal{Y}}(\mathbf{v})$ fulfill $\check{\mathcal{Y}}(\mathbf{v}^*) = 0$ and let condition (43) be satisfied. Then for each $0 \leq \mathbf{r} \leq \mathbf{r}^*$, on a set of probability greater $1 - e^{-\mathbf{x}}$

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

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

Remark 43 Note that the entropy of the original set $\mathcal{Y}_\circ(\mathbf{r}) \subset \mathbb{R}^{p^*}$ is equal to $2p^*$. So in order to control the norm $\|\check{\mathcal{Y}}(\mathbf{v})\|$ one only pays with the additional summand $2p$.

Proof In what follows, we use the representation

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

This implies

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

Due to Lemma 44 the process $\mathcal{U}(\mathbf{r}, \mathbf{v}, \mathbf{u}) \stackrel{\text{def}}{=} \frac{1}{\epsilon \|\mathcal{D}(\mathbf{v}-\mathbf{v}^*)\|} \mathbf{u}^\top \check{\mathcal{Y}}(\mathbf{v})$ satisfies condition $(\mathcal{E}d)$ (see (41)) as process on $U(\mathbf{r}^*)$ where

$$U(\mathbf{r}) \stackrel{\text{def}}{=} \mathcal{Y}_\circ(\mathbf{r}) \times B_{\mathbf{r}}(0). \quad (44)$$

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 42 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 $1 - e^{-x}$

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

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 Spokoiny (2012), Lemma 2.10). This gives $\mathbb{Q}(U) = 2p^* + 2p$. Finally observe that

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

■

Lemma 44 *Suppose that $\check{\mathcal{Y}}(\mathbf{v})$ satisfies for each $\|\mathbf{u}\| \leq 1$ and $|\lambda| \leq \mathfrak{g}$ the inequality (43). Then the process $\mathcal{U}(\mathbf{v}, \mathbf{u}) = \frac{1}{2\epsilon \|\mathcal{D}(\mathbf{v}-\mathbf{v}^*)\|} \check{\mathcal{Y}}(\mathbf{v})^\top \mathbf{u}_1$ satisfies $(\mathcal{E}d)$ from (41) with $|\lambda| \leq \mathfrak{g}/2$, $d((\mathbf{v}, \mathbf{u}), (\mathbf{v}^\circ, \mathbf{u}^\circ))^2 = \|\mathcal{D}(\mathbf{v}-\mathbf{v}^*)\|^2 + \|\mathbf{u}-\mathbf{u}^\circ\|^2$, $\nu = 2\nu_0$ and $U \subset \mathbb{R}^{p^*+p}$ defined in (44), i.e. for any $(\mathbf{v}, \mathbf{u}_1), (\mathbf{v}^\circ, \mathbf{u}_2) \in U$*

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

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 (43), we find

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

■

Appendix C. 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

$$\mathbb{P} \left(\sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \left\{ \|\check{\mathcal{Y}}(\mathbf{v})\| \right\} \geq \mathbf{C} \epsilon_2 \delta_1(\mathbf{x}, p^*) \mathbf{r} \right) \leq e^{-\mathbf{x}}.$$

We derive such a bound in a very similar manner to Theorem E.1 of Andresen and Spokoiny (2014).

We want to apply Corollary 2.2 of the supplement of Spokoiny (2012). Again we slightly generalized the formulation but the proof remains the same.

Corollary 45 *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 $\mathbf{v}, \mathbf{v}^\circ \in U(\mathbf{r})$*

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)}{d(\mathbf{v}, \mathbf{v}^\circ)} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathbf{g}. \quad (45)$$

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

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

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 $\mathbb{Q} > 0$

$$\mathfrak{z}_1(\mathbf{x}, \mathbb{Q}) \stackrel{\text{def}}{=} \begin{cases} \sqrt{2(\mathbf{x} + \mathbb{Q})} & \text{if } \sqrt{2(\mathbf{x} + \mathbb{Q})} \leq \mathfrak{g}_0, \\ \mathfrak{g}_0^{-1}(\mathbf{x} + \mathbb{Q}) + \mathfrak{g}_0/2 & \text{otherwise.} \end{cases} \quad (46)$$

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 $\mathbf{v}, \mathbf{v}^\circ \in \mathcal{Y}_o(\mathbf{r}) = \{\mathbf{v} \in \mathcal{Y} : \|\mathbf{v} - \mathbf{v}^*\|_{\mathcal{Y}} \leq \mathbf{r}\} \subset \mathbb{R}^{p^*}$

$$\sup_{\|\mathbf{u}_1\| \leq 1} \sup_{\|\mathbf{u}_2\| \leq 1} \log \mathbb{E} \exp \left\{ \lambda \frac{\mathbf{u}_1^\top (\mathcal{Y}(\mathbf{v}) - \mathcal{Y}(\mathbf{v}^\circ)) \mathbf{u}_2}{\epsilon_2 \|\mathcal{D}(\mathbf{v} - \mathbf{v}^\circ)\|} \right\} \leq \frac{\nu_2^2 \lambda^2}{2}. \quad (47)$$

Remark 46 In the setting of Theorem 14 we have $\|\mathbf{v} - \mathbf{v}^\circ\|_{\mathcal{Y}} = \|\mathcal{D}(\mathbf{v} - \mathbf{v}^\circ)\|$ and

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

and condition (47) becomes $(\mathcal{E}\mathcal{D}_2)$ from 2.1.

Theorem 47 Let a random process $\mathcal{Y}(\mathbf{v}) \in \mathbb{R}^{p^* \times p^*}$ fulfill $\mathcal{Y}(\mathbf{v}^*) = 0$ and let condition (47) be satisfied. Then for each $0 \leq \mathbf{r} \leq \mathbf{r}^*$, on a set of probability greater than $1 - e^{-x}$

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| \leq 9\epsilon_2 \nu_2 \mathfrak{z}_1(\mathbf{x}, 6p^*) \mathbf{r},$$

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

Remark 48 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

$$\|\mathcal{Y}(\mathbf{v})\| = \epsilon_2 \sup_{\|\mathbf{u}_1\| \leq \mathbf{r}} \sup_{\|\mathbf{u}_2\| \leq \mathbf{r}} \frac{1}{\epsilon_2 \mathbf{r}^2} \mathbf{u}_1^\top \check{\mathcal{Y}}(\mathbf{v}) \mathbf{u}_2.$$

This implies

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

Due to Lemma 49 the process $\mathcal{U}(\mathbf{v}) \stackrel{\text{def}}{=} \frac{1}{\epsilon \mathbf{r}^2} \mathbf{u}_1^\top \mathcal{Y}(\mathbf{v}) \mathbf{u}_2$ satisfies condition $(\mathcal{E}d)$ (see (45)) as process on

$$U(\mathbf{r}) \stackrel{\text{def}}{=} \mathcal{Y}_\circ(\mathbf{r}) \times B_{\mathbf{r}}(0) \times B_{\mathbf{r}}(0) \subset \mathbb{R}^{3p^*}. \quad (48)$$

This allows to apply Corollary 45 to obtain the desired result. We get on a set of probability greater $1 - e^{-x}$

$$\sup_{\mathbf{v} \in \mathcal{Y}_\circ(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| \leq \sup_{(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2) \in U(\mathbf{r})} \left\{ \frac{1}{\mathbf{r}^2} \mathbf{u}_1^\top \mathcal{Y}(\mathbf{v}) \mathbf{u}_2 \right\} \leq 9\epsilon_2 \nu_2 \delta_1 \left(\mathbf{x}, \mathbb{Q}(U(\mathbf{r}^*)) \right) \mathbf{r}.$$

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 Spokoiny (2012), Lemma 2.10). This gives the claim. \blacksquare

Lemma 49 *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 $|\lambda| \leq \mathbf{g}$ the inequality (47). Then the process*

$$\mathcal{U}(\mathbf{v}, \mathbf{u}_1, \mathbf{u}_2) = \frac{1}{2\epsilon_2 \mathbf{r}^2} \mathbf{u}_1^\top \mathcal{Y}(\mathbf{v}) \mathbf{u}_2$$

satisfies $(\mathcal{E}d)$ from (45) with $U \subset \mathbb{R}^{3p^}$ defined in (48), with $|\lambda| \leq \mathbf{g}/3$ and with*

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

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

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

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 (47), 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 \left(\frac{1}{r^2} \check{\mathcal{Y}}(\mathbf{v}) - \frac{1}{r^2} \check{\mathcal{Y}}(\mathbf{v}^\circ) \right) \mathbf{u}_2}{\epsilon_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}{\epsilon_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)}{\epsilon_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}{\epsilon_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}{\epsilon_2 \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\|} \right\} \\
 &\leq \frac{9\nu_2^2 \lambda^2}{2}, \quad \lambda \leq \mathfrak{g}/3.
 \end{aligned}$$

■

References

- A. Andresen. Finite sample analysis of profile m-estimation in the single index model. *Electronic Journal of Statistics*, 9(2):2528–2641, 2015.
- A. Andresen and V. Spokoiny. Critical dimension in profile semiparametric estimation. *Electron. J. Statist.*, 8(2):3077–3125, 2014. ISSN 1935-7524. doi: 10.1214/14-EJS982.
- S. Balakrishnan, M. J. Wainwright, and B. Yu. Statistical guarantees for the em algorithm: From population to sample-based analysis. *arXiv: 1408.2156*, 2014.
- Xiaohong Chen. Large sample sieve estimation of semi-nonparametric models. *Handbook of Econometrics*, 6:55495632, 2007.
- G. Cheng. How many iterations are sufficient for efficient semiparametric estimation? *Scandinavian Journal of Statistics*, 40(3):592–618, 2013.
- M. Delecroix., W. Haerdle, and M. Hristache. Efficient estimation in single-index regression. Technical report, SFB 373, Humboldt Univ. Berlin, 1997.

- A.P Dempster, N.M. Laird, and D.B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society*, Series B, 39:1–38, 1977.
- U. Grenander. *Abstract Inference*. Wiley, New York, 1981.
- W. Haerdle, P. Hall, and H. Ichimura. Optimal smoothing in single-index models. *Ann. Statist.*, 21:157–178, 1993.
- I.A. Ibragimov and R.Z. Khas'minskij. "Statistical estimation. Asymptotic theory. Transl. from the Russian by Samuel Kotz.". "New York - Heidelberg -Berlin: Springer-Verlag", 1981.
- P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. *STOC*, pages 665–674, 2013.
- R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from few entries. *IEEE Transactions on Information Theory*, 56(6):2980–2998, 2010.
- M.R. Kosorok. *Introduction to Empirical Processes and Semiparametric Inference*. Springer in Statistics, 2005.
- G.J. McLachlan and T Krishnan. *The EM Algorithm and Extensions*. Wiley, New York, 1997.
- S. A. Murphy and A. W. van der Vaart. On profile likelihood. *Journal of the American Statistical Association*, 95(450):449–465, 2000.
- P. Netrapalli, P. Jain, and S. Sanghavi. Phase retrieval using alternating minimization. *NIPS*, pages 2796–2804, 2013.
- Vladimir Spokoiny. Parametric estimation. Finite sample theory. *Ann. Statist.*, 40(6): 2877–2909, 2012.
- J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12:389–434, 2012.
- C.F.J Wu. On the convergence properties of the em algorithm. *Annals of Statistics*, 11: 95–103, 1983.
- X. Yi, C. Caramanis, and S. Sanghavi. Alternating minimization for mixed linear regression. *arXiv: 1310.3745*, 2013.